

Final exam

PARALLEL COMPUTING

Andrea Boido

March 25, 2019

Abstract

I introduce briefly the main aspects and advantages of the parallelization in computational science, in particular referring to the OpenCL framework. Then I explain and analyse an application of this technique to the eigenproblem and eventually to the numerical solution of the 2D Schrödinger equation.

1 Theory

1.1 Parallel computing

In the last years the increasing difficulties that emerged in designing and building faster and faster (in terms of frequency) processors forced us to go towards another way to improve the computational power of computers: the parallelization. It mainly consists in the splitting of the CPUs in subpieces called *compute units* (or *cores*) that have the capability of working semi-independently (it will be clear what I mean with “semi” afterwards) and simultaneously. In this way different parts of a given algorithm can be executed by different compute units at the same time (if the algorithm allows such a splitting), which results typically in a shorter time to complete the task. This is the same kind of approach that have been used for many years and on a much bigger scale in supercomputers, even though there are big differences in the concrete implementation of the idea. Nevertheless the basic features that we will consider are mostly common.

Among this scenario a key role has been played recently also by the development of the GPGPU, standing for *General-purpose computing on graphics processing units*, which is basically the practice of exploiting the characteristic features of a GPU to perform computing tasks that are traditionally carried on by a CPU. The main advantage is that GPUs are built specifically to be as fast as possible in graphics processing and thus if we are able to translate our data of interest into graphical form and make the GPU handle them, it is likely that we achieve a great speed-up. Moreover GPUs generally have a higher number of compute units compared to CPUs, hence we can furtherly take advantage of the parallelization paradigm. Compared to CPUs, they also provide an impressive power efficiency.

It is clear that these techniques can be very useful in the field of scientific computing and, for what concerns us, in particular in computational physics where the models are becoming increasingly complex and the simulations resources-demanding. Especially the use of GPUs is becoming more and more widespread even inside the supercomputers, so that GPUs producers like AMD and NVIDIA have built graphics cards designed specifically for the purpose of doing computations (for instance the series *Tesla* by NVIDIA).

Obviously the parallelization carries also some downsides; the greater is probably the higher complexity both of the algorithms, which have to be optimized properly for taking advantage of this approach, and also of the coding phase, which requires the developer to manually handle many more tasks compared to traditional programs (called *single-thread* in the following), for instance the synchronization between the compute units, memory management and data transfers. Writing a good parallelized program is much more difficult than writing a good single-thread one. However I will speak more concretely about some of these aspects in the following.

The most famous and used mainstream software for parallel computing are CUDA, a platform for GPGPU on NVIDIA graphics cards designed to work with C, C++, and Fortran, OpenMP an API¹ for multiprocessing programming again in C, C++, and Fortran (only on CPUs), and OpenCL, which is the main subject of this report. In the next section I will briefly introduce some basic aspects of the latter one.

1.2 Brief introduction to OpenCL

OpenCL is an open source standard for parallel programming developed initially by Apple and maintained by the consortium Khronos Group. Its main characteristic is that the same OpenCL code can run across *heterogeneous* platforms: CPUs, GPUs but also other kinds of data processors (for instance FPGAs²); the work of interfacing with the different hardware is done mostly by the OpenCL library, leaving the programmer the only job of writing a good code without worrying too much about the device it will run on. However in the following I will focus on standard computers.

The OpenCL standard defines a set of specifications for a C/C++ library that have to be implemented concretely by hardware vendors. The various implementations of OpenCL are called *platforms* and OpenCL allows one to use even more than one platform inside the same program, so that one can take full advantage of all the computational power he has in its computer. For instance if a computer has an Intel CPU and an AMD GPU, one can write a program using both of them to do computations simultaneously! The only limit to such a broad parallelization is the algorithm itself.

Each OpenCL program is composed by two main parts: the *host program* and the *kernels*. The former is a code written in C or C++ that runs on the CPU and takes care of initializing the OpenCL environment and managing the overall execution of the algorithm. Kernels are kinds of functions written in the OpenCL language, which is actually quite similar to C, that run on a chosen target device which can be a GPU or even the CPU itself; they are dispatched to the device (or devices) by the host together with the data of interest to be elaborated.

On the host side there are a couple of fundamental structures defined in the OpenCL standard that need to be understood in order to learn how to code. They are listed below; refer to figure 1 for a pictorial representation.

- *Context* - A logical container inside which the host puts a set of devices selected to work together. It is a structure that becomes useful when the code must run on complex machines with distributed processors, while when using a standard computer one typically initializes a single context. Note that devices running on different platforms cannot be in the same context, so it may be useful to define more than one context also in the latter case if one wants to run OpenCL on a CPU and a GPU from different vendors.

¹Application Programming Interface.

²Field-programmable gate array; they are integrated circuits designed for specific applications.

- *Command queue* - The name is rather self-explanatory: it is the structure that allows the host to tell devices what to do. Each device has his own command queue and the instructions that are enqueued by the host are executed sequentially, unless otherwise specified by the host. Through the command queue the host mainly deploys kernels to devices and transfers data from/to them.
- *Program* - A container storing a set of kernels. This structure is necessary since kernels have to be compiled run-time during the execution of the code and thus many of them are grouped inside a single program that is built just once. Clearly the compiling process depends on the particular platform and device one is using, so the compiling function requires to have gained informations about them in advance, which is done with some simple functions included in the OpenCL library. Also we can pass optional arguments to the compiling function through flags.

In addition to this, the OpenCL standard defines additional data types and data structures that must be used to interface properly with its functions.

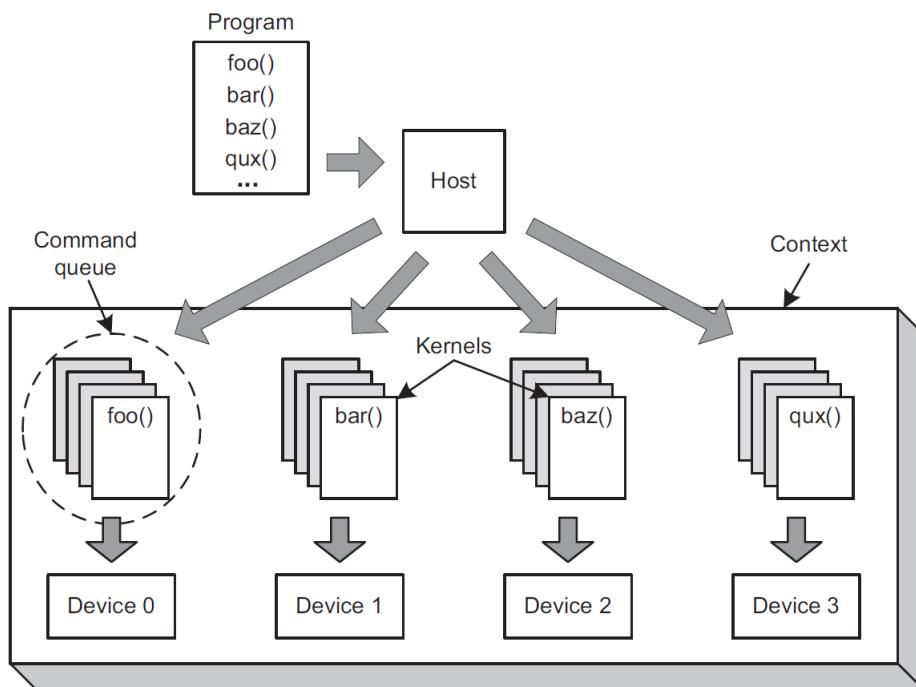


Figure 1: *Host managing the execution* (picture taken from the book “OpenCL in action”).

Switching to the kernel side, one has to understand how kernels are executed on the devices and this requires to understand how a device is structured within the OpenCL framework. Refer to figure 2 for a pictorial representation of the following discussion. As I already said, physically each device is divided into compute units but also each of them has multiple processing elements³. The data being processed is stored in various different memory location, mainly RAM, cache, VRAM, registers, depending also on what kind of device is being used, and each memory has a different capacity and speed (typically the higher the capacity, the lower the speed and viceversa). From the more abstract point of

³This is thanks to the SIMD architecture, meaning “Single instruction, multiple data” that allows a processor to execute simultaneously the same task with different input data. Here one should distinguish between CPUs and GPUs since they are quite different, but I will not go into further details of hardware-related topics.

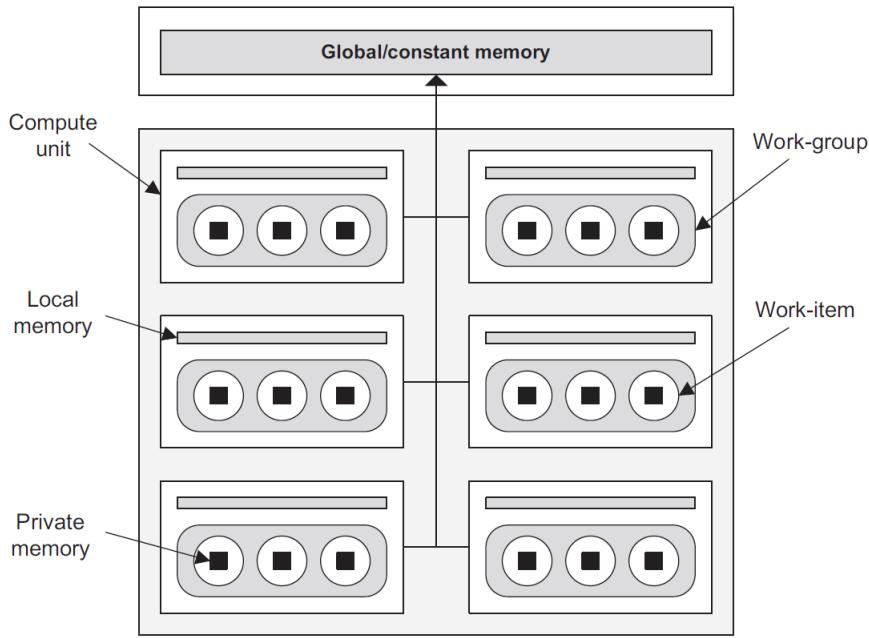


Figure 2: *Device model (picture taken from the book “OpenCL in action”)*.

view of OpenCL what happens is that, at the time of deploying a kernel to a device, the host fixes also the number of its instances that have to be executed; this number is called *global size* and each instance is called *work item*. So each work item executes the kernel once for a given set of input parameters, in general different from those of the other work items. Physically a work item is executed by a single processing element lying inside a compute unit of the device. Moreover a bunch of work items can be grouped together to form a *work group*; a given work group will be executed by a single compute unit and the number of work items inside it is called *local size*. The maximum possible local size is established by the device producer, while the number of work groups and the global size can be arbitrarily big, though one has to keep in mind that each compute unit can execute a single work group at a time. A work item is thus identified by both a *global id* and a *local id* whereas a work group is identified by the *work group id*. For what concerns the memory, from the OpenCL point of view we have three kinds of memory:

- *Global memory* - A memory shared among all the work items, which is usually the biggest but also the slowest. Most of the data the host transfers to the device goes into the global memory (exceptions can be only single variables that go directly into the private memory).
- *Local memory* - A memory common to all work items inside a work group. It should be sensibly faster than global memory but with less capacity and the host cannot access it directly. It is useful to store intermediate results.
- *Private memory* - The own memory of each work item. It should be the fastest of all but it is even smaller than local memory and can be used only by the single work item.

The physical partition of such memories is not easy to understand in general and it highly depends on the device producer; anyway to fix ideas for example we can think of GPUs as storing global memory inside VRAM and local/private memory inside the cache and the registers. The last important feature that has to be understand is that an uncontrolled

access to shared memory (global and private) could result in bugs inside the code. Indeed the programmer has almost no control on the timing of the events that happens on a device, thus for instance it may be possible that a work item accesses a value in the local memory that is not the correct one since another work item is still processing it. In order to avoid such problems one can implement *barriers* inside a kernel which are instructions that tell the device to wait until all the work items inside a work group have done their job; this way one can synchronize the execution within work groups. Unfortunately up to now it is impossible to synchronize work items belonging to different work groups.

Once one has understood the basic principles of working of the OpenCL framework, it is pretty clear what is the direction to follow in writing a code. Whenever an iterative cycle of definite size occurs inside an algorithm, it can be replaced by a kernel executing (hopefully simultaneously) a work item for each iteration, so that the execution time will be dramatically reduced, at least in principle. Obviously this approach cannot be always followed since a wide class of iterative algorithms are structured to use at a given iteration data computed during the previous one (we will see examples of this situation in the following).

Finally I should mention that the OpenCL kernel language provides the programmer with a set of data structures and built-in functions (mainly arithmetic functions) that are highly optimized for parallelization at a lower level and should help in improving furtherly the computations speed, for examples vectors, atomic operations etc. Some of them like vectors have been slightly used in my code, nevertheless these are quite advanced tools that should not be used until having achieved a minimum of experience with OpenCL, thus I will not speak about them.

1.3 Environment specifications and settings

All the programs analysed in the following have runned on my desktop PC. It has an Intel Core i5 4670 as CPU with 8 GB of DDR3 RAM, an AMD Radeon HD 7870 with 2 GB of GDDR5 VRAM as GPU and Windows 10 version 1809 as operative system. I installed on the system both the AMD platform⁴ and the Intel platform⁵ so as to compare the performance of OpenCL on different devices. The Intel CPU has also an integrated GPU but unfortunately it does not support the OpenCL extension for the double type, so I could not include it in the comparison. See table 1 for more detailed technical specifications of the two OpenCL devices. Note that though the latest version of OpenCL is the 2.2, my hardware is capable of running only the 1.2.

Device	Intel Core i5 4670	AMD Radeon HD 7870
Clock	3400 MHz	1000 MHz
OpenCL version	1.2	1.2
# (compute units)	4	20
Global memory	2047 MB	2048 MB
Local memory	32 kB	32 kB
Max work group size	8192	256

Table 1: *Technical specification of the two devices used (taken from GPU Caps Viewer).*

As a compiler for C/C++ I used the Windows porting of GCC 8.1.0 provided by MinGW-W64⁶.

⁴AMD Accelerated Parallel Processing

⁵Intel(R) OpenCL

⁶<https://mingw-w64.org>

I also installed *LAPACK* and its C interface *LAPACKE*⁷ with the main aim of checking the results of my algorithms, and the C++ library for sparse linear algebra *ViennaCL*⁸ that is capable of taking advantage of OpenCL and that I used for solving numerically the Schrödinger equation. Unfortunately I did not find any similar library for C, therefore I was forced to use also C++.

2 Matrix diagonalization

The main purpose of the following application is to write an original code using OpenCL that finds the lowest eigenpair of a given matrix (assumed to be hermitian) and to test its performance, in particular comparing it with the corresponding single-thread algorithm.

2.1 Algorithm

Despite the appearance of a rather basic task, the matrix diagonalization is not easy at all in computational science and includes lots of subtleties that prevent to identify a single method as the universally best one. Provided that the main purpose of this application was to evaluate OpenCL performances and not to write an excellent diagonalization algorithm, I ignored most of these issues and I tried to proceed straightforwardly with a relatively easy implementable solution. I individuated in the QR algorithm the best option.

The QR algorithm relies on the linear algebra statement that every real square matrix A can be written as $A = Q_0 R_0$ where Q_0 is an orthogonal matrix and R_0 is an upper triangular matrix; this is called QR decomposition. Provided that one is able to perform such a decomposition on a computer (it is not trivial how to do it, I will discuss it in the following), then he can build a succession of matrices $A_k = R_{k-1} Q_{k-1}$ with $k \in \mathbb{N}$. Then I can write:

$$A_{k+1} = R_k Q_k = Q_k^{-1} A_k Q_k = Q_k^T A_k Q_k = Q_k^T \dots Q_0^T A Q_0 \dots Q_k = (Q_0 \dots Q_k)^T A (Q_0 \dots Q_k)$$

It can be proved that this succession converges to a diagonal matrix, thus defining

$$D \equiv \lim_{k \rightarrow +\infty} A_k \quad P \equiv \lim_{k \rightarrow +\infty} (Q_0 \dots Q_k)$$

it is clear that $D = P^T A P$ with D containing the eigenvalues and P the eigenvectors of A . Clearly in a real implementation the limit has to be intended as a convergence in the computational sense i.e. when the quantities of interest do not vary between an iteration and the other up to a given precision.

However this algorithm carries some fundamental problems:

- the convergence is quite slow, especially if eigenvalues are close to each other;
- it is computationally expensive: each iteration step requires to compute two products of $N \times N$ matrices, thus even without considering the QR decomposition each single iteration step has a complexity of $\mathcal{O}(N^3)$; assuming that the number of steps is proportional to N , which is not even assured because of the above issue, one would get an $\mathcal{O}(N^4)$ complexity;
- it requires the initial matrix to be real, which is not the case in general if one wants to use the algorithm for solving quantum mechanics problems;

⁷<https://icl.cs.utk.edu/lapack-for-windows/lapack/>

⁸<http://viennacl.sourceforge.net/>

- without including further adjustments (that necessarily increase the coding complexity) the algorithm is quite unstable and could lead to completely wrong results depending on the initial matrix.

