

LIBXSMM

Library for small matrix-matrix multiplications targeting Intel Architecture (x86). The library generates code for the following instruction set extensions: Intel SSE3, Intel AVX, Intel AVX2, IMCI (KNCni) for Intel Xeon Phi coprocessors (“KNC”), and Intel AVX-512 as found in the Intel Xeon Phi processor family (“KNL”) and future Intel Xeon processors. Historically the library was solely targeting the Intel Many Integrated Core Architecture “MIC”) using intrinsic functions, however meanwhile optimized assembly code is generated for the fore mentioned instruction set extensions.

What is a small matrix-matrix multiplication? When characterizing the problem size using the M, N, and K parameters, a problem size suitable for LIBXSMM falls approximately within $(M \ N \ K)^{(1/3)} \leq 60$ (which illustrates that non-square matrices or even “tall and skinny” shapes are covered as well). However the code generator only generates code up to the specified threshold. Raising the threshold may not only generate excessive amounts of code (due to unrolling), but also missing to implement a tiling scheme to effectively utilize the L2 cache. For problem sizes above the configurable threshold, LIBXSMM is falling back to BLAS.

How to determine whether an application can benefit from using LIBXSMM or not? Given the application uses BLAS to carry out matrix multiplications, one may link against Intel MKL 11.2 (or higher), set the environment variable MKL_VERBOSE=1, and run the application using a representative workload (env MKL_VERBOSE=1 ./workload > verbose.txt). The collected output is the starting point for evaluating the problem sizes as imposed by the workload (grep -a “MKL_VERBOSE DGEMM” verbose.txt | cut -d, -f3-5).

Interface

The interface of the library is *generated* according to the Build Instructions (therefore the header file ‘include/libxsmm.h’ is **not** stored in the code repository). The generated interface also defines certain preprocessor symbols to store the properties the library was built for.

To perform the matrix-matrix multiplication $cm \times n = cm \times n + am \times k * bk \times n$, one of the following interfaces can be used:

```
/** If non-zero function pointer is returned, call (*function)(M, N, K). */
libxsmm_smm_function libxsmm_smm_dispatch(int m, int n, int k);
libxsmm_dmm_function libxsmm_dmm_dispatch(int m, int n, int k);
/** Automatically dispatched matrix-matrix multiplication. */
void libxsmm_smm(int m, int n, int k, const float* a, const float* b, float* c);
void libxsmm_dmm(int m, int n, int k, const double* a, const double* b, double* c);
/** Non-dispatched matrix-matrix multiplication using inline code. */
void libxsmm_simm(int m, int n, int k, const float* a, const float* b, float* c);
void libxsmm_dimm(int m, int n, int k, const double* a, const double* b, double* c);
/** Matrix-matrix multiplication using BLAS. */
void libxsmm_sblasmm(int m, int n, int k, const float* a, const float* b, float* c);
void libxsmm_dblasmm(int m, int n, int k, const double* a, const double* b, double* c);
```

With C++ function overloading, the library allows to omit the ‘s’ and ‘d’ denoting the numeric type in the above C interface. Further, a type ‘libxsmm_mm_dispatch<type>’ can be used to instantiate a functor rather than making a distinction for the numeric type in ‘libxsmm_?mm_dispatch’.

Build Instructions

To compile the library run:

```
make
```

The interface is produced inside of the ‘include’ directory. The library archives are produced inside of the ‘lib’ directory with the ‘mic’ subdirectory containing the native library and the ‘intel64’ folder storing the hybrid archive containing host and MIC code.

To remove intermediate files use:

```
make clean
```

or to remove all generated files including the interface and library archive files:

```
make realclean
```

The usual `make install` is simply a shortcut for `make; make clean`.

