

Designing a high-performance boundary element library with OpenCL and Numba

T. Betcke

Department of Mathematics, University College London

M. W. Scroggs

Department of Engineering, University of Cambridge

Abstract—The Bempp boundary element library is a well known library for the simulation of a range of electrostatic, acoustic and electromagnetic problems in homogeneous bounded and unbounded domains. It originally started as a traditional C++ library with a Python interface. Over the last two years we have completely redesigned Bempp as a native Python library, called Bempp-cl, that provides computational backends for OpenCL (using PyOpenCL) and Numba. The OpenCL backend implements kernels for GPUs and SIMD optimised for CPUs. In this paper we will discuss the design of Bempp-cl, provide performance comparisons on different compute devices and discuss the advantages and disadvantages of OpenCL as compared to Numba.

■ **THE BEMPP BOUNDARY ELEMENT LIBRARY** (originally named BEM++) started in 2011 as a project to develop an open-source C++ library for the fast solution of boundary integral equations. The original release came with a simple Python wrapper to the C++ interface. Over time more and more functionality was moved into the Python interface. However, computationally intensive routines and the main data structures remained in C++. This proved a burden when efforts started to modernise the library to be able to make better use of Single-Instruction-Multiple-Data (SIMD) optimisation on CPUs and offload computation to GPUs. In 2018 we decided to rewrite the

computational core from scratch. The aims were to support explicit SIMD optimisation on CPUs with various instruction lengths, be able to offload computations to AMD, Intel, and Nvidia GPUs, and to base the complete codebase on Python. These aims naturally lead to the choice of building a Python library based around OpenCL (using the PyOpenCL interface) and Numba. The name of the new library was Bempp-cl. At the end of 2019 we released an initial version 0.1 that we used to test the library in a number of real-world scenarios. This was followed later in 2020 by version 0.2, the first release that we considered feature complete and mature for application use.

Since then we have used Bempp-cl in a number of practical application and many of our users are slowly migrating to it from the old C++ based Bempp. In this article we want to discuss the design choices behind Bempp-cl and provide a number of performance benchmarks on different compute devices, including CPUs, AMD, Intel, and Nvidia GPUs.

[To be continued with some more details about content and overview of the Sections]

Boundary element methods with Bempp

In this section we want to provide a brief introduction to boundary element methods and describe the necessary steps for their numerical discretisation and solution.

The most simple boundary integral equation is of the form

$$\int_{\Gamma} g(\mathbf{x}, \mathbf{y}) \phi(\mathbf{y}) ds(\mathbf{y}) = f(\mathbf{x}), \quad \mathbf{x} \in \Gamma. \quad (1)$$

The function $g(\mathbf{x}, \mathbf{y})$ is a Green's function, f is a given right-hand side and ϕ is an unknown surface density over the boundary Γ of a bounded three dimensional domain $\Omega \subset \mathbb{R}^3$. To provide a concrete example. To compute the electrostatic capacity of an object Ω one solves the above equation with $f(\mathbf{x})$ the constant function 1 and $g(\mathbf{x}, \mathbf{y}) = \frac{1}{4\pi|\mathbf{x}-\mathbf{y}|}$. The normalized capacity is then obtained as $c = \frac{1}{4\pi} \int_{\Gamma} \phi(\mathbf{x}) ds(\mathbf{x})$. Many practical problems have a significantly more complex structure and can involve block systems of integral equations. Nevertheless, the fundamental structure of what Bempp-cl does is well described by the above problem and will be described in the following.

The first step is the **definition of the surface** Γ . Surfaces are represented in Bempp-cl as a triangulation into flat triangles (see 1). The triangulation is internally represented as an array of node coordinates and an associated connectivity array of node indices that define each triangle. In this step also topology data is computed. In particular, for each triangle the neighboring triangles and the type of intersection (i.e. via joint node or edge) is computed.