The latter issue was ignored because I was mostly interested in computing the lowest eigen-pair which proved to be quite accurate most of the times. For what concerns the first issue I tried to mitigate it by using the shifted QR algorithm which consists in applying the QR factorization to the matrices $\tilde{A}_k = A_k - \lambda \mathbb{I}$ where λ is the element (N, N) of A_k ; this modification was suggested by some references found on internet but I have not observed any tangible speed improvement in the convergence. The second and the third problem can be partially addressed by doing a preliminary step. Indeed any complex matrix can be brought to a real tridiagonal matrix by means of a unitary transformation i.e. given M complex we can write $M = U^\dagger T U$ with T real tridiagonal and U unitary. Also if M is hermitian, then T will be symmetric. So if one performs this transformation on the starting hermitian matrix he solves completely the third problem. Moreover the QR decomposition of a tridiagonal matrix can be computed quite easily and with a good speed by means of Givens' rotations. These exploit the fact that any 2-component vector $v = (v_1, v_2)$ can be rotated so that its second entry vanishes:

$$\begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = \begin{pmatrix} \|v\| \\ 0 \end{pmatrix} \quad \text{where } c = \frac{v_1}{\|v\|}, s = \frac{v_2}{\|v\|}$$

By subsequently multiplying T with suitable block matrices each of them containing such a rotation and the identity elsewhere, one can cancel its element below the diagonal ending up with an upper triangular matrix which is exactly the R of the QR decomposition. Then Q will be simply the transpose of the product of the block diagonal matrices. This method has a complexity $\mathcal{O}(N^2)$ (N rotations and for each of them a product between block diagonal matrices that accounts for another N in this particular case) which is negligible when compared to those of matrix products. Note that since the tridiagonal form is preserved at each iteration of the QR algorithm, I could have simplified a bit also the matrix products but I decided not to do it in order to keep the code not too much complex.

The only part that is missing is how to find the proper unitary transformation that makes M tridiagonal. This job is done through the Lanczos algorithm. I followed precisely the procedure illustrated on Wikipedia⁹, so I will not report here the entire algorithm. The key points are that each iteration step requires the computation of a matrix-vector multiplication, of two scalar products, and some algebraic componentwise operations on vectors. The complexity of this algorithm for a generic hermitian matrix is $\mathcal{O}(N^3)$ so, considering that it has to be used only once unlike the matrix products in the QR algorithm, it does not affect the overall complexity that much.

Now that we have the full picture we can consider how the parallelization paradigm can be implemented. Unfortunately this particular algorithm is not very suitable for a parallel approach since both the Lanczos algorithm and the QR algorithm include that kind of loops where at each iteration all the previous data is needed. The best I could do was to implement in OpenCL the single linear algebra operation taking place inside these loops, such as matrix-matrix and matrix-vector multiplications. However in section 2.3 I will show that this is enough to achieve already a consistent speed-up, while a more in-depth discussion about the OpenCL implementation of the above described algorithm is carried on in the following section.

⁹https://en.wikipedia.org/wiki/Lanczos_algorithm

2.2 Code development

Firstly I implemented the algorithm described in section 2.1 in a simple single-thread C code. Initially I designed it for computing the complete eigendecomposition but unfortunately the instability of the QR algorithm did not allow me to get correct results. In particular for matrices larger than some hundreds of elements some eigenvalues were counted as double eigenvalues, therefore the algorithm missed much of them. Anyhow the lowest eigenvalue (and eigenvector) was the correct one most of the times so I slightly modified the code to keep only this one. In particular I set the convergence condition only on the lowest element of the diagonal of A_k at each iteration as:

$$\left| \lambda_{min}^{(k)} - \lambda_{min}^{(k-1)} \right| < 10^{-12}$$

Once the base algorithm was ready I converted it to run parallelized. Firstly I had to initialize OpenCL and compile the program. See the code below.

```

1 //opencl initialization
2
3 FILE *f;
4 int i,j;
5 char str[40];
6 char name[40];
7 char *source_str;
8 size_t source_size;
9 sprintf(str, "-cl-std=CL1.2 -D N=%d",N);
10
11 int dev=1;
12 cl_platform_id *platforms;
13 cl_uint num_platforms;
14 cl_device_id device;
15 cl_uint num_devices;
16 cl_context context;
17 cl_command_queue queue;
18 cl_program program;
19 size_t glob_size;
20 size_t glob_size2;
21 size_t loc_size=N;
22
23 clGetPlatformIDs(5,NULL,&num_platforms);
24 platforms=(cl_platform_id*) malloc(sizeof(cl_platform_id)*num_platforms
   );
25 clGetPlatformIDs(num_platforms,platforms,NULL);
26 //printf("Number of Platforms detected: %d\n",num_platforms);
27 for(i=0;i<num_platforms;i++){
28     clGetPlatformInfo(platforms[i],CL_PLATFORM_NAME,sizeof(name),&name,
29         NULL);
30     //printf("%s\n",name);
31 }
32
33 if(dev==1){
34     clGetDeviceIDs(platforms[0],CL_DEVICE_TYPE_GPU,1,&device,&
35         num_devices);
36     for(i=1;loc_size>256;i++){
37         if(loc_size%i==0)
38             loc_size=N/i;
39         //printf("\nAMD GPU selected\n");
40     }
}

```

```

41     else if(dev==2){
42         clGetDeviceIDs(platforms[1], CL_DEVICE_TYPE_CPU, 1, &device, &
43             num_devices);
44         for(i=1; loc_size>8192; i++){
45             if(loc_size%i==0)
46                 loc_size=N/i;
47             }
48             //printf("\nIntel CPU selected\n");
49         }
50
51         context=clCreateContext( NULL, 1,&device, NULL, NULL, NULL);
52         queue=clCreateCommandQueueWithProperties(context, device, NULL, NULL);
53
54         f=fopen("kernels.cl", "r");
55         source_str = (char*)malloc(MAX_SOURCE_SIZE);
56         source_size = fread( source_str, 1, MAX_SOURCE_SIZE, f);
57         fclose(f);
58         program=clCreateProgramWithSource(context, 1,(const char **)&source_str
59             ,(const size_t *)&source_size,NULL);
60         clBuildProgram(program,1,&device,str,NULL,NULL);
61         free(source_str);
62
63         cl_kernel lanczos1 = clCreateKernel(program, "lanczos1", NULL);
64         cl_kernel lanczos2 = clCreateKernel(program, "lanczos2", NULL);
65         cl_kernel transpose = clCreateKernel(program, "transpose", NULL);
66         cl_kernel mat_mul = clCreateKernel(program, "mat_mul", NULL);
67         cl_kernel transfer = clCreateKernel(program, "transfer", NULL);

```

The main aspects to notice are the following.

- At line 9 I specify the compilation flags; in particular I tell the compiler to use OpenCL 1.2, which is the latest supported by my hardware, and I pass the input dimension N as a parameter to the kernels.
- At lines 23-30 I get informations about the installed platforms.
- At lines 33-48 I get informations about the device running in a given platform and I set the dimension of the work group depending on its capabilities and on the input parameter N ; in particular I will use both the CPU and the GPU (separately) in the tests.
- At line 50 I create the OpenCL context for the selected device.
- At line 51 I create the command queue for the selected device.
- At lines 53-59 I read the source of the kernels from the file “kernels.cl” into a string and I build the program.
- At lines 61-65 I create the five kernels that I will use in the rest of the program.

The five kernels I have written are functions that replace some parts of the original single-thread C code and are executed in parallel on the selected device. In particular their jobs are the following:

- *lanczos1*: computes the matrix-vector multiplication and some componentwise operation on the vectors; N work items, each of them computing one component of the resulting vector;

- *lanczos2*: similar to the previous but it does another part of the algorithm and it does not compute any matrix-vector product;
- *transpose*: computes the transpose of a matrix; N^2 work items, one for each element of the resulting matrix;
- *mat_mul*: computes the product of two matrices; N^2 work items, one for each element of the resulting matrix;
- *transfer*: copies the diagonal and the subdiagonal of a matrix inside two arrays; N work items, one for each element of the diagonal.

Notice that I had to split the Lanczos algorithm in two kernels since they are separated by the computation of a scalar product which is difficult to parallelize. I tried to implement it but I ended up with numerical instability and I was not able to understand where it came from, thus I preferred to let the host compute these scalar products. Also in these two kernels I had to deal with complex numbers that are not built-in types in the OpenCL language, so I exploited a simple set of functions that I found on internet that defines the *cdouble* custom type simply as a vector *double2* and defines properly the operations on complex numbers. The *double2* type is a vector data type included in the kernel language that is basically a couple of double values; the advantage of such type is that the operations involving it should be furtherly optimized for the parallelization, even though I used it with a different aim.

Now as an example I will briefly explain the kernel *mat_mul* which is quite instructive. Below it is reported the code.

```

1  __kernel void mat_mul(__global double *A,__global double *B,__global
2    double *C)
3  {
4    int id=get_global_id(0);
5
6    int i,j,k;
7    double temp=0;
8    i=id/N;
9    j=id-i*N;
10   for(k=0;k<N;k++){
11     temp+=A[i*N+k]*B[k*N+j];
12   }
13   C[id]=temp;
14 }
```

As you can see from line 1 the kernel is defined pretty much as any standard C function. The important differences are that kernels are forced to return *void* and must include the specifier *_kernel* before the definition. This particular function has three pointers to double as arguments and they are preceded by the specifier *_global*, which indicates that they will be stored in the global memory of the device. The first two will contain the matrices to be multiplied and the third the result of the product. Then at line 3 each work item executing the kernel collects its own global id (which is a number between 0 and N^2) and stores it inside the variable *id*. At lines 7-8 it computes the matricial indices corresponding to its global id and then in the *for* cycle at lines 9-11 it computes the scalar product entering in the correct position of the resulting matrix. Finally this value is stored in the right position inside the memory addressed by the pointer *C*.

On the host side, in order to call a kernel one has to create *buffer* objects and copy inside them the data he wants to pass as kernel arguments. The copying operation has to be enqueued in the device command queue just like the call of the kernel. Then, after the kernel execution, data can be copied back from the buffer to some data structure of the host

by enqueueing another operation. See the code below for an example referred to the *mat_mul* kernel.

```

1   A_buff=clCreateBuffer(context,CL_MEM_READ_WRITE,N*N*sizeof(double),NULL
2   ,NULL);
3   B_buff=clCreateBuffer(context,CL_MEM_READ_WRITE,N*N*sizeof(double),NULL
4   ,NULL);
5   C_buff=clCreateBuffer(context,CL_MEM_READ_WRITE,N*N*sizeof(double),NULL
6   ,NULL);

7   //temp=A*Q
8   clEnqueueWriteBuffer(queue, A_buff, CL_TRUE, 0,N*N*sizeof(double), A
9   , 0, NULL, NULL);
10  clEnqueueWriteBuffer(queue, B_buff, CL_TRUE, 0,N*N*sizeof(double), Q
11  , 0, NULL, NULL);
12  clSetKernelArg(mat_mul, 0, sizeof(cl_mem), &A_buff);
13  clSetKernelArg(mat_mul, 1, sizeof(cl_mem), &B_buff);
14  clSetKernelArg(mat_mul, 2, sizeof(cl_mem), &C_buff);
15  clEnqueueNDRangeKernel(queue,mat_mul,1,NULL,&glob_size,NULL,0,NULL,
16  NULL);

```

Once they proved to work as expected, I put together in a single program the single-thread and the OpenCL codes in order to test their performances. The full host code can be found in appendix A.1 and the kernels in appendix A.2¹⁰. The former is designed to create some basic statistics out of a sample of execution times that is collected by running the diagonalization algorithms a certain number of times for each given N . Both N and the dimension of the sample are taken as input from the file “inp.txt”. To collect each element of the sample the code initializes a random hermitian matrix whose entries have a real part and an imaginary part uniformly distributed between -100 and 100 (function *init* - lines 146-178) and then it finds its lowest eigenvalue using four different methods: the single-thread code (function *stCPU* - lines 292-419), the OpenCL code running on the GPU (function *openclGPU* - lines 452-752), the OpenCL code running on the CPU (function *openclCPU* - lines 757-1056), and also the LAPACK routine *zheev* (function *lapack* - lines 425-446) for checking purposes. The execution times of each of these functions are collected using the function *clock* of the C time library but I keep in the sample only those corresponding to a correct eigenvalue determination by comparing them with the value found by LAPACK up to a precision of 10^{-7} . Then the program computes the mean and the standard deviation of the execution times sample and also the rate of correct results, and writes them in the file “out.txt” in a readable format (main of the code - lines 58-137).

I have to say that I did not deal much with error handling and pre/post conditions since the program was not intended to be rock solid for future applications but rather to test OpenCL performances, and its current behavior is enough for this purpose.

The code execution has been managed by a simple Python script that compiles the host code and launches it for each couple of input parameters (matrix dimension and cardinality of the statistical sample) that I have chosen.

2.3 Results

The program runned on my computer for a few days mapping the execution times and the rate of correct results for different values of N between 10 and 1000. As the input dimension raised I had to reduce the cardinality of the statistical sample otherwise the program would have taken too much time to finish the execution. For the same reason I also stopped the

¹⁰In the version reported here the eigenvector is not computed but in the end it is just a matter of multiplying two matrices.

		LAPACK		Single-thread			OpenCL - GPU			OpenCL - CPU		
N	#	t [s]	σ_t [s]	t [s]	σ_t [s]	η	t [s]	σ_t [s]	η	t [s]	σ_t [s]	η
10	100	0	0	0.002	0.005	99%	0.2	0.1	99%	0.08	0.01	99%
20	100	0.000	0.002	0.02	0.04	97%	0.3	0.3	97%	0.13	0.05	97%
30	100	0.000	0.003	0.2	0.7	93%	0.6	1.5	93%	0.2	0.4	93%
40	100	0.001	0.003	0.2	0.2	89%	0.4	0.3	89%	0.2	0.1	84%
50	100	0.002	0.005	0.4	0.3	67%	0.4	0.2	67%	0.2	0.1	76%
60	100	0.003	0.006	1.0	1.1	75%	0.6	0.6	80%	0.3	0.2	78%
70	100	0.004	0.007	1.8	1.2	82%	0.6	0.3	82%	0.4	0.3	87%
80	100	0.005	0.007	3.1	2.7	93%	0.7	0.5	93%	0.6	0.8	89%
90	100	0.006	0.008	7.2	21.4	93%	1.1	2.9	94%	0.9	0.8	91%
100	100	0.008	0.008	8.2	8.6	90%	1.0	0.8	90%	1.6	1.9	93%
120	50	0.013	0.007	13	9	90%	1.1	0.6	90%	2.3	1.9	90%
140	50	0.021	0.007	23	15	98%	1.5	0.7	94%	6.7	14.9	94%
160	50	0.030	0.005	48	31	98%	2.2	1.3	98%	7.2	5.1	96%
180	50	0.040	0.008	69	45	98%	2.5	1.4	98%	13	12	98%
200	50	0.053	0.008	91	67	98%	2.8	1.8	98%	17	14	98%
220	50	0.073	0.007	174	119	100%	3.5	2.2	100%	23	22	98%
240	50	0.092	0.008	252	483	100%	5.2	3.1	100%	30	21	100%
260	50	0.117	0.008	282	208	100%	6.9	4.6	100%	50	32	100%
280	50	0.143	0.007	347	243	100%	8.4	5.7	100%	57	48	100%
300	50	0.170	0.006	567	414	100%	11	7	100%	75	49	100%
320	50	0.206	0.008	-	-	-	14	14	100%	106	68	100%
340	50	0.245	0.009	-	-	-	14	10	100%	113	76	96%
360	50	0.290	0.009	-	-	-	18	11	98%	176	177	100%
380	50	0.332	0.008	-	-	-	25	23	100%	160	114	98%
400	50	0.39	0.01	-	-	-	32	22	100%	239	205	100%
450	20	0.55	0.01	-	-	-	41	30	100%	430	290	100%
500	20	0.75	0.02	-	-	-	60	42	100%	615	408	100%
550	20	0.99	0.02	-	-	-	66	59	100%	703	751	100%
600	20	1.27	0.02	-	-	-	171	114	95%	1034	856	100%
700	10	2.01	0.04	-	-	-	249	176	100%	1411	1245	100%
800	10	2.98	0.03	-	-	-	303	307	100%	2840	2333	100%
900	10	4.1	0.3	-	-	-	792	645	100%	10166	6108	100%
1000	10	5.6	0.3	-	-	-	774	656	100%	5431	5378	100%

Table 2: Results of the diagonalization algorithm performance test.

single-thread code at $N = 300$. The results of this analysis are reported in table 2 (# is the dimension of the statistical sample and η is the rate i.e. the number of correct eigenvalues obtained divided by #). The plot of the execution times vs input dimension is report is reported in figure 3 (the y-axis is in logarithmic scale).