The library can be configured to accept row-major (default) or column-major order matrices. This is accomplished using the variable `ROW_MAJOR` (0 for column-major, and row-major order otherwise). The following configures the column-major format:

```
make ROW_MAJOR=0
```

By default, LIBXSMM is not optimized for particular matrix sizes (M, N, and K values). Specializing the library for certain matrix sizes (and therefore optimizing the performance) can be achieved in the following way:

```
make M="2 4" N="1" K="$(echo $(seq 2 5))"
```

The above example is generating the following set of (M,N,K) values:

```
(2,1,2), (2,1,3), (2,1,4), (2,1,5),  
(4,1,2), (4,1,3), (4,1,4), (4,1,5)
```

The index sets are in a loop-nest relationship (M(N(K))) when generating the indices. Moreover, an empty index set resolves to the next non-empty outer index set of the loop nest (including to wrap around from the M to K set). An empty index set is not participating anymore in the loop-nest relationship. Here is an example of generating multiplication routines which are “squares” with respect to M and N (since N inherits the current value of the “M loop”):

```
make M="$(echo $(seq 2 5))" K="$(echo $(seq 2 5))"
```

An even more flexible specialization is possible by using the MNK variable when building the library. It takes a list of indices which are eventually grouped (using commas):

```
make MNK="2 3, 23"
```

Each group of indexes is combined into all possible triplets generating the following set of (M,N,K) values:

```
(2,2,2), (2,2,3), (2,3,2), (2,3,3),  
(3,2,2), (3,2,3), (3,3,2), (3,3,3), (23,23,23)
```

Testing the generated cases means capturing the console output of the cp2k code sample:

```
make MNK="2 3, 23" test
```

The recorded output file can be further evaluated. For example:

```
grep "diff" samples/cp2k/cp2k-perf.txt | grep -v "diff=0.000"
```

Performance

Tuning

The build system allows to conveniently select the target system using an AVX flag when invoking “make”. The default is to generate code according to the feature bits of the host system running the compiler. The device-side defaults to “MIC” targeting the Intel Xeon Phi family of coprocessors (“KNC”). However beside of generating all supported host code path (and letting the compiler pick the one representing the host), specifying a particular code path will not only save some time when “printing the code” but also enabling cross-compilation for a target different that is different from the compiler’s host: SSE=3 (in fact !=0), AVX=1, AVX=2 (with FMA), and AVX=3 are supported. The latter is targeting the Intel Knights Landing processor family (“KNL”) and future Intel Xeon processors using Intel AVX-512 foundational instructions (AVX-512F):

```
make AVX=3
```

The library supports generating code using an “implicitly aligned leading dimension” for the destination matrix of a multiplication. The latter is enabling aligned store instructions, and also hints the inlinable code accordingly. The client code may be arranged accordingly by checking the build parameters of the library (at compile-time using the preprocessor). Aligned store instructions imply a leading dimension which is a multiple of the default alignment:

```
make ALIGNED_STORES=1
```

The default alignment (ALIGNMENT=64) as well as a non-default alignment for the store instructions (ALIGNED_STORES=n) can be specified when invoking “make”. The “implicitly aligned leading dimension” optimization is not expected to have a big impact due to the relatively low amount of store instructions in the instruction mix. In contrast, supporting an “implicitly aligned leading dimension” for loading the input matrices is supposed to make a bigger impact, however this is not yet exposed by the build system because: (1) aligning a batch of input matrices implies usually larger code changes for the client code whereas accumulating into a local temporary destination matrix is a relatively minor change, and (2) today’s Advanced Vector Extensions including the AVX-512 capable hardware supports unaligned load/store instructions.

Auto-dispatch

The function ‘libxsmm_?mm_dispatch’ helps amortizing the cost of the dispatch when multiple calls with the same M , N , and K are needed. In contrast, the automatic code dispatch uses three levels:

1. Specialized routine,
2. Inlinable C code, and
3. BLAS library call.

All three levels are accessible directly (see Interface) in order to allow a customized code dispatch. The level 2 and 3 may be supplied by the Intel Math Kernel Library (Intel MKL) 11.2 DIRECT CALL feature. Beside of the generic interface, one can call a specific kernel e.g., ‘libxsmm_dmm_4_4_4’ multiplying 4x4 matrices. For the latter, the code generator includes prototypes for all specialized functions into the generated interface (header file).

