

TensorFlow™ with LIBXSMM

Getting Started

This document is an early recipe for building and running TensorFlow with LIBXSMM. The amount of covered code paths as well as the performance of these code paths will be improved as the integration progresses. This means that executing TensorFlow for any real workload (or benchmark) beside of the cases shown below **may not use LIBXSMM at all** (same is/remains true for any system without Intel AVX2 instruction set extension). Relying on the master revision of TensorFlow is perfectly fine since the integration work is merged on a frequent basis. However, to capture the latest revision of the integration one may rely on:

```
git clone https://github.com/benoitsteiner/tensorflow-xsmm.git
wget https://github.com/hfp/libxsmm/archive/master.zip
sha256sum libxsmm-master.zip
```

To try LIBXSMM's master revision, the file `tensorflow/workspace.bzl` can be adjusted by using the SHA256 sum from above:

```
native.new_http_archive(
    name = "libxsmm_archive",
    urls = [
        "https://github.com/hfp/libxsmm/archive/master.zip",
    ],
    sha256 = "<sha256sum libxsmm-master.zip>",
    strip_prefix = "libxsmm-master",
    build_file = str(Label("//third_party:libxsmm.BUILD")),
)
```

Beside of the regular prerequisites, nothing else is needed to use TensorFlow with LIBXSMM. Prior to bazel 0.4.5, it was helpful to change TensorFlow's configure script by replacing `bazel clean --expunge` with `bazel clean --expunge_async` (at least when NFS-hosted).

```
cd tensorflow-xsmm
./configure
```

When behind a HTTP-proxy, the environment variables should carry `https://` or `http://`:

```
export https_proxy=https://proxy.domain.com:912
export http_proxy=http://proxy.domain.com:911
```

Please note that certain content under `tensorflow/models` may cause an error during the configuration. In such a case, simply configure with “models” temporarily not present. In general, if the build step of any of the bazel commands goes wrong, `-s --verbose_failures` can be added to the command line (`-s` shows the full command of each of the build steps). The flags `--define tensorflow_xsmm=1`, `--define eigen_xsmm=1`, and `--define tensorflow_xsmm_backward=1` are not actually needed for all the cases, but are supplied for consistency.

More important, one line of below target flags should be added to any “bazel” command line:

- AVX2/HSW/BDW: `--copt=-mfma --copt=-mavx2`
- AVX-512/SKX: `--copt=-mfma --copt=-mavx512f --copt=-mavx512cd --copt=-mavx512bw --copt=-mavx512vl` (and `--copt=-mavx512dq` depending on certain fixes being already present in TensorFlow)
- AVX-512/KNL: `--copt=-mfma --copt=-mavx512f --copt=-mavx512cd --copt=-mavx512pf --copt=-mavx512er`

For Intel AVX-512 support, GCC 5.x should be used (see section Non-default Compiler). LIBXSMM supports Intel AVX2 as the baseline code path for all JIT-generated DNN-code (SMM domain also supports AVX). For Intel AVX-512 (on top of AVX2), the foundational instructions are sufficient in many cases, but for the sparse domain the Core-flavor is a prerequisite (“Skylake server” or SKX), and VNNI/QFMA instructions are honored on Intel Xeon Phi code-named “Knights Mill” (KNM).

Generally, please follow the guide to build TensorFlow from the sources. Please invoke the following commands to build and install a pip-package, which represents your build of TensorFlow:

```
bazel build -c opt --copt=-O3 --linkopt=-pthread \
    --define tensorflow_xsmm=1 --define eigen_xsmm=1 --define tensorflow_xsmm_backward=1 \
    <line-of-target-flags-from-above> \
    //tensorflow/tools/pip_package:build_pip_package

bazel-bin/tensorflow/tools/pip_package/build_pip_package /tmp/tensorflow_pkg
```

```
sudo pip install \
  --proxy proxy.domain.com:912 \
  -I /tmp/tensorflow_pkg/<package-name-build-above.whl>
```

Benchmarks

This document is an early recipe for building and running TensorFlow with LIBXSMM. Please do not expect any performance advantage (at this point) when comparing to TensorFlow without LIBXSMM!

```
git clone https://github.com/soumith/convnet-benchmarks.git
cd tensorflow-xsmm
mkdir -p tensorflow/models
ln -s /path/to/convnet-benchmarks/tensorflow tensorflow/models/convnetbenchmarks
bazel build -c opt --copt=-O3 --linkopt=-pthread \
  --define tensorflow_xsmm=1 --define eigen_xsmm=1 --define tensorflow_xsmm_backward=1 \
  <line-of-target-flags-from-above> \
  //tensorflow/models/convnetbenchmarks:benchmark_alexnet \
  //tensorflow/models/convnetbenchmarks:benchmark_overfeat \
  //tensorflow/models/convnetbenchmarks:benchmark_vgg \
  //tensorflow/models/convnetbenchmarks:benchmark_googlenet
```

