



NVIDIA NVSHMEM

DG-08911-001_v1.0.1 | June 2020

Installation Guide



TABLE OF CONTENTS

Chapter 1. Overview.....	1
Chapter 2. Hardware And Software Requirements.....	2
2.1. Hardware Requirements.....	2
2.2. Software Requirements.....	2
Chapter 3. Installation.....	3
3.1. Downloading NVSHMEM.....	3
3.2. Building And Installing NVSHMEM.....	3
3.3. Using NVSHMEM In Your Applications.....	4
3.3.1. Using NVSHMEM With Your C Or C++ Program.....	4
3.3.2. Using NVSHMEM With Your MPI Program.....	4
3.4. Running Performance Tests.....	5
3.5. "Hello World" Example.....	5
Chapter 4. Support.....	7

Chapter 1.

OVERVIEW

NVIDIA® NVSHMEM™ is a programming interface that implements a Partitioned Global Address Space (PGAS) model across a cluster of NVIDIA GPUs. NVSHMEM provides an easy-to-use interface to allocate memory that is symmetrically distributed across the GPUs. In addition to a CPU-side interface, NVSHMEM also provides a CUDA kernel-side interface that allows CUDA® threads to access any location in the symmetrically-distributed memory.

Chapter 2.

HARDWARE AND SOFTWARE REQUIREMENTS

NVIDIA® NVSHMEM™ has the following hardware and software requirements.

2.1. Hardware Requirements

NVSHMEM requires the following hardware:

- ▶ x86_64 CPU architecture
- ▶ NVIDIA® Tesla® GPUs, NVIDIA® Volta® GPU architecture or newer. A full list is available at the following website:
<https://developer.nvidia.com/cuda-gpus>
- ▶ All GPUs must be P2P connected (directly via NVLink or PCIe) or must have GPUDirect RDMA enabled with a Mellanox InfiniBand adapter.

2.2. Software Requirements

NVSHMEM requires the following software:

- ▶ 64-bit Linux. For a complete compatibility matrix, see the [NVIDIA CUDA Installation Guide](#) for Linux.
- ▶ CUDA 10.1 or newer
- ▶ [Mellanox OFED](#)
- ▶ [nv_peer_mem for GPUDirect RDMA](#)
- ▶ [GDRCopy v2.0 or newer](#)

Chapter 3.

INSTALLATION

3.1. Downloading NVSHMEM

1. Ensure you are a registered partner and have access to [NVONLINE](#).
2. Download and extract the NVSHMEM `txz` archive. The extracted directory contains the following files and subdirectories:

File or Directory	Description
<code>src/</code>	Contains NVSHMEM sources and headers.
<code>perftest/</code>	Contains tests showing use of NVSHMEM APIs with performance reporting.
<code>examples/</code>	Contains examples showing use of some common use cases of NVSHMEM.
<code>scripts/</code>	Contains helper scripts, for example, script to download, build and install Hydra.
<code>COPYRIGHT.txt</code>	Copyright information.
<code>NVSHMEM-SLA.txt</code>	NVSHMEM service level agreement (SLA).

3.2. Building And Installing NVSHMEM

1. Set the environment variable `CUDA_HOME` to point to the CUDA Toolkit.
2. Set the environment variable `GDRCOPY_HOME` to point to the GDRCopy installation.
3. If MPI and/or SHMEM support is required, set `NVSHMEM_MPI_SUPPORT=1` and/or `NVSHMEM_SHMEM_SUPPORT=1`. Set the environment variables `MPI_HOME` and `SHMEM_HOME` to point to the MPI and OpenSHMEM installations, respectively.

4. If the MPI library is neither OpenMPI nor its derivative, set `NVSHMEM_MPI_IS_OMPI=0`.



When using OpenMPI and OSHMEM, the paths are the same. To use OSHMEM, OpenMPI needs to be built with UCX support. NVSHMEM has been tested with OpenMPI 4.0.1 and UCX 1.8.0. Other MPI and OpenSHMEM installations are expected to work. By default, MPI support is enabled and OpenSHMEM support is disabled.

5. Set `NVSHMEM_PREFIX` to specify the location where NVSHMEM will be installed.
6. Run `make -j8 install` to build and install the library.

3.3. Using NVSHMEM In Your Applications

3.3.1. Using NVSHMEM With Your C Or C++ Program

1. Include `nvshmem.h` and `nvshmemx.h` from `include/`.
2. Point to the `include/` and `lib/` paths.
3. NVSHMEM users: If your C or C++ program only uses NVSHMEM, install Hydra Process Manager using the bash script `install_hydra.sh` under the `scripts/` directory. Provide the download and install location as arguments, for example:

```
./install_hydra.sh <download_path> <install_path>
```

Use `nvshmrn` launcher under `bin/` (of the Hydra install path) to run the NVSHMEM job.

3.3.2. Using NVSHMEM With Your MPI Program



The only MPI library currently tested is OpenMPI, however, derivatives of OpenMPI such as SpectrumMPI as well as MPICH derivatives are expected to work.