Further, a preprocessor symbol denotes the largest problem size ($M \times N \times K$) that belongs to level (1) and (2), and therefore determines if a matrix-matrix multiplication falls back to level (3) of calling the LAPACK/BLAS library linked with the LIBXSMM. This threshold can be configured using for example:

```
make THRESHOLD=$((24 * 24 * 24))
```

The maximum of the given threshold and the largest requested specialization refines the value of the threshold. If a problem size falls below the threshold, dispatching the code requires to figure out whether a specialized routine exists or not. This can be implemented by bisecting a table of all specialized functions (binary search). At the expense of storing function pointers for the entire problem space below the threshold, a direct lookup can be used instead. The actual behavior can be configured using for example:

```
make SPARSITY=2
```

A binary search is implemented when a sparsity (calculated at construction time of the library) is above the given SPARSITY value. Raising the given value prevents generating a binary search (and generates a direct lookup) whereas a value below or equal one is generating the binary search. Furthermore, the size of the direct lookup table is limited to 512 KB (currently hardcoded). The overhead of auto-dispatched multiplications based on the binary search becomes negligible with reasonable problem sizes (above ~20x20 matrices), but may be significant for very small auto-dispatched matrix-matrix multiplication.

Results

The library does not claim to be “optimal” or “best-performing”, and the presented results are modeling certain applications which are not representative “in general”. Instead, information on how to reproduce the results are given for each of the reported cases (see figure 1-3).

Implementation

Limitations

Beside of the inlinable code path, the library is currently limited to a single code path which is selected at build time of the library. Without a specific flag (SSE=1, AVX=1|2|3), the assembly code generator emits code for all supported instruction set extensions whereas the Intrinsic code generator (GENASM=0) is actually covering only IMCI (KNCni) and Intel AVX-512F. However, the compiler is picking only one of the generated code paths according to its code generation flags (or according to what is native with respect to the compiler-host). A future version of the library may be including all code paths at build time and allow for runtime-dynamic dispatch of the most suitable code path.

A future version of the library may support an auto-tuning stage when generating the code (to find M,N,K -combinations for specialized routines which are beneficial compared to the code generated from the inlinable C code path). Auto-tuning the compiler code generation using a profile-guided optimization may be another option to be incorporated into the build system (Makefile). However, auto-tuning also imposes a reasonable burden (similar to profile guided optimization) and the current library is supposed to generate extremely efficient code for everything that does not rely on L2 cache blocking (below THRESHOLD), and except for very small problem sizes where a dedicated kernel may beat our implementation (e.g., $MNK=4 \times 4 \times 4$).

Roadmap

Although the library is under development, the published interface is rather stable and may only be extended in future revisions. The following issues are being addressed in upcoming revisions:

- Full xGEMM interface, and native FORTRAN interface
- Just-in-Time (JIT) runtime dynamic code generation
- API supporting sparse matrices and other cases

