

# LIBXSMM

LIBXSMM is a library for small dense and small sparse matrix-matrix multiplications as well as for deep learning primitives such as small convolutions targeting Intel Architecture (x86). The library is generating code for the following instruction set extensions: Intel SSE, 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 Intel Xeon processors (Skylake-EP “SKX”). Small convolutions are currently only optimized for Intel AVX-512. Historically the library was solely targeting the Intel Many Integrated Core Architecture “MIC”) using intrinsic functions, meanwhile optimized assembly code is targeting all aforementioned instruction set extensions (static code generation), and Just-In-Time (JIT) code generation is targeting Intel AVX and beyond.

**What is the background of the name “LIBXSMM”?** The “MM” stands for Matrix Multiplication, and the “S” clarifies the working domain i.e., Small Matrix Multiplication. The latter also means the name is neither a variation of “MXM” nor an eXtreme Small Matrix Multiplication but rather about Intel Architecture (x86) - and no, the library is 64-bit only. The spelling of the name might follow the syllables of libx\smm, libx’smm, or libx-smm.

**What is a small matrix multiplication?** When characterizing the problem size using the M, N, and K parameters, a problem size suitable for LIBXSMM falls approximately within  $(M N K)^{1/3} \leq 80$  (which illustrates that non-square matrices or even “tall and skinny” shapes are covered as well). The library is typically used to generate code up to the specified threshold. Raising the threshold may not only generate excessive amounts of code (due to unrolling in M and K dimension), but also miss to implement a tiling scheme to effectively utilize the cache hierarchy. For auto-dispatched problem sizes above the configurable threshold, LIBXSMM is falling back to BLAS.

**What about “medium-sized” matrix multiplication?** A more recent addition are GEMM routines which are parallelized using OpenMP (`libxsmm_?gemm_omp`). These routines leverage the same specialized kernel routines as the small matrix multiplications, in-memory code generation (JIT), and automatic code/parameter dispatch but they are implementing a tile-based multiplication scheme i.e., a scheme suitable for larger problem sizes.

**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 use the Call Wrapper, and measure the application performance e.g., time to solution. However, the latter can significantly improve when using LIBXSMM’s API directly. To check whether there are applicable GEMM-calls, the Verbose Mode can help to collect an insight. Further, when an application uses Intel MKL 11.2 (or higher), then running the application with the environment variable `MKL_VERBOSE=1` (`env MKL_VERBOSE=1 ./workload > verbose.txt`) can collect a similar insight (`grep -a "MKL_VERBOSE DGEMM(N,N" verbose.txt | cut -d'(' -f2 | cut -d, -f3-5`).

**What is a small convolution?** In the last years, new workloads such as deep learning and more specifically convolutional neural networks (CNN) emerged, and are pushing the limits of today’s hardware. One of the expensive kernels is a small convolution with certain kernel sizes (3, 5, or 7) such that calculations in the frequency space is not the most efficient method when compared with direct convolutions. LIBXSMM’s current support for convolutions aims for an easy to use invocation of small (direct) convolutions, which are intended for CNN training and classification. The Interface is currently ramping up, and the functionality increases quickly towards a broader set of use cases.

## Interface for Matrix Multiplication

The interface of the library is *generated* according to the Build Instructions, and it is therefore **not** stored in the code repository. Instead, one may have a look at the code generation template files for C/C++ and FORTRAN.

In order to initialize the dispatch-table or other internal resources, an explicit initialization routine helps to avoid lazy initialization overhead when calling LIBXSMM for the first time. The library deallocates internal resources at program exit, but also provides a companion to the aforementioned initialization (`finalize`).

```
/** Initialize the library; pay for setup cost at a specific point. */
void libxsmm_init(void);
/** De-initialize the library and free internal memory (optional). */
void libxsmm_finalize(void);
```

To perform the dense matrix-matrix multiplication  $C_{m \times n} = \alpha \cdot A_{m \times k} \cdot B_{k \times n} + \beta \cdot C_{m \times n}$ , the full-blown GEMM interface can be treated with “default arguments” (which is deviating from the BLAS standard, however without compromising the binary compatibility).

```
/** Automatically dispatched dense matrix multiplication (single/double-precision, C code). */
libxsmm_?gemm(NULL/*transa*/, NULL/*transb*/, &m/*required*/, &n/*required*/, &k/*required*/,
             NULL/*alpha*/, a/*required*/, NULL/*lda*/, b/*required*/, NULL/*ldb*/,
             NULL/*beta*/, c/*required*/, NULL/*ldc*/);
```

```

/** Automatically dispatched dense matrix multiplication (C++ code). */
libxsmm_gemm(NULL/*transa*/, NULL/*transb*/, m/*required*/, n/*required*/, k/*required*/,
  NULL/*alpha*/, a/*required*/, NULL/*lda*/, b/*required*/, NULL/*ldb*/,
  NULL/*beta*/, c/*required*/, NULL/*ldc*/);

```

For the C interface (with type prefix ‘s’ or ‘d’), all arguments and in particular m, n, and k are passed by pointer. This is needed for binary compatibility with the original GEMM/BLAS interface. The C++ interface is also supplying overloaded versions where m, n, and k are allowed to be passed by-value (making it clearer that m, n, and k are non-optional arguments).