Once a triangulation is given we need to define the necessary data structures for the discretisation. Bempp-cl uses a Galerkin discretisation.

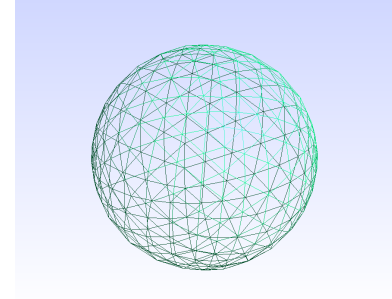


Figure 1. Discretisation of a sphere into flat surface triangles.

Consider (1). We represent the solution ϕ as $\phi = \sum_{j=1}^N x_j \phi_j$, where the ϕ_i are basis functions of a **trial space**, defined in the most simple case as constant 1 on the triangle τ_j and 0 everywhere also. Several other types of discretisation spaces are available in Bempp. Moreover, we require a **test (or dual) space** for the discretisation. Let $\psi = \sum_{i=1}^N y_i \psi_i$, where each ψ_i is a basis function of a space of test functions. For simplicity, here we will assume that the ψ_i are also piecewise constant. The discrete representation of the above problem then takes the form

$$Ax = b$$

with

$$A_{ij} = \int_{\Gamma} \psi_i(\mathbf{x}) \int_{\Gamma} g(\mathbf{x}, \mathbf{y}) \phi_j(\mathbf{x}) ds(\mathbf{y}) ds(\mathbf{x})$$

and $f_i = \int_{\Gamma} \psi_i(\mathbf{x}) f(\mathbf{x}) ds(\mathbf{x})$. In the case of piecewise constant trial and test functions the definition of A_{ij} simplifies to $A_{ij} = \int_{\Gamma} \int_{\Gamma} g(\mathbf{x}, \mathbf{y}) ds(\mathbf{y}) ds(\mathbf{x})$.

In Bempp-cl an operator definition consists of the type of the operator (e.g. Laplace single-layer in the above example), the definition of the trial and test spaces, and the definition of the range space. The range space is required for operator products and not relevant for the purpose of this paper. The function f is represented as a **Grid function object**, which consists of either the dual representation in the forms of the integrals $b_i = \int_{\Gamma} \psi_i(\mathbf{x}) f(\mathbf{x}) ds(\mathbf{x})$ or directly through its coefficients f_j in the representation $f = \sum_{j=1}^N f_j \phi_j$.

Once the grid, the space objects, and the operator are defined, the main computational step

is performed, namely the **computation of the discrete matrix entries** A_{ij} . For pairs of triangles τ_i and τ_j that do not share a joint edge or vertex this is done through a simple numerical quadrature rule that approximates $A_{ij} \approx \sum_{\ell} \sum_q g(\mathbf{x}_{\ell}, \mathbf{y}_q) \psi_i(\mathbf{x}_{\ell}) \phi_j(\mathbf{y}_q) \omega_{\ell} \omega_q$, where the \mathbf{x}_{ℓ} and \mathbf{y}_q are quadrature points in the corresponding triangles τ_i and τ_j , and the values ω_i and ω_j are the quadrature weights. In the case that two triangles share a joint vertex/edge or the triangles τ_i and τ_j are identical, corresponding singularity adapted quadrature numerical quadrature rules are used that are based on singularity removing coordinate transformations.

The values b_i of the right-hand side vector b are similarly computed through a numerical quadrature rule.

In the final step, **Bempp solves the underlying linear system of equations** either through a direct LU decomposition or through iterative solvers such as Gmres. The solution can then be evaluated away from the surface Γ through domain potential operators and exported in various formats for visualization.

In summary, to solve a boundary integral equation problem, the following steps are performed by Bempp

- 1) Import of the surface description as triangulation data.
- 2) Definition of the trial space, test space and range space.
- 3) Discretization into a matrix problem $Ax = b$.
- 4) Solution of the matrix problem by either a direct or iterative solver.
- 5) Evaluation of domain potential operators for visualization and post-processing.