There are some quite interesting observation I can point out by looking at this data. Firstly we can see how LAPACK is always much faster than my algorithm, but this is a behavior I was expecting as I have said in the previous section. Yet we can clearly see a remarkable boost in the performance both when passing from a single-thread to a parallelized algorithm and when passing from the CPU to the GPU implementation of OpenCL. This is undoubtedly a success! Apart from the lowest dimensions (approximately up to $N = 100$) the GPU is about 25-50x faster than the single thread code and about 5-10x faster than the CPU running OpenCL. For what concerns the performances at dimensions lower than $N = 100$ we can see that the situation is less clear, with the single-thread algorithm that is the fastest up to $N = 40$ while OpenCL running on the CPU is the fastest from $N = 50$ to $N = 90$. I think that this is mainly due to the fact that at low dimensions the data transfer from the host memory to the devices memory becomes a bottleneck for the program. Indeed this is not a particularly fast operation and thus when the actual computational time becomes very

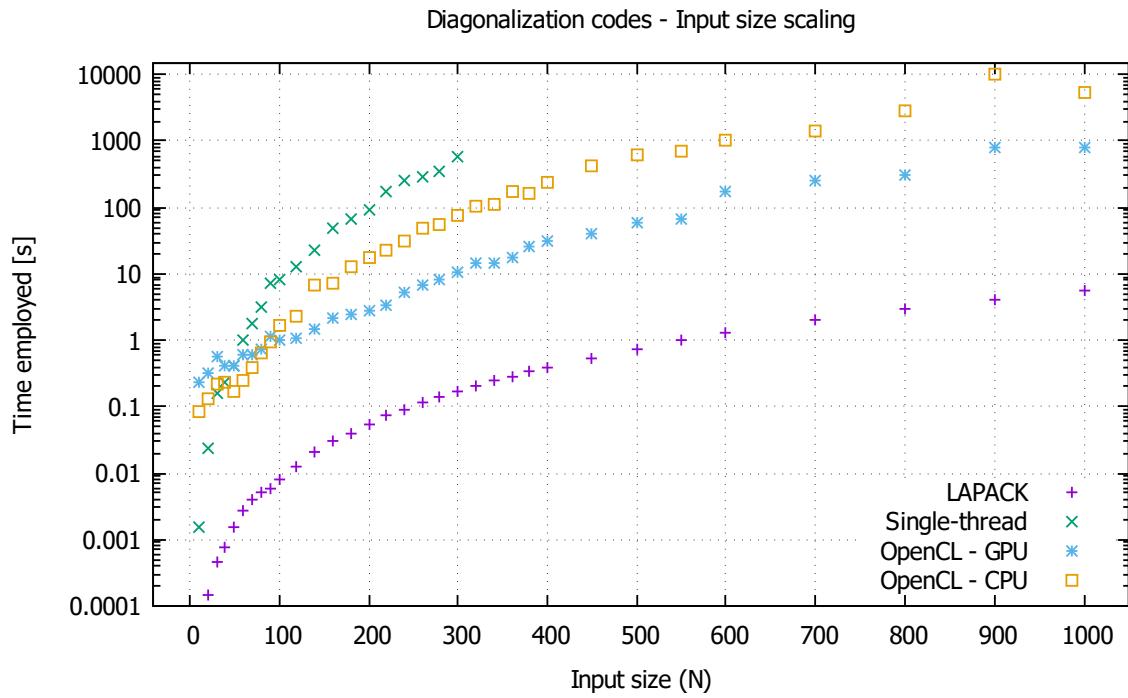


Figure 3: Scaling of the four codes used for solving the eigenproblem with respect to the size of the input matrix.

short this contribution to the total time becomes significant.

Next if one looks at the standard deviations he certainly will notice that in most cases it is very big, almost always of the same order of magnitude of the mean itself. This means that the collected times are spread in a wide range. A possible explanation of this great variability could be the algorithm itself, in particular the fact that the convergence speed depends highly on the initial matrix and thus the number of iterations required by the QR algorithm to get to the result is rather casual.

The last interesting comment I can do is about the rate of success. The instability of the QR algorithm lead to a certain number of wrong eigenvalues, but we can see how this issue becomes less and less important as the input dimension N increases. For $N > 200$ we get almost 100% of the times the correct eigenvalue, which is not bad in this context. Instead we cannot trust the program much for lower N and this is a further reason to prefer the parallelized approach.

Finally I performed a more quantitative analysis of the scaling of these algorithms. I linearised the data in the former plot and I fitted them with a line. The results are reported in figure 4. From the plot one can see even more clearly the anomalous behavior of the execution times for low values of N when using OpenCL. For this reason I excluded these data from the linear fit (the fitted points are those darker in the plot). The obtained scaling are:

$$\text{LAPACK} \sim \mathcal{O}(N^{2.69})$$

$$\text{Single-thread} \sim \mathcal{O}(N^{3.70})$$

$$\text{OpenCL on GPU} \sim \mathcal{O}(N^{3.62})$$

$$\text{OpenCL on CPU} \sim \mathcal{O}(N^{3.69})$$

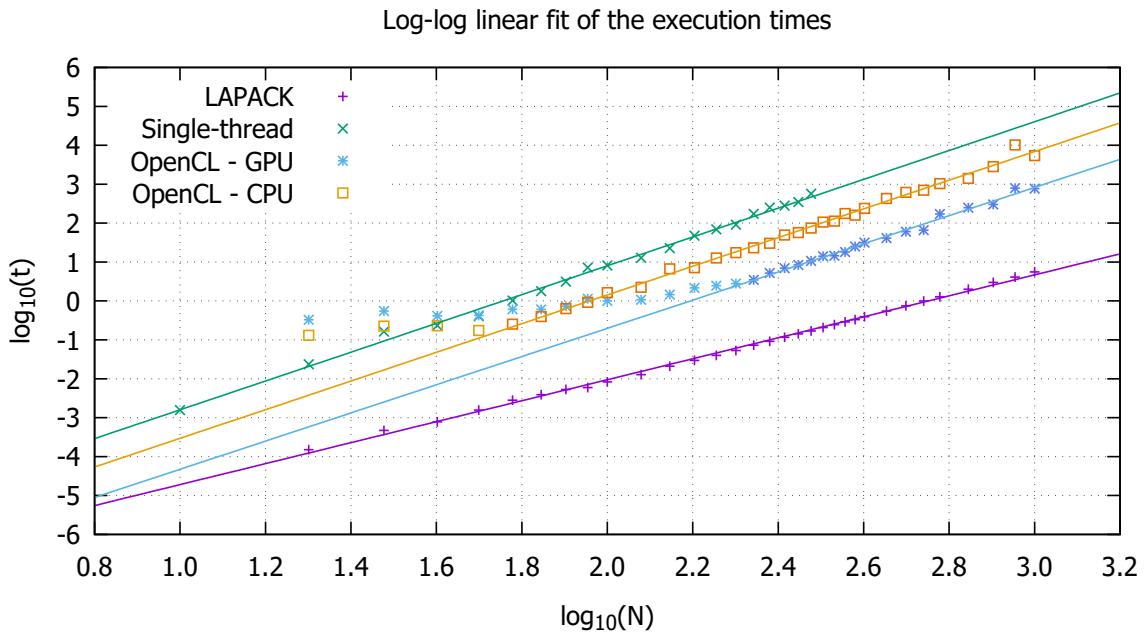


Figure 4: Log-log plot of the executions times vs input sizes and linear fit.

Despite the expectations there is not a great difference in the scaling of the three codes I realised (while LAPACK is much faster). This can be seen also on the plot: the three fitting lines corresponding to my algorithms are pretty much parallel unlike that corresponding to LAPACK. What is really different is the y-intercept which is the parameter related to the coefficient in front of the power of N in the scaling. So apparently the parallelization is not able to improve significantly the scaling of an algorithm but only to speedup its execution affecting only the multiplicative coefficient. Nevertheless more analysis and data would be required to state anything conclusive, also because the parallelization algorithm has still optimization possibilities.

3 2D Schrödinger equation

After having learned the basic working principles of OpenCL, I exploited its potentialities to solve the 2D time-independent Schrödinger equation for different potentials (in particular I found the ground state and first excited state).

3.1 Recall of FDM

I exploited the FDM¹¹ as already done in exercise 6, but this time I used its 2D version. They are quite similar, the only difference is that now the matrix to be diagonalized is no more tridiagonal as in 1D. For simplicity I consider a normalized Schrödinger equation (i.e. $m = 1/2$, $\hbar = 1$):

$$[-\nabla^2 + V(x, y)] \psi(x, y) = E \psi(x, y)$$

After an homogeneous discretization of the domain into intervals of width Δ (the same for both the x and the y coordinates), exploiting the usual approximation for the second

¹¹Finite Difference Method.

derivative one gets the expression:

$$-\frac{1}{\Delta^2} (\psi_{n+1,m} + \psi_{n-1,m} + \psi_{n,m+1} + \psi_{n,m-1}) + \left(V_{n,m} + \frac{4}{\Delta^2} \right) \psi_{n,m} = E \psi_{n,m}$$

where $\psi_{n,m} \equiv \psi(x_n, y_m)$, with (x_n, y_m) being a point of the lattice, and analogously for $V_{n,m}$. This equation is not in the form we desire yet (i.e. that of an eigenvalue problem) since we have too many indices. Being N the number of lattice points in the x coordinate we apply the following redefinitions:

$$\psi_{n,m} \longrightarrow \psi_{n+mN} \quad V_{n,m} \longrightarrow V_{n+mN}$$

Then the above equation can be rewritten as:

$$-\frac{1}{\Delta^2} (\psi_{n+1+mN} + \psi_{n-1+mN} + \psi_{n+(m+1)N} + \psi_{n+(m-1)N}) + \left(V_{n+mN} + \frac{4}{\Delta^2} \right) \psi_{n+mN} = E \psi_{n+mN}$$

So at the price of having increased the dimensionality of the problem, we reached the desired expression. The eigenvalues of this operator are the energies of the eigenstates and the eigenvectors contain the values of the corresponding wavefunction evaluated in the points of the lattice. In the following we will assume that $n = 1, \dots, N$ and $m = 1, \dots, N$ that is the domain is symmetric in the two coordinates.

The interesting thing is that it is much more difficult to solve analytically a 2D problem than a 1D problem, whereas the numerical solution is essentially the same, only with bigger matrices. Therefore with enough time/computational resources we will always be able to find the solution. Actually this is true for a generic number of dimensions: the only limits are the physical means.

For what follows it is also important to recall the following fact valid in any dimension: though the wavefunction is complex in general one can always choose the eigenfunctions of a given system to be real¹². Indeed if $\psi(x)$ is a solution of the time-independent Schrödinger equation, then also $\psi^*(x)$ is a solution of the same equation; now since the Schrödinger equation is linear, the function $\bar{\psi}(x) \equiv \psi(x) + \psi^*(x)$ is still a solution and it is real by construction.

In particular the potentials I will consider in the following are:

(A) Decoupled harmonic oscillator: $V(x, y) = \frac{1}{2} (x^2 + y^2)$

(B) Regularized electrostatic potential¹³: $V(x, y) = \begin{cases} -\frac{3}{\sqrt{x^2+y^2}} & \text{for } x^2 - y^2 > \frac{9}{100} \\ -10 & \text{for } x^2 - y^2 \leq \frac{9}{100} \end{cases}$

(C) Central mexican hat: $V(x, y) = -(x^2 + y^2) + \frac{1}{8} (x^2 + y^2)^2$

(D) Double mexican hat: $V(x, y) = -x^2 + \frac{1}{8} x^4 - y^2 + \frac{1}{8} y^4$

¹²This is true if the potential $V(x)$ is real which is the case for any physical interesting system.

¹³I had to regularize the Coulomb potential since otherwise its divergence at $x = y = 0$ would have broken the numerical solution.

Clearly the potential (A) has an analytical solution, and indeed I used it for checking the correctness of the numerical results. The energy eigenvalues are given by:

$$E_{\vec{n}} = \hbar\omega(n_1 + n_2 + 1)$$

where $n_1, n_2 \in \mathbb{N}$. With the above defined normalization $\hbar\omega = \sqrt{2}$ so that the ground state and the first excited level have energies:

$$E_0 = \sqrt{2} \approx 1.141214 \quad E_1 = 2\sqrt{2} \approx 2.828427 \quad (1)$$

Also note that potentials (B) and (C) could be simplified to a 1D problem exploiting the central symmetry, yet the aim of this program was to solve general 2D problems therefore I kept them 2D.

3.2 Code development

By looking empirically at the potentials, I set the domain of the algorithm to be $[-4, 4] \times [-4, 4]$ and I chose the lattice to have 401 points for each axis, which results in a total of 160801 points separated by a distance of $\Delta = 0.02$ in each coordinate.

Considering that for this application I had to deal with matrices much bigger than 1000×1000 (they have about $160801 \times 160801 \approx 2.6 \cdot 10^{10}$ elements) and that are sparse instead of dense, I could not use the program described in section 2.2. I searched on internet for some OpenCL library including functions for computing the lowest eigenvalues of a sparse matrix. The only one that I was able to find is the C++ library *ViennaCL*. In particular it contains also an example code that uses Lanczos algorithm and then the bisection method to find the greater k eigenpairs (an approximation of them actually), where k is specified as an input parameter. So once computed the matrix to diagonalize I simply passed its opposite to the ViennaCL function, since I was interested in the lowest eigepairs. The input format required for the matrix to be diagonalized is the Matrix Market exchange format (extension *.mtx*) which is essentially a text file that contains at each row the coordinates and the value of an element of the matrix. An element is assumed to be 0 if not specified in the file. It is clear that this format is very convenient when it is necessary to store a sparse matrix since one can avoid to write a bunch of useless zeros. Another important parameter of this function is the size of the Krylov subspace used in the computations. The behavior of the output when varying this parameter was rather strange, thus I had to use a bit of fine-tuning to find the optimal value, comparing every time the result for the potential (A) with its analytical solution. Eventually I set it to 400 for the computation of the ground state and to 900 for that of the first excited state.

Given that, the number of points used to partition the domain was the highest that allowed my computer to find the first two eigenstates of the quantum systems. In fact, despite the code was still very fast (less than a minute for computing the second lowest eigenpair), a higher number of points resulted in crashes of the code due to VRAM saturation. This limitation is also what prevented me from solving also some 3D problem.

The resulting eigenfunctions are written in a suitable form into the file “psi.txt”. The full code I used is reported in appendix B.

3.3 Results

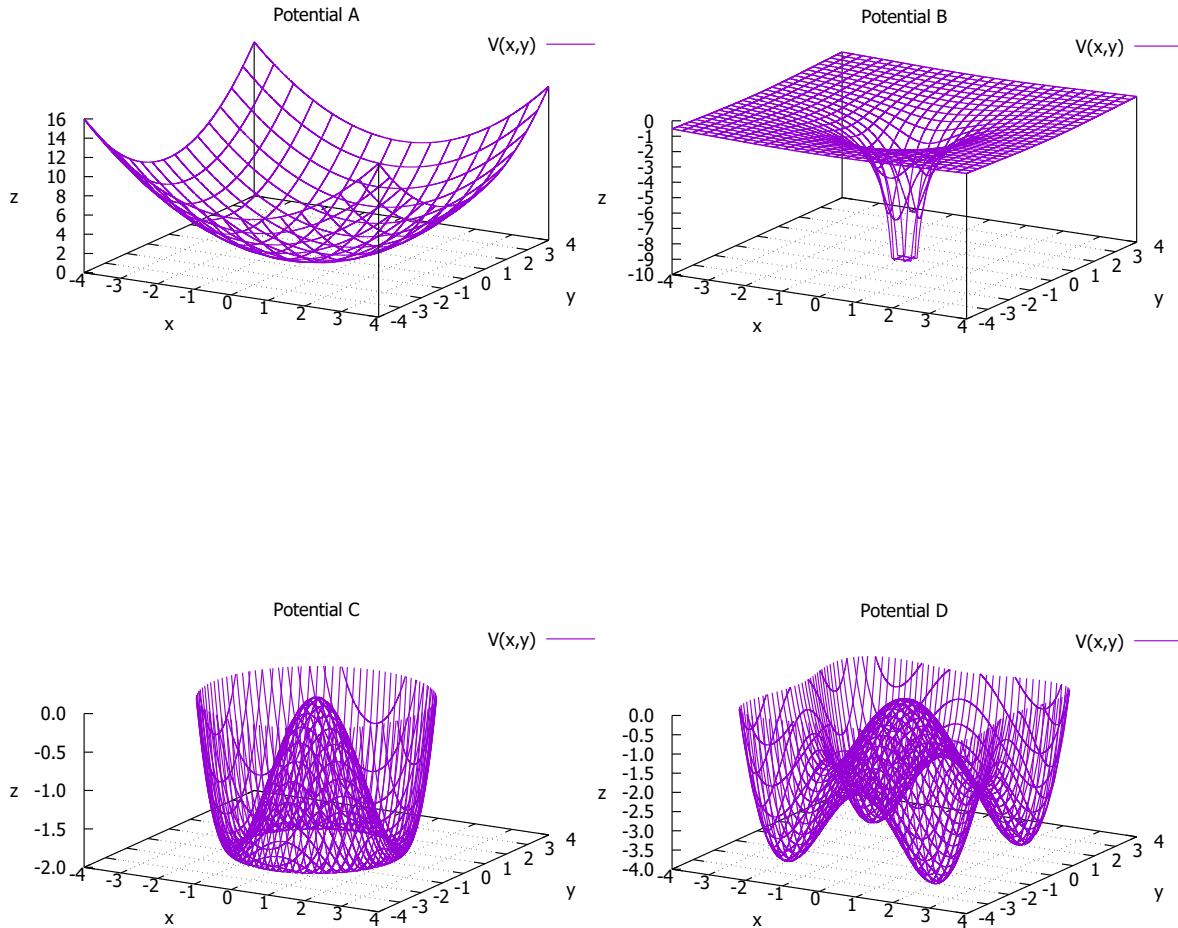
The resulting eigenvalues of the four above reported potentials can be found in table 3. Also the figures below contain the plots of the four potentials and those of the modulus squared of the eigenfunctions found¹⁴. A positive aspect is that they are quite well defined, thus it