The FORTRAN interface supports optional arguments (without affecting the binary compatibility with the original BLAS interface) by allowing to omit arguments where the C/C++ interface allows for NULL to be passed.

```

! Automatically dispatched dense matrix multiplication (single/double-precision).
CALL libxsmm_?gemm(m=m, n=n, k=k, a=a, b=b, c=c)
! Automatically dispatched dense matrix multiplication (generic interface).
CALL libxsmm_gemm(m=m, n=n, k=k, a=a, b=b, c=c)

```

For convenience, a BLAS-based dense matrix multiplication (`libxsmm_blas_gemm`) is provided for all supported languages which is simply re-exposing the underlying GEMM/BLAS implementation. The BLAS-based GEMM might be useful for validation/benchmark purposes, and more important as a fallback when building an application-specific dispatch mechanism.

```

/** Automatically dispatched dense matrix multiplication (single/double-precision). */
libxsmm_blas_?gemm(NULL/*transa*/, NULL/*transb*/, &m/*required*/, &n/*required*/, &k/*required*/,
  NULL/*alpha*/, a/*required*/, NULL/*lda*/, b/*required*/, NULL/*ldb*/,
  NULL/*beta*/, c/*required*/, NULL/*ldc*/);

```

A more recently added variant of matrix multiplication is parallelized based on the OpenMP standard. The associated routines will open an internal parallel region by default, however participating on an already opened parallel region (without relying on nested parallelism) is also possible by using the environment variable `LIBXSMM_MT` (0: small-sized, 1: sequential, and 2: parallelized/default). The actual parallelism is based on “classic” OpenMP by default (thread-based), but can be adjusted to OpenMP 3.0 task-based parallelism (environment variable `LIBXSMM_TASKS=1`). At least the latter parallelization is dynamically scheduled. Please note that these routines are hosted by the extension library (`libxsmmext`) keeping the main library agnostic with respect to a particular threading runtime.

```

/** OpenMP parallelized dense matrix multiplication (single/double-precision). */
libxsmm_?gemm_omp(&transa, &transb, &m, &n, &k, &alpha, a, &lda, b, &ldb, &beta, c, &ldc);

```

Successively calling a particular kernel (i.e., multiple times) allows for amortizing the cost of the code dispatch. Moreover, in order to customize the dispatch mechanism, one can rely on the following interface.

```

/** If non-zero function pointer is returned, call (*function_ptr)(a, b, c). */
libxsmm_smmfunction libxsmm_smmdispatch(int m, int n, int k,
  const int* lda, const int* ldb, const int* ldc,
  const float* alpha, const float* beta,
  const int* flags, const int* prefetch);
/** If non-zero function pointer is returned, call (*function_ptr)(a, b, c). */
libxsmm_dmmfunction libxsmm_dmmdispatch(int m, int n, int k,
  const int* lda, const int* ldb, const int* ldc,
  const double* alpha, const double* beta,
  const int* flags, const int* prefetch);

```

A variety of overloaded function signatures is provided allowing to omit arguments not deviating from the configured defaults. In C++, a type `libxsmm_mmmfunction<type>` can be used to instantiate a functor rather than making a distinction for the numeric type in `libxsmm_?mmdispatch`. Similarly in FORTRAN, when calling the generic interface (`libxsmm_mmdispatch`) the given `LIBXSMM_?MMFUNCTION` is dispatched such that `libxsmm_call` can be used to actually perform the function call using the PROCEDURE POINTER wrapped by `LIBXSMM_?MMFUNCTION`. Beside of dispatching code, one can also call a specific kernel (e.g., `libxsmm_dmm_4_4_4`) using the prototype functions included for statically generated kernels.

## Interface for Convolutions

In order to achieve best performance with small convolutions for CNN on SIMD architectures, a specific data layout has to be used. As this layout depends on several architectural parameters, the goal of LIBXSMM interface is to hide

this complexity from the user by providing copy-in and copy-out routines. These happen on custom datatype which themselves are later bound to a convolution operation. The interface is available for C.

The main concept in LIBXSMM's frontend is that everything is circled around `libxsmm_dnn_conv_handle` which will define all properties of a layer operation. A handle can be created by describing the convolutional layer and calling a create function:

```
/** simplified LIBXSMM types which are needed to create a handle */
/** Structure which describes the input and output of data (DNN). */
typedef struct libxsmm_dnn_conv_desc {
    int N; /* number of images in mini-batch */
    int C; /* number of input feature maps */
    int H; /* height of input image */
    int W; /* width of input image */
    int K; /* number of output feature maps */
    int R; /* height of filter kernel */
    int S; /* width of filter kernel */
    int u; /* vertical stride */
    int v; /* horizontal stride */
    int pad_h_in; /* height of zero-padding in input buffer, ignored */
    int pad_w_in; /* width of zero-padding in input buffer, ignored */
    int pad_h_out; /* height of zero-padding in output buffer */
    int pad_w_out; /* width of zero-padding in output buffer */
    libxsmm_dnn_conv_algo algo; /* convolution algorithm used */
    libxsmm_dnn_conv_format buffer_format; /* format which is for buffer buffers */
    libxsmm_dnn_conv_format filter_format; /* format which is for filter buffers */
    libxsmm_dnn_conv_fuse_ops fuse_ops; /* used ops into convolutions */
    libxsmm_dnn_conv_option options; /* additional options */
    libxsmm_dnn_datatype datatype_in; /* datatypes use for all buffers */
    libxsmm_dnn_datatype datatype_out; /* datatypes use for all buffers */
} libxsmm_dnn_conv_desc;

/** Type of algorithm used for convolutions. */
typedef enum libxsmm_dnn_conv_algo {
    /** direct convolution. */
    LIBXSMM_DNN_CONV_ALGO_DIRECT
} libxsmm_dnn_conv_algo;

/** Denotes the element/pixel type of an image/channel. */
typedef enum libxsmm_dnn_conv_datatype {
    LIBXSMM_DNN_DATATYPE_F32
} libxsmm_dnn_datatype;

libxsmm_dnn_conv_handle* libxsmm_dnn_create_conv_handle_check(
    libxsmm_dnn_conv_desc conv_desc,
    libxsmm_dnn_datatype conv_datatype,
    libxsmm_dnn_conv_algo conv_algo,
    libxsmm_dnn_err_t* status);
```