All of these steps are accelerated through the use of either Numba or OpenCL. In the following section we provide a high-level overview of the library and how these acceleration techniques are deployed before we deep dive into the design of the computational kernels.

A high level overview of Bempp-cl

The main user visible component of Bempp-cl is the module **bempp.api**, which defines all user interface functions and other high-level routines. In particular, it contains the definition of the main

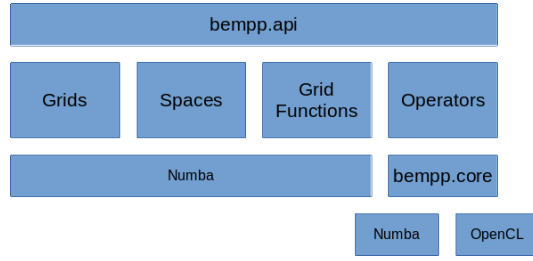


Figure 2. The layout of Bempp-cl with its computational backends.

object types: **Grids, Spaces, Gridfunctions, Operators**. All computational backend routines are part of **bempp.core**. We are currently supporting a full Numba backend and an OpenCL backend. An overview of this structure is provided in Figure 2.

Outside the operator discretisation Numba is used in the following contexts.

- Topological computations for the grid. Finding out all neighbour relationships between triangles requires parses through the grid data structure with linear complexity, which is accelerated with Numba.
- Definition of local-to-global maps for function spaces. Again, this requires traversal through the grid and assigning relationships between global and local indices.
- Grid Functions. A right-hand side function f can be defined as a Python callable. This is just-in-time compiled via Numba and then the product with the corresponding basis functions integrated in each triangle via numerical quadrature, again via Numba accelerated routines.
- Sparse matrices. Sparse surface matrices are assembled through Numba accelerated routines.

The main computational complexity in discretising the boundary integral equation (1) is coming from the discretisation of the left-hand side integral operator. Using dense methods this has quadratic complexity in the number of surface triangles.

Once the user defines an operator with its associated discretisation spaces and calls the **weak_form** method to discretize it the code calls a regular integrator to assemble all the

interactions between non-adjacent elements and a singular integrator to compute the interactions between adjacent triangles (only necessary if trial and test space are defined on the same grid). The result is stored internally as a dense Numpy array. The corresponding discretisation routines are proxies that forward to computational routines in either Numba or OpenCL, depending on the user preferences. For OpenCL assembly the code checks additional parameters, such as the default vector length for SIMD operators (e.g. 4 for double precision and 8 for single precision in Intel AVX2, or 1 if a GPU is used), and whether the discretisation should proceed in single or double precision. The OpenCL kernel is then compiled for the underlying compute device using PyOpenCL and executed. If the computational backend is Numba, the call is forwarded to the corresponding Numba discretisation routines and executed. For simple piecewise constant function spaces or other spaces, where the support of each basis function is localized to a single triangle, only one call to the computational routines is necessary. If the support of basis functions is larger than a single triangle different threads may need to sum into the same global degree of freedom. We discuss this in more detail in Section ??.

Assembling boundary integral operators with OpenCL

In this section we discuss in more detail the assembly of boundary integral operators with OpenCL and how we integrated this into our Python workflow. We start with a brief introduction to OpenCL and then dive into how we use OpenCL as part of Bempp-cl.

What is OpenCL?

OpenCL (<https://www.khronos.org/opencl/>) is a heterogeneous compute standard for CPUs, GPUs, FPGAs, and other types of devices that provide conformant drivers. At its core OpenCL executes compute kernels that can be written in C99, or more recently, in C++. The current version of OpenCL is 3.0, though the most widely implemented standard is OpenCL 1.2, which Bempp-cl uses. OpenCL splits the computational tasks into work-items, each of which represents a single execution of a compute kernel. Work items

are grouped together into work-groups, which share local memory. Barrier synchronization is only allowed within a work-group. All work-items are uniquely indexed by a one, two, or three dimensional index space, called NDRange. Kernels are launched onto a compute device (e.g. CPU or GPU) from a host device. OpenCL allows kernels to be loaded as strings and compiled on the fly for a given device, making it well suited for launching from high-productivity languages. To launch an OpenCL kernel the following steps need to be performed.

