

---

# **RCCL Documentation**

*Release 2.29.7*

**Advanced Micro Devices, Inc.**

**Jun 11, 2026**



# CONTENTS

<b>1</b>	<b>What is RCCL?</b>	<b>3</b>
<b>2</b>	<b>Install RCCL</b>	<b>5</b>
2.1	Install the ROCm Core SDK . . . . .	5
2.2	Install RCCL on Linux . . . . .	5
2.3	Install a nightly build . . . . .	6
<b>3</b>	<b>Building RCCL from source</b>	<b>7</b>
3.1	Requirements . . . . .	7
3.2	Build and install RCCL using the install script . . . . .	7
3.3	Building the library using CMake . . . . .	10
3.4	Testing RCCL . . . . .	10
<b>4</b>	<b>Running RCCL using Docker</b>	<b>13</b>
<b>5</b>	<b>Using the RCCL Tuner plugin API</b>	<b>15</b>
5.1	API description . . . . .	16
5.2	Build and usage instructions . . . . .	17
<b>6</b>	<b>Using the NCCL Net plugin API</b>	<b>19</b>
6.1	Plugin architecture . . . . .	19
6.2	API (v6) . . . . .	20
<b>7</b>	<b>Troubleshooting RCCL</b>	<b>25</b>
7.1	Collecting system information . . . . .	25
7.2	Collecting RCCL information . . . . .	26
7.3	Analyzing performance issues . . . . .	27
7.4	RCCL and NCCL comparisons . . . . .	28
<b>8</b>	<b>RCCL usage tips</b>	<b>29</b>
8.1	NPKit . . . . .	29
8.2	Enabling peer-to-peer transport . . . . .	29
8.3	Ignoring CPU affinity with multi-node . . . . .	30
8.4	Improving performance on the MI300X . . . . .	30
<b>9</b>	<b>RCCL library specification</b>	<b>35</b>
9.1	Communicator functions . . . . .	35
9.2	Collective communication operations . . . . .	37
9.3	Group semantics . . . . .	41
9.4	Library functions . . . . .	42
9.5	Types . . . . .	42

9.6 Enumerations . . . . .	43
<b>10 API library</b>	<b>47</b>
10.1 Operations . . . . .	47
<b>11 RCCL environment variables</b>	<b>77</b>
11.1 Configuration and setup . . . . .	77
11.2 Logging and debugging . . . . .	78
11.3 Algorithm and protocol control . . . . .	80
11.4 Network and topology . . . . .	80
11.5 Development and testing (advanced) . . . . .	82
11.6 Multi-communicator ordering . . . . .	82
<b>12 License</b>	<b>83</b>
<b>13 Attributions</b>	<b>85</b>
<b>Index</b>	<b>87</b>

The ROCm Communication Collectives Library (RCCL) is a stand-alone library that provides multi-GPU and multi-node collective communication primitives optimized for AMD GPUs. It uses PCIe and xGMI high-speed interconnects. To learn more, see [What is RCCL?](#)

The RCCL public repository is located within the rocm-systems repo at <https://github.com/ROCm/rocm-systems/tree/develop/projects/rccl>.

#### Install

- [Install RCCL](#)
- [Build from source](#)
- [Running RCCL using Docker](#)

#### How to

- [Using the RCCL Tuner plugin](#)
- [Using the NCCL Net plugin](#)
- [Troubleshoot RCCL](#)
- [RCCL usage tips](#)

#### Examples

- [RCCL Tuner plugin examples](#)
- [NCCL Net plugin examples](#)

#### API reference

- [Library specification](#)
- [API library](#)
- [Environment variables](#)

To contribute to the documentation, see [Contributing to ROCm](#).

You can find licensing information on the [Licensing](#) page.



## WHAT IS RCCL?

The ROCm Communication Collectives Library (RCCL) includes multi-GPU and multi-node collective communication primitives optimized for AMD GPUs. It implements routines such as `all-reduce`, `all-gather`, `reduce`, `broadcast`, `reduce-scatter`, `gather`, `scatter`, `all-to-allv`, and `all-to-all`, as well as direct point-to-point (GPU-to-GPU) send and receive operations. It is optimized to achieve high bandwidth on platforms using PCIe and xGMI and networking using InfiniBand Verbs or TCP/IP sockets. RCCL supports an arbitrary number of GPUs installed in a single node or multiple nodes and can be used in either single- or multi-process (for example, MPI) applications.

The collective operations are implemented using ring and tree algorithms and have been optimized for throughput and latency by leveraging topology awareness, high-speed interconnects, and RDMA-based collectives. For best performance, small operations can be either batched into larger operations or aggregated through the API.

RCCL uses PCIe and xGMI high-speed interconnects for intra-node communication as well as InfiniBand, RoCE, and TCP/IP for inter-node communication. It supports an arbitrary number of GPUs installed in a single-node or multi-node platform and can easily integrate into single- or multi-process (for example, MPI) applications.



## INSTALL RCCL

Before you begin, verify that your system is supported. For more information, see ROCm Core SDK components.

For advanced workflows, source builds, or custom configurations, see *Building RCCL from source*.

### 2.1 Install the ROCm Core SDK

RCCL is included with the ROCm Core SDK on Linux. For the most complete installation, we recommend that developers use the `amdrocm-core-sdk` meta package.

For instructions, see Install AMD ROCm. Use the selector panel on that page to view instructions appropriate for your system environment.

### 2.2 Install RCCL on Linux

Alternatively, if you want to install RCCL without additional ROCm libraries and tools, install the `amdrocm-rccl` package. This package also includes base ROCm packages.

1. Complete the ROCm installation prerequisites to install dependencies and configure GPU access permissions.
2. Install the RCCL package that matches your ROCm version and development package needs. Package names use the following format:

```
amdrocm-rccl<-dev/devel><rocm_version>
```

Where:

- `<-dev/devel>` specifies whether to install library files and headers. Omit this suffix to only install runtime packages.
  - `-dev` is used on Debian-based distributions, including Ubuntu.
  - `-devel` is used on RPM-based distributions, including RHEL and SLES.
- `<rocm_version>` is the ROCm Core SDK version to install. Omit this suffix to install the latest available version.

For example: `amdrocm-rccl-dev7.13-gfx950`

Use the following command to install the latest RCCL development package release for supported GPU architectures:

### Debian-based distros

```
sudo apt install amdrocm-rccl-dev
```

### RHEL-based distros

```
sudo dnf install amdrocm-rccl-devel
```

### SLES

```
sudo zypper install amdrocm-rccl-devel
```

## 2.3 Install a nightly build

The [TheRock](#) build system also publishes nightly builds for the ROCm Core SDK and its components, including RCCL. See [Nightly release status](#) for details.

## BUILDING RCCL FROM SOURCE

To build RCCL as part of the ROCm Core SDK, see [TheRock build instructions](#). TheRock is the recommended way to build ROCm components from source.

Alternatively, you can build RCCL standalone using the following instructions.

To build RCCL directly from the source code, follow these steps. This guide also includes instructions explaining how to test the build. For information on using the quick start install script to build RCCL, see [Install RCCL](#).

### 3.1 Requirements

The following prerequisites are required to build RCCL:

1. ROCm-supported GPUs
2. Having the ROCm stack installed on the system, including the [HIP runtime](#) and the HIP-Clang compiler.

### 3.2 Build and install RCCL using the install script

To quickly install RCCL using the install script, follow these steps. For additional tips, see [RCCL usage tips](#).

#### 3.2.1 Quick start RCCL build

RCCL directly depends on the HIP runtime plus the HIP-Clang compiler, which are part of the ROCm software stack. For ROCm installation instructions, see the `rocm:install/rocm`.

Use the [install.sh helper script](#), located in the root directory of the RCCL repository, to build and install RCCL with a single command. It uses hard-coded configurations that can be specified directly when using `cmake`. However, it's a great way to get started quickly and provides an example of how to build and install RCCL.

### 3.2.2 Run the install script

To build the library using the install script, use this command:

```
./install.sh
```

To build for only for your local machine's AMD GPU architecture, reducing build time, add the `-l` flag.

```
./install.sh -l
```

For more information on the build options and flags for the install script, run the following command:

```
./install.sh --help
```

The RCCL build and installation helper script options are as follows:

```
RCCL build & installation helper script
Options:
  --address-sanitizer      Build with address sanitizer enabled
  --amdgpu_targets        Only compile for specified GPU architecture(s). For
↳ multiple targets, separate by ';' (builds for all supported GPU architectures by
↳ default)
  --cmake-options          Pass additional CMake options (e.g. --cmake-options "--
↳ DFOO=BAR -DBAZ=ON")
  --debug                  Build debug library
  --debug-fast            Build debug library with lto optimization disabled (fast
↳ build times)
  -d|--dependencies      Install RCCL dependencies
  --device-linker         Build with assembly-extract device linker (default)
  --disable-colltrace     Build without collective trace
  --disable-roctx         Build without ROCTX logging
  --disable-warp-speed    Disable WARP_SPEED kernel optimizations
  --dump-asm              Disassemble code and dump assembly with inline code
  -c|--enable-code-coverage Enable code coverage
  --enable_backtrace      Build with custom backtrace support
  --enable-mpi-tests      Enable MPI-based tests (requires --debug and MPI
↳ installation; set MPI_PATH if not in /opt/ompi)
  -f|--fast               Quick-build RCCL (local gpu arch only, no backtrace, and
↳ collective trace support)
  --force-reduce-pipeline Force reduce_copy sw pipeline to be used for every reduce
↳ based collectives and datatypes
  --generate-sym-kernels  Generate symmetric memory kernels (default: OFF)
  -h|--help               Prints this help message
  -i|--install            Install RCCL library (see --prefix argument below)
  -j|--jobs               Specify how many parallel compilation jobs to run (16 by
↳ default)
  --kernel-resource-use   Dump GPU kernel resource usage (e.g., VGPRs, scratch,
↳ spill) at link stage
  -l|--local_gpu_only     Only compile for local GPU architecture
  --log-trace             Build with log trace enabled (i.e. NCCL_DEBUG=TRACE)
  --no_clean              Don't delete files if they already exist
  --no-device-linker      Disable device linker, use standard -fgpu-rcd
  --npkit-enable          Compile with npkit enabled
  --openmp-test-enable    Enable OpenMP in rccl unit tests
```

(continues on next page)

(continued from previous page)

```

-p|--package_build      Build RCCL package
--prefix                Specify custom directory to install RCCL to (default: `/  

↳opt/rocm`)
-q|--quiet-warnings     Suppress majority of compiler warnings (not recommended)
--rocshmem              Build with rocSHMEM support
--run_tests_all         Run all rccl unit tests (must be built already)
-r|--run_tests_quick    Run small subset of rccl unit tests (must be built
↳already)
--static                Build RCCL as a static library instead of shared library
-t|--tests_build        Build rccl unit tests, but do not run
--time-trace            Plot the build time of RCCL (requires `ninja-build`
↳package installed on the system)
--verbose                Show compile commands

Available RCCL-specific CMake options for --cmake-options:
-DBUILD_EXT_EXAMPLES=ON      Build ext-  

↳(default: OFF)
-DDWORDX4_INTRINSICS=OFF     Disable dwordx4 intrinsics (default: ON)
-DENABLE_COMPRESS=OFF        Disable GPU code compression (default: ON)
-DENABLE_IFC=ON              Enable indirect function call (default: OFF)
-DDEFAULT_INJECTION=OFF      Disable fault injection (default: ON)
-DPROFILE=ON                  Enable profiling (default: OFF)
-DRCCCL_ROCPROFILER_REGISTER=OFF Disable rocprofiler-register support (default:
↳ON)
-DTIMETRACE=ON                Enable time-trace during compilation (default:
↳OFF)

Environment variables:
ONLY_FUNCS                  Build only specified collective functions (debug builds
↳only).
                             Restricts GPU kernel generation to the listed collectives,
↳ significantly
                             reducing build time during development. Use '|' to
↳separate multiple functions.
                             Example: ONLY_FUNCS="AllReduce|SendRecv" ./install.sh --
↳debug -t
                             Available: AllReduce, Broadcast, Reduce, AllGather,
↳ReduceScatter,
                             AlltoAllPivot, SendRecv, AlltoAllGda,
↳AlltoAllvGda
                             Advanced: Specify algo, protocol, redop, and type per
↳collective.
                             ONLY_FUNCS="AllReduce RING SIMPLE Sum f32|SendRecv"
ROCSHMEM_INSTALL_DIR       Path to a pre-built rocSHMEM installation (skips building
↳from source)

```

**Tip:** By default, the RCCL install script builds all the GPU targets that are defined in `DEFAULT_GPUS` in `CMakeLists.txt`. To target specific GPUs and potentially reduce the build time, use `--amdgpu_targets` along with a semi-colon (;) separated string list of the GPU targets.