To run a Hybrid MPI + NVSHMEM program, use the `mpirun` launcher available in the MPI installation.

Similarly, NVSHMEM can be used from OpenSHMEM programs, you cannot use the launchers provided in the NVSHMEM package. The only OpenSHMEM version currently tested is OSHMEM in OpenMPI. Other OpenSHMEM implementations such as Sandia OpenSHMEM (SOS) are expected to work. Use the `oshrun` launcher available in the OpenMPI installation (or follow the launcher specification of your OpenSHMEM library) to run the hybrid OpenSHMEM/NVSHMEM job.

3.4. Running Performance Tests

Before you can run performance tests, you first must build them.

1. Set the `CUDA_HOME`, `NVSHMEM_HOME` and `MPI_HOME` (if the NVSHMEM library was built with `NVSHMEM_MPI_SUPPORT=1`) environment variables to build NVSHMEM performance tests:

```
CUDA_HOME=<path to supported CUDA installation>
NVSHMEM_HOME=<path to directory where NVSHMEM is installed>
MPI_HOME=<path to MPI installation>
```

Configuring OpenMPI using the `-with-ucx` option is required for OpenMPI/OSHMEM interoperability. If you have built NVSHMEM with MPI and OpenSHMEM support (`NVSHMEM_MPI_SUPPORT=1` and `NVSHMEM_SHMEM_SUPPORT=1`), building `perftest/` also requires MPI and OpenSHMEM support to be enabled.

Build without SHMEM interoperability: To build NVSHMEM performance tests without SHMEM interoperability, set the environment variable `NVSHMEM_SHMEM_SUPPORT` to 0. By default, performance tests are installed under `perftest/perftest_install`. To install to a different path, set `PERFTEST_INSTALL` to point to the desired path.

2. Update `LD_LIBRARY_PATH` to point to `$CUDA_HOME/lib64` and `$MPI_HOME/lib`.
3. Run performance tests as NVSHMEM jobs (assuming Hydra is installed under `HYDRA_HOME`), hybrid MPI+NVSHMEM jobs, or hybrid OpenSHMEM+NVSHMEM jobs with the following commands (using `perftest/device/pt-to-pt/put.cu` as an example):

NVSHMEM job:

```
$HYDRA_HOME/bin/nvshmrn -n <up to number of P2P or InfiniBand NIC
accessible GPUs> $PERFTEST_INSTALL/device/pt-to-pt/shmem_put_bw
```

Hybrid MPI/NVSHMEM job:

```
$MPI_HOME/bin/mpirun -n <up to number of GPUs accessible by P2P or
InfiniBand NIC> -x NVSHMEMTEST_USE_MPI_LAUNCHER=1 $PERFTEST_INSTALL/device/
pt-to-pt/shmem_put_bw
```

Hybrid OpenSHMEM/NVSHMEM job:

```
$MPI_HOME/bin/oshrun -n <up to number of GPUs accessible by P2P or
InfiniBand NIC> -x USE_SHMEM_IN_TEST=1 $PERFTEST_INSTALL/device/pt-to-pt/
shmem_put_bw
```

3.5. "Hello World" Example

1. Save the following code as `nvshmemHelloWorld.cu`:

```
#include <stdio.h>
#include <cuda.h>
```

```

#include <nvshmem.h>
#include <nvshmemx.h>

__global__ void simple_shift(int *destination) {
    int mype = nvshmem_my_pe();
    int npes = nvshmem_n_pes();
    int peer = (mype + 1) % npes;

    nvshmem_int_p(destination, mype, peer);
}

int main(void) {
    int mype_node, msg;
    cudaStream_t stream;

    nvshmem_init();
    mype_node = nvshmem_team_my_pe(NVSHMEMX_TEAM_NODE);
    cudaSetDevice(mype_node);
    cudaStreamCreate(&stream);

    int *destination = (int *) nvshmem_malloc(sizeof(int));

    simple_shift<<<1, 1, 0, stream>>>(destination);
    nvshmemx_barrier_all_on_stream(stream);
    cudaMemcpyAsync(&msg, destination, sizeof(int), cudaMemcpyDeviceToHost,
stream);

    cudaStreamSynchronize(stream);
    printf("%d: received message %d\n", nvshmem_my_pe(), msg);

    nvshmem_free(destination);
    nvshmem_finalize();
    return 0;
}

```

2. Build `nvshmemHelloWorld.cu` with the following command:

```

nvcc -rdc=true -cubin g++ -gencode=$NVCC_GENCODE -I $NVSHMEM_HOME/include
nvshmemHelloWorld.cu -o nvshmemHelloWorld.out -L $NVSHMEM_HOME/lib -
lnvshmem
-lcuda

```

Where `arch=compute_70`, `code=sm_70` is the value of `NVCC_GENCODE` for V100 GPUs.