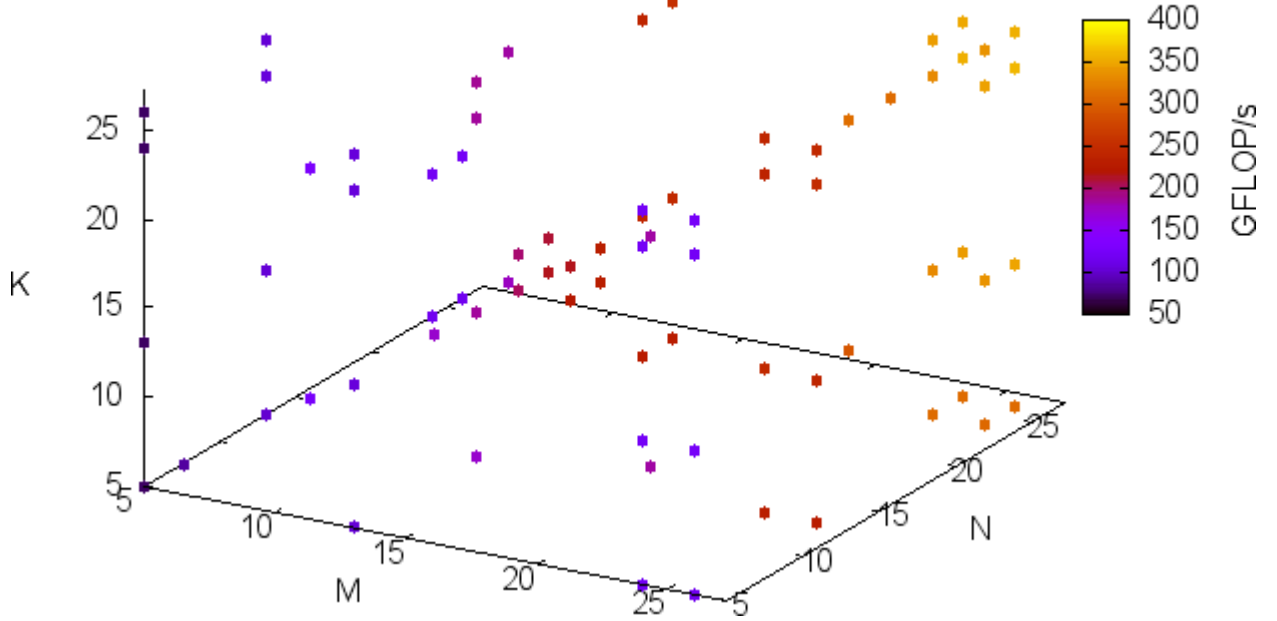


Figure 1: This performance plot for a dual-socket Intel Xeon E5-2697v3 (“Haswell”) shows a “compact selection” (to make the plot visually more appealing) out of 328 specializations as useful for CP2K Open Source Molecular Dynamics [1]. The code has been generated and built by running “./make.sh -cp2k AVX=2 -j”. This and below plots were generated by running “cd samples/cp2k ; ./cp2k-plot.sh specialized cp2k-specialized.png -1”. Please note that larger problem sizes (MNK) carry a higher arithmetic intensity which usually leads to higher performance (less bottlenecked by memory bandwidth).