Therefore, a sample call looks like:

```
/** Macro to check for an error. */
#define CHKERR_LIBXSMM_DNN(A) if (A != LIBXSMM_DNN_SUCCESS) \
    fprintf(stderr, "%s\n", libxsmm_dnn_get_error(A));
/** declare LIBXSMM variables */
libxsmm_dnn_conv_desc conv_desc;
libxsmm_dnn_err_t status;
libxsmm_dnn_conv_handle* libxsmm_handle;
/** setting conv_desc values.... */
conv_desc.N = ...
/** create handle */
libxsmm_handle = libxsmm_dnn_create_conv_handle_check(conv_desc, &status);
CHKERR_LIBXSMM_DNN(status);
```

Next activation and filter buffers need to be created, initialized and bound to the handle. Afterwards the convolution could be executed by a threading environment of choice:

```
libxsmm_dnn_buffer* libxsmm_input;
libxsmm_dnn_buffer* libxsmm_output;
libxsmm_dnn_filter* libxsmm_filter;
```

```

/* setup LIBXSMM layer information */
libxsmm_input = libxsmm_dnn_create_input_buffer_check(libxsmm_handle, &status);
CHKERR_LIBXSMM_DNN(status);
libxsmm_output = libxsmm_dnn_create_output_buffer_check(libxsmm_handle, &status);
CHKERR_LIBXSMM_DNN(status);
libxsmm_filter = libxsmm_dnn_create_filter_check(libxsmm_handle, &status);
CHKERR_LIBXSMM_DNN(status);

/* copy in data to LIBXSMM format: naive format is: */
/* (mini-batch)(number-featuremaps)(featuremap-height)(featuremap-width) for layers, */
/* and the naive format for filters is: */
/* (number-output-featuremaps)(number-input-featuremaps)(kernel-height)(kernel-width) */
CHKERR_LIBXSMM_DNN(libxsmm_dnn_copyin_buffer(libxsmm_input, (void*)naive_input));
CHKERR_LIBXSMM_DNN(libxsmm_dnn_zero_buffer(libxsmm_output));
CHKERR_LIBXSMM_DNN(libxsmm_dnn_copyin_filter(libxsmm_filter, (void*)naive_filter));

/* bind layer to handle */
CHKERR_LIBXSMM_DNN(libxsmm_dnn_bind_input_buffer(libxsmm_handle, libxsmm_input));
CHKERR_LIBXSMM_DNN(libxsmm_dnn_bind_output_buffer(libxsmm_handle, libxsmm_output));
CHKERR_LIBXSMM_DNN(libxsmm_dnn_bind_filter(libxsmm_handle, libxsmm_filter));

/* run the convolution */
#pragma omp parallel
{
    CHKERR_LIBXSMM_DNN(libxsmm_dnn_convolve_st(libxsmm_handle, LIBXSMM_DNN_CONV_KIND_FWD, 0,
        omp_get_thread_num(), omp_get_num_threads()));
}

/* copy out data */
CHKERR_LIBXSMM_DNN(libxsmm_dnn_copyout_buffer(libxsmm_output, (void*)naive_libxsmm_output));

```

## Service Functions

For convenient operation of the library and to ease integration, a number of service routines are available. They do not exactly belong to the core functionality of LIBXSMM (SMM or DNN domain), but users are encouraged to rely on these routines of the API. There are two categories: (1) routines which are available for C and Fortran, and (2) routines which are only available with the C interface.

## Getting and Setting the Target Architecture

There are ID based and string based functions to query the code path (as determined by the CPUID), or to set the code path regardless of the presented CPUID features. The latter may degrade performance (if a lower set of instruction set extensions is requested), which can be still useful for studying the performance impact of different instruction set extensions. This functionality is available for the C and Fortran interface, and there is an environment variable which corresponds to `libxsmm_set_target_arch` (LIBXSMM\_TARGET).

**NOTE:** There is no additional check performed if an unsupported instruction set extension is requested, and incompatible JIT-generated code may be executed (unknown instruction signaled).

```

int libxsmm_get_target_archid(void);
void libxsmm_set_target_archid(int id);

const char* libxsmm_get_target_arch(void);
void libxsmm_set_target_arch(const char* arch);

```

## Getting and Setting the Verbosity

The Verbose Mode (level of verbosity) can be controlled using the C or Fortran API, and there is an environment variable which corresponds to `libxsmm_set_verbosity` (LIBXSMM\_VERBOSE).

```

int libxsmm_get_verbosity(void);
void libxsmm_set_verbosity(int level);

```

## Timer Facility

Due to the performance oriented nature of LIBXSMM, timer-related functionality is available for the C and Fortran interface. This is used for instance by the code samples, which measure the duration of executing various code regions. Both “tick” functions (`libxsmm_timer_[x]tick`) are based on monotonic timer sources, which use a platform-specific resolution. The `xtick`-counter attempts to directly rely on the time stamp counter instruction (RDTSC), but it is not necessarily counting real CPU cycles due to varying CPU clock speed (Turbo Boost), different clock domains

(e.g., depending on the instructions executed), and other reasons (which are out of scope in this context).

**NOTE:** `libxsmm_timer_xtick` is not directly suitable for `libxsmm_timer_duration` (seconds).

```
unsigned long long libxsmm_timer_tick(void);
unsigned long long libxsmm_timer_xtick(void);
double libxsmm_timer_duration(unsigned long long tick0, unsigned long long tick1);
```

## Memory Allocation

Without further claims on the properties of the memory allocation (e.g., thread scalability), there are C functions that allocate aligned memory one of which allows to specify the alignment (or to specify an automatically chosen alignment). The automatic alignment is also exposed by a `malloc` compatible signature. The size of the automatic alignment depends on a heuristic, which uses the size of the requested buffer.

**NOTE:** only `libxsmm_free` is supported in order to deallocate the memory.

```
void* libxsmm_aligned_malloc(size_t size, int alignment);
void* libxsmm_malloc(size_t size);
void libxsmm_free(const volatile void* memory);
```

## Build Instructions

The build system relies on GNU Make (typically associated with the `make` command, but e.g. FreeBSD is calling it `gmake`). The build can be customized by using key-value pairs. Key-value pairs can be supplied in two ways: (1) after the “make” command, or (2) prior to the “make” command (`env`) which is effectively the same as exporting the key-value pair as an environment variable (`export`, or `setenv`). Of course both methods can be mixed, however the second method may require to supply the `-e` flag. Please note that the `CXX`, `CC`, and `FC` keys are handled such that they are taken into account in any case.