## 3.3 Building the library using CMake

To build the library from source, follow these steps:

```
git clone --recursive https://github.com/ROCm/rccl.git
cd rccl
mkdir build
cd build
cmake ..
make -j 16      # Or some other suitable number of parallel jobs
```

If you have already cloned the repository, you can checkout the external submodules manually.

```
git submodule update --init --recursive --depth=1
```

You can substitute a different installation path by providing the path as a parameter to `CMAKE_INSTALL_PREFIX`, for example:

```
cmake -DCMAKE_INSTALL_PREFIX=$PWD/rccl-install -DCMAKE_BUILD_TYPE=Release ..
```

---

**Note:** Ensure ROCm CMake is installed using the command `apt install rocm-cmake`. By default, CMake builds the component in debug mode unless `DCMAKE_BUILD_TYPE` is specified.

---

### 3.3.1 Building the RCCL package and install package:

After you have cloned the repository and built the library as described in the previous section, use this command to build the package on Debian-based distros:

```
cd rccl/build
make package
sudo dpkg -i *.deb
```

---

**Note:** The RCCL package install process requires `sudo` or root access because it creates a directory named `rccl` in `/opt/rocm/`. This is an optional step. RCCL can be used directly by including the path containing `librccl.so`.

---

## 3.4 Testing RCCL

The RCCL unit tests are implemented using the Googletest framework in RCCL. These unit tests require Googletest 1.10 or higher to build and run (this dependency can be installed using the `-d` option for `install.sh`). To run the RCCL unit tests, go to the `build` folder and the `test` subfolder, then run the appropriate RCCL unit test executables.

The RCCL unit test names follow this format:

```
CollectiveCall.[Type of test]
```

Filtering of the RCCL unit tests can be done using environment variables and by passing the `--gtest_filter` command line flag:

```
UT_DATATYPES=ncclBfloat16 UT_REDOPS=prod ./rccl-UnitTests --gtest_filter="AllReduce.C*"
```

This command runs only the AllReduce correctness tests with the float16 datatype. A list of the available environment variables for filtering appears at the top of every run. See the [Googletest documentation](#) for more information on how to form advanced filters.

There are also other performance and error-checking tests for RCCL. They are maintained separately at <https://github.com/ROCm/rccl-tests>.

---

**Note:** For more information on how to build and run rccl-tests, see the [rccl-tests README file](#) .

---



## RUNNING RCCL USING DOCKER

To use Docker to run RCCL, Docker must already be installed on the system. To build the Docker image and run the container, follow these steps.

1. Build the Docker image

By default, the Dockerfile uses `docker.io/rocm/dev-ubuntu-22.04:latest` as the base Docker image. It then installs RCCL and `rccl-tests` (in both cases, it uses the version from the `develop` branch).

Use this command to build the Docker image:

```
docker build -t rccl-tests -f Dockerfile.ubuntu --pull .
```

The base Docker image, `rccl` repository, `rccl-tests` repository, and GPU targets can be modified by using `--build-args` in the `docker build` command above. For example, to use a different base Docker image for the MI250 GPU, use this command:

```
docker build -t rccl-tests -f Dockerfile.ubuntu --build-arg="ROCM_IMAGE_NAME=rocm/  
↪dev-ubuntu-20.04" --build-arg="ROCM_IMAGE_TAG=6.2" --build-arg="GPU_TARGETS=gfx90a  
↪" --pull .
```

2. Launch an interactive Docker container on a system with AMD GPUs:

```
docker run --rm --device=/dev/kfd --device=/dev/dri --group-add video --ipc=host --  
↪network=host --cap-add=SYS_PTRACE --security-opt seccomp=unconfined -it rccl-  
↪tests /bin/bash
```

To run, for example, the `all_reduce_perf` test from `rccl-tests` on 8 AMD GPUs from inside the Docker container, use this command for ROCm 6.4.2 or later:

```
mpirun --allow-run-as-root -np 8 --mca pml ucx --mca btl ^openib -x NCCL_  
↪DEBUG=VERSION /workspace/rccl-tests/build/all_reduce_perf -b 1 -e 16G -f 2 -g 1
```

For ROCm 6.4.1 or earlier:

```
mpirun --allow-run-as-root -np 8 --mca pml ucx --mca btl ^openib -x NCCL_  
↪DEBUG=VERSION -x HSA_NO_SCRATCH_RECLAIM=1 /workspace/rccl-tests/build/all_reduce_  
↪perf -b 1 -e 16G -f 2 -g 1
```

For more information on the `rccl-tests` options, see the [Usage guidelines](#) in the GitHub repository.



## USING THE RCCL TUNER PLUGIN API

An external plugin enables users to hand-tailor the selection of an algorithm, protocol, and number of channels (thread blocks) based on an input configuration specifying the message size, number of nodes and GPUs, and link types (for instance, PCIe, XGMI, or NET). One advantage of this plugin is that each user can create and maintain their own hand-tailored tuner without relying on RCCL to develop and maintain it. This topic describes the API required to implement an external tuner plugin for RCCL.

The following usage notes are relevant when using the RCCL Tuner plugin API:

- The API allows partial outputs: tuners can set only the algorithm and protocol and let RCCL set the remaining fields, such as the number of channels.
- If `getCollInfo()` fails, RCCL uses its default internal mechanisms to determine the best collective configuration.
- `getCollInfo` is called for each collective invocation per communicator, so special care must be taken to avoid introducing excessive latency.
- The supported RCCL algorithms are `NCCL_ALGO_TREE`, and `NCCL_ALGO_RING`.
- The supported RCCL protocols are `NCCL_PROTO_SIMPLE`, `NCCL_PROTO_LL` and `NCCL_PROTO_LL128`.
  - Until support is present for network collectives, use the example in the `pluginGetCollInfo` API implementation to ignore other algorithms as follows:

```
if ((a == NCCL_ALGO_COLLNET_DIRECT || a == NCCL_ALGO_COLLNET_CHAIN) &&
    collNetSupport != 1) continue;
if ((a == NCCL_ALGO_NVLS || a == NCCL_ALGO_NVLS_TREE) && nvlsSupport != 1)
    continue;
if (a == NCCL_ALGO_NVLS && collNetSupport != 1) continue;
```

---

**Note:** The `example plugin` uses math models to approximate the bandwidth and latency of the available selection of algorithms and protocols and select the one with the lowest calculated latency. It is customized for the AMD Instinct MI300 accelerators and RoCEv2 networks on a limited number of nodes. This example, which is intended for demonstration purposes only, is not meant to be inclusive of all potential AMD GPUs and network configuration.

---

## 5.1 API description

To build a custom tuner, implement the `ncclTuner_v1_t` structure.

### 5.1.1 Structure: `ncclTuner_v1_t`

#### Fields

- `name`
  - **Type:** `const char*`
  - **Description:** The name of the tuner, which can be used for logging purposes when `NCCL_DEBUG=info` and `NCCL_DEBUG_SUBSYS=tune` are set.

#### Functions

- `init` (called upon communicator initialization with `ncclCommInitRank`)

Initializes the tuner states. Each communicator initializes its tuner. `nNodes x nRanks` = the total number of GPUs participating in the collective communication.

- **Parameters:**

- \* `nRanks` (`size_t`): The number of devices (GPUs).
- \* `nNodes` (`size_t`): The number of operating system nodes (physical nodes or VMs).
- \* `logFunction` (`ncclDebugLogger_t`): A log function for certain debugging info.

- **Return:**

- \* **Type:** `ncclResult_t`
- \* **Description:** The result of the initialization.

- `getCollInfo` (called for each collective call per communicator)

Retrieves information about the collective algorithm, protocol, and number of channels for the given input parameters.

- **Parameters:**

- \* `collType` (`ncclFunc_t`): The collective type, for example, `allreduce`, `allgather`, etc.
- \* `nBytes` (`size_t`): The size of the collective in bytes.
- \* `collNetSupport` (`int`): Whether `collNet` supports this type.
- \* `nvlsSupport` (`int`): Whether NVLink SHARP supports this type.
- \* `numPipeOps` (`int`): The number of operations in the group.

- **Outputs:**

- \* `algorithm` (`int*`): The selected algorithm to be used for the given collective.
- \* `protocol` (`int*`): The selected protocol to be used for the given collective.
- \* `nChannels` (`int*`): The number of channels (and SMs) to be used.

- **Description:**

If `getCollInfo()` does not return `ncclSuccess`, RCCL falls back to its default tuning for the given collective. The tuner is allowed to leave fields unset, in which case RCCL automatically sets those fields.

– **Return:**

\* **Type:** `ncclResult_t`

\* **Description:** The result of the operation.

- `destroy` (called upon communicator finalization with `ncclCommFinalize`)

Terminates the plugin and cleans up any resources allocated by the tuner.

– **Return:**

\* **Type:** `ncclResult_t`

\* **Description:** The result of the cleanup process.

## 5.2 Build and usage instructions

To use the external plugin, implement the desired algorithm and protocol selection technique using the API described above. As a reference, the following example is based on the MI300 tuning table by default.

### 5.2.1 Building and using the example `libnccl-tuner.so` file

1. Build the `libnccl-tuner.so` file following the example Makefile.

```
cd $RCCL_HOME/plugins/tuner/example/  
make
```

2. Tell RCCL to use the custom `libnccl-tuner.so` file by setting the following environment variable to the file path:

```
export NCCL_TUNER_PLUGIN=$RCCL_HOME/plugins/tuner/example/libnccl-tuner.so
```



## USING THE NCCL NET PLUGIN API

NCCL provides a way to use external plugins to let NCCL run on many network types. This topic describes the NCCL Net plugin API and explains how to implement a network plugin for NCCL.

Plugins implement the NCCL network API and decouple NCCL binary builds, which are built against a particular version of the GPU stack (such as NVIDIA CUDA), from the network code, which is built against a particular version of the networking stack. Using this method, you can easily integrate any CUDA version with any network stack version.

NCCL network plugins are packaged as a shared library called `librccl-net.so`. The shared library contains one or more implementations of the NCCL Net API in the form of versioned structs, which are filled with pointers to all required functions.

### 6.1 Plugin architecture

When NCCL is initialized, it searches for a `librccl-net.so` library and dynamically loads it, then searches for symbols inside the library.

The `NCCL_NET_PLUGIN` environment variable allows multiple plugins to coexist. If it's set, NCCL looks for a library named `librccl-net- $\{NCCL\_NET\_PLUGIN\}$ .so`. It is therefore recommended that you name the library according to that pattern, with a symlink pointing from `librccl-net.so` to `librccl-net- $\{NCCL\_NET\_PLUGIN\}$ .so`. This lets users select the correct plugin if there are multiple plugins in the path.