The above bazel command may be combined to build `//tensorflow/tools/pip_package:build_pip_package` as well. When completed (wheel installed!) e.g., run the “Alexnet” benchmark:

```
LIBXSMM_VERBOSE=2 \
bazel-bin/tensorflow/models/convnetbenchmarks/benchmark_alexnet \
  --data_format=NHWC --forward_only=true --batch_size=256 2>&1 \
| tee output_alexnet.log
```

In case of an `ImportError: No module named builtins` one can resolve the problem with `sudo pip install future --upgrade`.

Non-default Compiler

LIBXSMM does not impose to build for a specific code path, and always exploits the most suitable instruction set extension for JIT-enabled code paths. However, LIBXSMM may also use non-JIT code paths which are CPUID-dispatched when the static code path has lower capabilities. This only works when using GCC 4.9 (or later) or the Intel Compiler. If TensorFlow does not match the highest possible CPU target (`march=native`), a performance penalty is possible.

It is recommended to rely on a pre-built compiler by using for instance the “devtools” package (RedHat) or similar (depends on the Linux distribution). It apparently helps to build `bazel` using the compiler intended for building TensorFlow. To use a custom-built compiler with TensorFlow may not only ask to source this compiler:

```
export LD_LIBRARY_PATH=/software/gnu/gcc-6.3.0/lib64:/software/gnu/gcc-6.3.0/lib:${LD_LIBRARY_PATH}
export PATH=/software/gnu/gcc-6.3.0/bin:${PATH}
```

but to further advertise the different compiler-runtime to the linker (`ld`).

```
echo "/software/gnu/gcc-6.3.0/lib64" > software-gnu-gcc630.conf
sudo mv software-gnu-gcc630.conf /etc/ld.so.conf.d/
sudo ldconfig
```

If there are still problems when using the custom compiler (mismatched GLIBC version), the symbolic link `libstdc++.so.6` (under `/usr/lib64`) may be adjusted to point to the latest update of the same major GLIBC version.

Performance Profiling

To gain insight into performance bottlenecks, one might source the Intel VTune Amplifier and run:

```
amplxe-cl -r result -data-limit 0 \
  -collect advanced-hotspots -knob collection-detail=stack-sampling -- \
  bazel-bin/tensorflow/models/convnetbenchmarks/benchmark_alexnet \
  --data_format=NHWC --forward_only=true --batch_size=64 --num_batches=50
```

Regression Tests

There are two aspects of LIBXSMM enabled within TensorFlow: (1) sparse CNN, and (2) CNN. To build and test the sparse routines:

```

bazel build -c opt --copt=-O3 --linkopt=-pthread \
  --define tensorflow_xsmm=1 --define eigen_xsmm=1 --define tensorflow_xsmm_backward=1 \
  <line-of-target-flags-from-above> \
  //tensorflow/core/kernels:sparse_matmul_op_test

bazel-bin/tensorflow/core/kernels/sparse_matmul_op_test --benchmarks=all
bazel-bin/tensorflow/core/kernels/sparse_matmul_op_test

bazel run -c opt --copt=-O3 --linkopt=-pthread \
  --define tensorflow_xsmm=1 --define eigen_xsmm=1 --define tensorflow_xsmm_backward=1 \
  <line-of-target-flags-from-above> \
  //tensorflow/python/kernel_tests:sparse_matmul_op_test

```

To build and test the regular CNN routines (note that below `bazel run...` may be deadlocking during the test):

```

bazel build -c opt --copt=-O3 --linkopt=-pthread \
  --define tensorflow_xsmm=1 --define eigen_xsmm=1 --define tensorflow_xsmm_backward=1 \
  <line-of-target-flags-from-above> \
  //tensorflow/core/kernels:conv_ops_test

bazel-bin/tensorflow/core/kernels/conv_ops_test

bazel run -c opt --copt=-O3 --linkopt=-pthread \
  --define tensorflow_xsmm=1 --define eigen_xsmm=1 --define tensorflow_xsmm_backward=1 \
  <line-of-target-flags-from-above> \
  //tensorflow/python/kernel_tests:conv_ops_test

```