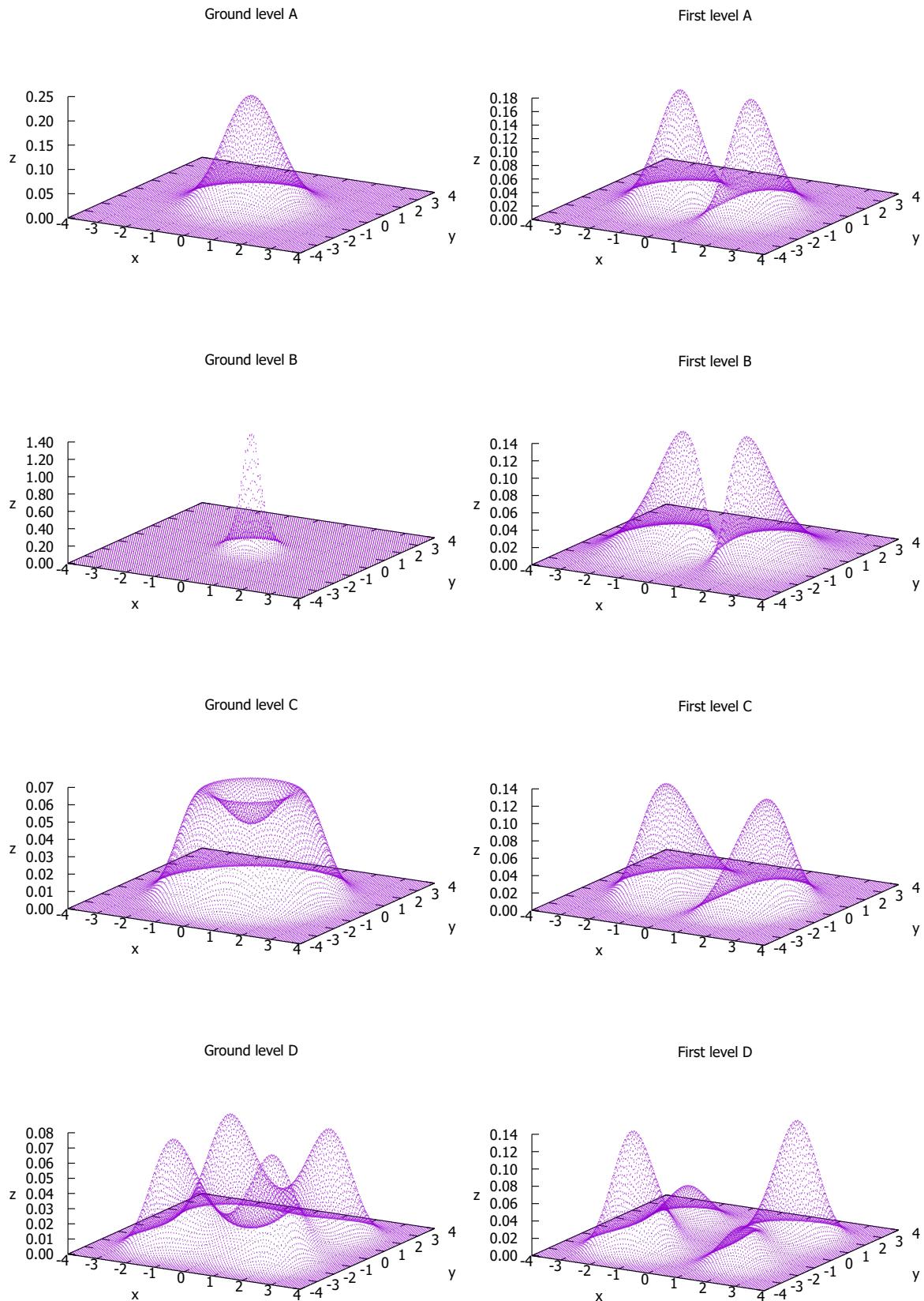
¹⁴I have plotted only a subset of the total 160801 points since reducing the size of the images was a typographic need.

Potential	E_0	E_1
A	1.413258	2.829293
B	-5.178028	-0.944380
C	-0.837434	-0.280717
D	-1.710583	-1.477471

Table 3: First two energy eigenvalues for the potentials above reported.

seems that the code has at least a good precision. By considering the numerical estimates of the energy eigenvalues for the potential (A) and comparing them with (1) we can state that the accuracy is also quite good, even though not excellent: both the energies differs from the exact ones by about 10^{-4} . However another consistency check comes from the fact that all the eigenvalues have the correct sign (plus for the harmonic oscillator since there $V(x, y) \geq 0$ always, minus for the other three since it is not true anymore). Moreover all the ground state found have a shape that matches with that one could have expected by looking at the potentials. So I conclude that the code works quite well.





4 Conclusion and self-evaluation

The main result of the present work is that I have learned the ideas behind the parallelization techniques in computational science and the basic working principles of OpenCL, which is a very powerful tool that provides a way to implement such techniques in almost any mainstream computer. Now I am able to code basic applications implementing the OpenCL framework and I have gained the capabilities of studying more advanced topics.

OpenCL itself has proven to yield a consistent gain in the performances of an algorithm even if implemented without spending a great effort in the optimization. Therefore it can certainly constitute a powerful ally to save time for heavy computations.

The two application that I realised were quite successful in their intents even though they did not reach perfect results. In particular the diagonalization algorithm hid much more difficulties than I thought and I had to devote the majority of the time employed in fixing issues that occurred while writing it. Especially the numerical instability bothered me for a long time and eventually I had to renounce to some further passages that I would have liked to implement in a parallelized way. On the other hand the application to the solution of the 2D time-independent Schrödinger equation was pretty straightforward thanks to the functions provided by the library *ViennaCL*. Here I learned how to generalize the FDM to higher dimensions and also how to write huge sparse matrices in a suitable form, that is with the Matrix Market exchange format. Also in this application OpenCL proved to have remarkable performances, finding the greatest eigenpairs of a matrix with some billions elements in less than a minute.

Overall I am quite satisfied with the work done and I am glad to have opened some doors for future studying possibilities.

A Matrix diagonalization code

A.1 Host code

```

1 ///////////////////////////////////////////////////////////////////
2 //
3 // Evaluates the capabilities of the hardware by computing some statistics
4 // about
5 // the computation of the lowest eigenvalues of a matrix in four different
6 // ways:
7 // 1) LAPACK library
8 // 2) Lanczos + QR algorithm - pure single-thread C
9 // 3) Lanczos + QR algorithm - OpenCL over GPU
10 // 4) Lanczos + QR algorithm - OpenCL over CPU
11 //
12 // libraries
13
14 #include <stdio.h>
15 #include <stdlib.h>
16 #include <math.h>
17 #include <complex.h>
18 #include <time.h>
19 #include <string.h>
20 #include <lapacke.h>
21
22 #ifdef __APPLE__
23 #include <OpenCL/opencl.h>
24 #else
25 #include <CL/cl.h>
26 #endif
27
28
29 #define MAX_SOURCE_SIZE (0x100000)
30
31
32 int N;    //dimension of the hermitian matrix
33 int it;   //number of iterations
34
35
36
37
38 //PROTOTYPES
39
40 void init(double complex *cmat,double *mat);
41 void prtcmpx(double complex z);
42 double complex dot(double complex *u, double complex *v);
43 double norm(double complex *u);
44 double complex* matvec(double complex *M, double complex *v);
45 double* QR(double *diag,double *subdiag);
46 double* matmul(double *A, double *B);
47 double* transp(double *A);
48
49 void stCPU(double *time ,double *eig0 ,double complex *M);
50 void openclGPU(double *time ,double *eig0 ,double *M);
51 void openclCPU(double *time ,double *eig0 ,double *M);
52 void lapack(double *time ,double *eig0 ,double complex *M);
53
54
55

```

```

56
57
58 //MAIN
59
60 int main(){
61
62     FILE *f;
63     f=fopen("inp.txt","r");
64     fscanf(f,"%d",&N);
65     fscanf(f,"%d",&it);
66     fclose(f);
67
68     int i,j,nhit[3];
69     double time[4],eig[4],mean[4],temp[4][it],dev[4];
70     double *M;
71     double complex *cM;
72
73     printf("DIMENSION: %d\n",N);
74     printf("NUMBER OF ITERATIONS: %d\n\n",it);
75
76     mean[0]=0; dev[0]=0; nhit[0]=0;
77     mean[1]=0; dev[1]=0; nhit[1]=0;
78     mean[2]=0; dev[2]=0; nhit[2]=0;
79     mean[3]=0; dev[3]=0;
80
81     for(i=0;i<it;i++){
82
83         printf("\nIteration %d\n\n",i+1);
84
85         M=(double*)malloc(sizeof(double)*2*N*N);
86         cM=malloc(sizeof(double complex)*N*N);
87         init(cM,M);
88
89         if(N<=300)
90             stCPU(&time[0],&eig[0],cM);
91         lapack(&time[3],&eig[3],cM);
92         openclGPU(&time[1],&eig[1],M);
93         openclCPU(&time[2],&eig[2],M);
94
95         free(M);
96         free(cM);
97
98         mean[3]+=time[3];
99         temp[3][i]=time[3];
100        for(j=0;j<3;j++){
101            if(fabs(eig[j]-eig[3])<1e-7){
102                mean[j]+=time[j];
103                temp[j][nhit[j]]=time[j];
104                nhit[j]++;
105            }
106        }
107    }
108
109    mean[0]=(double)mean[0]/nhit[0];
110    mean[1]=(double)mean[1]/nhit[1];
111    mean[2]=(double)mean[2]/nhit[2];
112    mean[3]=(double)mean[3]/it;
113
114    for(i=0;i<nhit[0];i++){
115        dev[0]+=(temp[0][i]-mean[0])*(temp[0][i]-mean[0]);
116    }

```

```

117     for(i=0;i<nhit[1];i++){
118         dev[1]+=(temp[1][i]-mean[1])*(temp[1][i]-mean[1]);
119     }
120     for(i=0;i<nhit[2];i++){
121         dev[2]+=(temp[2][i]-mean[2])*(temp[2][i]-mean[2]);
122     }
123     for(i=0;i<it;i++){
124         dev[3]+=(temp[3][i]-mean[3])*(temp[3][i]-mean[3]);
125     }
126     dev[0]=sqrt((double)dev[0]/(nhit[0]-1));
127     dev[1]=sqrt((double)dev[1]/(nhit[1]-1));
128     dev[2]=sqrt((double)dev[2]/(nhit[2]-1));
129     dev[3]=sqrt((double)dev[3]/(it-1));
130
131     f=fopen("out.txt","a");
132     fprintf(f,"%d %d %f %f %f %f %f %f %f %f %f %f\n",N,it,mean[3],dev
133     [3],mean[0],dev[0],
134     (double)nhit[0]/it,mean[1],dev[1],(double)nhit[1]/it,mean[2],dev
135     [2],(double)nhit[2]/it);
136     fclose(f);
137
138
139
140
141
142
143 //FUNCTIONS
144
145
146 //random hermitian matrix generation
147
148 void init(double complex *cmat,double *mat){
149     double x,y;
150     int i,j;
151     srand((unsigned int)time(NULL));
152     for(i=0;i<N;i++){
153         for(j=i;j<N;j++){
154             if(i==j){
155                 x = ((double)rand()/(double)(RAND_MAX)) * 200-100;
156                 y=0;
157                 cmat[i*N+j]=x+y*I;
158                 mat[i*2*N+2*j]=x;
159                 mat[i*2*N+2*j+1]=y;
160             }
161             else{
162                 x = ((double)rand()/(double)(RAND_MAX)) * 200-100;
163                 y = ((double)rand()/(double)(RAND_MAX)) * 200-100;
164                 cmat[i*N+j]=x+y*I;
165                 cmat[j*N+i]=x-y*I;
166                 mat[i*2*N+2*j]=x;
167                 mat[i*2*N+2*j+1]=y;
168                 mat[j*2*N+2*i]=x;
169                 mat[j*2*N+2*i+1]=-y;
170             }
171         }
172     }
173     return;
174 }
175

```

```

176
177 //print complex number
178
179 void prtcmpx(double complex z){
180     printf("(%.f,%.f)\n",creal(z),cimag(z));
181     return;
182 }
183
184
185 //dot product
186
187 double complex dot(double complex *u, double complex *v){
188     int i;
189     double complex z=0+0*I;
190     for(i=0;i<N;i++)
191         z+=conj(u[i])*v[i];
192     return z;
193 }
194
195
196 //euclidean norm
197
198 double norm(double complex *u){
199     return sqrt(cabs(dot(u,u)));
200 }
201
202
203 //matrix-vector multiplication
204
205 double complex* matvec(double complex *M, double complex *v){
206     int i,j;
207     double complex *z;
208     z=malloc(sizeof(double complex)*N);
209     for(i=0;i<N;i++){
210         z[i]=0+0*I;
211         for(j=0;j<N;j++)
212             z[i]+=M[i*N+j]*v[j];
213     }
214     return z;
215 }
216
217
218 //QR decomposition
219
220 double* QR(double *diag,double *subdiag){
221     int i,j;
222     double norm,gam,sig,temp1,temp2,temp3;
223     double *Q;
224
225     Q=(double*)calloc(N*N,sizeof(double));
226     norm=sqrt(diag[0]*diag[0]+subdiag[0]*subdiag[0]);
227     gam=diag[0]/norm;
228     sig=subdiag[0]/norm;
229     Q[0]=gam;
230     Q[1]=-sig;
231     Q[N]=sig;
232     Q[N+1]=gam;
233     temp1=-sig*subdiag[0]+gam*diag[1];
234     temp2=gam*subdiag[1];
235
236

```

```

237     for(i=1;i<N-1;i++){
238         norm=sqrt(temp1*temp1+subdiag[i]*subdiag[i]);
239         gam=temp1/norm;
240         sig=subdiag[i]/norm;
241         for(j=0;j<i+1;j++){
242             temp3=Q[j*N+i];
243             Q[j*N+i]=temp3*gam;
244             Q[j*N+i+1]=-temp3*sig;
245         }
246         Q[(i+1)*N+i]=sig;
247         Q[(i+1)*N+i+1]=gam;
248         temp1=-sig*temp2+gam*diag[i+1];
249         temp2=gam*subdiag[i+1];
250     }
251     return Q;
252 }
253
254
255 //matrix multiplication
256
257 double* matmul(double *A, double *B){
258     int i,j,k;
259     double temp;
260     double *M;
261     M=malloc(sizeof(double)*N*N);
262     for(i=0;i<N;i++){
263         for(j=0;j<N;j++){
264             temp=0;
265             for(k=0;k<N;k++){
266                 temp+=A[i*N+k]*B[k*N+j];
267             }
268             M[i*N+j]=temp;
269         }
270     }
271     return M;
272 }
273
274
275 //matrix transpose
276
277 double* transp(double *A){
278     int i,j;
279     double *M;
280     M=malloc(sizeof(double)*N*N);
281     for(i=0;i<N;i++)
282         for(j=0;j<N;j++)
283             M[j*N+i]=A[i*N+j];
284     return M;
285 }
286
287
288
289
290
291
292 //SINGLE-THREAD CPU
293
294
295 void stCPU(double *time ,double *eig0 ,double complex *M){
296     printf("Single-thread CPU algorithm started!\n");