To generate the interface of the library inside of the ‘include’ directory and to build the static library (by default, `STATIC=1` is activated), simply run the following command:

```
make
```

If the build process is not successful, it may help to avoid more advanced GCC flags. This is useful with a tool chain, which pretends to be GCC-compatible (or is treated as such) but actually fails to consume the aforementioned flags. In such a case (`CCE`, etc.) one may raise the compatibility:

```
make COMPATIBLE=1
```

By default, only the non-coprocessor targets are built (`OFFLOAD=0` and `KNC=0`). In general, the subfolders of the ‘lib’ directory are separating the build targets where the ‘mic’ folder is containing the native library (`KNC=1`) targeting the Intel Xeon Phi coprocessor (“KNC”), and the ‘intel64’ folder is storing either the hybrid archive made of CPU and coprocessor code (`OFFLOAD=1`), or an archive which is only containing the CPU code. By default, an `OFFLOAD=1` implies `KNC=1`.

To remove intermediate files, or to remove all generated files and folders (including the interface and the library archives), run one of the following commands:

```
make clean
make realclean
```

By default, LIBXSMM uses the JIT backend which is automatically building optimized code. However, one can also statically specialize for particular matrix sizes (M, N, and K values), for convolutions the options below can be ignored:

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

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

```
(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 (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 indexes which are eventually grouped (using commas):

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

Each group of the above 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)
```

Of course, both mechanisms (M/N/K and MNK based) can be combined using the same command line (make). Static optimization and JIT can also be combined (no need to turn off the JIT backend). Testing the library is supported by a variety of targets with “test” and “test-all” being the most prominent for this matter.

Functionality of LIBXSMM, which is unrelated to GEMM can be used without introducing a dependency to BLAS. This can be achieved in two ways: (1) building a special library with `make BLAS=0`, or (2) linking the application against the ‘libxsmmnoblas’ library. Some care must be taken with any matrix multiplication which does not appear to require BLAS for the given test arguments. However, it may fall back to BLAS (at runtime of the application), if an unforeseen input is given (problem size, or unsupported GEMM arguments).

**NOTE:** by default, a C/C++ and a FORTRAN compiler is needed (some sample code is written in C++). Beside of specifying the compilers (`make CXX=g++ CC=gcc FC=gfortran` and maybe `AR=ar`), the need for a FORTRAN compiler can be relaxed (`make FC=` or `make FORTRAN=0`). The latter affects the availability of the MODULE file and the corresponding ‘libxsmmf’ library (the interface ‘libxsmmf.f’ is still generated). FORTRAN code can make use of LIBXSMM in three different ways:

- By relying on the module file, and by linking against ‘libxsmmf’, ‘libxsmm’, and (optionally) ‘libxsmmext’,
- By including the interface ‘libxsmmf.f’ and linking against ‘libxsmm’, and (optionally) ‘libxsmmext’, or
- By declaring e.g., `libxsmm_?gemm` (BLAS signature) and linking ‘libxsmm’ (and ‘libxsmmext’ if needed).

At the expense of a limited set of functionality (`libxsmm_?gemm[_omp]`, `libxsmm_blas_?gemm`, and `libxsmm_[s|d]otrans[_omp]`), the latter method also works with FORTRAN 77 (otherwise the FORTRAN 2003 standard is necessary). For the “omp” functionality, the ‘libxsmmext’ library needs to be present at the link line. For no code change at all, the Call Wrapper might be of interest.

## Link Instructions

The library is agnostic with respect to the threading-runtime, and therefore an application is free to use any threading runtime (e.g., OpenMP). The library is also thread-safe, and multiple application threads can call LIBXSMM’s routines concurrently. Forcing OpenMP (`OMP=1`) for the entire build of LIBXSMM is not supported and untested (‘libxsmmext’ is automatically built with OpenMP enabled).

Similarly, an application is free to choose any BLAS or LAPACK library (if the link model available on the OS supports this), and therefore linking GEMM routines when linking LIBXSMM itself (by supplying `BLAS=1|2`) may prevent a user from making this decision at the time of linking the actual application.

**NOTE:** LIBXSMM does not support to dynamically link against ‘libxsmm’ or ‘libxsmmext’ (“so”), when BLAS is linked statically (“a”).

## Installation

Installing LIBXSMM makes possibly the most sense when combining the JIT backend (enabled by default) with a collection of statically generated SSE kernels (by specifying M, N, K, or MNK). If the JIT backend is not disabled, statically generated kernels are only registered for dispatch if the CPUID flags at runtime are not supporting a more specific instruction set extension (code path). Since the JIT backend does not support or generate SSE code by itself, the library is compiled by selecting SSE code generation if not specified otherwise (`AVX=1|2|3`, or with `SSE=0` falling back to an “arch-native” approach). Limiting the static code path to SSE4.2 allows to practically target any deployed system, however using `SSE=0` and `AVX=0` together is falling back to generic code, and any static kernels are not specialized using the assembly code generator.