#### 6.1.1 Struct versioning

After a library is found, NCCL looks for a symbol named `ncc1Net_vX`, with `X` increasing over time. This versioning pattern ensures that the plugin and the NCCL core are compatible.

Plugins are encouraged to provide a number of these symbols, implementing many versions of the NCCL Net API. This is so the same plugin can be compiled for and support a wide range of NCCL versions.

Conversely, and to ease transition, NCCL can choose to support different plugin versions. It can look for the latest `ncc1Net` struct version but also search for older versions, so that older plugins still work.

## 6.1.2 In-network collective operations (collNet)

In addition to the `ncclNet` structure, network plugins can provide a `collNet` structure which implements any supported in-network collective operations. This is an optional structure provided by the network plugin, but its versioning is tied to the `ncclNet` structure and many functions are common between the two to ease implementation. The `collNet` structure can be used by the NCCL `collNet` algorithm to accelerate inter-node reductions in `allReduce`.

## 6.1.3 Header management

To help users effortlessly build plugins, plugins should copy the `ncclNet_vX` definitions they support to their list of internal includes. An example is shown in `plugins/net/example/`, which stores all headers in the `nccl/` directory and provides thin layers to implement old versions on top of newer ones.

The `nccl/` directory is populated with `net_vX.h` files, which extract all relevant definitions from the old API versions. It also provides error codes in `err.h`.

## 6.2 API (v6)

Here is the main `ncclNet_v6` struct. Each function is explained in later sections.

```
typedef struct {
// Name of the network (mainly for logs)
const char* name;
// Initialize the network.
ncclResult_t (*init)(ncclDebugLogger_t logFunction);
// Return the number of adapters.
ncclResult_t (*devices)(int* ndev);
// Get various device properties.
ncclResult_t (*getProperties)(int dev, ncclNetProperties_v6_t* props);
// Create a receiving object and provide a handle to connect to it. The
// handle can be up to NCCL_NET_HANDLE_MAXSIZE bytes and will be exchanged
// between ranks to create a connection.
ncclResult_t (*listen)(int dev, void* handle, void** listenComm);
// Connect to a handle and return a sending comm object for that peer.
// This call must not block for the connection to be established, and instead
// should return successfully with sendComm == NULL with the expectation that
// it will be called again until sendComm != NULL.
ncclResult_t (*connect)(int dev, void* handle, void** sendComm);
// Finalize connection establishment after remote peer has called connect.
// This call must not block for the connection to be established, and instead
// should return successfully with recvComm == NULL with the expectation that
// it will be called again until recvComm != NULL.
ncclResult_t (*accept)(void* listenComm, void** recvComm);
// Register/Deregister memory. Comm can be either a sendComm or a recvComm.
// Type is either NCCL_PTR_HOST or NCCL_PTR_CUDA.
ncclResult_t (*regMr)(void* comm, void* data, int size, int type, void** mhandle);
/* DMA-BUF support */
ncclResult_t (*regMrDmaBuf)(void* comm, void* data, size_t size, int type, uint64_t
↪offset, int fd, void** mhandle);
ncclResult_t (*deregMr)(void* comm, void* mhandle);
// Asynchronous send to a peer.
```

(continues on next page)

(continued from previous page)

```

// May return request == NULL if the call cannot be performed (or would block)
ncclResult_t (*isend)(void* sendComm, void* data, int size, int tag, void* mhandle,
↳void** request);
// Asynchronous recv from a peer.
// May return request == NULL if the call cannot be performed (or would block)
ncclResult_t (*irecv)(void* recvComm, int n, void** data, int* sizes, int* tags, void**
↳mhandles, void** request);
// Perform a flush/fence to make sure all data received with NCCL_PTR_CUDA is
// visible to the GPU
ncclResult_t (*iflush)(void* recvComm, int n, void** data, int* sizes, void** mhandles,
↳void** request);
// Test whether a request is complete. If size is not NULL, it returns the
// number of bytes sent/received.
ncclResult_t (*test)(void* request, int* done, int* sizes);
// Close and free send/recv comm objects
ncclResult_t (*closeSend)(void* sendComm);
ncclResult_t (*closeRecv)(void* recvComm);
ncclResult_t (*closeListen)(void* listenComm);
} ncclNet_v6_t;

```

## 6.2.1 Error codes

All plugins functions use NCCL error codes as their return value. `ncclSuccess` should be returned upon success. Otherwise, plugins can return one of the following codes:

- `ncclSystemError` is the most common error for network plugins. It should be returned when a call to the Linux kernel or a system library fails. This typically includes all network and hardware errors.
- `ncclInternalError` is returned when the NCCL core code is using the network plugin in an incorrect way, for example, allocating more requests than it should or passing an invalid argument in API calls.
- `ncclInvalidUsage` should be returned when the error is most likely due to user error. This can include mis-configuration, but also size mismatches.
- `ncclInvalidArgument` should not typically be used by plugins because arguments should be checked by the NCCL core layer.
- `ncclUnhandledCudaError` is returned when an error is received from NVIDIA CUDA. Network plugins should not need to rely on CUDA, so this error should not be common.

## 6.2.2 Operational overview

NCCL first calls the `init` function, queries the number of network devices with the `devices` function, and retrieves the properties from each network device using `getProperties`.

To establish a connection between two network devices, NCCL first calls `listen` on the receiving side. It passes the returned handle to the sender side of the connection, and uses it to call `connect`. Finally, `accept` is called on the receiving side to finalize the establishment of the connection.

After the connection is established, communication is performed using the functions `isend`, `irecv`, and `test`. Prior to calling `isend` or `irecv`, NCCL calls the `regMr` function on all buffers to allow RDMA NICs to prepare the buffers. `deregMr` is used to unregister buffers.

In certain conditions, `iflush` is called after a receive call completes to allow the network plugin to flush data and ensure the GPU processes the newly written data.

To close the connections, NCCL calls `closeListen` to close the object returned by `listen`, `closeSend` to close the object returned by `connect`, and `closeRecv` to close the object returned by `accept`.

## 6.2.3 API Functions

The RCCL Tuner plugin API provides the following interface for initialization, connection management, and communications.

### Initialization

- `name` - The name field should point to a character string with the name of the network plugin. This name is used for all logging, especially when `NCCL_DEBUG=INFO` is set.

---

**Note:** Setting `NCCL_NET=<plugin name>` ensures a specific network implementation is used, with a matching name. This is not to be confused with `NCCL_NET_PLUGIN` which defines a suffix for the `librccl-net.so` library name to load.

---

- `init` - As soon as NCCL finds the plugin and the correct `ncclNet` symbol, it calls the `init` function. This allows the plugin to discover network devices and ensure they are usable. If the `init` function does not return `ncclSuccess`, then NCCL does not use the plugin and falls back to internal ones.

To allow the plugin logs to seamlessly integrate into the NCCL logs, NCCL provides a logging function to `init`. This function is typically used to allow `INFO` and `WARN` macros within the plugin code by adding the following definitions:

```
#define WARN(...) logFunction(NCCL_LOG_WARN, NCCL_ALL, __FILE__, __LINE__, __VA_
↪ARGS__)
#define INFO(FLAGS, ...) logFunction(NCCL_LOG_INFO, (FLAGS), __func__, __LINE__, __
↪VA_ARGS__)
```

- `devices` - After the plugin is initialized, NCCL queries the number of devices available. This should not be zero. Otherwise, NCCL initialization will fail. If no device is present or usable, the `init` function should not return `ncclSuccess`.
- `getProperties` - Right after retrieving the number of devices, NCCL queries the properties for each available network device. These properties are necessary when multiple adapters are present to ensure NCCL uses each adapter in the optimal way.
  - The `name` is only used for logging.
  - The `pciPath` is the base for all topology detection and should point to the PCI device directory in `/sys`. This is typically the directory pointed to by `/sys/class/net/eth0/device` or `/sys/class/infiniband/mlx5_0/device`. If the network interface is virtual, then `pciPath` should be `NULL`.
  - The `guid` field is used to determine whether network adapters are connected to multiple PCI endpoints. For normal cases, this is set to the device number. If multiple network devices have the same `guid`, then NCCL understands them to be sharing the same network port to the fabric. In this case, it will not use the port multiple times.
  - The `ptrSupport` field indicates whether or not CUDA pointers are supported. If so, it should be set to `NCCL_PTR_HOST|NCCL_PTR_CUDA`. Otherwise, it should be set to `NCCL_PTR_HOST`. If the plugin supports `dmabuf`, it should set `ptrSupport` to `NCCL_PTR_HOST|NCCL_PTR_CUDA|NCCL_PTR_DMABUF` and provide a `regMrDmaBuf` function.

- The `regIsGlobal` field allows NCCL to register buffers in advance, for example, using a loopback connection. Later, it also lets NCCL expect that a subsequent registration on a buffer from a previous registration will happen nearly immediately, because the buffer is already known by the network adapter. A typical implementation maintains a registration cache, with the call to `ncclCommRegister` creating the initial entry in the cache using `regMr()` on a loopback connection. Any later call to the NCCL system can call `regMr()` again on the real connection, with the real buffer (which could be at a different offset within the original buffer, with a smaller size, for example). It could then call `deregMr()` immediately afterwards. The `ncclCommDeregister` call should issue the final call to `deregMr()` and effectively remove the mapping on the network adapter.
- The `speed` field indicates the speed of the network port in Mbps ( $10^6$  bits per second). This ensures proper optimization of flows within the node.
- The `port` field indicates the port number. This is important for topology detection and flow optimization within the node when a NIC with a single PCI connection is connected to the fabric through multiple ports.
- The `latency` field indicates the network latency in microseconds. This can be useful to improve the NCCL tuning and ensure NCCL switches from tree to ring at the correct size.
- The `maxComms` field indicates the maximum number of connections that can be created.
- The `maxRecvs` field indicates the maximum number for grouped receive operations (see grouped receive).

## Connection establishment

Connections are used in an unidirectional manner, with a sender side and a receiver side.

- `listen` - To create a connection, NCCL calls `listen` on the receiver side. This function accepts a device number as an input argument and returns a local `listenComm` object and a `handle` to pass to the other side of the connection, so that the sender can connect to the receiver. The `handle` is a buffer of size `NCCL_NET_HANDLE_MAXSIZE` and is provided by NCCL. This call should never block, but unlike `connect` and `accept`, `listenComm` should never be `NULL` if the call succeeds.
- `connect` - NCCL uses its bootstrap infrastructure to provide the `handle` to the sender side, then calls `connect` on the sender side on a given device index `dev` and provides the `handle`. `connect` should not block either. Instead, it sets `sendComm` to `NULL` and returns `ncclSuccess`. In that case, NCCL will keep calling `accept` again until it succeeds.
- `accept` - To finalize the connection, the receiver side calls `accept` on the `listenComm` object previously returned by the `listen` call. If the sender did not connect yet, `accept` should not block. It should return `ncclSuccess`, setting `recvComm` to `NULL`. NCCL will keep calling `accept` again until it succeeds.
- `closeListen` / `closeSend` / `closeRecv` - When a `listenComm`, `sendComm`, or `recvComm` object is no longer needed, NCCL calls `closeListen`, `closeSend`, or `closeRecv` to free the associated resources.

## Communication

Communication is handled using the asynchronous send and receive operations: `isend`, `irecv`, and `test`. To support RDMA capabilities, buffer registration and flush functions are provided.

To keep track of asynchronous send, receive, and flush operations, requests are returned to NCCL, then queried using `test`. Each `sendComm` or `recvComm` must be able to handle `NCCL_NET_MAX_REQUESTS` requests in parallel.

---

**Note:** This value should be multiplied by the multi-receive capability of the plugin for the sender side, so the plugin can effectively have `NCCL_NET_MAX_REQUESTS` multi-receive operations happening in parallel. If `maxRecvs` is 8

and `NCCL_NET_MAX_REQUESTS` is 8, then each `sendComm` must be able to handle up to 64 (8x8) concurrent `isend` operations.