```

```

298     clock_t start, end;
299     start = clock();
300
301
302     //lanczos algorithm
303
304     int i,j;
305     double a,b;
306     double complex *P,*w,*v;
307     double *diag,*subdiag;
308
309     v=malloc(sizeof(double complex)*N);
310     P=malloc(sizeof(double complex)*N*N);
311     diag=malloc(sizeof(double)*N);
312     subdiag=malloc(sizeof(double)*(N-1));
313
314     P[0]=1+0*I;
315     v[0]=P[0];
316     for(i=1;i<N;i++){
317         P[i*N]=0+0*I;
318         v[i]=P[i*N];
319     }
320     w=matvec(M,v);
321     a=creal(dot(w,v));
322     diag[0]=a;
323     for(i=0;i<N;i++)
324         w[i]=(a+0*I)*v[i];
325
326     for(i=1;i<N;i++){
327         b=norm(w);
328         subdiag[i-1]=b;
329         for(j=0;j<N;j++){
330             v[j]=w[j]/(b+0*I);
331             P[j*N+i]=v[j];
332         }
333         free(w);
334         w=matvec(M,v);
335         a=dot(w,v);
336         diag[i]=a;
337         for(j=0;j<N;j++)
338             w[j]=w[j]-(a+0*I)*v[j]-(b+0*I)*P[j*N+i-1];
339     }
340     free(w);
341     free(v);
342     free(P);
343
344
345     //QR algorithm
346
347     int chk=0;
348     double mu,conf,conf2=0;
349     double *Q,*A;
350     mu=diag[N-1];
351     A=malloc(sizeof(double)*N*N);
352     for(i=0;i<N;i++){
353         for(j=0;j<N;j++){
354             if(i==j){
355                 A[i*N+j]=diag[i];
356             }
357             else if(i==j+1){
358                 A[i*N+j]=subdiag[i-1];
359             }
360         }
361     }
362
363     //orthogonalization
364     for(i=1;i<N;i++){
365         for(j=0;j<N;j++){
366             if(i>j)
367                 Q[i*N+j]=0;
368             else if(i==j)
369                 Q[i*N+j]=1;
370             else
371                 Q[i*N+j]=-subdiag[i-1];
372         }
373     }
374
375     //QR decomposition
376     for(i=0;i<N;i++){
377         for(j=i;j<N;j++){
378             if(i>j)
379                 A[i*N+j]=0;
380             else if(i==j)
381                 A[i*N+j]=1;
382             else
383                 A[i*N+j]=-subdiag[i-1];
384         }
385     }
386
387     //free memory
388     free(diag);
389     free(subdiag);
390     free(Q);
391     free(A);
392
393
394     //print results
395     for(i=0;i<N;i++){
396         for(j=0;j<N;j++)
397             printf("%f ",A[i*N+j]);
398         printf("\n");
399     }
400
401
402     //free memory
403     free(w);
404     free(v);
405     free(P);
406
407
408     //print results
409     for(i=0;i<N;i++){
410         for(j=0;j<N;j++)
411             printf("%f ",Q[i*N+j]);
412         printf("\n");
413     }
414
415
416     //free memory
417     free(w);
418     free(v);
419     free(P);
420
421
422     //print results
423     for(i=0;i<N;i++){
424         for(j=0;j<N;j++)
425             printf("%f ",A[i*N+j]);
426         printf("\n");
427     }
428
429
430     //free memory
431     free(w);
432     free(v);
433     free(P);
434
435
436     //print results
437     for(i=0;i<N;i++){
438         for(j=0;j<N;j++)
439             printf("%f ",Q[i*N+j]);
440         printf("\n");
441     }
442
443
444     //free memory
445     free(w);
446     free(v);
447     free(P);
448
449
450     //print results
451     for(i=0;i<N;i++){
452         for(j=0;j<N;j++)
453             printf("%f ",A[i*N+j]);
454         printf("\n");
455     }
456
457
458     //free memory
459     free(w);
460     free(v);
461     free(P);
462
463
464     //print results
465     for(i=0;i<N;i++){
466         for(j=0;j<N;j++)
467             printf("%f ",Q[i*N+j]);
468         printf("\n");
469     }
470
471
472     //free memory
473     free(w);
474     free(v);
475     free(P);
476
477
478     //print results
479     for(i=0;i<N;i++){
480         for(j=0;j<N;j++)
481             printf("%f ",A[i*N+j]);
482         printf("\n");
483     }
484
485
486     //free memory
487     free(w);
488     free(v);
489     free(P);
490
491
492     //print results
493     for(i=0;i<N;i++){
494         for(j=0;j<N;j++)
495             printf("%f ",Q[i*N+j]);
496         printf("\n");
497     }
498
499
500     //free memory
501     free(w);
502     free(v);
503     free(P);
504
505
506     //print results
507     for(i=0;i<N;i++){
508         for(j=0;j<N;j++)
509             printf("%f ",A[i*N+j]);
510         printf("\n");
511     }
512
513
514     //free memory
515     free(w);
516     free(v);
517     free(P);
518
519
520     //print results
521     for(i=0;i<N;i++){
522         for(j=0;j<N;j++)
523             printf("%f ",Q[i*N+j]);
524         printf("\n");
525     }
526
527
528     //free memory
529     free(w);
530     free(v);
531     free(P);
532
533
534     //print results
535     for(i=0;i<N;i++){
536         for(j=0;j<N;j++)
537             printf("%f ",A[i*N+j]);
538         printf("\n");
539     }
540
541
542     //free memory
543     free(w);
544     free(v);
545     free(P);
546
547
548     //print results
549     for(i=0;i<N;i++){
550         for(j=0;j<N;j++)
551             printf("%f ",Q[i*N+j]);
552         printf("\n");
553     }
554
555
556     //free memory
557     free(w);
558     free(v);
559     free(P);
560
561
562     //print results
563     for(i=0;i<N;i++){
564         for(j=0;j<N;j++)
565             printf("%f ",A[i*N+j]);
566         printf("\n");
567     }
568
569
570     //free memory
571     free(w);
572     free(v);
573     free(P);
574
575
576     //print results
577     for(i=0;i<N;i++){
578         for(j=0;j<N;j++)
579             printf("%f ",Q[i*N+j]);
580         printf("\n");
581     }
582
583
584     //free memory
585     free(w);
586     free(v);
587     free(P);
588
589
590     //print results
591     for(i=0;i<N;i++){
592         for(j=0;j<N;j++)
593             printf("%f ",A[i*N+j]);
594         printf("\n");
595     }
596
597
598     //free memory
599     free(w);
600     free(v);
601     free(P);
602
603
604     //print results
605     for(i=0;i<N;i++){
606         for(j=0;j<N;j++)
607             printf("%f ",Q[i*N+j]);
608         printf("\n");
609     }
610
611
612     //free memory
613     free(w);
614     free(v);
615     free(P);
616
617
618     //print results
619     for(i=0;i<N;i++){
620         for(j=0;j<N;j++)
621             printf("%f ",A[i*N+j]);
622         printf("\n");
623     }
624
625
626     //free memory
627     free(w);
628     free(v);
629     free(P);
630
631
632     //print results
633     for(i=0;i<N;i++){
634         for(j=0;j<N;j++)
635             printf("%f ",Q[i*N+j]);
636         printf("\n");
637     }
638
639
640     //free memory
641     free(w);
642     free(v);
643     free(P);
644
645
646     //print results
647     for(i=0;i<N;i++){
648         for(j=0;j<N;j++)
649             printf("%f ",A[i*N+j]);
650         printf("\n");
651     }
652
653
654     //free memory
655     free(w);
656     free(v);
657     free(P);
658
659
660     //print results
661     for(i=0;i<N;i++){
662         for(j=0;j<N;j++)
663             printf("%f ",Q[i*N+j]);
664         printf("\n");
665     }
666
667
668     //free memory
669     free(w);
670     free(v);
671     free(P);
672
673
674     //print results
675     for(i=0;i<N;i++){
676         for(j=0;j<N;j++)
677             printf("%f ",A[i*N+j]);
678         printf("\n");
679     }
680
681
682     //free memory
683     free(w);
684     free(v);
685     free(P);
686
687
688     //print results
689     for(i=0;i<N;i++){
690         for(j=0;j<N;j++)
691             printf("%f ",Q[i*N+j]);
692         printf("\n");
693     }
694
695
696     //free memory
697     free(w);
698     free(v);
699     free(P);
700
701
702     //print results
703     for(i=0;i<N;i++){
704         for(j=0;j<N;j++)
705             printf("%f ",A[i*N+j]);
706         printf("\n");
707     }
708
709
710     //free memory
711     free(w);
712     free(v);
713     free(P);
714
715
716     //print results
717     for(i=0;i<N;i++){
718         for(j=0;j<N;j++)
719             printf("%f ",Q[i*N+j]);
720         printf("\n");
721     }
722
723
724     //free memory
725     free(w);
726     free(v);
727     free(P);
728
729
730     //print results
731     for(i=0;i<N;i++){
732         for(j=0;j<N;j++)
733             printf("%f ",A[i*N+j]);
734         printf("\n");
735     }
736
737
738     //free memory
739     free(w);
740     free(v);
741     free(P);
742
743
744     //print results
745     for(i=0;i<N;i++){
746         for(j=0;j<N;j++)
747             printf("%f ",Q[i*N+j]);
748         printf("\n");
749     }
750
751
752     //free memory
753     free(w);
754     free(v);
755     free(P);
756
757
758     //print results
759     for(i=0;i<N;i++){
760         for(j=0;j<N;j++)
761             printf("%f ",A[i*N+j]);
762         printf("\n");
763     }
764
765
766     //free memory
767     free(w);
768     free(v);
769     free(P);
770
771
772     //print results
773     for(i=0;i<N;i++){
774         for(j=0;j<N;j++)
775             printf("%f ",Q[i*N+j]);
776         printf("\n");
777     }
778
779
780     //free memory
781     free(w);
782     free(v);
783     free(P);
784
785
786     //print results
787     for(i=0;i<N;i++){
788         for(j=0;j<N;j++)
789             printf("%f ",A[i*N+j]);
790         printf("\n");
791     }
792
793
794     //free memory
795     free(w);
796     free(v);
797     free(P);
798
799
800     //print results
801     for(i=0;i<N;i++){
802         for(j=0;j<N;j++)
803             printf("%f ",Q[i*N+j]);
804         printf("\n");
805     }
806
807
808     //free memory
809     free(w);
810     free(v);
811     free(P);
812
813
814     //print results
815     for(i=0;i<N;i++){
816         for(j=0;j<N;j++)
817             printf("%f ",A[i*N+j]);
818         printf("\n");
819     }
820
821
822     //free memory
823     free(w);
824     free(v);
825     free(P);
826
827
828     //print results
829     for(i=0;i<N;i++){
830         for(j=0;j<N;j++)
831             printf("%f ",Q[i*N+j]);
832         printf("\n");
833     }
834
835
836     //free memory
837     free(w);
838     free(v);
839     free(P);
840
841
842     //print results
843     for(i=0;i<N;i++){
844         for(j=0;j<N;j++)
845             printf("%f ",A[i*N+j]);
846         printf("\n");
847     }
848
849
850     //free memory
851     free(w);
852     free(v);
853     free(P);
854
855
856     //print results
857     for(i=0;i<N;i++){
858         for(j=0;j<N;j++)
859             printf("%f ",Q[i*N+j]);
860         printf("\n");
861     }
862
863
864     //free memory
865     free(w);
866     free(v);
867     free(P);
868
869
870     //print results
871     for(i=0;i<N;i++){
872         for(j=0;j<N;j++)
873             printf("%f ",A[i*N+j]);
874         printf("\n");
875     }
876
877
878     //free memory
879     free(w);
880     free(v);
881     free(P);
882
883
884     //print results
885     for(i=0;i<N;i++){
886         for(j=0;j<N;j++)
887             printf("%f ",Q[i*N+j]);
888         printf("\n");
889     }
890
891
892     //free memory
893     free(w);
894     free(v);
895     free(P);
896
897
898     //print results
899     for(i=0;i<N;i++){
900         for(j=0;j<N;j++)
901             printf("%f ",A[i*N+j]);
902         printf("\n");
903     }
904
905
906     //free memory
907     free(w);
908     free(v);
909     free(P);
910
911
912     //print results
913     for(i=0;i<N;i++){
914         for(j=0;j<N;j++)
915             printf("%f ",Q[i*N+j]);
916         printf("\n");
917     }
918
919
920     //free memory
921     free(w);
922     free(v);
923     free(P);
924
925
926     //print results
927     for(i=0;i<N;i++){
928         for(j=0;j<N;j++)
929             printf("%f ",A[i*N+j]);
930         printf("\n");
931     }
932
933
934     //free memory
935     free(w);
936     free(v);
937     free(P);
938
939
940     //print results
941     for(i=0;i<N;i++){
942         for(j=0;j<N;j++)
943             printf("%f ",Q[i*N+j]);
944         printf("\n");
945     }
946
947
948     //free memory
949     free(w);
950     free(v);
951     free(P);
952
953
954     //print results
955     for(i=0;i<N;i++){
956         for(j=0;j<N;j++)
957             printf("%f ",A[i*N+j]);
958         printf("\n");
959     }
960
961
962     //free memory
963     free(w);
964     free(v);
965     free(P);
966
967
968     //print results
969     for(i=0;i<N;i++){
970         for(j=0;j<N;j++)
971             printf("%f ",Q[i*N+j]);
972         printf("\n");
973     }
974
975
976     //free memory
977     free(w);
978     free(v);
979     free(P);
980
981
982     //print results
983     for(i=0;i<N;i++){
984         for(j=0;j<N;j++)
985             printf("%f ",A[i*N+j]);
986         printf("\n");
987     }
988
989
990     //free memory
991     free(w);
992     free(v);
993     free(P);
994
995
996     //print results
997     for(i=0;i<N;i++){
998         for(j=0;j<N;j++)
999             printf("%f ",Q[i*N+j]);
1000        printf("\n");
1001    }
1002
1003
1004     //free memory
1005     free(w);
1006     free(v);
1007     free(P);
1008
1009
1010     //print results
1011     for(i=0;i<N;i++){
1012         for(j=0;j<N;j++)
1013             printf("%f ",A[i*N+j]);
1014         printf("\n");
1015     }
1016
1017
1018     //free memory
1019     free(w);
1020     free(v);
1021     free(P);
1022
1023
1024     //print results
1025     for(i=0;i<N;i++){
1026         for(j=0;j<N;j++)
1027             printf("%f ",Q[i*N+j]);
1028         printf("\n");
1029     }
1030
1031
1032     //free memory
1033     free(w);
1034     free(v);
1035     free(P);
1036
1037
1038     //print results
1039     for(i=0;i<N;i++){
1040         for(j=0;j<N;j++)
1041             printf("%f ",A[i*N+j]);
1042         printf("\n");
1043     }
1044
1045
1046     //free memory
1047     free(w);
1048     free(v);
1049     free(P);
1050
1051
1052     //print results
1053     for(i=0;i<N;i++){
1054         for(j=0;j<N;j++)
1055             printf("%f ",Q[i*N+j]);
1056         printf("\n");
1057     }
1058
1059
1060     //free memory
1061     free(w);
1062     free(v);
1063     free(P);
1064
1065
1066     //print results
1067     for(i=0;i<N;i++){
1068         for(j=0;j<N;j++)
1069             printf("%f ",A[i*N+j]);
1070         printf("\n");
1071     }
1072
1073
1074     //free memory
1075     free(w);
1076     free(v);
1077     free(P);
1078
1079
1080     //print results
1081     for(i=0;i<N;i++){
1082         for(j=0;j<N;j++)
1083             printf("%f ",Q[i*N+j]);
1084         printf("\n");
1085     }
1086
1087
1088     //free memory
1089     free(w);
1090     free(v);
1091     free(P);
1092
1093
1094     //print results
1095     for(i=0;i<N;i++){
1096         for(j=0;j<N;j++)
1097             printf("%f ",A[i*N+j]);
1098         printf("\n");
1099     }
1100
1101
1102     //free memory
1103     free(w);
1104     free(v);
1105     free(P);
1106
1107
1108     //print results
1109     for(i=0;i<N;i++){
1110         for(j=0;j<N;j++)
1111             printf("%f ",Q[i*N+j]);
1112         printf("\n");
1113     }
1114
1115
1116     //free memory
1117     free(w);
1118     free(v);
1119     free(P);
1120
1121
1122     //print results
1123     for(i=0;i<N;i++){
1124         for(j=0;j<N;j++)
1125             printf("%f ",A[i*N+j]);
1126         printf("\n");
1127     }
1128
1129
1130     //free memory
1131     free(w);
1132     free(v);
1133     free(P);
1134
1135
1136     //print results
1137     for(i=0;i<N;i++){
1138         for(j=0;j<N;j++)
1139             printf("%f ",Q[i*N+j]);
1140         printf("\n");
1141     }
1142
1143
1144     //free memory
1145     free(w);
1146     free(v);
1147     free(P);
1148
1149
1150     //print results
1151     for(i=0;i<N;i++){
1152         for(j=0;j<N;j++)
1153             printf("%f ",A[i*N+j]);
1154         printf("\n");
1155     }
1156
1157
1158     //free memory
1159     free(w);
1160     free(v);
1161     free(P);
1162
1163
1164     //print results
1165     for(i=0;i<N;i++){
1166         for(j=0;j<N;j++)
1167             printf("%f ",Q[i*N+j]);
1168         printf("\n");
1169     }
1170
1171
1172     //free memory
1173     free(w);
1174     free(v);
1175     free(P);
1176
1177
1178     //print results
1179     for(i=0;i<N;i++){
1180         for(j=0;j<N;j++)
1181             printf("%f ",A[i*N+j]);
1182         printf("\n");
1183     }
1184
1185
1186     //free memory
1187     free(w);
1188     free(v);
1189     free(P);
1190
1191
1192     //print results
1193     for(i=0;i<N;i++){
1194         for(j=0;j<N;j++)
1195             printf("%f ",Q[i*N+j]);
1196         printf("\n");
1197     }
1198
1199
1200     //free memory
1201     free(w);
1202     free(v);
1203     free(P);
1204
1205
1206     //print results
1207     for(i=0;i<N;i++){
1208         for(j=0;j<N;j++)
1209             printf("%f ",A[i*N+j]);
1210         printf("\n");
1211     }
1212
1213
1214     //free memory
1215     free(w);
1216     free(v);
1217     free(P);
1218
1219
1220     //print results
1221     for(i=0;i<N;i++){
1222         for(j=0;j<N;j++)
1223             printf("%f ",Q[i*N+j]);
1224         printf("\n");
1225     }
1226
1227
1228     //free memory
1229     free(w);
1230     free(v);
1231     free(P);
1232
1233
1234     //print results
1235     for(i=0;i<N;i++){
1236         for(j=0;j<N;j++)
1237             printf("%f ",A[i*N+j]);
1238         printf("\n");
1239     }
1240
1241
1242     //free memory
1243     free(w);
1244     free(v);
1245     free(P);
1246
1247
1248     //print results
1249     for(i=0;i<N;i++){
1250         for(j=0;j<N;j++)
1251             printf("%f ",Q[i*N+j]);
1252         printf("\n");
1253     }
1254
1255
1256     //free memory
1257     free(w);
1258     free(v);
1259     free(P);
1260
1261
1262     //print results
1263     for(i=0;i<N;i++){
1264         for(j=0;j<N;j++)
1265             printf("%f ",A[i*N+j]);
1266         printf("\n");
1267     }
1268
1269
1270     //free memory
1271     free(w);
1272     free(v);
1273     free(P);
1274
1275
1276     //print results
1277     for(i=0;i<N;i++){
1278         for(j=0;j<N;j++)
1279             printf("%f ",Q[i*N+j]);
1280         printf("\n");
1281     }
1282
1283
1284     //free memory
1285     free(w);
1286     free(v);
1287     free(P);
1288
1289
1290     //print results
1291     for(i=0;i<N;i++){
1292         for(j=0;j<N;j++)
1293             printf("%f ",A[i*N+j]);
1294         printf("\n");
1295     }
1296
1297
1298     //free memory
1299     free(w);
1300     free(v);
1301     free(P);
1302
1303
1304     //print results
1305     for(i=0;i<N;i++){
1306         for(j=0;j<N;j++)
1307             printf("%f ",Q[i*N+j]);
1308         printf("\n");
1309     }
1310
1311
1312     //free memory
1313     free(w);
1314     free(v);
1315     free(P);
1316
1317
1318     //print results
1319     for(i=0;i<N;i++){
1320         for(j=0;j<N;j++)
1321             printf("%f ",A[i*N+j]);
1322         printf("\n");
1323     }
1324
1325
1326     //free memory
1327     free(w);
1328     free(v);
1329     free(P);
1330
1331
1332     //print results
1333     for(i=0;i<N;i++){
1334         for(j=0;j<N;j++)
1335             printf("%f ",Q[i*N+j]);
1336         printf("\n");
1337     }
1338
1339
1340     //free memory
1341     free(w);
1342     free(v);
1343     free(P);
1344
1345
1346     //print results
1347     for(i=0;i<N;i++){
1348         for(j=0;j<N;j++)
1349             printf("%f ",A[i*N+j]);
1350         printf("\n");
1351     }
1352
1353
1354     //free memory
1355     free(w);
1356     free(v);
1357     free(P);
1358
1359
1360     //print results
1361     for(i=0;i<N;i++){
1362         for(j=0;j<N;j++)
1363             printf("%f ",Q[i*N+j]);
1364         printf("\n");
1365     }
1366
1367
1368     //free memory
1369     free(w);
1370     free(v);
1371     free(P);
1372
1373
1374     //print results
1375     for(i=0;i<N;i++){
1376         for(j=0;j<N;j++)
1377             printf("%f ",A[i*N+j]);
1378         printf("\n");
1379     }
1380
1381
1382     //free memory
1383     free(w);
1384     free(v);
1385     free(P);
1386
1387
1388     //print results
1389     for(i=0;i<N;i++){
1390         for(j=0;j<N;j++)
1391             printf("%f ",Q[i*N+j]);
1392         printf("\n");
1393     }
1394
1395
1396     //free memory
1397     free(w);
1398     free(v);
1399     free(P);
1400
1401
1402     //print results
1403     for(i=0;i<N;i++){
1404         for(j=0;j<N;j++)
1405             printf("%f ",A[i*N+j]);
1406         printf("\n");
1407     }
1408
1409
1410     //free memory
1411     free(w);
1412     free(v);
1413     free(P);
1414
1415
1416     //print results
1417     for(i=0;i<N;i++){
1418         for(j=0;j<N;j++)
1419             printf("%f ",Q[i*N+j]);
1420         printf("\n");
1421     }
1422
1423
1424     //free memory
1425     free(w);
1426     free(v);
1427     free(P);
1428
1429
1430     //print results
1431     for(i=0;i<N;i++){
1432         for(j=0;j<N;j++)
1433             printf("%f ",A[i*N+j]);
1434         printf("\n");
1435     }
1436
1437
1438     //free memory
1439     free(w);
1440     free(v);
1441     free(P);
1442
1443
1444     //print results
1445     for(i=0;i<N;i++){
1446         for(j=0;j<N;j++)
1447             printf("%f ",Q[i*N+j]);
1448         printf("\n");
1449     }
1450
1451
1452     //free memory
1453     free(w);
1454     free(v);
1455     free(P);
1456
1457
1458     //print results
1459     for(i=0;i<N;i++){
1460         for(j=0;j<N;j++)
1461             printf("%f ",A[i*N+j]);
1462         printf("\n");
1463     }
1464
1465
1466     //free memory
1467     free(w);
1468     free(v);
1469     free(P);
1470
1471
1472     //print results
1473     for(i=0;i<N;i++){
1474         for(j=0;j<N;j++)
1475             printf("%f ",Q[i*N+j]);
1476         printf("\n");
1477     }
1478
1479
1480     //free memory
1481     free(w);
1482     free(v);
1483     free(P);
1484
1485
1486     //print results
1487     for(i=0;i<N;i++){
1488         for(j=0;j<N;j++)
1489             printf("%f ",A[i*N+j]);
1490         printf("\n");
1491     }
1492
1493
1494     //free memory
1495     free(w);
1496     free(v);
1497     free(P);
1498
1499
1500     //print results
1501     for(i=0;i<N;i++){
1502         for(j=0;j<N;j++)
1503             printf("%f ",Q[i*N+j]);
1504         printf("\n");
1505     }
1506
1507
1508     //free memory
1509     free(w);
1510     free(v);
1511     free(P);
1512
1513
1514     //print results
1515     for(i=0;i<N;i++){
1516         for(j=0;j<N;j++)
1517             printf("%f ",A[i*N+j]);
1518         printf("\n");
1519     }
1520
1521
1522     //free memory
1523     free(w);
1524     free(v);
1525     free(P);
1526
1527
1528     //print results
1529     for(i=0;i<N;i++){
1530         for(j=0;j<N;j++)
1531             printf("%f ",Q[i*N+j]);
1532         printf("\n");
1533     }
1534
1535
1536     //free memory
1537     free(w);
1538     free(v);
1539     free(P);
1540
1541
1542     //print results
1543     for(i=0;i<N;i++){
1544         for(j=0;j<N;j++)
1545             printf("%f ",A[i*N+j]);
1546         printf("\n");
1547     }
1548
1549
1550     //free memory
1551     free(w);
1552     free(v);
1553     free(P);
1554
1555
1556     //print results
1557     for(i=0;i<N;i++){
1558         for(j=0;j<N;j++)
1559             printf("%f ",Q[i*N+j]);
1560         printf("\n");
1561     }
1562
1563
1564     //free memory
1565     free(w);
1566     free(v);
1567     free(P);
1568
1569
1570     //print results
1571     for(i=0;i<N;i++){
1572         for(j=0;j<N;j++)
1573             printf("%f ",A[i*N+j]);
1574         printf("\n");
1575     }
1576
1577
1578     //free memory
1579     free(w);
1580     free(v);
1581     free(P);
1582
1583
1584     //print results
1585     for(i=0;i<N;i++){
1586         for(j=0;j<N;j++)
1587             printf("%f ",Q[i*N+j]);
1588         printf("\n");
1589     }
1590
1591
1592     //free memory
1593     free(w);
1594     free(v);
1595     free(P);
1596
1597
1598     //print results
1599     for(i=0;i<N;i++){
1600         for(j=0;j<N;j++)
1601             printf("%f ",A[i*N+j]);
1602         printf("\n");
1603     }
1604
1605
1606     //free memory
1607     free(w);
1608     free(v);
1609     free(P);
1610
1611
1612     //print results
1613     for(i=0;i<N;i++){
1614         for(j=0;j<N;j++)
1615             printf("%f ",Q[i*N+j]);
1616         printf("\n");
1617     }
1618
1619
1620     //free memory
1621     free(w);
1622     free(v);
1623     free(P);
1624
1625
1626     //print results
1627     for(i=0;i<N;i++){
1628         for(j=0;j<N;j++)
1629             printf("%f ",A[i*N+j]);
1630         printf("\n");
1631     }
1632
1633
1634     //free memory
1635     free(w);
1636     free(v);
1637     free(P);
1638
1639
1640
```

```

359         }
360         else if(i==j-1){
361             A[i*N+j]=subdiag[i];
362         }
363         else{
364             A[i*N+j]=0;
365         }
366     }
367     diag[i]-=mu;
368 }
369
370 Q=QR(diag,subdiag);
371 while(chk==0){
372     double *temp,*Qt,*tdiag,*tsubdiag,*Qnew,*Ak;
373     temp=matmul(A,Q);
374     Qt=transp(Q);
375     Ak=matmul(Qt,temp);
376     free(temp);
377     free(Qt);
378     mu=Ak[N*N-1];
379     tdiag=malloc(sizeof(double)*N);
380     tsubdiag=malloc(sizeof(double)*(N-1));
381     for(i=0;i<N-1;i++){
382         tdiag[i]=Ak[i*N+i]-mu;
383         tsubdiag[i]=Ak[i*N+i+1];
384     }
385     free(Ak);
386     tdiag[N-1]=0;
387     Qnew=QR(tdiag,tsubdiag);
388     temp=matmul(Q,Qnew);
389     free(Q);
390     free(Qnew);
391     Q=temp;
392
393     conf=tdiag[0];
394     for(i=1;i<N;i++){
395         if(tdiag[i]<conf)
396             conf=tdiag[i];
397     }
398     conf+=mu;
399     if(fabs(conf-conf2)<1e-12)
400         chk=1;
401     conf2=conf;
402     //printf("%f\n",conf2);
403     free(tdiag);
404     free(tsubdiag);
405 }
406
407 free(diag);
408 free(subdiag);
409 free(A);
410 free(Q);
411
412 *eig0=conf2;
413
414 end = clock();
415 *time=((double) (end - start)) / CLOCKS_PER_SEC;
416 printf("Single-thread CPU algorithm finished!\n");
417 printf("Execution time: %f\n",*time);
418 printf("Lowest eigenvalue: %f\n\n",*eig0);
419 }