3. Run the `nvshmemHelloWorld` sample with one of the following commands:
When running on a single host with 2 GPUs (connected by PCI-E, NVLink or Infiniband):

```

$HYDRA_HOME/bin/nvshmrn -n 2 -ppn 2 ./nvshmemHelloWorld.out

```

When running on two hosts with 1 GPU per host (connected by InfiniBand):

```

$HYDRA_HOME/bin/nvshmrn -n 2 -ppn 1 --hosts hostname1,hostname2 ./
nvshmemHelloWorld.out

```


Chapter 4. SUPPORT

Report bugs and submit feature requests using [NVONLINE](#) or by emailing nvshmem@nvidia.com.

Notice

This document is provided for information purposes only and shall not be regarded as a warranty of a certain functionality, condition, or quality of a product. NVIDIA Corporation (“NVIDIA”) makes no representations or warranties, expressed or implied, as to the accuracy or completeness of the information contained in this document and assumes no responsibility for any errors contained herein. NVIDIA shall have no liability for the consequences or use of such information or for any infringement of patents or other rights of third parties that may result from its use. This document is not a commitment to develop, release, or deliver any Material (defined below), code, or functionality.

NVIDIA reserves the right to make corrections, modifications, enhancements, improvements, and any other changes to this document, at any time without notice.

Customer should obtain the latest relevant information before placing orders and should verify that such information is current and complete.

NVIDIA products are sold subject to the NVIDIA standard terms and conditions of sale supplied at the time of order acknowledgement, unless otherwise agreed in an individual sales agreement signed by authorized representatives of NVIDIA and customer (“Terms of Sale”). NVIDIA hereby expressly objects to applying any customer general terms and conditions with regards to the purchase of the NVIDIA product referenced in this document. No contractual obligations are formed either directly or indirectly by this document.

NVIDIA products are not designed, authorized, or warranted to be suitable for use in medical, military, aircraft, space, or life support equipment, nor in applications where failure or malfunction of the NVIDIA product can reasonably be expected to result in personal injury, death, or property or environmental damage. NVIDIA accepts no liability for inclusion and/or use of NVIDIA products in such equipment or applications and therefore such inclusion and/or use is at customer’s own risk.

NVIDIA makes no representation or warranty that products based on this document will be suitable for any specified use. Testing of all parameters of each product is not necessarily performed by NVIDIA. It is customer’s sole responsibility to evaluate and determine the applicability of any information contained in this document, ensure the product is suitable and fit for the application planned by customer, and perform the necessary testing for the application in order to avoid a default of the application or the product. Weaknesses in customer’s product designs may affect the quality and reliability of the NVIDIA product and may result in additional or different conditions and/or requirements beyond those contained in this document. NVIDIA accepts no liability related to any default, damage, costs, or problem which may be based on or attributable to: (i) the use of the NVIDIA product in any manner that is contrary to this document or (ii) customer product designs.

No license, either expressed or implied, is granted under any NVIDIA patent right, copyright, or other NVIDIA intellectual property right under this document. Information published by NVIDIA regarding third-party products or services does not constitute a license from NVIDIA to use such products or services or a warranty or endorsement thereof. Use of such information may require a license from a third party under the patents or other intellectual property rights of the third party, or a license from NVIDIA under the patents or other intellectual property rights of NVIDIA.

Reproduction of information in this document is permissible only if approved in advance by NVIDIA in writing, reproduced without alteration and in full compliance with all applicable export laws and regulations, and accompanied by all associated conditions, limitations, and notices.

THIS DOCUMENT AND ALL NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, "MATERIALS") ARE BEING PROVIDED "AS IS." NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE. TO THE EXTENT NOT PROHIBITED BY LAW, IN NO EVENT WILL NVIDIA BE LIABLE FOR ANY DAMAGES, INCLUDING WITHOUT LIMITATION ANY DIRECT, INDIRECT, SPECIAL, INCIDENTAL, PUNITIVE, OR CONSEQUENTIAL DAMAGES, HOWEVER CAUSED AND REGARDLESS OF THE THEORY OF LIABILITY, ARISING OUT OF ANY USE OF THIS DOCUMENT, EVEN IF NVIDIA HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES. Notwithstanding any damages that customer might incur for any reason whatsoever, NVIDIA's aggregate and cumulative liability towards customer for the products described herein shall be limited in accordance with the Terms of Sale for the product.

VESA DisplayPort

DisplayPort and DisplayPort Compliance Logo, DisplayPort Compliance Logo for Dual-mode Sources, and DisplayPort Compliance Logo for Active Cables are trademarks owned by the Video Electronics Standards Association in the United States and other countries.

HDMI

HDMI, the HDMI logo, and High-Definition Multimedia Interface are trademarks or registered trademarks of HDMI Licensing LLC.

OpenCL

OpenCL is a trademark of Apple Inc. used under license to the Khronos Group Inc.

Trademarks

NVIDIA, the NVIDIA logo, and CUDA, CUDA Toolkit, GPU, Kepler, Mellanox, NVLink, NVSHMEM, and Tesla are trademarks and/or registered trademarks of NVIDIA Corporation in the United States and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

Copyright

© 2020 NVIDIA Corporation. All rights reserved.