---

- `regMr` - Prior to sending or receiving data, NCCL calls `regMr` with any buffers later used for communication. It provides a `sendComm` or `recvComm` object for the `comm` argument, the buffer pointer `data`, the `size`, and the `type`. The `type` is either `NCCL_PTR_HOST` or `NCCL_PTR_CUDA` if the network supports CUDA pointers.

The network plugin can use the output argument `mhandle` to store any reference to the memory registration, because `mhandle` is returned for all `isend`, `irecv`, `iflush`, and `deregMr` calls.

- `regMrDmaBuf` - If the plugin has set the `NCCL_PTR_DMABUF` property in `ptrSupport`, NCCL uses `regMrDmaBuf` instead of `regMr`. If the property was not set, `regMrDmaBuf` can be set to `NULL`.
- `deregMr` - When buffers are no longer used for communication, NCCL calls `deregMr` to let the plugin free resources. This function is used to deregister handles returned by `regMr` and `regMrDmaBuf`.
- `isend` - Data is sent through the connection using `isend`, passing the `sendComm` object previously created by `connect`, the buffer described by `data`, `size`, and `mhandle`. A `tag` must be used if the network supports multi-recv operations (see `irecv`) to distinguish between different send requests matching the same multi-recv. Otherwise it can be set to `0`.

The `isend` operation returns a handle in the `request` argument for further calls to `test`. If the `isend` operation cannot be initiated, `request` is set to `NULL`. NCCL will call `isend` again later.

- `irecv` - To receive data, NCCL calls `irecv` with the `recvComm` returned by `accept`. The argument `n` configures NCCL for multi-recv, to allow grouping of multiple sends through a single network connection. Each buffer can be described by the `data`, `sizes`, and `mhandles` arrays. `tags` specify a tag for each receive so that each of the `n` independent `isend` operations is received into the right buffer.

If all receive operations can be initiated, `irecv` returns a handle in the `request` pointer. Otherwise, it sets the pointer to `NULL`. In the case of multi-recv, all `n` receive operations are handled by a single request handle.

The sizes provided to `irecv` can (and will) be larger than the size of the `isend` operation. However, it is an error if the receive size is smaller than the send size.

---

**Note:** For a given connection, send and receive operations should always match in the order they were posted. Tags provided for receive operations are only used to assign a given send operation to one of the buffers of the first (multi-)receive operation in the queue, not to allow for out-of-order tag matching on any receive operation posted.

---

- `test` - After an `isend` or `irecv` operation is initiated, NCCL calls `test` on the request handles until the operation completes. When that happens, `done` is set to 1 and `sizes` is set to the real size sent or received, the latter could potentially be lower than the size passed to `irecv`.

In the case of a multi-recv, all receives are considered as part of a single operation, the goal being to allow aggregation. Therefore, they share a single request and a single `done` status. However, they can have different sizes, so if `done` is non-zero, the `sizes` array should contain the `n` sizes corresponding to the buffers passed to `irecv`.

After `test` returns 1 in `done`, the request handle can be freed. This means that NCCL will never call `test` again on that request, unless it is reallocated by another call to `isend` or `irecv`.

- `iflush` - After a receive operation completes, if the operation was targeting GPU memory and received a non-zero number of bytes, NCCL calls `iflush`. This lets the network flush any buffer to ensure the GPU can read it immediately without seeing stale data. This flush operation is decoupled from the `test` code to improve the latency of LL\* protocols, because those are capable of determining when data is valid or not.

`iflush` returns a request which must be queried using `test` until it completes.

## TROUBLESHOOTING RCCL

This topic explains the steps to troubleshoot functional and performance issues with RCCL. While debugging, collect the output from the commands in this guide. This data can be used as supporting information when submitting an issue report to AMD.

### 7.1 Collecting system information

Collect this information about the ROCm version, GPU/accelerator, platform, and configuration.

- Verify the ROCm version. This might be a release version or a mainline or staging version. Use this command to display the version:

```
cat /opt/rocm/.info/version
```

Run the following command and collect the output:

```
rocm_agent_enumerator
```

Also, collect the name of the GPU or accelerator:

```
rocminfo
```

- Run these `amd-smi` commands to display the system topology.

```
amd-smi
amd-smi topology
amd-smi static --driver
amd-smi firmware
amd-smi xgmi
```

- Determine the values of the `PATH` and `LD_LIBRARY_PATH` environment variables.

```
echo $PATH
echo $LD_LIBRARY_PATH
```

- Collect the HIP configuration.

```
/opt/rocm/bin/hipconfig --full
```

- Verify the network settings and setup. Use the `ibv_devinfo` command to display information about the available RDMA devices and determine whether they are installed and functioning properly. Run `rdma link` to print a summary of the network links.

```
ibv_devinfo
rdma link
```

### 7.1.1 Isolating the issue

The problem might be a general issue or specific to the architecture or system. To narrow down the issue, collect information about the GPU or accelerator and other details about the platform and system. Some issues to consider include:

- Is ROCm running on:
  - A bare-metal setup
  - In a Docker container (determine the name of the Docker image)
  - In an SR-IOV virtualized
  - Some combination of these configurations
- Is the problem only seen on a specific GPU architecture?
- Is it only seen on a specific system type?
- Is it happening on a single node or multinode setup?
- Use the following troubleshooting techniques to attempt to isolate the issue.
  - Build or run the develop branch version of RCCL and see if the problem persists.
  - Try an earlier RCCL version (minor or major).
  - If you recently changed the ROCm runtime configuration, AMD Kernel-mode GPU Driver (KMD), or compiler, rerun the test with the previous configuration.

## 7.2 Collecting RCCL information

Collect the following information about the RCCL installation and configuration.

- Run the `ldd` command to list any dynamic dependencies for RCCL.

```
ldd <specify-path-to-librccl.so>
```

- Determine the RCCL version. This might be the pre-packaged component in `/opt/rocm/lib` or a version that was built from source. To verify the RCCL version, enter the following command, then run either `rccl-tests` or an e2e application.

```
export NCCL_DEBUG=VERSION
```

- Run `rccl-tests` and collect the results. For information on how to build and run `rccl-tests`, see the [rccl-tests GitHub](#).
- Collect the RCCL logging information. Enable the debug logs, then run `rccl-tests` or any e2e workload to collect the logs. Use the following command to enable the logs.

```
export NCCL_DEBUG=INFO
```

## 7.2.1 Using the RCCL Replayer

The RCCL Replayer is a debugging tool designed to analyze and replay the collective logs obtained from RCCL runs. It can be helpful when trying to reproduce problems, because it uses dummy data and doesn't have any dependencies on non-RCCL calls. For more information, see [RCCL Replayer GitHub documentation](#).

You must build the RCCL Replayer before you can use it. To build it, run these commands. Ensure `MPI_DIR` is set to the path where MPI is installed.

```
cd rccl/tools/rccl_replayer
MPI_DIR=/path/to/mpi make
```

To use the RCCL Replayer, follow these steps:

1. Collect the per-rank logs from the RCCL run by adding the following environment variables. This prevents any race conditions that might cause ranks to interrupt the output from other ranks.

```
NCCL_DEBUG=INFO NCCL_DEBUG_SUBSYS=COLL NCCL_DEBUG_FILE=some_name_here.%h.%p.log
```

2. Combine all the logs into a single file. This will become the input to the RCCL Replayer.

```
cat some_name_here_*.log > some_name_here.log
```

3. Run the RCCL Replayer using the following command. Replace `<numProcesses>` with the number of MPI processes to run, `</path/to/logfile>` with the path to the collective log file generated during the RCCL runs, and `<numGpusPerMpiRank>` with the number of GPUs per MPI rank used in the application.

```
mpirun -np <numProcesses> ./rcclReplayer </path/to/logfile> <numGpusPerMpiRank>
```

In a multi-node application environment, you can replay the collective logs on multiple nodes using the following command:

```
mpirun --hostfile <path/to/hostfile.txt> -np <numProcesses> ./rcclReplayer </path/
↳to/logfile> <numGpusPerMpiRank>
```

---

**Note:** Depending on the MPI library you're using, you might need to modify the `mpirun` command.

---

## 7.3 Analyzing performance issues

If the issues involve performance issues in an e2e workload, try the following microbenchmarks and collect the results. Follow the instructions in the subsequent sections to run these benchmarks and provide the results to the support team.

- TransferBench
- RCCL Unit Tests
- rccl-tests

### 7.3.1 Collect the TransferBench data

TransferBench allows you to benchmark simultaneous copies between user-specified devices. For more information, see the [TransferBench documentation](#).

To collect the TransferBench data, follow these steps:

1. Clone the TransferBench Git repository.

```
git clone https://github.com/ROCm/TransferBench.git
```

2. Change to the new directory and build the component.

```
cd TransferBench
make
```

3. Run the TransferBench utility with the following parameters and save the results.

```
USE_FINE_GRAIN=1 GFX_UNROLL=2 ./TransferBench a2a 64M 8
```

### 7.3.2 Collect the RCCL microbenchmark data

To use the RCCL tests to collect the RCCL benchmark data, follow these steps:

1. Disable NUMA auto-balancing using the following command:

```
sudo sysctl kernel.numa_balancing=0
```

Run the following command to verify the setting. The expected output is 0.

```
cat /proc/sys/kernel/numa_balancing
```

2. Build MPI, RCCL, and rccl-tests. To download and install MPI, see either [OpenMPI](#) or [MPICH](#). To learn how to build and run rccl-tests, see the [rccl-tests GitHub](#).
3. Run rccl-tests with MPI and collect the performance numbers.

## 7.4 RCCL and NCCL comparisons

If you are also using NVIDIA hardware or NCCL and notice a performance gap between the two systems, collect the system and performance data on the NVIDIA platform. Provide both sets of data to the support team.

## RCCL USAGE TIPS

This topic describes some of the more common RCCL extensions, such as NPKit, and provides tips on how to configure and customize the application.

### 8.1 NPKit

RCCL integrates [NPKit](#), a profiler framework that enables the collection of fine-grained trace events in RCCL components, especially in giant collective GPU kernels. See the [NPKit sample workflow for RCCL](#) for a fully-automated usage example. It also provides useful templates for the following manual instructions.

To manually build RCCL with NPKit enabled, pass `-DNPKIT_FLAGS="-DENABLE_NPKIT -DENABLE_NPKIT_... (other NPKit compile-time switches)"` to the `cmake` command. All NPKit compile-time switches are declared in the RCCL code base as macros with the prefix `ENABLE_NPKIT_`. These switches control the information that is collected.

---

**Note:** NPKit only supports the collection of non-overlapped events on the GPU. The `-DNPKIT_FLAGS` settings must follow this rule.

---

To manually run RCCL with NPKit enabled, set the environment variable `NPKIT_DUMP_DIR` to the NPKit event dump directory. NPKit only supports one GPU per process. To manually analyze the NPKit dump results, use [npkit\\_trace\\_generator.py](#).

MSCCL and MSCCL++ integration has been removed from RCCL. The legacy API symbols `mscclLoadAlgo`, `mscclRunAlgo`, and `mscclUnloadAlgo` remain as no-ops for link compatibility.

### 8.2 Enabling peer-to-peer transport

To enable peer-to-peer access on machines with PCIe-connected GPUs, set the HSA environment variable as follows:

```
HSA_FORCE_FINE_GRAIN_PCIE=1
```

This feature requires GPUs that support peer-to-peer access along with proper large BAR addressing support.

## 8.3 Ignoring CPU affinity with multi-node

Depending on the job launcher and the requirements of your workload, performance as the communication workload scales can be improved by setting `NCCL_IGNORE_CPU_AFFINITY`. This allows the RCCL communication library to ignore the job's supplied CPU affinity and use the GPU affinity only.

```
NCCL_IGNORE_CPU_AFFINITY=1
```

For general usage, this environment variable is not set so it doesn't interfere with the user or launcher supplied preferences.