```

```

420
421
422
423
424
425 //LAPACK
426
427 void lapack(double *time ,double *eig0 ,double complex *M){
428
429     printf("Lapack started!\n");
430     clock_t start , end;
431     start = clock();
432
433     lapack_int n;
434     n=N;
435     double *eig;
436     eig=malloc(sizeof(double)*N);
437     LAPACKE_zheev(LAPACK_ROW_MAJOR , 'V' , 'U' , n , M , n , eig);
438     *eig0=eig [0];
439     free(eig);
440
441     end = clock();
442     *time=((double) (end - start)) / CLOCKS_PER_SEC;
443     printf("Lapack finished!\n");
444     printf("Execution time: %f\n",*time);
445     printf("Lowest eigenvalue: %f\n\n",*eig0);
446 }
447
448
449
450
451
452 //OpenCL GPU
453
454 void openclGPU(double *time ,double *eig0 ,double *M){
455
456     printf("OpenCL (GPU) started!\n");
457     clock_t start , end;
458     start = clock();
459
460     //opencl initialization
461
462     FILE *f;
463     int i,j;
464     char str[40];
465     char name[40];
466     char *source_str;
467     size_t source_size;
468     sprintf(str,"-cl-std=CL1.2 -D N=%d",N);
469
470     int dev=1;
471     cl_platform_id *platforms;
472     cl_uint num_platforms;
473     cl_device_id device;
474     cl_uint num_devices;
475     cl_context context;
476     cl_command_queue queue;
477     cl_program program;
478     size_t glob_size;
479     size_t glob_size2;
480     size_t loc_size=N;

```

```

481
482     clGetPlatformIDs(5, NULL, &num_platforms);
483     platforms=(cl_platform_id*) malloc(sizeof(cl_platform_id)*num_platforms
484         );
484     clGetPlatformIDs(num_platforms, platforms, NULL);
485     //printf("Number of Platforms detected: %d\n", num_platforms);
486     for(i=0;i<num_platforms;i++){
487         clGetPlatformInfo(platforms[i], CL_PLATFORM_NAME, sizeof(name), &name,
488             NULL);
489         //printf("%s\n", name);
490     }
491
492     if(dev==1){
493         clGetDeviceIDs(platforms[0], CL_DEVICE_TYPE_GPU, 1, &device, &
494             num_devices);
495         for(i=1;loc_size>256;i++){
496             if(loc_size%i==0)
497                 loc_size=N/i;
498         }
499         //printf("\nAMD GPU selected\n");
500     }
501     else if(dev==2){
502         clGetDeviceIDs(platforms[1], CL_DEVICE_TYPE_CPU, 1, &device, &
503             num_devices);
504         for(i=1;loc_size>8192;i++){
505             if(loc_size%i==0)
506                 loc_size=N/i;
507         }
508         //printf("\nIntel CPU selected\n");
509     }
510     context=clCreateContext( NULL, 1,&device, NULL, NULL, NULL);
511     queue=clCreateCommandQueueWithProperties(context,device,NULL,NULL);
512
513     f=fopen("kernels.cl", "r");
514     source_str = (char*)malloc(MAX_SOURCE_SIZE);
515     source_size = fread( source_str, 1, MAX_SOURCE_SIZE, f);
516     fclose(f);
517     program=clCreateProgramWithSource(context,1,(const char **)&source_str
518         ,(const size_t *)&source_size,NULL);
519     clBuildProgram(program,1,&device,str,NULL,NULL);
520     free(source_str);
521
522     cl_kernel lanczos1 = clCreateKernel(program, "lanczos1", NULL);
523     cl_kernel lanczos2 = clCreateKernel(program, "lanczos2", NULL);
524     cl_kernel transpose = clCreateKernel(program, "transpose", NULL);
525     cl_kernel mat_mul = clCreateKernel(program, "mat_mul", NULL);
526     cl_kernel transfer = clCreateKernel(program, "transfer", NULL);
527
528
529     //lanczos algorithm
530
531     double a,b=1;
532     double *P, *v, *vold,*w;
533     double *diag,*subdiag,*ab;
534
535     cl_double alpha;
536     cl_double beta;

```

```

537     cl_mem M_buff;
538     cl_mem v_buff;
539     cl_mem vold_buff;
540     cl_mem w_buff;
541     cl_mem ab_buff;
542
543     diag=malloc(sizeof(double)*N);
544     subdiag=malloc(sizeof(double)*N);
545     w=calloc(2*N,sizeof(double));
546     v=malloc(sizeof(double)*2*N);
547     vold=calloc(2*N,sizeof(double));
548     P=calloc(2*N*N,sizeof(double));
549     ab=malloc(N*sizeof(double));
550     P[0]=1;
551     P[1]=0;
552     w[0]=P[0];
553     w[1]=P[1];
554
555     glob_size=N;
556     M_buff=clCreateBuffer(context ,CL_MEM_READ_ONLY ,2*N*N*sizeof(double) ,
557                           NULL ,NULL );
557     v_buff=clCreateBuffer(context ,CL_MEM_READ_WRITE ,2*N*sizeof(double) ,NULL
558                           ,NULL );
558     vold_buff=clCreateBuffer(context ,CL_MEM_READ_WRITE ,2*N*sizeof(double) ,
559                           NULL ,NULL );
559     w_buff=clCreateBuffer(context ,CL_MEM_READ_WRITE ,2*N*sizeof(double) ,NULL
560                           ,NULL );
560     ab_buff=clCreateBuffer(context ,CL_MEM_WRITE_ONLY ,N*sizeof(double) ,NULL ,
561                           NULL );
561
562     for(i=0;i<N;i++){
563
564         beta=b;
565
566         clEnqueueWriteBuffer(queue , M_buff , CL_TRUE , 0,2*N*N*sizeof(double) ,
567                               M , 0 , NULL , NULL );
567         clEnqueueWriteBuffer(queue , w_buff , CL_TRUE , 0,2*N*sizeof(double) , w
568                               , 0 , NULL , NULL );
568         clSetKernelArg(lanczos1 , 0 , sizeof(cl_double) , &beta);
569         clSetKernelArg(lanczos1 , 1 , sizeof(cl_mem) , &M_buff);
570         clSetKernelArg(lanczos1 , 2 , sizeof(cl_mem) , &v_buff);
571         clSetKernelArg(lanczos1 , 3 , sizeof(cl_mem) , &w_buff);
572         clSetKernelArg(lanczos1 , 4 , sizeof(cl_mem) , &ab_buff);
573         clEnqueueNDRangeKernel(queue , lanczos1 ,1 ,NULL ,&glob_size ,&loc_size ,0 ,
574                               NULL ,NULL );
574         clEnqueueReadBuffer(queue , v_buff , CL_TRUE , 0,2*N*sizeof(double) , v ,
575                               0 , NULL , NULL );
575         clEnqueueReadBuffer(queue , ab_buff , CL_TRUE , 0,N*sizeof(double) , ab ,
576                               0 , NULL , NULL );
576
577         a=ab[0];
578         for(j=1;j<N;j++){
579             a+=ab[j];
580         }
581         diag[i]=a;
582
583         if(i==0)
584             beta=0;
585
586         alpha=a;
587