- Initialize buffers on the compute device.
- Compute relevant data from host to the compute device.
- Load a kernel string and just-in-time compile it for the device.
- Execute the kernel for the given buffers and chosen index space.
- Copy results back to the host.

An excellent feature of OpenCL is its very good support for vectorized operations. OpenCL provides vector data types, such as **double4**, that can hold four double values in a SIMD register and defines a number of standard operations for these vector types. This makes it easy to explicitly target modern SIMD execution in a portable way while avoiding difficult compiler intrinsics and keeping kernel code readable.

Python has excellent OpenCL support through the PyOpenCL library by Andreas Kloeckner, which automates much of the initialization of the OpenCL environment and makes it easy to create buffers and launch OpenCL kernels from Python.

OpenCL Assembly in Bempp-cl

Bempp-cl has OpenCL kernels for all its boundary operators. All operators have the same interface and are launched in the same way. In the first step the relevant data will need to be copied to the compute device. This data consists of:

- Test and trial indices of triangles over which to be integrated.
- Signs of the normal directions for the spaces.
- Test and trial grid as flat floating point array, defining each triangle through nine floating point numbers, specifying the (x, y, z) coordinates of each of the three nodes of a triangle.

```

#include "bempp_base_types.h"
#include "bempp_helpers.h"
#include "bempp_spaces.h"
#include "kernels.h"

__kernel void kernel_function(
    __global uint* testIndices, __global uint* trialIndices,
    __global int *testNormalSigns, __global int *trialNormalSigns,
    __global REALTYPE* testGrid, __global REALTYPE* trialGrid,
    __global uint* testConnectivity, __global uint* trialConnectivity,
    __global uint* testLocal2Global, __global uint* trialLocal2Global,
    __global REALTYPE* testLocalMultipliers,
    __global REALTYPE* trialLocalMultipliers, __constant REALTYPE* quadPoints,
    __constant REALTYPE* quadWeights, __global REALTYPE* globalResult,
    __global REALTYPE* kernel_parameters,
    int nTest, int nTrial, char gridsAreDisjoint) {
    /* Variable declarations */

```

Figure 3. Definition of the OpenCL compute kernel for scalar integral equations.

- Test and trial connectivity information that stores the indices of each node for each triangle.
- Test and trial mappings of local triangle degrees of freedom to global degrees of freedom.
- Test and trial basis function multipliers for each triangle.
- Quadrature points and quadrature weights.
- A buffer that contains the global assembled matrix.
- Additional kernel parameters, such as the wavenumber for Helmholtz problems.
- Number of test and trial degrees of freedom.
- A single byte that is set to one if the test and trial grids are disjoint.

The actual kernel definition is shown in Figure 3. Before the kernel can be launched, it needs to be configured and just-in-time compiled. Kernel configuration happens through C-style preprocessor definitions that are passed through the just-in-time compiler. These include the names of the test and trial space, the name of the function that evaluates the integral kernel, choosing between single and double precision types, and for SIMD enhanced kernel the vector length of the SIMD types. For example, in Figure 3 all floating point types have the name `REALTYPE`. This is substituted with either `float` or `double` during just-in-time compilation.