## 8.4 Improving performance on the MI300X

This section outlines ways to improve RCCL performance on MI300X systems, including guidelines for systems with fewer than eight GPUs and the most efficient GPU partition modes.

### 8.4.1 Configuration with fewer than eight GPUs

On a system with eight MI300X accelerators, each pair of accelerators is connected with dedicated Infinity Fabric™ links in a fully connected topology. For collective operations, this can achieve good performance when all eight accelerators (and all Infinity Fabric links) are used. When fewer than eight GPUs are used, however, this can only achieve a fraction of the potential bandwidth on the system. However, if your workload warrants using fewer than eight MI300X accelerators on a system, you can set the run-time variable `NCCL_MIN_NCHANNELS` to increase the number of channels. For example:

```
export NCCL_MIN_NCHANNELS=32
```

Increasing the number of channels can benefit performance, but it also increases GPU utilization for collective operations. Additionally, RCCL pre-defines a higher number of channels when only two or four accelerators are in use on a 8\*MI300X system. In this situation, RCCL uses 32 channels with two MI300X accelerators and 24 channels for four MI300X accelerators.

### 8.4.2 NPS4 and CPX partition modes

The term compute partitioning modes, or Modular Chiplet Platform (MCP), refers to the logical partitioning of XCDs into devices in the ROCm stack. The names are derived from the number of logical partitions that are created out of the eight XCDs. In the default mode, SPX (Single Partition X-celerator), all eight XCDs are viewed as a single logical compute element, meaning that the `amd-smi` utility will show a single MI300X device. In CPX (Core Partitioned X-celerator) mode, each XCD appears as a separate logical GPU, for example, as eight separate GPUs in `amd-smi` per MI300X. CPX mode can be viewed as having explicit scheduling privileges for each individual compute element (XCD).

While compute partitioning modes change the space on which you can assign work to compute units, the memory partitioning modes (known as Non-Uniform Memory Access (NUMA) Per Socket (NPS)) change the number of NUMA domains that a device exposes. In other words, it changes the number of HBM stacks which are accessible to a compute unit, and therefore the size of its memory space. However, for the MI300X, the number of memory partitions must be less than or equal to the number of compute partitions. NPS4 (viewing pairs of HBM stacks as a disparate element), for example, is only enabled when in CPX mode (viewing each XCD as a disparate element).

- Compute partition modes

- In SPX mode, workgroups launched to the device are distributed round-robin to the XCDs in the device, meaning that the programmer cannot have explicit control over which XCD a workgroup is assigned to.
  - In CPX mode, workgroups are launched to a single XCD, meaning the programmer has explicit control over work placement onto the XCDs.
- Memory partition modes
    - In NPS1 mode (compatible with CPX and SPX), the entire memory is accessible to all XCDs.
    - In NPS4 mode (compatible with CPX), each memory quadrant of the memory is directly visible to the logical devices in its quadrant. An XCD can still access all portions of memory through multi-GPU programming techniques.

The MI300 CPX mode can be accessed using the following `amdsmi:index` commands.

```
amd-smi set --gpu all --compute-partition CPX
amd-smi set --gpu all --memory-partition NPS4
```

### RCCL performance with CPX and NPS4

To run RCCL allreduce on 64 GPUs with CPX+NPS4 mode on the MI300X, use this example:

```
mpirun -np 64 --bind-to numa rccl-tests/build/all_reduce_perf -b 8 -e 1G -f 2 -g 1
```

To run RCCL allreduce on 8 GPUs in the same OAM with CPX+NPS4 mode on the MI300X, use this example:

```
export ROCR_VISIBLE_DEVICES=0,1,2,3,4,5,6,7
mpirun -np 8 --bind-to numa rccl-tests/build/all_reduce_perf -b 8 -e 1G -f 2 -g 1
```

RCCL delivers improved allreduce performance in CPX mode for TP=8 (8 GPUs in the same OAM) on the MI300X.

```
export HIP_FORCE_DEV_KERNARG=1
export ROCR_VISIBLE_DEVICES=0,1,2,3,4,5,6,7
mpirun -np 8 --bind-to numa rccl-tests/build/all_reduce_perf -b 32 -e 1G -f 2 -g 1 -G 2 -
↳w 20 -n 50
```

Here are the benchmark results for in-place (where the output buffer is used as the input buffer) and out-of-place allreduce bus bandwidth.

A significant performance improvement is achievable with optimized CPX mode, which peaks at ~340 GB/s with a single OAM. The difference in bus bandwidth between the unoptimized and optimized modes increases as the buffer size grows.

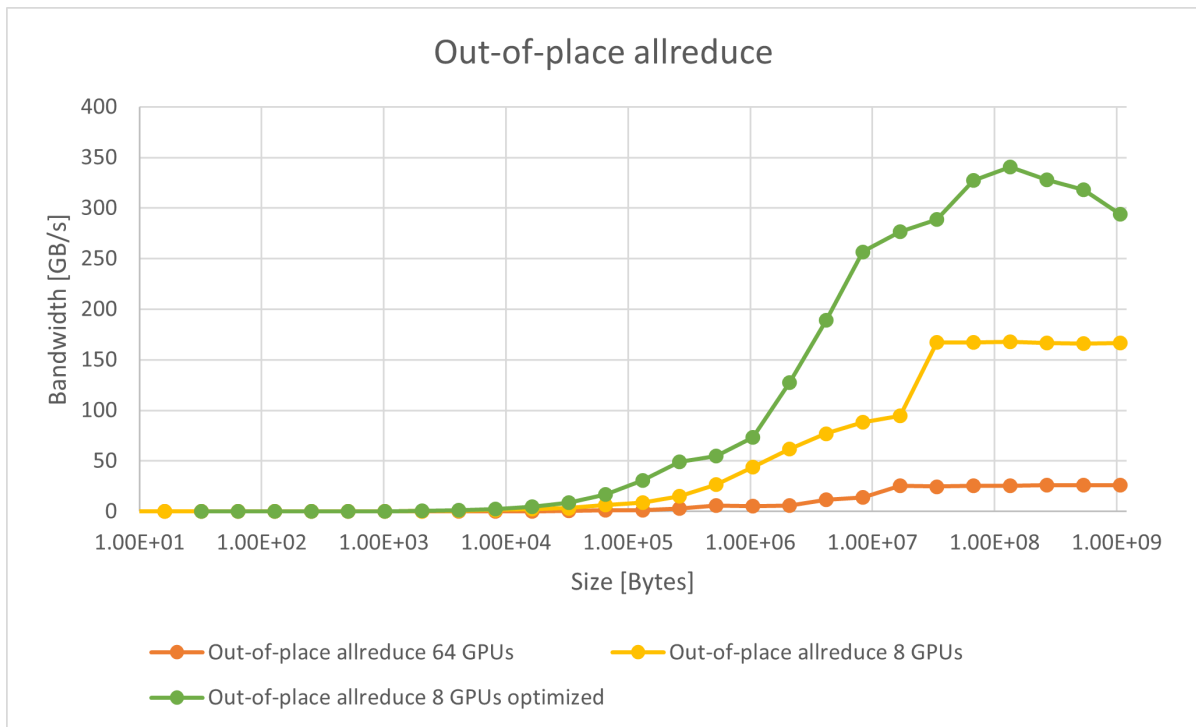
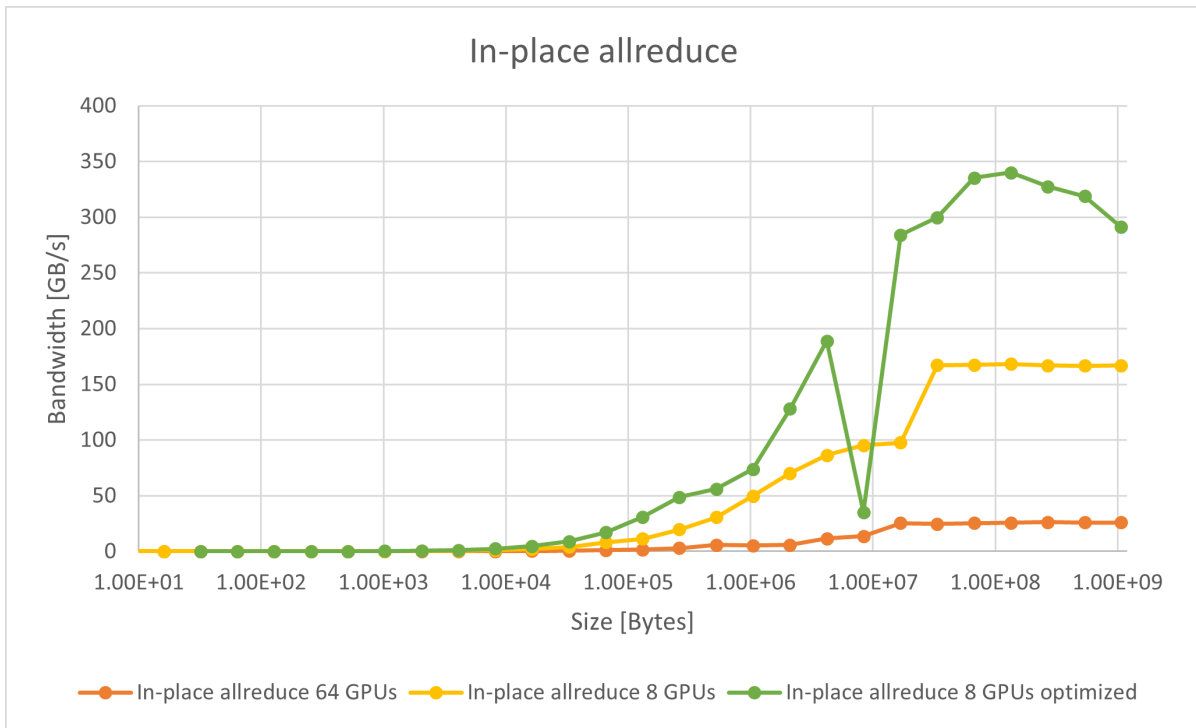
### Using RCCL and CPX in PyTorch

The PyTorch `all_reduce` benchmark is used to reproduce the performance reported by RCCL-Tests with the RCCL and CPX optimizations.

---

**Note:** To use RCCL with CPX mode in PyTorch, check the RCCL version used by PyTorch.

For a virtualenv with a .whl-based PyTorch setup (such as `nightly/rocm6.2`), this would be in `<path-to-your-venv>/lib/<python-version>/site-packages/torch/lib/librccl.so`. This is the version of RCCL that is packaged as part of ROCm version 6.2.



RCCL for CPX mode was enabled in ROCm 6.3.0. To use the CPX features, replace the existing `librccl.so` with one from ROCm 6.3.0 or newer or from a local build of the RCCL develop branch.

To test the effects of RCCL on PyTorch, the [stas00 all reduce benchmark](#) was used. The following command is used to run a single OAM allreduce benchmark:

```
export ROCR_VISIBLE_DEVICES=0,1,2,3,4,5,6,7
python -u -m torch.distributed.run --nproc_per_node=8 --rdzv_endpoint localhost:6000 --
↳rdzv_backend c10d all_reduce_bench.py
```

For better performance, the `HIP_FORCE_DEV_KERNARG` and `TORCH_NCCL_USE_TENSOR_REGISTER_ALLOCATOR_HOOK` environment variables are set during the benchmark in the following manner:

```
export TORCH_NCCL_USE_TENSOR_REGISTER_ALLOCATOR_HOOK=1
export HIP_FORCE_DEV_KERNARG=1
export ROCR_VISIBLE_DEVICES=0,1,2,3,4,5,6,7
python -u -m torch.distributed.run --nproc_per_node=8 --rdzv_endpoint localhost:6000 --
↳rdzv_backend c10d all_reduce_bench.py
```

The default allreduce PyTorch benchmark peak bus bandwidth performance is ~170 GB/s on a single OAM with ROCm 6.2.4, while the optimized run for CPX on a single OAM peaks at ~315 GB/s.