```

```

588     clEnqueueWriteBuffer(queue, vold_buff, CL_TRUE, 0, 2*N*sizeof(double)
589     , vold, 0, NULL, NULL);
590     clSetKernelArg(lanczos2, 0, sizeof(cl_double), &alpha);
591     clSetKernelArg(lanczos2, 1, sizeof(cl_double), &beta);
592     clSetKernelArg(lanczos2, 2, sizeof(cl_mem), &v_buff);
593     clSetKernelArg(lanczos2, 3, sizeof(cl_mem), &vold_buff);
594     clSetKernelArg(lanczos2, 4, sizeof(cl_mem), &w_buff);
595     clSetKernelArg(lanczos2, 5, sizeof(cl_mem), &ab_buff);
596     clEnqueueNDRangeKernel(queue, lanczos2, 1, NULL, &glob_size, &loc_size, 0,
597     NULL, NULL);
598     clEnqueueReadBuffer(queue, w_buff, CL_TRUE, 0, 2*N*sizeof(double), w,
599     0, NULL, NULL);
600     clEnqueueReadBuffer(queue, ab_buff, CL_TRUE, 0, N*sizeof(double), ab,
601     0, NULL, NULL);
602
603     b=ab[0];
604     for(j=1;j<N;j++){
605         b+=ab[j];
606     }
607     b=sqrt(b);
608     subdiag[i]=b;
609
610     for(j=0;j<N;j++){
611         P[2*j*N+2*i]=v[2*j];
612         P[2*j*N+2*i+1]=v[2*j+1];
613         vold[2*j]=v[2*j];
614         vold[2*j+1]=v[2*j+1];
615     }
616 }
617
618 clReleaseMemObject(M_buff);
619 clReleaseMemObject(v_buff);
620 clReleaseMemObject(vold_buff);
621 clReleaseMemObject(w_buff);
622 clReleaseMemObject(ab_buff);
623 clReleaseKernel(lanczos1);
624 clReleaseKernel(lanczos2);
625
626
627
628 //QR algorithm
629
630 int chk=0;
631 double mu,conf,conf2=0;
632 double *Q,*A;
633
634 cl_mem A_buff;
635 cl_mem B_buff;
636 cl_mem C_buff;
637
638 glob_size=N*N;
639 glob_size2=N;
640 A_buff=clCreateBuffer(context,CL_MEM_READ_WRITE,N*N*sizeof(double),NULL
641 ,NULL);
642 B_buff=clCreateBuffer(context,CL_MEM_READ_WRITE,N*N*sizeof(double),NULL
643 ,NULL);
644 C_buff=clCreateBuffer(context,CL_MEM_READ_WRITE,N*N*sizeof(double),NULL
645 ,NULL);

```

```

        ,NULL);
643 v_buff=clCreateBuffer(context,CL_MEM_WRITE_ONLY,N*sizeof(double),NULL,
644             NULL);
645 w_buff=clCreateBuffer(context,CL_MEM_WRITE_ONLY,N*sizeof(double),NULL,
646             NULL);
647
648 mu=diag[N-1];
649 A=malloc(sizeof(double)*N*N);
650 for(i=0;i<N;i++){
651     for(j=0;j<N;j++){
652         if(i==j){
653             A[i*N+j]=diag[i];
654         }
655         else if(i==j+1){
656             A[i*N+j]=subdiag[i-1];
657         }
658         else if(i==j-1){
659             A[i*N+j]=subdiag[i];
660         }
661         else{
662             A[i*N+j]=0;
663         }
664     }
665     diag[i]-=mu;
666 }
667
668 Q=QR(diag,subdiag);
669 while(chk==0){
670     double *tdiag,*tsubdiag,*Qnew;
671
672     //temp=A*Q
673     clEnqueueWriteBuffer(queue, A_buff, CL_TRUE, 0,N*N*sizeof(double), A
674             , 0, NULL, NULL);
675     clEnqueueWriteBuffer(queue, B_buff, CL_TRUE, 0,N*N*sizeof(double), Q
676             , 0, NULL, NULL);
677     clSetKernelArg(mat_mul, 0, sizeof(cl_mem), &A_buff);
678     clSetKernelArg(mat_mul, 1, sizeof(cl_mem), &B_buff);
679     clSetKernelArg(mat_mul, 2, sizeof(cl_mem), &C_buff);
680     clEnqueueNDRangeKernel(queue,mat_mul,1,NULL,&glob_size,NULL,0,NULL,
681             NULL);
682
683     //Qt=trasp(Q)
684     clEnqueueWriteBuffer(queue, B_buff, CL_TRUE, 0,N*N*sizeof(double), Q
685             , 0, NULL, NULL);
686     clSetKernelArg(transpose, 0, sizeof(cl_mem), &B_buff);
687     clSetKernelArg(transpose, 1, sizeof(cl_mem), &A_buff);
688     clEnqueueNDRangeKernel(queue,transpose,1,NULL,&glob_size,NULL,0,NULL,
689             NULL);
690
691     tdiag=malloc(sizeof(double)*N);
692     tsubdiag=malloc(sizeof(double)*N);
693
694     clSetKernelArg(transfer, 0, sizeof(cl_mem), &B_buff);
695     clSetKernelArg(transfer, 1, sizeof(cl_mem), &v_buff);

```

```

695     clSetKernelArg(transfer, 2, sizeof(cl_mem), &w_buff);
696     clEnqueueNDRangeKernel(queue, transfer, 1, NULL, &glob_size2, NULL, 0, NULL
697                             , NULL);
698     clEnqueueReadBuffer(queue, v_buff, CL_TRUE, 0, N*sizeof(double), tdiag
699                          , 0, NULL, NULL);
700     clEnqueueReadBuffer(queue, w_buff, CL_TRUE, 0, N*sizeof(double),
701                          tsubdiag, 0, NULL, NULL);
702
703     mu=tsubdiag[N-1];
704     tdiag[N-1]=0;
705     Qnew=QR(tdiag,tsubdiag);
706
707     //Q=Q*Qnew
708     clEnqueueWriteBuffer(queue, A_buff, CL_TRUE, 0, N*N*sizeof(double), Q
709                           , 0, NULL, NULL);
710     clEnqueueWriteBuffer(queue, B_buff, CL_TRUE, 0, N*N*sizeof(double),
711                           Qnew, 0, NULL, NULL);
712     clSetKernelArg(mat_mul, 0, sizeof(cl_mem), &A_buff);
713     clSetKernelArg(mat_mul, 1, sizeof(cl_mem), &B_buff);
714     clSetKernelArg(mat_mul, 2, sizeof(cl_mem), &C_buff);
715     clEnqueueNDRangeKernel(queue, mat_mul, 1, NULL, &glob_size, NULL, 0, NULL,
716                            NULL);
717     clEnqueueReadBuffer(queue, C_buff, CL_TRUE, 0, N*N*sizeof(double), Q,
718                          0, NULL, NULL);
719     free(Qnew);
720
721     conf=tdiag[0];
722     for(i=0;i<N;i++){
723         if(tdiag[i]<conf)
724             conf=tdiag[i];
725     }
726     free(tdiag);
727     free(tsubdiag);
728     conf+=mu;
729     if(fabs(conf-conf2)< 1e-13)
730         chk=1;
731     conf2=conf;
732     //printf("%f\n", conf2);
733 }
734 *eig0=conf2;
735
736 //free memory
737 free(P);
738 free(diag);
739 free(subdiag);
740 free(A);
741 free(Q);
742
743 //Deallocate resources
744 clReleaseMemObject(A_buff);
745 clReleaseMemObject(B_buff);
746 clReleaseMemObject(C_buff);
747 clReleaseKernel(mat_mul);
748 clReleaseKernel(transpose);
749 clReleaseKernel(transfer);
750 clReleaseCommandQueue(queue);
751 clReleaseProgram(program);
752 clReleaseContext(context);
753
754 end = clock();
755 *time=((double) (end - start)) /CLOCKS_PER_SEC;

```

```

749     printf("OpenCL (GPU) finished!\n");
750     printf("Execution time: %f\n", *time);
751     printf("Lowest eigenvalue: %f\n\n", *eig0);
752 }
753
754
755
756
757 //OpenCL CPU
758
759 void openclCPU(double *time ,double *eig0 ,double *M){
760
761     printf("OpenCL (CPU) started!\n");
762     clock_t start, end;
763     start = clock();
764
765     //opencl initialization
766
767     FILE *f;
768     int i,j;
769     char str[40];
770     char name[40];
771     char *source_str;
772     size_t source_size;
773     sprintf(str,"-cl-std=CL1.2 -D N=%d",N);
774
775     int dev=2;
776     cl_platform_id *platforms;
777     cl_uint num_platforms;
778     cl_device_id device;
779     cl_uint num_devices;
780     cl_context context;
781     cl_command_queue queue;
782     cl_program program;
783     size_t glob_size;
784     size_t glob_size2;
785     size_t loc_size=N;
786
787     clGetPlatformIDs(5,NULL,&num_platforms);
788     platforms=(cl_platform_id*) malloc(sizeof(cl_platform_id)*num_platforms
    );
789     clGetPlatformIDs(num_platforms,platforms,NULL);
//printf("Number of Platforms detected: %d\n",num_platforms);
791     for(i=0;i<num_platforms;i++){
792         clGetPlatformInfo(platforms[i],CL_PLATFORM_NAME,sizeof(name),&name,
    NULL);
//printf("%s\n",name);
793     }
794
795
796
797     if(dev==1){
798         clGetDeviceIDs(platforms[0],CL_DEVICE_TYPE_GPU,1,&device,&
    num_devices);
799         for(i=1;loc_size>256;i++){
800             if(loc_size%i==0)
801                 loc_size=N/i;
802         }
//printf("\nAMD GPU selected\n");
803     }
804     else if(dev==2){
805         clGetDeviceIDs(platforms[1],CL_DEVICE_TYPE_CPU,1,&device,&

```

```

        num_devices);
807    for(i=1; loc_size>8192; i++){
808        if(loc_size%i==0)
809            loc_size=N/i;
810    }
811    //printf("\nIntel CPU selected\n");
812 }
813
814 context=clCreateContext( NULL , 1,&device , NULL , NULL , NULL );
815 queue=clCreateCommandQueueWithProperties( context , device ,NULL ,NULL );
816
817 f=fopen("kernels.cl" , "r");
818 source_str = (char*)malloc(MAX_SOURCE_SIZE);
819 source_size = fread( source_str , 1 , MAX_SOURCE_SIZE , f );
820 fclose(f);
821 program=clCreateProgramWithSource(context ,1,(const char **)&source_str
822     ,(const size_t *)&source_size ,NULL );
823 clBuildProgram(program ,1,&device ,str ,NULL ,NULL );
824 free(source_str);
825
826 cl_kernel lanczos1 = clCreateKernel(program , "lanczos1" , NULL );
827 cl_kernel lanczos2 = clCreateKernel(program , "lanczos2" , NULL );
828 cl_kernel transpose = clCreateKernel(program , "transpose" , NULL );
829 cl_kernel mat_mul = clCreateKernel(program , "mat_mul" , NULL );
830 cl_kernel transfer = clCreateKernel(program , "transfer" , NULL );
831
832
833 //lanczos algorithm
834
835 double a,b=1;
836 double *P , *v , *vold,*w;
837 double *diag,*subdiag,*ab;
838
839 cl_double alpha;
840 cl_double beta;
841 cl_mem M_buff;
842 cl_mem v_buff;
843 cl_mem vold_buff;
844 cl_mem w_buff;
845 cl_mem ab_buff;
846
847 diag=malloc(sizeof(double)*N);
848 subdiag=malloc(sizeof(double)*N);
849 w=calloc(2*N,sizeof(double));
850 v=malloc(sizeof(double)*2*N);
851 vold=calloc(2*N,sizeof(double));
852 P=calloc(2*N*N,sizeof(double));
853 ab=malloc(N*sizeof(double));
854 P[0]=1;
855 P[1]=0;
856 w[0]=P[0];
857 w[1]=P[1];
858
859 glob_size=N;
860 M_buff=clCreateBuffer(context ,CL_MEM_READ_ONLY ,2*N*N*sizeof(double) ,
861 NULL ,NULL );
862 v_buff=clCreateBuffer(context ,CL_MEM_READ_WRITE ,2*N*sizeof(double) ,NULL
863 ,NULL );
864 vold_buff=clCreateBuffer(context ,CL_MEM_READ_WRITE ,2*N*sizeof(double) ,
865 NULL ,NULL );

```

```

863     w_buff=clCreateBuffer(context ,CL_MEM_READ_WRITE ,2*N*sizeof(double) ,NULL
864         ,NULL);
864     ab_buff=clCreateBuffer(context ,CL_MEM_WRITE_ONLY ,N*sizeof(double) ,NULL ,
865         NULL);
865
866     for(i=0;i<N;i++){
867
868         beta=b;
869
870         clEnqueueWriteBuffer(queue , M_buff , CL_TRUE , 0,2*N*N*sizeof(double) ,
871             M , 0 , NULL , NULL);
871         clEnqueueWriteBuffer(queue , w_buff , CL_TRUE , 0,2*N*sizeof(double) , w
872             , 0 , NULL , NULL);
872         clSetKernelArg(lanczos1 , 0 , sizeof(cl_double) , &beta);
873         clSetKernelArg(lanczos1 , 1 , sizeof(cl_mem) , &M_buff);
874         clSetKernelArg(lanczos1 , 2 , sizeof(cl_mem) , &v_buff);
875         clSetKernelArg(lanczos1 , 3 , sizeof(cl_mem) , &w_buff);
876         clSetKernelArg(lanczos1 , 4 , sizeof(cl_mem) , &ab_buff);
877         clEnqueueNDRangeKernel(queue ,lanczos1 ,1 ,NULL ,&glob_size ,&loc_size ,0 ,
878             NULL ,NULL);
878         clEnqueueReadBuffer(queue ,v_buff , CL_TRUE , 0,2*N*sizeof(double) , v ,
879             0 , NULL , NULL);
879         clEnqueueReadBuffer(queue ,ab_buff , CL_TRUE , 0,N*sizeof(double) , ab ,
880             0 , NULL , NULL);
880
881         a=ab [0];
882         for(j=1;j<N;j++){
883             a+=ab [j];
884         }
885         diag [i]=a;
886
887         if(i==0)
888             beta=0;
889
890         alpha=a;
891
892         clEnqueueWriteBuffer(queue , vold_buff , CL_TRUE , 0,2*N*sizeof(double)
893             , vold , 0 , NULL , NULL);
893         clSetKernelArg(lanczos2 , 0 , sizeof(cl_double) , &alpha);
894         clSetKernelArg(lanczos2 , 1 , sizeof(cl_double) , &beta);
895         clSetKernelArg(lanczos2 , 2 , sizeof(cl_mem) , &v_buff);
896         clSetKernelArg(lanczos2 , 3 , sizeof(cl_mem) , &vold_buff);
897         clSetKernelArg(lanczos2 , 4 , sizeof(cl_mem) , &w_buff);
898         clSetKernelArg(lanczos2 , 5 , sizeof(cl_mem) , &ab_buff);
899         clEnqueueNDRangeKernel(queue ,lanczos2 ,1 ,NULL ,&glob_size ,&loc_size ,0 ,
900             NULL ,NULL);
900         clEnqueueReadBuffer(queue ,w_buff , CL_TRUE , 0,2*N*sizeof(double) , w ,
901             0 , NULL , NULL);
901         clEnqueueReadBuffer(queue ,ab_buff , CL_TRUE , 0,N*sizeof(double) , ab ,
902             0 , NULL , NULL);
902
903         b=ab [0];
904         for(j=1;j<N;j++){
905             b+=ab [j];
906         }
907         b=sqrt(b);
908         subdiag [i]=b;
909
910         for(j=0;j<N;j++){
911             P [2*j*N+2*i]=v [2*j];
912             P [2*j*N+2*i+1]=v [2*j+1];