Each work-item computes via numerical quadrature all interactions via the integral kernel of basis functions on the trial element with basis functions on the test element. Before summing

the result into the global result buffer, the kernel checks via the connectivity information if the test and trial triangle are adjoint or identical. For this it simply checks if at least one of the node indices of the test triangle is identical to one of the node indices of the trial triangle. If this is true and the grids are not disjoint, the result of the kernel is discarded and not summed back into the global result buffer. The reason is that for adjoint triangles a separate singular quadrature rule is used. The effect is that a few work-items do work that is discarded. However, in a grid with N elements the number of triangle pairs requiring a singular quadrature rule is $\mathcal{O}(N)$, while the total number of interactions is N^2 . Hence, only a tiny fraction of work-items are discarded.

SIMD optimized kernels

The above description was valid for kernels that do not take advantage of CPU SIMD operations are kernels that run directly on a GPU device. For SIMD optimization on CPUs we proceed slightly differently. The corresponding kernel in principle works as described above, but batches together the interaction on one test triangle with X trial triangles, where X is either 4, 8, or 16, depending on the number of available SIMD lanes. This strategy allows us to optimize almost all floating point operations within a kernel run for SIMD operation. If the number of trial elements is not divisible by 4, 8, or 16, then the few remaining trial elements are assembled with the standard-non vectorized kernel.

Each kernel definition is stored in two vari-

ants, one with the ending `_novec.cl` and another one with the ending `_vec.cl`. The vectorized variant is configured via preprocessor directives for the desired number of vector lanes. Having to develop two OpenCL kernel codes for each operator is a certain amount of overhead. If we write a new kernel, we first implement the non-vectorized version. Then, with the help of preprocessor directives and a number of helper functions that do the actual implementation of operations depending on whether the kernel is vectorized or not, it is usually a matter of an hour or two to convert the non-vectorized kernel into a vectorized version.

Alternatively, some CPU OpenCL runtime environments have the option to try and auto-vectorize kernels by batching together work items on SIMD lanes, similarly to what we do manually. In our experience this works well for very simple kernels but often fails for more complex OpenCL kernels. This is why we decided to implement this strategy manually.

A completely different SIMD strategy is based on batching together quadrature evaluations within a single test/triangle pair. There are two disadvantages to this approach. First, it only works well if the number of quadrature points is a multiple of the available SIMD lanes. Second, other operations such as the geometry calculations for each element then can not be SIMD optimized as these are only performed once per test/triangle pair.

Assembling the singular part of integral operators

The assembly of the singular part of an integral operator works a bit differently. Remember that the singular part consists of triangle parts which are adjacent to each other or identical. There are $\mathcal{O}(N)$ of such pairs. We are using fully numerical quadrature rules for these integrals that are based on subdividing the four-dimensional integration domain and using transformation techniques to remove the singularities. This gives highly accurate evaluations of these integration pairs but requires a large number (typically over 1000) quadrature points per triangle pair. Hence, we create one workgroup for each triangle pair and inside this workgroup have a number of work-items that are evaluating the quadrature

rules and then at the end sum up the result together. Depending on how two triangles are in relationship to each other different types of singular quadrature rules are needed. We solve this pre-loading all possible quadrature rules onto the device and also store for each triangle pair an index into the required quadrature rule so that the kernel function can select the correct rule to evaluate. For the singular quadrature rules we did not implement separate SIMD optimized kernels as the proposed implementation is already highly efficient and requires only a fraction of the computational time of the regular quadrature rules described above. At the end, the singular integral values are either summed into the overall result matrix, or if desired by the user, stored as separate sparse matrix.

Numba assembly of integral operators

The Numba assembly of integral operators is structured differently from the OpenCL assembly. In OpenCL we create a large number of work-items that are executed in parallel on the compute device by the runtime system with the details of the parallel execution left to the device driver. For Numba we use classical loop based parallelism. Each loop iteration is the assembly of one test element with all trial elements. Within each loop we partially optimized for SIMD instruction, though this depends on the autovectorization of the just-in-time compiler and cannot be explicitly targeted in Numba. Again, we separate regular and singular interactions. In the next section we give some more details on SIMD optimization in Numba at the example of the Laplace kernel.

SIMD optimized evaluation of the Laplace kernel with Numba