### 8.4.3 Context tracking on GPUs

Context tracking is disabled by default for optimal performance. However, enabling of context tracking can significantly improve performance in certain scenarios. To enable context tracking, set the following environment variable:

```
export RCCL_ENABLE_CONTEXT_TRACKING=1
```



## RCCL LIBRARY SPECIFICATION

This document provides details of the API library.

### 9.1 Communicator functions

*ncclResult\_t* **ncclGetUniqueId**(*ncclUniqueId* \*uniqueId)

Generates an ID for `ncclCommInitRank`.

Generates an ID to be used in `ncclCommInitRank`. `ncclGetUniqueId` should be called once by a single rank and the ID should be distributed to all ranks in the communicator before using it as a parameter for `ncclCommInitRank`.

#### Parameters

**uniqueId** – [out] Pointer to where `uniqueId` will be stored

#### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommInitRank**(*ncclComm\_t* \*comm, int nranks, *ncclUniqueId* commId, int rank)

Creates a new communicator (multi thread/process version).

Rank must be between 0 and `nranks-1` and unique within a communicator clique. Each rank is associated to a CUDA device, which has to be set before calling `ncclCommInitRank`. `ncclCommInitRank` implicitly synchronizes with other ranks, so it must be called by different threads/processes or use `ncclGroupStart/ncclGroupEnd`.

#### Parameters

- **comm** – [out] Pointer to created communicator
- **nranks** – [in] Total number of ranks participating in this communicator
- **commId** – [in] UniqueId required for initialization
- **rank** – [in] Current rank to create communicator for

#### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommInitAll**(*ncclComm\_t* \*comm, int ndev, const int \*devlist)

Creates a clique of communicators (single process version).

This is a convenience function to create a single-process communicator clique. Returns an array of `ndev` newly initialized communicators in `comm`. `comm` should be pre-allocated with size at least `ndev*sizeof(ncclComm_t)`. If `devlist` is NULL, the first `ndev` HIP devices are used. Order of `devlist` defines user-order of processors within the communicator.

#### Parameters

- **comm** – [out] Pointer to array of created communicators
- **ndev** – [in] Total number of ranks participating in this communicator
- **devlist** – [in] Array of GPU device indices to create for

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommDestroy**(*ncclComm\_t* comm)

Frees local resources associated with communicator object.

Destroy all local resources associated with the passed in communicator object

**Parameters**

**comm** – [in] Communicator to destroy

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommAbort**(*ncclComm\_t* comm)

Abort any in-progress calls and destroy the communicator object.

Frees resources associated with communicator object and aborts any operations that might still be running on the device.

**Parameters**

**comm** – [in] Communicator to abort and destroy

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommCount**(const *ncclComm\_t* comm, int \*count)

Gets the number of ranks in the communicator clique.

Returns the number of ranks in the communicator clique (as set during initialization)

**Parameters**

- **comm** – [in] Communicator to query
- **count** – [out] Pointer to where number of ranks will be stored

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommCuDevice**(const *ncclComm\_t* comm, int \*device)

Get the ROCm device index associated with a communicator.

Returns the ROCm device number associated with the provided communicator.

**Parameters**

- **comm** – [in] Communicator to query
- **device** – [out] Pointer to where the associated ROCm device index will be stored

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommUserRank**(const *ncclComm\_t* comm, int \*rank)

Get the rank associated with a communicator.

Returns the user-ordered “rank” associated with the provided communicator.

**Parameters**

- **comm** – [in] Communicator to query
- **rank** – [out] Pointer to where the associated rank will be stored

#### Returns

Result code. See *Result Codes* for more details.

## 9.2 Collective communication operations

Collective communication operations must be called separately for each communicator in a communicator clique.

They return when operations have been enqueued on the hipstream.

Since they may perform inter-CPU synchronization, each call has to be done from a different thread or process, or need to use Group Semantics (see below).

*ncclResult\_t* **ncclReduce**(const void \*sendbuff, void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, *ncclRedOp\_t* op, int root, *ncclComm\_t* comm, hipStream\_t stream)

Reduce.

Reduces data arrays of length *count* in *sendbuff* into *recvbuff* using *op* operation. *recvbuff\** may be NULL on all calls except for root device. *root\** is the rank (not the HIP device) where data will reside after the operation is complete. In-place operation will happen if *sendbuff* == *recvbuff*.

#### Parameters

- **sendbuff** – [in] Local device data buffer to be reduced
- **recvbuff** – [inout] Data buffer where result is stored (only for *root* rank). May be null for other ranks.
- **count** – [in] Number of elements in every send buffer
- **datatype** – [in] Data buffer element datatype
- **op** – [in] Reduction operator type
- **root** – [in] Rank where result data array will be stored
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

#### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclBcast**(void \*buff, size\_t count, *ncclDataType\_t* datatype, int root, *ncclComm\_t* comm, hipStream\_t stream)

(Deprecated) Broadcast (in-place)

Copies *count* values from *root* to all other devices. *root* is the rank (not the CUDA device) where data resides before the operation is started. This operation is implicitly in-place.

#### Parameters

- **buff** – [inout] Input array on *root* to be copied to other ranks. Output array for all ranks.
- **count** – [in] Number of elements in data buffer
- **datatype** – [in] Data buffer element datatype
- **root** – [in] Rank owning buffer to be copied to others
- **comm** – [in] Communicator group object to execute on

- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclBroadcast**(const void \*sendbuff, void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, int root, *ncclComm\_t* comm, hipStream\_t stream)

Broadcast.

Copies *count* values from *sendbuff* on *root* to *recvbuff* on all devices. *root\** is the rank (not the HIP device) where data resides before the operation is started. *sendbuff\** may be NULL on ranks other than *root*. In-place operation will happen if *sendbuff* == *recvbuff*.

**Parameters**

- **sendbuff** – [in] Data array to copy (if *root*). May be NULL for other ranks
- **recvbuff** – [inout] Data array to store received array
- **count** – [in] Number of elements in data buffer
- **datatype** – [in] Data buffer element datatype
- **root** – [in] Rank of broadcast root
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclAllReduce**(const void \*sendbuff, void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, *ncclRedOp\_t* op, *ncclComm\_t* comm, hipStream\_t stream)

All-Reduce.

Reduces data arrays of length *count* in *sendbuff* using *op* operation, and leaves identical copies of result on each *recvbuff*. In-place operation will happen if *sendbuff* == *recvbuff*.

**Parameters**

- **sendbuff** – [in] Input data array to reduce
- **recvbuff** – [inout] Data array to store reduced result array
- **count** – [in] Number of elements in data buffer
- **datatype** – [in] Data buffer element datatype
- **op** – [in] Reduction operator
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclReduceScatter**(const void \*sendbuff, void \*recvbuff, size\_t recvcount, *ncclDataType\_t* datatype, *ncclRedOp\_t* op, *ncclComm\_t* comm, hipStream\_t stream)

Reduce-Scatter.

Reduces data in *sendbuff* using *op* operation and leaves reduced result scattered over the devices so that *recvbuff* on rank *i* will contain the *i*-th block of the result. Assumes *sendcount* is equal to *n ranks*\**recvcount*, which means that *sendbuff* should have a size of at least *n ranks*\**recvcount* elements. In-place operations will happen if *recvbuff* == *sendbuff* + *rank* \* *recvcount*.

**Parameters**

- **sendbuff** – [in] Input data array to reduce
- **recvbuff** – [inout] Data array to store reduced result subarray
- **recvcount** – [in] Number of elements each rank receives
- **datatype** – [in] Data buffer element datatype
- **op** – [in] Reduction operator
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclAllGather**(const void \*sendbuff, void \*recvbuff, size\_t sendcount, *ncclDataType\_t* datatype, *ncclComm\_t* comm, hipStream\_t stream)

All-Gather.

Each device gathers *sendcount* values from other GPUs into *recvbuff*, receiving data from rank *i* at offset  $i * \text{sendcount}$ . Assumes *recvcount* is equal to  $n\text{ranks} * \text{sendcount}$ , which means that *recvbuff* should have a size of at least  $n\text{ranks} * \text{sendcount}$  elements. In-place operations will happen if  $\text{sendbuff} == \text{recvbuff} + \text{rank} * \text{sendcount}$ .

**Parameters**

- **sendbuff** – [in] Input data array to send
- **recvbuff** – [inout] Data array to store the gathered result
- **sendcount** – [in] Number of elements each rank sends
- **datatype** – [in] Data buffer element datatype
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclSend**(const void \*sendbuff, size\_t count, *ncclDataType\_t* datatype, int peer, *ncclComm\_t* comm, hipStream\_t stream)

Send.

Send data from *sendbuff* to rank *peer*. Rank *peer* needs to call *ncclRecv* with the same *datatype* and the same *count* as this rank. This operation is blocking for the GPU. If multiple *ncclSend* and *ncclRecv* operations need to progress concurrently to complete, they must be fused within a *ncclGroupStart* / *ncclGroupEnd* section.

**Parameters**

- **sendbuff** – [in] Data array to send
- **count** – [in] Number of elements to send
- **datatype** – [in] Data buffer element datatype
- **peer** – [in] Peer rank to send to
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclRecv**(void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, int peer, *ncclComm\_t* comm, hipStream\_t stream)

Receive.

Receive data from rank *peer* into *recvbuff*. Rank *peer* needs to call `ncclSend` with the same datatype and the same count as this rank. This operation is blocking for the GPU. If multiple `ncclSend` and `ncclRecv` operations need to progress concurrently to complete, they must be fused within a `ncclGroupStart/ncclGroupEnd` section.

**Parameters**

- **recvbuff** – [inout] Data array to receive
- **count** – [in] Number of elements to receive
- **datatype** – [in] Data buffer element datatype
- **peer** – [in] Peer rank to send to
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclGather**(const void \*sendbuff, void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, int root, *ncclComm\_t* comm, hipStream\_t stream)

Gather.

Each rank sends count elements from sendbuff to the root rank. On the root rank, data from rank *i* is placed at `recvbuff + i*count`. On non-root ranks, `recvbuff` is not used. `root` is the rank where data will be gathered.

In-place operations will happen if `sendbuff == recvbuff + root * count`.

**Parameters**

- **sendbuff** – [in] Data array to send
- **recvbuff** – [inout] Data array to recv
- **count** – [in] Number of elements
- **datatype** – [in] Data buffer element datatype
- **root** – [in] Rank of gather root
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclScatter**(const void \*sendbuff, void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, int root, *ncclComm\_t* comm, hipStream\_t stream)

Scatter.

On the root rank, count elements from `sendbuff+i*count` are sent to rank *i*. On non-root ranks, `sendbuff` is not used. Each rank receives count elements into `recvbuff`. `root` is the rank that will distribute the data.

In-place operations will happen if `recvbuff == sendbuff + root * count`.

**Parameters**

- **sendbuff** – [in] Data array to send
- **recvbuff** – [inout] Data array to recv
- **count** – [in] Number of elements
- **datatype** – [in] Data buffer element datatype
- **root** – [in] Rank of scatter root
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

#### Returns

Result code. See *Result Codes* for more details.

`ncclResult_t ncclAllToAll` (const void \*sendbuff, void \*recvbuff, size\_t count, `ncclDataType_t` datatype, `ncclComm_t` comm, `hipStream_t` stream) \_\_attribute\_\_((deprecated("ncclAllToAll is replaced with ncclAlltoAll and will be removed in the future")))

All-To-All.