There are two main mechanisms to install LIBXSMM (both mechanisms can be combined): (1) building the library in an out-of-tree fashion, and (2) installing into a certain location. Building in an out-of-tree fashion looks like:

```
cd libxsmm-install  
make -f /path/to/libxsmm/Makefile
```

For example, installing into a specific location (incl. a selection of statically generated Intel SSE kernels) looks like:

```
make MNK="1 2 3 4 5" PREFIX=/path/to/libxsmm-install install
```

Performing `make install-minimal` omits the documentation (default: `'PREFIX/share/libxsmm'`). Moreover, `PINCDIR`, `POUTDIR`, `PBINDIR`, and `PDOCDIR` allow to customize the locations underneath of the `PREFIX` location. To build a general package for an unpredictable audience (Linux distribution, or similar), it is advised to not over-specify or customize the build step i.e., `JIT`, `SSE`, `AVX`, `OMP`, `BLAS`, etc. should not be used. The following is building and installing a complete set of libraries where the generated interface matches both the static and the shared libraries:

```
make PREFIX=/path/to/libxsmm-install STATIC=0 install
make PREFIX=/path/to/libxsmm-install install
```

## Running

### Call Wrapper

Since the library is binary compatible with existing GEMM calls (BLAS), these calls can be replaced at link-time or intercepted at runtime of an application such that LIBXSMM is used instead of the original BLAS library. There are two cases to consider:

- An application which is linked statically against BLAS requires to wrap the `'sgemm_'` and the `'dgemm_'` symbol (an alternative is to wrap only `'dgemm_'`), and a special build of the `libxsmm(ext)` library is required (`make WRAP=1` to wrap SGEMM and DGEMM, or `make WRAP=2` to wrap only DGEMM):

```
gcc [...] -Wl,--wrap=sgemm_ --wrap=dgemm_ /path/to/libxsmmext.a /path/to/libxsmm.a /path/to/your_regular_blas.a
```

Relinking the application as shown above can often be accomplished by copying, pasting, modifying the linker command, and then re-invoking the modified link step. This linker command may appear as console output of the application's "make" command (or a similar build system).

The static link-time wrapper technique may only work with a GCC tool chain (GNU Binutils: `1a`, or `1a` via compiler-driver), and it has been tested with GNU GCC, Intel Compiler, and Clang. However, this does not work under Microsoft Windows (even when using the GNU tool chain), and it may not work under OS X (Compiler 6.1 or earlier, later versions have not been tested).

- An application which is dynamically linked against BLAS allows for intercepting the GEMM calls at startup time (runtime) of the unmodified executable by using the `LD_PRELOAD` mechanism. The shared library of LIBXSMM (`make STATIC=0`) allows to intercept the GEMM calls of the application:

```
LD_PRELOAD=/path/to/libxsmmext.so ./myapplication
```

The behavior of the intercepted GEMM routines (statically wrapped or via `LD_PRELOAD`) can be controlled with the environment variable `LIBXSMM_MT` i.e., 0: calling sequential below-threshold routines without OpenMP (default when only linking `'libxsmm'`), 1: OpenMP-parallelized behavior but without an internal parallel region, and 2: OpenMP-parallelized routines with internal parallel region (default when linking `'libxsmmext'`). In any case, the wrapper mechanism also supports to fall back to BLAS.

```
LIBXSMM_MT=0 ./myapplication
```

**NOTE:** Using the same multiplication kernel in a consecutive fashion (batch-processing) allows to extract higher performance, when using LIBXSMM's native programming interface.

### Verbose Mode

The verbose mode allows for an insight into the code dispatch mechanism by receiving a small tabulated statistic as soon as the library terminates. The design point for this functionality is to not impact the performance of any critical code path i.e., verbose mode is always enabled and does not require symbols (`SYM=1`) or debug code (`DBG=1`). The statistics appears (`stderr`) when the environment variable `LIBXSMM_VERBOSE` is set to a non-zero value. For example:

```
LIBXSMM_VERBOSE=1 ./myapplication
[... application output]
```

HSW/SP	TRY	JIT	STA	COL
0..13	7	7	0	0
14..23	0	0	0	0
24..80	3	3	0	0

The tables are distinct between single-precision and double-precision, but either table is pruned if all counters are zero. If both tables are pruned, the library shows the code path which would have been used for JIT'ing the code: `LIBXSMM_TARGET=hsw` (otherwise the code path is shown in the table's header). The actual counters are collected for

three buckets: small kernels ( $MNK^{1/3} \leq 13$ ), medium-sized kernels ( $13 < MNK^{1/3} \leq 23$ ), and larger kernels ( $23 < MNK^{1/3} \leq 80$ ; the actual upper bound depends on `LIBXSMM_MAX_MNK` as selected at compile-time). Keep in mind, that “larger” is supposedly still fairly small in terms of arithmetic intensity (which grows linearly with the kernel size). Unfortunately, the arithmetic intensity depends on the way a kernel is used (which operands are loaded/stored into main memory) and it is not performance-neutral to collect this information.

The TRY counter represents all attempts to register statically generated kernels, and all attempts to dynamically generate and register kernels. The TRY counter includes rejected JIT requests due to unsupported GEMM arguments. The JIT and STA counters distinct the successful cases of the aforementioned event (TRY) into dynamically (JIT) and statically (STA) generated code. In case the capacity ( $O(n) = 10^5$ ) of the code registry is exhausted, no more kernels can be registered although further attempts are not prevented. Registering many kernels ( $O(n) = 10^3$ ) may ramp the number of hash key collisions (COL), which can degrade performance. The latter is prevented if the small thread-local cache is utilized effectively.

Since explicitly JIT-generated code (`libxsmm_mmdispatch`) does not fall under the THRESHOLD criterion, the above table is extended by one line if large kernels have been requested. This indicates a missing threshold-criterion (customized dispatch), or asks for cache-blocking the matrix multiplication. The latter is already implemented by LIBXSMM’s “medium-sized” GEMM routines (`libxsmm_gemm_omp`), which perform a tiled multiplication.

**NOTE:** setting `LIBXSMM_VERBOSE` to a negative value will dump each generated JIT kernel to a file with each file being named similar to the function name shown in Intel VTune.