```

```

913         vold[2*j]=v[2*j];
914         vold[2*j+1]=v[2*j+1];
915     }
916 }
917
918 clReleaseMemObject(M_buff);
919 clReleaseMemObject(v_buff);
920 clReleaseMemObject(vold_buff);
921 clReleaseMemObject(w_buff);
922 clReleaseMemObject(ab_buff);
923 clReleaseKernel(lanczos1);
924 clReleaseKernel(lanczos2);
925
926 free(w);
927 free(v);
928 free(vold);
929
930
931
932 //QR algorithm
933
934 int chk=0;
935 double mu,conf,conf2=0;
936 double *Q,*A;
937
938 cl_mem A_buff;
939 cl_mem B_buff;
940 cl_mem C_buff;
941
942 glob_size=N*N;
943 glob_size2=N;
944 A_buff=clCreateBuffer(context,CL_MEM_READ_WRITE,N*N*sizeof(double),NULL,
945 ,NULL);
946 B_buff=clCreateBuffer(context,CL_MEM_READ_WRITE,N*N*sizeof(double),NULL,
947 ,NULL);
948 C_buff=clCreateBuffer(context,CL_MEM_READ_WRITE,N*N*sizeof(double),NULL,
949 ,NULL);
950 v_buff=clCreateBuffer(context,CL_MEM_WRITE_ONLY,N*sizeof(double),NULL,
951 ,NULL);
952 w_buff=clCreateBuffer(context,CL_MEM_WRITE_ONLY,N*sizeof(double),NULL,
953 ,NULL);
954
955 mu=diag[N-1];
956 A=malloc(sizeof(double)*N*N);
957 for(i=0;i<N;i++){
958     for(j=0;j<N;j++){
959         if(i==j){
960             A[i*N+j]=diag[i];
961         }
962         else if(i==j+1){
963             A[i*N+j]=subdiag[i-1];
964         }
965         else if(i==j-1){
966             A[i*N+j]=subdiag[i];
967         }
968         else{
969             A[i*N+j]=0;
970         }
971     }
972     diag[i]-=mu;
973 }

```

```

969
970     Q=QR(diag,subdiag);
971     while(chk==0){
972         double *tdiag,*tsubdiag,*Qnew;
973
974         //temp=A*Q
975         clEnqueueWriteBuffer(queue, A_buff, CL_TRUE, 0,N*N*sizeof(double), A
976             , 0, NULL, NULL);
977         clEnqueueWriteBuffer(queue, B_buff, CL_TRUE, 0,N*N*sizeof(double), Q
978             , 0, NULL, NULL);
979         clSetKernelArg(mat_mul, 0, sizeof(cl_mem), &A_buff);
980         clSetKernelArg(mat_mul, 1, sizeof(cl_mem), &B_buff);
981         clSetKernelArg(mat_mul, 2, sizeof(cl_mem), &C_buff);
982         clEnqueueNDRangeKernel(queue,mat_mul,1,NULL,&glob_size,NULL,0,NULL,
983             NULL);
984
985         //Qt=trasp(Q)
986         clEnqueueWriteBuffer(queue, B_buff, CL_TRUE, 0,N*N*sizeof(double), Q
987             , 0, NULL, NULL);
988         clSetKernelArg(transpose, 0, sizeof(cl_mem), &B_buff);
989         clSetKernelArg(transpose, 1, sizeof(cl_mem), &A_buff);
990         clEnqueueNDRangeKernel(queue,transpose,1,NULL,&glob_size,NULL,0,NULL
991             ,NULL);
992
993         //Ak=Qt*temp
994         clSetKernelArg(mat_mul, 0, sizeof(cl_mem), &A_buff);
995         clSetKernelArg(mat_mul, 1, sizeof(cl_mem), &C_buff);
996         clSetKernelArg(mat_mul, 2, sizeof(cl_mem), &B_buff);
997         clEnqueueNDRangeKernel(queue,mat_mul,1,NULL,&glob_size,NULL,0,NULL,
998             NULL);
999
1000         tdiag=malloc(sizeof(double)*N);
1001         tsubdiag=malloc(sizeof(double)*N);
1002
1003         clSetKernelArg(transfer, 0, sizeof(cl_mem), &B_buff);
1004         clSetKernelArg(transfer, 1, sizeof(cl_mem), &v_buff);
1005         clSetKernelArg(transfer, 2, sizeof(cl_mem), &w_buff);
1006         clEnqueueNDRangeKernel(queue,transfer,1,NULL,&glob_size2,NULL,0,NULL
1007             ,NULL);
1008         clEnqueueReadBuffer(queue,v_buff, CL_TRUE, 0,N*sizeof(double), tdiag
1009             , 0, NULL, NULL);
1010         clEnqueueReadBuffer(queue,w_buff, CL_TRUE, 0,N*sizeof(double),
1011             tsubdiag, 0, NULL, NULL);
1012
1013         mu=tsubdiag[N-1];
1014         tdiag[N-1]=0;
1015         Qnew=QR(tdiag,tsubdiag);
1016
1017         //Q=Q*Qnew
1018         clEnqueueWriteBuffer(queue, A_buff, CL_TRUE, 0,N*N*sizeof(double), Q
1019             , 0, NULL, NULL);
1020         clEnqueueWriteBuffer(queue, B_buff, CL_TRUE, 0,N*N*sizeof(double),
1021             Qnew, 0, NULL, NULL);
1022         clSetKernelArg(mat_mul, 0, sizeof(cl_mem), &A_buff);
1023         clSetKernelArg(mat_mul, 1, sizeof(cl_mem), &B_buff);
1024         clSetKernelArg(mat_mul, 2, sizeof(cl_mem), &C_buff);
1025         clEnqueueNDRangeKernel(queue,mat_mul,1,NULL,&glob_size,NULL,0,NULL,
1026             NULL);
1027         clEnqueueReadBuffer(queue,C_buff, CL_TRUE, 0,N*N*sizeof(double), Q,
1028             0, NULL, NULL);
1029         free(Qnew);

```

```

1017
1018     conf=tdiag[0];
1019     for(i=0;i<N;i++){
1020         if(tdiag[i]<conf)
1021             conf=tdiag[i];
1022     }
1023     free(tdiag);
1024     free(tsubdiag);
1025     conf+=mu;
1026     if(fabs(conf-conf2)< 1e-13)
1027         chk=1;
1028     conf2=conf;
1029     //printf("%f\n", conf2);
1030 }
1031 *eig0=conf2;
1032
1033 //free memory
1034 free(P);
1035 free(diag);
1036 free(subdiag);
1037 free(A);
1038 free(Q);
1039
1040 //Deallocate resources
1041 clReleaseMemObject(A_buff);
1042 clReleaseMemObject(B_buff);
1043 clReleaseMemObject(C_buff);
1044 clReleaseKernel(mat_mul);
1045 clReleaseKernel(transpose);
1046 clReleaseKernel(transfer);
1047 clReleaseCommandQueue(queue);
1048 clReleaseProgram(program);
1049 clReleaseContext(context);
1050
1051 end = clock();
1052 *time=((double) (end - start)) / CLOCKS_PER_SEC;
1053 printf("OpenCL (CPU) finished!\n");
1054 printf("Execution time: %f\n",*time);
1055 printf("Lowest eigenvalue: %f\n\n",*eig0);
1056 }
```

A.2 Kernels

```

1 //enable support to double type
2 #pragma OPENCL EXTENSION cl_khr_fp64: enable
3
4
5 //2 component vector to hold the real and imaginary parts of a complex
6 //number:
7 typedef double2 cdouble;
8
9
10
11 /*
12 * Return Real (Imaginary) component of complex number:
13 */
14 inline double real(cdouble a){
15     return a.x;
```

```

16 }
17 inline double imag(cdouble a){
18     return a.y;
19 }
20
21 inline cdouble cc(cdouble z){
22     z.s1=-z.s1;
23     return z;
24 }
25
26
27 /*
28 * Multiply two complex numbers:
29 *
30 * a = (aReal + I*aImag)
31 * b = (bReal + I*bImag)
32 * a * b = (aReal + I*aImag) * (bReal + I*bImag)
33 *          = aReal*bReal +I*aReal*bImag +I*aImag*bReal +I^2*aImag*bImag
34 *          = (aReal*bReal - aImag*bImag) + I*(aReal*bImag + aImag*bReal)
35 */
36 inline cdouble cmult(cdouble a, cdouble b){
37     return (cdouble)( a.s0*b.s0 - a.s1*b.s1, a.s0*b.s1 + a.s1*b.s0);
38 }
39
40
41
42
43 // first part of the Lanczos algorithm iteration
44
45 __kernel void lanczos1(const double beta, __global cdouble *A, __global
46     cdouble *v, __global cdouble *w, __global double *alpha)
47 {
48     int id=get_global_id(0);
49     int gid=get_group_id(0);
50     int ls=get_local_size(0);
51     int lid=get_local_id(0);
52     int ngr=get_num_groups(0);
53
54     cdouble temp1,temp2;
55     temp2=0;
56     temp1=w[id]/beta;
57     v[id]=temp1;
58
59     int i;
60     barrier(CLK_LOCAL_MEM_FENCE);
61
62     for(i=0;i<N;i++){
63         temp2+=cmult(A[N*id+i],v[i]);
64     }
65     w[id]=temp2;
66
67     alpha[id]=real(cmult(cc(temp2),temp1));
68 }
69
70 // second part of the Lanczos algorithm iteration
71
72 __kernel void lanczos2(const double alpha,const double beta, __global
73     cdouble *v, __global cdouble *v2, __global cdouble *w, __global double
    *ab)
74 {

```

```

74
75     int id=get_global_id(0);
76     int gid=get_group_id(0);
77     int ls=get_local_size(0);
78     int lid=get_local_id(0);
79     int ngr=get_num_groups(0);
80
81     cdouble temp1;
82     temp1=w[id]-alpha*v[id]-beta*v2[id];
83     w[id]=temp1;
84
85     ab[id]=real(cmult(cc(temp1),temp1));
86 }
87
88
89
90 // transpose NxN matrix
91
92 __kernel void transpose(__global double *A,__global double *At)
93 {
94     int id=get_global_id(0);
95
96     int i,j;
97     i=id/N;
98     j=id-i*N;
99     At[j*N+i]=A[id];
100}
101
102
103 // multiply two NxN matrices
104
105 __kernel void mat_mul(__global double *A,__global double *B,__global
106     double *C)
107 {
108     int id=get_global_id(0);
109
110     int i,j,k;
111     double temp=0;
112     i=id/N;
113     j=id-i*N;
114     for(k=0;k<N;k++){
115         temp+=A[i*N+k]*B[k*N+j];
116     }
117     C[id]=temp;
118 }
119
120 // auxiliary function for QR iteration
121
122 __kernel void transfer(__global double *A,__global double *diag,__global
123     double *subdiag)
124 {
125     int id=get_global_id(0);
126
127     double temp,mu;
128     mu=A[N*N-1];
129     temp=A[id*N+id]-mu;
130     diag[id]=temp;
131     temp=A[id*N+id+1];
132     subdiag[id]=temp;
133     if(id==N-1)

```

```
133     subdiag[id]=mu;
134 }
```

B Schrödinger equation code

```

1 //////////////////////////////////////////////////////////////////
2 //
3 // The code generate the matrix corresponding to the 2D schrodinger
4 // equation with any potential.
5 // The matrix is written in the file mh.mtx in the matrix market
6 // format. Then the Lanczos algorithm from ViennaCL is called and
7 // the lowest eigenpair is found. The eigenvector is written in a
8 // suitable form to the file "psi.txt" so as to be plotted.
9 //
10 //////////////////////////////////////////////////////////////////
11 //
12 //
13 //
14 // include necessary system headers
15 #include <iostream>
16 #include <string>
17 #include <iomanip>
18 #include <fstream>
19 #include <cmath>
20
21 // tell viennaCL to use OpenCL
22 #define VIENNACL_WITH_OPENCL
23
24 // include basic scalar and vector types of ViennaCL
25 #include "viennacl/scalar.hpp"
26 #include "viennacl/vector.hpp"
27 #include "viennacl/matrix.hpp"
28 #include "viennacl/matrix_proxy.hpp"
29 #include "viennacl/compressed_matrix.hpp"
30 #include "viennacl/linalg/lanczos.hpp"
31 #include "viennacl/io/matrix_market.hpp"
32
33 // If you GPU does not support double precision, use 'float' instead of 'double':
34 typedef double ScalarType;
35
36
37 viennacl::vector<ScalarType> ground_state();
38
39
40 int main()
41 {
42     //Schrodinger equation generator
43
44     int N,i,j;
45     double delta, lim, temp1, V, x, y;
46     std::ofstream pot,sch;
47
48     N=401;
49     lim=4;
50     delta=2*(double)lim/(N-1);
51
52     pot.open("pot.txt");

```

```

53     sch.open("sch.mtx");
54     sch << N*N << "    " << N*N << "    " << N*N*5-2*(N-1)-4 << std::endl;
55
56     for(i=0;i<N;i++){
57         for(j=0;j<N;j++){
58
59             x=-lim+i*delta;
60             y=-lim+j*delta;
61             temp1=(double)1/(delta*delta);
62
63             //defining the potential energy
64
65             //V=0.5*x*x+0.5*y*y;                                //2D decoupled harmonic
66             //oscillator
67             V=-(x*x+y*y)+0.125*pow(x*x+y*y,2);           //central mexican hat
68             //V=-x*x+0.125*pow(x,4)-y*y+0.125*pow(y,4);      //double mexican
69             //hat
70             /*if(sqrt(x*x+y*y)>0.1)                         //regularised
71                 electrostatic potential
72                 V=-(double)3/sqrt(x*x+y*y);
73             else
74                 V=-10; */
75
76
77             if(i+1+j*N+1<=N*N)
78                 sch << i+j*N+1 << "    " << i+1+j*N+1 << "    " << temp1 << std::
79                 endl;
80             if(i+j*N>0)
81                 sch << i+j*N+1 << "    " << i-1+j*N+1 << "    " << temp1 << std::
82                 endl;
83             if(i+(j+1)*N+1<=N*N)
84                 sch << i+j*N+1 << "    " << i+(j+1)*N+1 << "    " << temp1 << std
85                 ::endl;
86             if(i+(j-1)*N+1>0)
87                 sch << i+j*N+1 << "    " << i+(j-1)*N+1 << "    " << temp1 << std
88                 ::endl;
89             sch << i+j*N+1 << "    " << i+j*N+1 << "    " << -V-4*temp1 << std::
90                 endl;
91
92             pot << x << "    " << y << "    " << V << std::endl;
93         }
94     }
95
96     sch.close();
97     pot.close();
98
99
100    //ground state computation
101    viennacl::vector<ScalarType> psi;
102    psi=ground_state();
103    std::cout << "Remember to change sign to the eigenvalue!" << std::endl;
104
105
106    //output
107    std::ofstream out;
108    out.open("psi.txt");
109    for(i=0;i<N;i++){
110        for(j=0;j<N;j++){
111            out << -lim+i*delta << "    " << -lim+j*delta << "    " << pow(psi[i
112                +j*N],2)/pow(delta,2) << std::endl;
113        }
114    }

```

```

105     }
106     out.close();
107
108     return 0;
109 }
110
111
112
113
114
115 /*
=====
116 Copyright (c) 2010-2016, Institute for Microelectronics,
117 Institute for Analysis and Scientific
118 Computing,
119 TU Wien.
120 Portions of this software are copyright by UChicago Argonne, LLC.
121 -----
122 ViennaCL - The Vienna Computing Library
123 -----
124 Project Head: Karl Rupp rupp@iue.tuwien.ac.at
125 (A list of authors and contributors can be found in the PDF manual)
126 License: MIT (X11), see file LICENSE in the base directory
127 -----
*/
128 viennacl::vector<ScalarType> ground_state(){
129
130     viennacl::vector<ScalarType> approx_eigenvector;
131     std::vector< std::map<unsigned int, ScalarType> > host_A;
132     if (!viennacl::io::read_matrix_market_file(host_A, "sch mtx"))
133     {
134         std::cout << "Error reading Matrix file" << std::endl;
135         return approx_eigenvector;
136     }
137
138     viennacl::compressed_matrix<ScalarType> A;
139     viennacl::copy(host_A, A);
140     viennacl::linalg::lanczos_tag ltag(0.75, // Select a power of 0.75
141                                         as the tolerance for the machine precision.
142                                         1, // Compute (approximations to) the 1
143                                         largest eigenvalues
144                                         1, // use full reorthogonalization
145                                         400); // Maximum size of the Krylov space
146
147     std::cout << "Running Lanczos algorithm (with eigenvectors). This might
148     take a while..." << std::endl;
149     viennacl::matrix<ScalarType> approx_eigenvectors_A(A.size1(), ltag.
150                                         num_eigenvalues());
151     std::vector<ScalarType> lanczos_eigenvalues = viennacl::linalg::eig(A,
152                                         approx_eigenvectors_A, ltag);
153
154     for (std::size_t i = 0; i < lanczos_eigenvalues.size(); i++)
155     {
156         std::cout << "Approx. eigenvalue " << std::setprecision(7) <<
157             lanczos_eigenvalues[i] << std::endl;
158         // test approximated eigenvector by computing A*v:
159         approx_eigenvector = viennacl::column(approx_eigenvectors_A,
160                                         static_cast<unsigned int>(i));
161         viennacl::vector<ScalarType> Aq = viennacl::linalg::prod(A,

```

```
    approx_eigenvector);
155   std::cout << " (" << viennacl::linalg::inner_prod(Aq,
156     approx_eigenvector) << " for <Av,v> with approx. eigenvector v"
157     << std::endl;
158 }
```