Device (i) send (j)th block of data to device (j) and be placed as (i)th block. Each block for sending/receiving has *count* elements, which means that *recvbuff* and *sendbuff* should have a size of *n ranks*\**count* elements. In-place operation is NOT supported. It is the user's responsibility to ensure that *sendbuff* and *recvbuff* are distinct.

*Deprecated:*

`ncclAllToAll` is replaced with `ncclAlltoAll` and will be removed in the future.

#### Parameters

- **sendbuff** – [in] Data array to send (contains blocks for each other rank)
- **recvbuff** – [inout] Data array to receive (contains blocks from each other rank)
- **count** – [in] Number of elements to send between each pair of ranks
- **datatype** – [in] Data buffer element datatype
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

#### Returns

Result code. See *Result Codes* for more details.

## 9.3 Group semantics

When managing multiple GPUs from a single thread, and since NCCL collective calls may perform inter-CPU synchronization, we need to “group” calls for different ranks/devices into a single call.

Grouping NCCL calls as being part of the same collective operation is done using `ncclGroupStart` and `ncclGroupEnd`. `ncclGroupStart` will enqueue all collective calls until the `ncclGroupEnd` call, which will wait for all calls to be complete. Note that for collective communication, `ncclGroupEnd` only guarantees that the operations are enqueued on the streams, not that the operation is effectively done.

Both collective communication and `ncclCommInitRank` can be used in conjunction of `ncclGroupStart/ncclGroupEnd`.

*ncclResult\_t* **ncclGroupStart**()

Group Start.

Start a group call. All calls to RCCL until `ncclGroupEnd` will be fused into a single RCCL operation. Nothing will be started on the HIP stream until `ncclGroupEnd`.

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclGroupEnd**()

Group End.

End a group call. Start a fused RCCL operation consisting of all calls since `ncclGroupStart`. Operations on the HIP stream depending on the RCCL operations need to be called after `ncclGroupEnd`.

**Returns**

Result code. See *Result Codes* for more details.

## 9.4 Library functions

*ncclResult\_t* **ncclGetVersion**(int \*version)

Return the `RCCL_VERSION_CODE` of RCCL in the supplied integer.

This integer is coded with the MAJOR, MINOR and PATCH level of RCCL.

**Parameters**

**version** – [out] Pointer to where version will be stored

**Returns**

Result code. See *Result Codes* for more details.

const char \***ncclGetErrorString**(*ncclResult\_t* result)

Returns a string for each result code.

Returns a human-readable string describing the given result code.

**Parameters**

**result** – [in] Result code to get description for

**Returns**

String containing description of result code.

## 9.5 Types

There are few data structures that are internal to the library. The pointer types to these structures are given below. The user would need to use these types to create handles and pass them between different library functions.

typedef struct ncclComm \***ncclComm\_t**

Opaque handle to communicator.

A communicator contains information required to facilitate collective communications calls

struct **ncclUniqueId**

Opaque unique id used to initialize communicators.

The *ncclUniqueId* must be passed to all participating ranks

## 9.6 Enumerations

This section provides all the enumerations used.

enum **ncclResult\_t**

Result type.

Return codes aside from `ncclSuccess` indicate that a call has failed

*Values:*

enumerator **ncclSuccess**

No error

enumerator **ncclUnhandledCudaError**

Unhandled HIP error

enumerator **ncclSystemError**

Unhandled system error

enumerator **ncclInternalError**

Internal Error - Please report to RCCL developers

enumerator **ncclInvalidArgument**

Invalid argument

enumerator **ncclInvalidUsage**

Invalid usage

enumerator **ncclRemoteError**

Remote process exited or there was a network error

enumerator **ncclInProgress**

RCCL operation in progress

enumerator **ncclNumResults**

Number of result types

enum **ncclRedOp\_t**

Reduction operation selector.

Enumeration used to specify the various reduction operations `ncclNumOps` is the number of built-in `ncclRedOp_t` values and serves as the least possible value for dynamic `ncclRedOp_t` values constructed by `ncclRedOpCreate` functions.

`ncclMaxRedOp` is the largest valid value for `ncclRedOp_t` and is defined to be the largest signed value (since compilers are permitted to use signed enums) that won't grow `sizeof(ncclRedOp_t)` when compared to previous RCCL versions to maintain ABI compatibility.

*Values:*

enumerator **ncclSum**

Sum

enumerator **ncclProd**

Product

enumerator **ncclMax**

Max

enumerator **ncclMin**

Min

enumerator **ncclAvg**

Average

enumerator **ncclNumOps**

Number of built-in reduction ops

enumerator **ncclMaxRedOp**

Largest value for ncclRedOp\_t

enum **ncclDataType\_t**

Data types.

Enumeration of the various supported datatype

*Values:*

enumerator **ncclInt8**

enumerator **ncclChar**

enumerator **ncclUInt8**

enumerator **ncclInt32**

enumerator **ncclInt**

enumerator **ncclUInt32**

enumerator **ncclInt64**

enumerator **ncclUInt64**

enumerator **ncclFloat16**

enumerator **ncclHalf**

enumerator **ncclFloat32**

enumerator **ncclFloat**

enumerator **ncclFloat64**

enumerator **ncclDouble**

enumerator **ncclBfloat16**

enumerator **ncclFloat8e4m3**

enumerator **ncclFloat8e5m2**

enumerator **ncclNumTypes**



## API LIBRARY

RCCL (pronounced “Rickle”) is a stand-alone library of standard collective communication routines for GPUs, implementing all-reduce, all-gather, reduce, broadcast, reduce-scatter, gather, scatter, and all-to-all. There is also initial support for direct GPU-to-GPU send and receive operations. It has been optimized to achieve high bandwidth on platforms using PCIe, xGMI as well as networking using InfiniBand Verbs or TCP/IP sockets. RCCL supports an arbitrary number of GPUs installed in a single node or multiple nodes, and can be used in either single- or multi-process (e.g., MPI) applications.

The collective operations are implemented using ring and tree algorithms and have been optimized for throughput and latency. For best performance, small operations can be either batched into larger operations or aggregated through the API.

### 10.1 Operations

struct **ncclConfig\_t**

Communicator configuration.

Users can assign value to attributes to specify the behavior of a communicator

#### Public Members

size\_t **size**

Should not be touched

unsigned int **magic**

Should not be touched

unsigned int **version**

Should not be touched

int **blocking**

Whether or not calls should block or not

int **cgaClusterSize**

Cooperative group array cluster size

int **minCTAs**  
Minimum number of cooperative thread arrays (blocks)

int **maxCTAs**  
Maximum number of cooperative thread arrays (blocks)

const char \***netName**  
Force NCCL to use a specific network

int **splitShare**  
Allow communicators to share resources

int **trafficClass**  
Traffic class

const char \***commName**  
Name of the communicator

int **collnetEnable**  
Check for collnet enablement

int **CTAPolicy**  
CTA Policy

int **shrinkShare**  
Shrink size

int **nvlsCTAs**  
Number of NVLS cooperative thread arrays (blocks)

int **nChannelsPerNetPeer**  
Number of channels per NET peer

int **nvlinkCentricSched**  
nvlinkCentricSched

int **graphUsageMode**  
CUDA Graph usage mode hint

int **numRmaCtx**  
Number of RMA contexts

struct **ncclSimInfo\_t**

**Public Members**

size\_t **size**

unsigned int **magic**

unsigned int **version**

float **estimatedTime**

struct **ncclUniqueId**

Opaque unique id used to initialize communicators.

The *ncclUniqueId* must be passed to all participating ranks

**Public Members**

char **internal**[NCCL\_UNIQUE\_ID\_BYTES]

Opaque array>

struct **ncclWaitSignalDesc\_t**

**Public Members**

int **opCnt**

int **peer**

int **sigIdx**

int **ctx**

file **nccl.h.in**

*#include <hip/hip\_runtime.h>#include <hip/hip\_fp16.h>#include <limits.h>#include <stdint.h>*

**Defines**

**NCCL\_H\_**

**NCCL\_MAJOR**

**NCCL\_MINOR**

NCCL\_PATCH

NCCL\_SUFFIX

NCCL\_VERSION\_CODE

NCCL\_VERSION(X, Y, Z)

RCCL\_BFLOAT16

RCCL\_FLOAT8

RCCL\_GATHER\_SCATTER

RCCL\_ALLTOALLV

RCCL\_ALLREDUCE\_WITH\_BIAS

RCCL\_SUSPEND\_RESUME

NCCL\_COMM\_NULL

NCCL\_UNIQUE\_ID\_BYTES

NCCL\_CONFIG\_UNDEF\_INT

NCCL\_CONFIG\_UNDEF\_PTR

NCCL\_SPLIT\_NOCOLOR

NCCL\_UNDEF\_FLOAT

NCCL\_API\_MAGIC

NCCL\_WIN\_DEFAULT

NCCL\_WIN\_COLL\_SYMMETRIC

NCCL\_WIN\_STRICT\_ORDERING

NCCL\_WIN\_REQUIRED\_ALIGNMENT

**NCCL\_CTA\_POLICY\_DEFAULT**

**NCCL\_CTA\_POLICY\_EFFICIENCY**

**NCCL\_CTA\_POLICY\_ZERO**

**NCCL\_SHRINK\_DEFAULT**

**NCCL\_SHRINK\_ABORT**

**NCCL\_REVOKE\_DEFAULT**

**NCCL\_CONFIG\_INITIALIZER**

**NCCL\_SIM\_INFO\_INITIALIZER**

**NCCL\_SUSPEND\_MEM**

Communicator suspend flag: release dynamic GPU memory allocations.

## Typedefs

typedef struct ncclComm **\*ncclComm\_t**

Opaque handle to communicator.

A communicator contains information required to facilitate collective communications calls

typedef struct ncclWindow\_vidmem **\*ncclWindow\_t**

typedef int **mscclAlgoHandle\_t**

Opaque handle to MSCCL algorithm.

## Enums

enum **ncclResult\_t**

Result type.

Return codes aside from `ncclSuccess` indicate that a call has failed

*Values:*

enumerator **ncclSuccess**

No error

enumerator **ncclUnhandledCudaError**

Unhandled HIP error

enumerator **ncclSystemError**

Unhandled system error

enumerator **ncclInternalError**

Internal Error - Please report to RCCL developers

enumerator **ncclInvalidArgument**

Invalid argument

enumerator **ncclInvalidUsage**

Invalid usage

enumerator **ncclRemoteError**

Remote process exited or there was a network error

enumerator **ncclInProgress**

RCCL operation in progress

enumerator **ncclNumResults**

Number of result types

enum **ncclCommMemStat\_t**

Communicator memory statistic selector.

Identifier passed to *ncclCommMemStats* to choose which memory counter to read.

*Values:*

enumerator **ncclStatGpuMemSuspend**

Allocated GPU memory that can be suspended (bytes)

enumerator **ncclStatGpuMemSuspended**

GPU memory suspended? (0=active, 1=suspended)

enumerator **ncclStatGpuMemPersist**

Allocated GPU memory that cannot be suspended (bytes)

enumerator **ncclStatGpuMemTotal**

Total allocated GPU memory tracked by NCCL (bytes)

enum **ncclRedOp\_dummy\_t**

Dummy reduction enumeration.

Dummy reduction enumeration used to determine value for `ncclMaxRedOp`

*Values:*

enumerator **ncclNumOps\_dummy**

**enum ncclRedOp\_t**

Reduction operation selector.

Enumeration used to specify the various reduction operations `ncclNumOps` is the number of built-in `ncclRedOp_t` values and serves as the least possible value for dynamic `ncclRedOp_t` values constructed by `ncclRedOpCreate` functions.

`ncclMaxRedOp` is the largest valid value for `ncclRedOp_t` and is defined to be the largest signed value (since compilers are permitted to use signed enums) that won't grow `sizeof(ncclRedOp_t)` when compared to previous RCCL versions to maintain ABI compatibility.