## Call Trace

During the initial steps of employing the LIBXSMM API, one may rely on a debug version of the library (`make DBG=1`). The latter also implies console output (`stderr`) in case of an error/warning condition inside of the library. It is also possible to print the execution flow (call trace) inside of LIBXSMM (can be combined with `DBG=1` or `OPT=0`):

```
make TRACE=1
```

Building an application which is able to trace calls (inside of the library) requires the shared library of LIBXSMM, alternatively the application is required to link the static library of LIBXSMM in a dynamic fashion (GNU tool chain: `-rdynamic`). Actually tracing calls (without debugger) can be accomplished by an environment variable called `LIBXSMM_TRACE`.

```
LIBXSMM_TRACE=1 ./myapplication
```

Syntactically up to three arguments separated by commas (which allows to omit arguments) are taken (*tid,i,n*): *tid* signifies the ID of the thread to be traced with `1...NTHREADS` being valid and where `LIBXSMM_TRACE=1` is filtering for the “main thread” (in fact the first thread running into the trace facility); grabbing all threads (no filter) can be achieved by supplying a negative id (which is also the default when omitted). The second argument is pruning higher levels of the call-tree with *i=1* being the default (level zero is the highest at the same level as the main function). The last argument is taking the number of inclusive call levels with *n=-1* being the default (signifying no filter).

Although the `ltrace` (Linux utility) provides similar insight, the trace facility might be useful due to the aforementioned filtering expressions. Please note that the trace facility is severely impacting the performance (even with `LIBXSMM_TRACE=0`), and this is not just because of console output but rather since inlining (internal) functions might be prevented along with additional call overhead on each function entry and exit. Therefore, debug symbols can be also enabled separately (`make SYM=1`; implied by `TRACE=1` or `DBG=1`) which might be useful when profiling an application. No facility of the library (other than `DBG` or `TRACE/LIBXSMM_TRACE`) is performing visible (console) or other non-private I/O (files).

## Performance

### Profiling

#### Intel VTune Amplifier

To analyze which kind of kernels have been called, and from where these kernels have been invoked (call stack), the library allows profiling its JIT code using Intel VTune Amplifier. To enable this support, VTune’s root directory needs to be set at build-time of the library. Enabling symbols (`SYM=1` or `DBG=1`) incorporates VTune’s JIT Profiling API:

```
source /path/to/vtune_amplifier/amplxe-vars.sh
make SYM=1
```

Above, the root directory is automatically determined from the environment (`VTUNE_AMPLIFIER_*_DIR`). This variable is present after source'ing the Intel VTune environment, but it can be manually provided as well (`make VTUNEROOT=/path/to/vtune_amplifier`). Symbols are actually not required to display kernel names for the dynamically generated code, however enabling symbols makes the analysis much more useful for the rest of the (static) code, and hence it has been made a prerequisite. For example, when “call stacks” are collected it is possible to find out where the JIT code has been invoked by the application:

```
amplxe-cl -r result-directory -data-limit 0 -collect advanced-hotspots \
          -knob collection-detail=stack-sampling -- ./myapplication
```

In case of an MPI-parallelized application, it might be useful to only collect results from a “representative” rank, and to also avoid running the event collector in every rank of the application. With Intel MPI both of the latter can be achieved by adding

```
-gtool 'amplxe-cl -r result-directory -data-limit 0 -collect advanced-hotspots \
       -knob collection-detail=stack-sampling:4=exclusive'
```

to the `mpirun` command line. Please notice the `:4=exclusive` (unrelated to VTune's command line syntax), which is related to `mpirun`'s `gtool` arguments; these arguments need to appear at the end of the `gtool`-string. For instance, the shown command line selects the 4th rank (otherwise all ranks are sampled) along with “exclusive” usage of the performance monitoring unit (PMU) such that only one event-collector runs for all ranks.

Intel VTune Amplifier presents invoked JIT code like functions, which belong to a module named “`libxsmm.jit`”. The function name as well as the module name are supplied by LIBXSMM using the aforementioned JIT Profiling API. For instance “`libxsmm_hsw_dnn_23x23x23_23_23_23_a1_b1_p0::mxm`” encodes an Intel AVX2 (“`hsw`”) double-precision kernel (“`d`”) for small dense matrix multiplications (“`mxm`”) which is multiplying matrices without transposing them (“`nn`”). The rest of the name encodes `M=N=K=LDA=LDB=LDC=23`, `Alpha=Beta=1.0` (all similar to GEMM), and no prefetch strategy (“`p0`”).

## Linux perf

With LIBXSMM, there is both basic (`perf map`) and extended support (`jitdump`) when profiling an application. To enable `perf` support at runtime, the environment `LIBXSMM_VERBOSE` needs to be set to a negative value.

- The basic support can be enabled at compile-time with `PERF=1` (implies `SYM=1`) using `make PERF=1`. At runtime of the application, a map-file (`jit-pid.map`) is generated (`/tmp` directory). This file is automatically read by Linux `perf`, and enriches the information about unknown code such as JIT'ted kernels.
- The support for “`jitdump`” can be enabled by supplying `JITDUMP=1` (implies `PERF=1`) or `PERF=2` (implies `JITDUMP=1`) when making the library: `make JITDUMP=1` or `make PERF=2`. At runtime of the application, a dump-file (`jit-pid.dump`) is generated (in `perf`'s debug directory, usually `$HOME/.debug/jit/`) which includes information about JIT'ted kernels (such as addresses, symbol names, code size, and the code itself). The dump file can be injected into `perf.data` (using `perf inject -j`), and it enables an annotated view of the assembly in `perf`'s report (requires a reasonably recent version of `perf`).