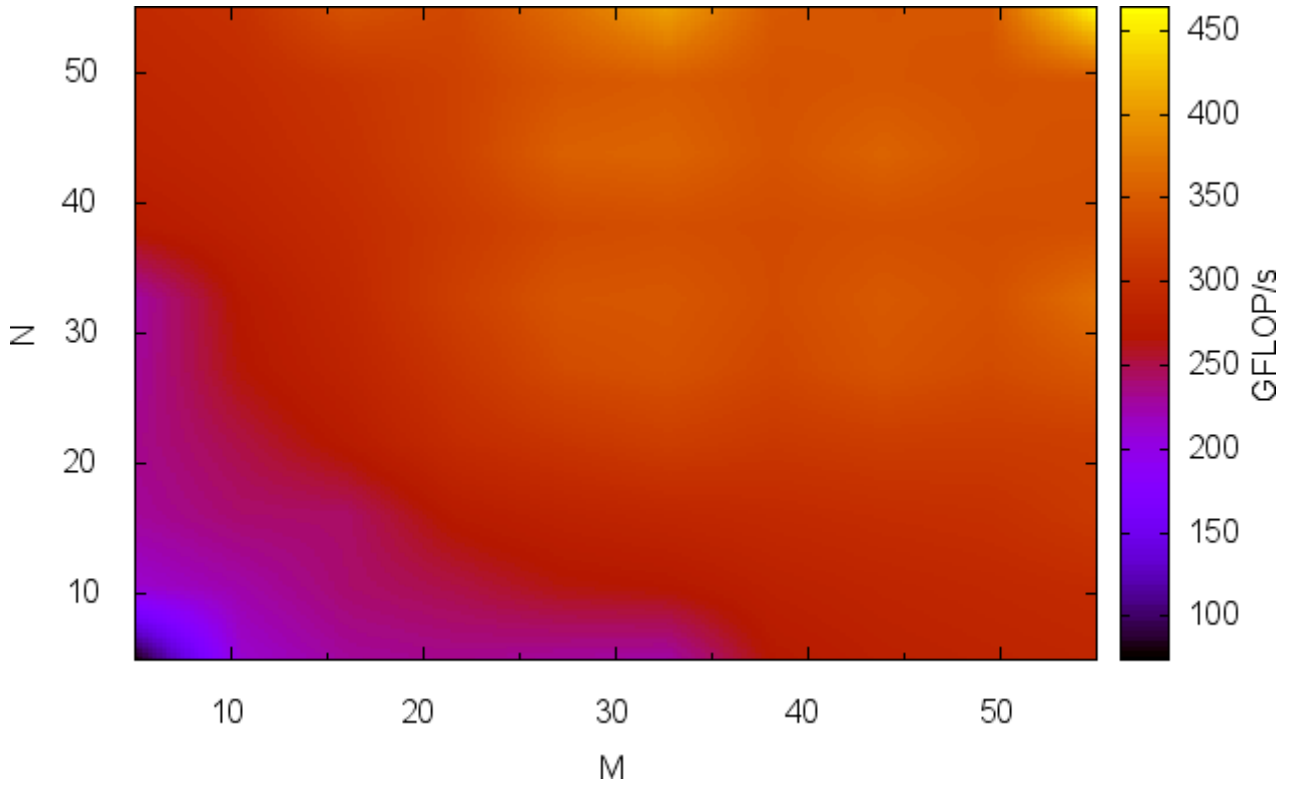


Figure 2: This plot summarizes the performance of the generated kernels by averaging the results over K (and therefore the bar on the right hand side may not show the same maximum when compared to other plots). The performance is well-tuned across the parameter space with no “cold islands”, and the lower left “cold” corner is fairly limited. Please refer to the first figure on how to reproduce the results.