## Tuning

Specifying a particular code path is not really necessary if the JIT backend is not disabled. However, disabling JIT compilation, statically generating a collection of kernels, and targeting a specific instruction set extension for the entire library looks like:

```
make JIT=0 AVX=3 MNK="1 2 3 4 5"
```

The above example builds a library which cannot be deployed to anything else but the Intel Knights Landing processor family (“`KNL`”) or future Intel Xeon processors supporting foundational Intel AVX-512 instructions (`AVX-512F`). The latter might be even more adjusted by supplying `MIC=1` (along with `AVX=3`), however this does not matter since critical code is in inline assembly (and not affected). Similarly, `SSE=0` (or `JIT=0` without `SSE` or `AVX` build flag) employs an “arch-native” approach whereas `AVX=1`, `AVX=2` (with `FMA`), and `AVX=3` are specifically selecting the kind of Intel AVX code. Moreover, controlling the target flags manually or adjusting the code optimizations is also possible. The following example is GCC-specific and corresponds to `OPT=3`, `AVX=3`, and `MIC=1`:

```
make OPT=3 TARGET="-mavx512f -mavx512cd -mavx512er -mavx512pf"
```

An extended interface can be generated which allows to perform software prefetches. Prefetching data might be helpful when processing batches of matrix multiplications where the next operands are farther away or otherwise

unpredictable in their memory location. The prefetch strategy can be specified similar as shown in the section Generator Driver i.e., by either using the number of the shown enumeration, or by exactly using the name of the prefetch strategy. The only exception is PREFETCH=1 which is automatically selecting a strategy according to an internal table (navigated by CPUID flags). The following example is requesting the “AL2jpst” strategy:

```
make PREFETCH=8
```

The prefetch interface is extending the signature of all kernels by three arguments (pa, pb, and pc). These additional arguments are specifying the locations of the operands of the next multiplication (the next a, b, and c matrices). Providing unnecessary arguments in case of the three-argument kernels is not big a problem (beside of some additional call-overhead), however running a kernel which is picking up more than three arguments and actually picking up garbage data is disabling the hardware prefetcher (due to software prefetches) followed by a misleading prefetch location plus an eventual page fault due to an out-of-bounds (garbage-)location.

Further, the generated configuration (template) of the library encodes the parameters for which the library was built for (static information). This helps optimizing client code related to the library’s functionality. For example, the LIBXSMM\_MAX\_\* and LIBXSMM\_AVG\_\* information can be used with the LIBXSMM\_PRAGMA\_LOOP\_COUNT macro in order to hint loop trip counts when handling matrices related to the problem domain of LIBXSMM.

### Auto-dispatch

The function `libxsmm_mmdispatch` helps amortizing the cost of the dispatch when multiple calls with the same M, N, and K are needed. The automatic code dispatch is orchestrating two levels:

1. Specialized routine (implemented in assembly code),
2. BLAS library call (fallback).

Both levels are accessible directly (see Interface section) allowing to customize the code dispatch. The fallback level may be supplied by the Intel Math Kernel Library (Intel MKL) 11.2 DIRECT CALL feature.

Further, a preprocessor symbol denotes the largest problem size ( $M \times N \times K$ ) that belongs to the first level, and therefore determines if a matrix multiplication falls back to BLAS. The problem size threshold can be configured by using for example:

```
make THRESHOLD=$((60 * 60 * 60))
```

The maximum of the given threshold and the largest requested specialization refines the value of the threshold. Please note that explicitly JIT’ing and executing a kernel is possible and independent of the threshold. If a problem size is below the threshold, dispatching the code requires to figure out whether a specialized routine exists or not.

In order to minimize the probability of key collisions (code cache), the preferred precision of the statically generated code can be selected:

```
make PRECISION=2
```

The default preference is to generate and register both single and double-precision code, and therefore no space in the dispatch table is saved (PRECISION=0). Specifying PRECISION=1|2 is only generating and registering either single-precision or double-precision code.

The automatic dispatch is highly convenient because existing GEMM calls can serve specialized kernels (even in a binary compatible fashion), however there is (and always will be) an overhead associated with looking up the code-registry and checking whether the code determined by the GEMM call is already JIT’ed or not. This lookup has been optimized using various techniques such as using specialized CPU instructions calculating a CRC32 checksum, avoiding costly synchronization (needed for thread-safety) until it is ultimately known that the requested kernel is not yet JIT’ed, and also a small thread-local cache of recently dispatched kernels. The latter of which can be adjusted in size (only power-of-two sizes) but also disabled:

```
make CACHE=0
```

Please note that measuring the relative cost of automatically dispatching a requested kernel depends on the kernel size (obviously smaller matrices are multiplied faster on an absolute basis), however smaller matrix multiplications are bottlenecked by memory bandwidth rather than arithmetic intensity. The latter implies the highest relative overhead when (artificially) benchmarking the very same multiplication out of the CPU-cache.

## JIT Backend

There might be situations in which it is up-front not clear which problem sizes will be needed when running an application. In order to leverage LIBXSMM's high-performance kernels, the library implements a JIT (Just-In-Time) code generation backend which generates the requested kernels on the fly (in-memory). This is accomplished by emitting the corresponding byte-code directly into an executable buffer. The actual JIT code is generated according to the CPUID flags, and therefore does not rely on the code path selected when building the library. In the current implementation, some limitations apply to the JIT backend specifically:

1. In order to stay agnostic to any threading model used, Pthread mutexes are guarding the updates of the JIT'ed code cache (link line with `-pthread` is required); building with `OMP=1` employs an OpenMP critical section as an alternative locking mechanism.
2. There is no support for the Intel SSE (Intel Xeon 5500/5600 series) and IMCI (Intel Xeon Phi coprocessor code-named Knights Corner) instruction set extensions. However, statically generated SSE-kernels can be leveraged without disabling support for JIT'ing AVX kernels.
3. There is no support for the Windows calling convention (only kernels with `PREFETCH=0` signature).

The JIT backend can also be disabled at build time (`make JIT=0`) as well as at runtime (`LIBXSMM_TARGET=0`, or anything prior to Intel AVX). The latter is an environment variable which allows to set a code path independent of the CPUID (`LIBXSMM_TARGET=0|1|sse|snb|hsw|knl|skx`). Please note that `LIBXSMM_TARGET` cannot enable the JIT backend if it was disabled at build time (`JIT=0`).

One can use the aforementioned `THRESHOLD` parameter to control the matrix sizes for which the JIT compilation will be automatically performed. However, explicitly requested kernels (by calling `libxsmm_?mmdispatch`) are not subject to a problem size threshold. In any case, JIT code generation can be used for accompanying statically generated code.