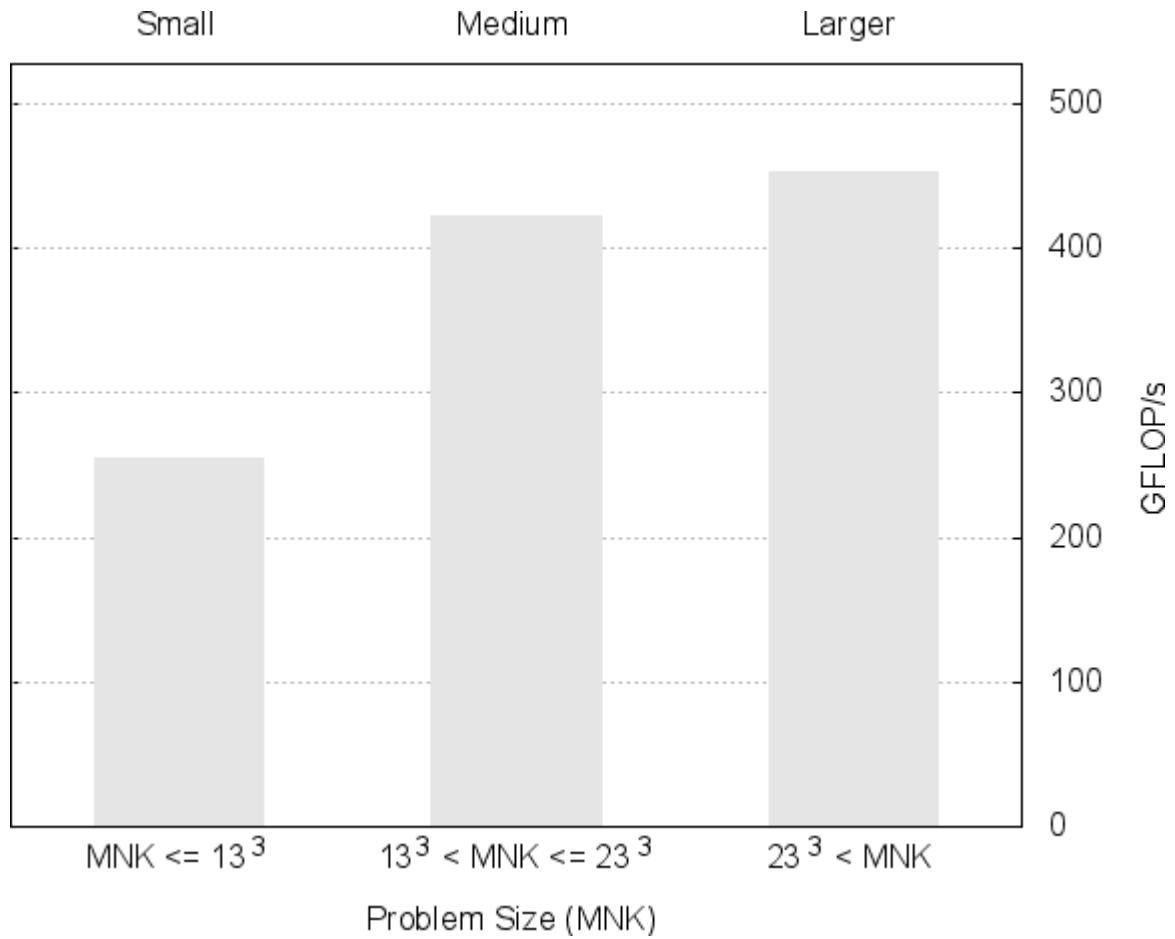
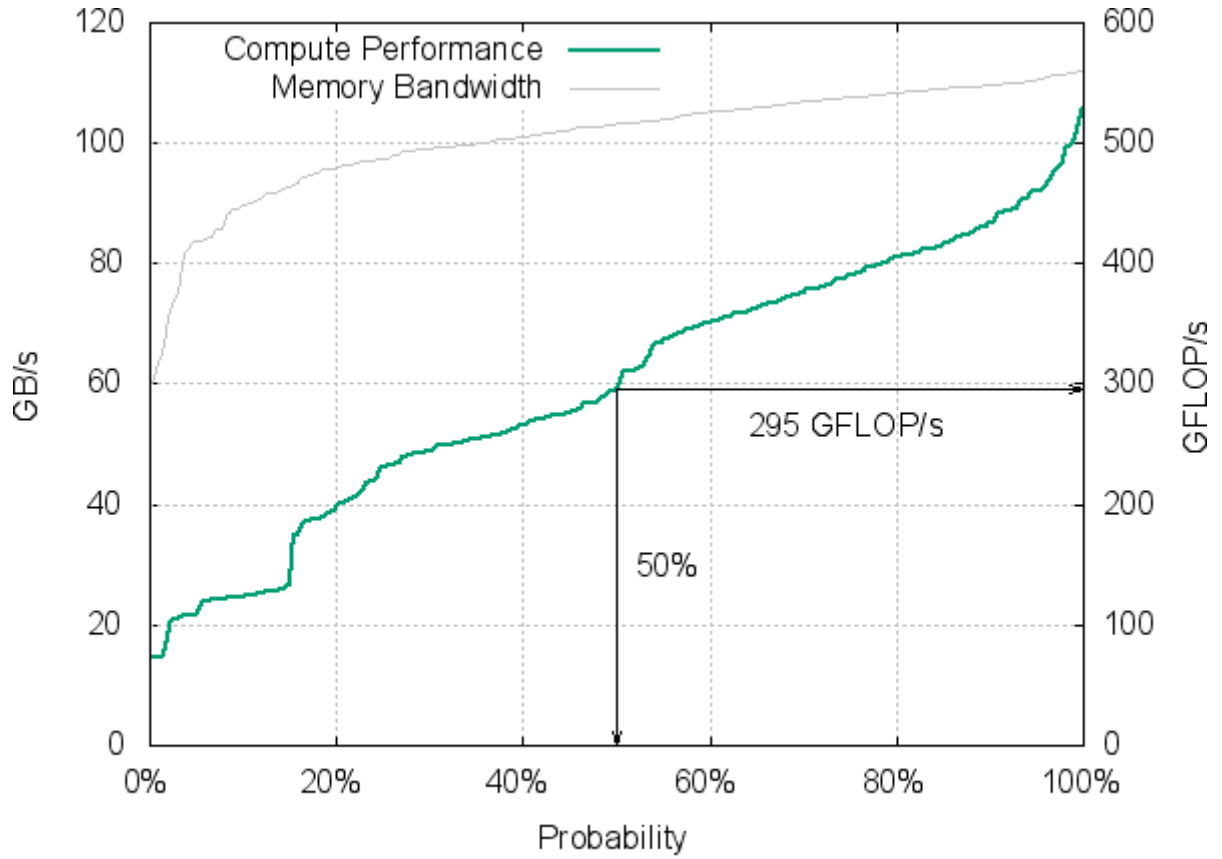


Figure 3: This plot shows the arithmetic average of the performance with respect to groups of problem sizes (MNK). The problem sizes are binned into three groups according to the shown intervals: “Small”, “Medium”, and “Larger” (notice that “larger” may still not be a large problem size). Please refer to the first figure on how to reproduce the results.



Minimum: 73 GFLOP/s Geo. Mean: 274 GFLOP/s Median: 296 GFLOP/s Maximum: 528 GFLOP/s

Figure 4: In order to further summarize the previous plots, this graph shows the cumulative distribution function (CDF) of the performance across all cases. Similar to the median value at 50%, one can read for example that 100% of the cases are yielding less or equal the largest discovered value. The value highlighted by the arrows is usually the median value, the plot script however attempts to highlight a single “fair performance value” representing all cases by linearly fitting the CDF, projecting onto the x-axis, and taking the midpoint of the projection (usually at 50%). Please note, that this diagram shows a statistical distribution and does not allow to identify any particular kernel. Moreover at any point of the x-axis (“Probability”), the “Compute Performance” and the “Memory Bandwidth” graph do not necessarily belong to the same kernel! Please refer to the first figure on how to reproduce the results.

Applications and References

- [1] <http://cp2k.org/>: Open Source Molecular Dynamics which (optionally) uses LIBXSMM. The application is generating batches of small matrix-matrix multiplications (“matrix stack”) out of a problem-specific distributed block-sparse matrix (see <https://github.com/cp2k/cp2k/tree/intel>).
- [2] <http://www.seissol.org/>: SeisSol is one of the leading codes for earthquake scenarios, in particular for simulating dynamic rupture processes. LIBXSMM provides highly optimized assembly kernels which form the computational back-bone of SeisSol (see https://github.com/TUM-I5/seissol_kernels/tree/lts_compressed).
- [3] <https://github.com/TUM-I5/GemmCodeGenerator>: Code generator for matrix-matrix multiplications used as an infrastructure to develop LIBXSMM’s assembly code generator.
- [4] <http://software.intel.com/xeonphicatalog>: Intel Xeon Phi Applications and Solutions Catalog.
- [5] <http://goo.gl/qsnOOf>: Intel 3rd Party Tools and Libraries.