Note: Modern Linux kernels are supporting transparent huge pages (THP). LIBXSMM is sanitizing this feature when setting the permissions for pages holding the executable code. However, we measured up to 30% slowdown when running JIT'ed code in cases where THP decided to deliver a huge page. For systems with Linux kernel 2.6.38 (or later) THP will be automatically disabled for the `mmap`'ed regions (using `madvise`).

## Generator Driver

In rare situations it might be useful to directly incorporate generated C code (with inline assembly regions). This is accomplished by invoking a driver program (with certain command line arguments). The driver program is built as part of LIBXSMM's build process (when requesting static code generation), but also available via a separate build target:

```
make generator
bin/libxsmm_gemm_generator
```

The code generator driver program accepts the following arguments:

1. `dense/dense_asm/sparse` (`dense` creates C code, `dense_asm` creates ASM)
2. Filename of a file to append to
3. Routine name to be created
4. M parameter
5. N parameter
6. K parameter
7. LDA (0 when 1. is "sparse" indicates A is sparse)
8. LDB (0 when 1. is "sparse" indicates B is sparse)
9. LDC parameter
10. alpha (1)
11. beta (0 or 1)
12. Alignment override for A (1 auto, 0 no alignment)
13. Alignment override for C (1 auto, 0 no alignment)
14. Architecture (`noarch`, `wsm`, `snb`, `hsw`, `knc`, `knl`, `skx`)
15. Prefetch strategy, see below enumeration (`dense/dense_asm` only)
16. single precision (SP), or double precision (DP)
17. CSC file (just required when 1. is "sparse"). Matrix market format.

The prefetch strategy can be:

1. "nopf": no prefetching at all, just 3 inputs (A, B, C)

2. “pfsigonly”: just prefetching signature, 6 inputs (A, B, C, A', B', C')
3. “BL2viaC”: uses accesses to C to prefetch B'
4. “curAL2”: prefetches current A ahead in the kernel
5. “curAL2\_BL2viaC”: combines curAL2 and BL2viaC
6. “AL2”: uses accesses to A to prefetch A'
7. “AL2\_BL2viaC”: combines AL2 and BL2viaC
8. “AL2jpst”: aggressive A' prefetch of first rows without any structure
9. “AL2jpst\_BL2viaC”: combines AL2jpst and BL2viaC
10. “AL2\_BL2viaC\_CL2”: combines AL2 and BL2viaC

Here are some examples of invoking the driver program:

```
bin/libxsmm_gemm_generator dense foo.c foo 16 16 16 32 32 32 1 1 1 1 hsw nopf DP
bin/libxsmm_gemm_generator dense_asm foo.c foo 16 16 16 32 32 32 1 1 1 1 knl AL2_BL2viaC DP
bin/libxsmm_gemm_generator sparse foo.c foo 16 16 16 32 0 32 1 1 1 1 hsw nopf DP bar.csc
```

Please note, there are additional examples given in `samples/generator` and `samples/seissol`.

## Results

The LIBXSMM repository provides an orphaned branch “results” which is collecting collateral material such as measured performance results along with explanatory figures. The results can be found at <https://github.com/hfp/libxsmm/tree/results#libxsmm-results>.

Please note that comparing performance results depends on whether or not streaming the operands of the matrix multiplication. For example, running a matrix multiplication code many time with all operands covered by the L1 cache may have an emphasis towards an implementation which actually performs worse for the real workload (if this real workload needs to stream some or all operands from the main memory).

## Contributions

Contributions are very welcome! Please visit <https://github.com/hfp/libxsmm/wiki/Contribute>.

## Applications

[1] <https://cp2k.org/>: Open Source Molecular Dynamics with its DBCSR component processing batches of small matrix multiplications (“matrix stacks”) out of a problem-specific distributed block-sparse matrix. Starting with CP2K 3.0, LIBXSMM can be used to substitute CP2K’s ‘libsmm’ library. Prior to CP2K 3.0, only the Intel-branch of CP2K was integrating LIBXSMM (see <https://github.com/hfp/libxsmm/raw/master/documentation/cp2k.pdf>).

[2] <https://github.com/SeisSol/SeisSol/>: 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/](https://github.com/TUM-I5/seissol_kernels/)).

[3] <https://github.com/Nek5000/NekBox>: NekBox is a version of the highly scalable and portable spectral element Nek5000 code which is specialized for box geometries, and intended for prototyping new methods as well as leveraging FORTRAN beyond the FORTRAN 77 standard. LIBXSMM provides optimized kernels aiming to conveniently substitute the MXM\_STD code.

[4] <https://github.com/Nek5000/Nek5000>: Nek5000 is the open-source, highly-scalable, always-portable spectral element code from <https://nek5000.mcs.anl.gov/>. The development branch of the Nek5000 code now incorporates LIBXSMM.

[5] <https://software.intel.com/en-us/articles/intel-xeon-phi-delivers-competitive-performance-for-deep-learning>: Intel Xeon Phi Delivers Competitive Performance For Deep Learning - And Getting Better Fast. Article mentioning LIBXSMM’s performance of convolution kernels with DeepBench. Intel Corporation, 2016.

## References

[1] <http://sc16.supercomputing.org/presentation/?id=pap364&sess=sess153>: LIBXSMM: Accelerating Small Matrix Multiplications by Runtime Code Generation (paper). SC’16: The International Conference for High Performance Computing, Networking, Storage and Analysis, Salt Lake City (Utah).

[2] [http://sc15.supercomputing.org/sites/all/themes/SC15images/tech\\_poster/tech\\_poster\\_pages/post137.html](http://sc15.supercomputing.org/sites/all/themes/SC15images/tech_poster/tech_poster_pages/post137.html): LIBXSMM: A High Performance Library for Small Matrix Multiplications (poster and abstract). SC’15: The International Conference for High Performance Computing, Networking, Storage and Analysis, Austin (Texas).