*Values:*

enumerator **ncclSum**

Sum

enumerator **ncclProd**

Product

enumerator **ncclMax**

Max

enumerator **ncclMin**

Min

enumerator **ncclAvg**

Average

enumerator **ncclNumOps**

Number of built-in reduction ops

enumerator **ncclMaxRedOp**

Largest value for `ncclRedOp_t`

**enum ncclDataType\_t**

Data types.

Enumeration of the various supported datatype

*Values:*

enumerator **ncclInt8**

enumerator **ncclChar**

enumerator **ncclUInt8**

enumerator **ncclInt32**

enumerator **ncclInt**

enumerator **ncclUint32**

enumerator **ncclInt64**

enumerator **ncclUint64**

enumerator **ncclFloat16**

enumerator **ncclHalf**

enumerator **ncclFloat32**

enumerator **ncclFloat**

enumerator **ncclFloat64**

enumerator **ncclDouble**

enumerator **ncclBfloat16**

enumerator **ncclFloat8e4m3**

enumerator **ncclFloat8e5m2**

enumerator **ncclNumTypes**

enum **ncclScalarResidence\_t**

Location and dereferencing logic for scalar arguments.

Enumeration specifying memory location of the scalar argument. Based on where the value is stored, the argument will be dereferenced either while the collective is running (if in device memory), or before the `ncclRedOpCreate()` function returns (if in host memory).

*Values:*

enumerator **ncclScalarDevice**

Scalar is in device-visible memory

enumerator **ncclScalarHostImmediate**

Scalar is in host-visible memory

## Functions

*ncclResult\_t* **ncclMemAlloc**(void \*\*ptr, size\_t size)

*ncclResult\_t* **pnccclMemAlloc**(void \*\*ptr, size\_t size)

*ncclResult\_t* **ncclMemFree**(void \*ptr)

*ncclResult\_t* **pnccclMemFree**(void \*ptr)

*ncclResult\_t* **ncclGetVersion**(int \*version)

Return the RCCL\_VERSION\_CODE of RCCL in the supplied integer.

This integer is coded with the MAJOR, MINOR and PATCH level of RCCL.

### Parameters

**version** – [out] Pointer to where version will be stored

### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclGetUniqueId**(*ncclUniqueId* \*uniqueId)

Generates an ID for `ncclCommInitRank`.

Generates an ID to be used in `ncclCommInitRank`. `ncclGetUniqueId` should be called once by a single rank and the ID should be distributed to all ranks in the communicator before using it as a parameter for `ncclCommInitRank`.

### Parameters

**uniqueId** – [out] Pointer to where uniqueId will be stored

### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommInitRankConfig**(*ncclComm\_t* \*comm, int nranks, *ncclUniqueId* commId, int rank, *ncclConfig\_t* \*config)

Create a new communicator with config.

Create a new communicator (multi thread/process version) with a configuration set by users. See *Communicator Configuration* for more details. Each rank is associated to a CUDA device, which has to be set before calling `ncclCommInitRank`.

### Parameters

- **comm** – [out] Pointer to created communicator
- **nranks** – [in] Total number of ranks participating in this communicator
- **commId** – [in] UniqueId required for initialization
- **rank** – [in] Current rank to create communicator for. [0 to nranks-1]
- **config** – [in] Pointer to communicator configuration

### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommInitRank**(*ncclComm\_t* \*comm, int nranks, *ncclUniqueId* commId, int rank)

Creates a new communicator (multi thread/process version).

Rank must be between 0 and n ranks-1 and unique within a communicator clique. Each rank is associated to a CUDA device, which has to be set before calling `ncclCommInitRank`. `ncclCommInitRank` implicitly synchronizes with other ranks, so it must be called by different threads/processes or use `ncclGroupStart/ncclGroupEnd`.

#### Parameters

- **comm** – [out] Pointer to created communicator
- **n ranks** – [in] Total number of ranks participating in this communicator
- **commId** – [in] UniqueId required for initialization
- **rank** – [in] Current rank to create communicator for

#### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommInitAll**(*ncclComm\_t* \*comm, int ndev, const int \*devlist)

Creates a clique of communicators (single process version).

This is a convenience function to create a single-process communicator clique. Returns an array of ndev newly initialized communicators in comm. comm should be pre-allocated with size at least ndev\*sizeof(ncclComm\_t). If devlist is NULL, the first ndev HIP devices are used. Order of devlist defines user-order of processors within the communicator.

#### Parameters

- **comm** – [out] Pointer to array of created communicators
- **ndev** – [in] Total number of ranks participating in this communicator
- **devlist** – [in] Array of GPU device indices to create for

#### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommFinalize**(*ncclComm\_t* comm)

Finalize a communicator.

`ncclCommFinalize` flushes all issued communications and marks communicator state as `ncclInProgress`. The state will change to `ncclSuccess` when the communicator is globally quiescent and related resources are freed; then, calling `ncclCommDestroy` can locally free the rest of the resources (e.g. communicator itself) without blocking.

#### Parameters

**comm** – [in] Communicator to finalize

#### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommDestroy**(*ncclComm\_t* comm)

Frees local resources associated with communicator object.

Destroy all local resources associated with the passed in communicator object

#### Parameters

**comm** – [in] Communicator to destroy

#### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommAbort**(*ncclComm\_t* comm)

Abort any in-progress calls and destroy the communicator object.

Frees resources associated with communicator object and aborts any operations that might still be running on the device.

**Parameters**

**comm** – [in] Communicator to abort and destroy

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommRevoke**(*ncclComm\_t* comm, int revokeFlags)

Revoke a communicator without destroying it.

Aborts in-flight collectives by raising the comm’s abort flag, stops the proxy service, and rejects subsequently enqueued collectives with `ncclInvalidUsage`. The abort flag is cleared once the asynchronous revoke job completes, so the communicator remains valid as a parent for `ncclCommSplit` / `ncclCommShrink` / `ncclCommGrow`, and may be torn down via `ncclCommDestroy` / `ncclCommAbort`. Because in-flight collectives are aborted (not drained), revoke can recover from peers that have failed mid-collective; output buffers of an aborted collective contain undefined data. `ncclCommFinalize` on a revoked communicator is invalid and returns `ncclInvalidUsage`. Pass `NCCL_REVOKE_DEFAULT` for `revokeFlags`; other values are rejected with `ncclInvalidArgument`.

**Parameters**

- **comm** – [in] Communicator to revoke
- **revokeFlags** – [in] Reserved; must be `NCCL_REVOKE_DEFAULT`

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommSplit**(*ncclComm\_t* comm, int color, int key, *ncclComm\_t* \*newcomm, *ncclConfig\_t* \*config)

Create one or more communicators from an existing one.

Creates one or more communicators from an existing one. Ranks with the same color will end up in the same communicator. Within the new communicator, key will be used to order ranks. `NCCL_SPLIT_NOCOLOR` as color will indicate the rank will not be part of any group and will therefore return a NULL communicator. If config is NULL, the new communicator will inherit the original communicator’s configuration

**Parameters**

- **comm** – [in] Original communicator object for this rank
- **color** – [in] Color to assign this rank
- **key** – [in] Key used to order ranks within the same new communicator
- **newcomm** – [out] Pointer to new communicator
- **config** – [in] Config file for new communicator. May be NULL to inherit from comm

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommShrink**(*ncclComm\_t* comm, int \*excludeRanksList, int excludeRanksCount, *ncclComm\_t* \*newcomm, *ncclConfig\_t* \*config, int shrinkFlags)

Shrink existing communicator.

Ranks in `excludeRanksList` will be removed from the existing communicator. Within the new communicator, ranks will be re-ordered to fill the gap of removed ones. If config is NULL, the new communicator will

inherit the original communicator's configuration. The flag enables NCCL to adapt to various states of the parent communicator, see `NCCL_SHRINK` flags.

#### Parameters

- **comm** – [in] Original communicator object for this rank
- **excludeRanksList** – [in] List of ranks to be excluded
- **excludeRanksCount** – [in] Number of ranks to be excluded
- **newcomm** – [out] Pointer to new communicator
- **config** – [in] Config file for new communicator. May be NULL to inherit from comm
- **shrinkFlags** – [in] Flag to adapt to various states of the parent communicator (see `NCCL_SHRINK` flags)

#### Returns

Result code. See *Result Codes* for more details.

```
ncclResult_t pnccCommShrink(ncclComm_t comm, int *excludeRanksList, int excludeRanksCount,  
                             ncclComm_t *newcomm, ncclConfig_t *config, int shrinkFlags)
```

```
ncclResult_t ncclCommGetUniqueId(ncclComm_t comm, ncclUniqueId *uniqueId)
```

```
ncclResult_t pnccCommGetUniqueId(ncclComm_t comm, ncclUniqueId *uniqueId)
```

```
ncclResult_t ncclCommGrow(ncclComm_t comm, int nRanks, const ncclUniqueId *uniqueId, int rank,  
                           ncclComm_t *newcomm, ncclConfig_t *config)
```

```
ncclResult_t pnccCommGrow(ncclComm_t comm, int nRanks, const ncclUniqueId *uniqueId, int rank,  
                           ncclComm_t *newcomm, ncclConfig_t *config)
```

```
ncclResult_t ncclCommInitRankScalable(ncclComm_t *newcomm, int nRanks, int myrank, int nId,  
                                       ncclUniqueId *commIds, ncclConfig_t *config)
```

Creates a new communicator (multi thread/process version), similar to `ncclCommInitRankConfig`.

Allows to use more than one *ncclUniqueId* (up to one per rank), indicated by `nId`, to accelerate the init operation. The number of *ncclUniqueIds* and their order must be the same for every rank.

#### Parameters

- **newcomm** – [out] Pointer to new communicator
- **nRanks** – [in] Total number of ranks participating in this communicator
- **myrank** – [in] Current rank
- **nId** – [in] Number of unique IDs
- **commIds** – [in] List of unique IDs
- **config** – [in] Config file for new communicator. May be NULL to inherit from comm

#### Returns

Result code. See *Result Codes* for more details.

```
const char *ncclGetErrorString(ncclResult_t result)
```

Returns a string for each result code.

Returns a human-readable string describing the given result code.

#### Parameters

**result** – [in] Result code to get description for

**Returns**

String containing description of result code.

```
const char *ncclGetLastError(ncclComm_t comm)
```

\_\_attribute\_\_((deprecated("ncclResetDebugInit is not supported as part of the NCCL API and will be removed in the future")))

```
ncclResult_t ncclCommGetAsyncError(ncclComm_t comm, ncclResult_t *asyncError)
```

Checks whether the comm has encountered any asynchronous errors.

Query whether the provided communicator has encountered any asynchronous errors

**Parameters**

- **comm** – [in] Communicator to query
- **asyncError** – [out] Pointer to where result code will be stored

**Returns**

Result code. See *Result Codes* for more details.

```
ncclResult_t ncclCommCount(const ncclComm_t comm, int *count)
```

Gets the number of ranks in the communicator clique.

Returns the number of ranks in the communicator clique (as set during initialization)

**Parameters**

- **comm** – [in] Communicator to query
- **count** – [out] Pointer to where number of ranks will be stored

**Returns**

Result code. See *Result Codes* for more details.

```
ncclResult_t ncclCommCuDevice(const ncclComm_t comm, int *device)
```

Get the ROCm device index associated with a communicator.

Returns the ROCm device number associated with the provided communicator.

**Parameters**

- **comm** – [in] Communicator to query
- **device** – [out] Pointer to where the associated ROCm device index will be stored

**Returns**

Result code. See *Result Codes* for more details.

```
ncclResult_t ncclCommUserRank(const ncclComm_t comm, int *rank)
```

Get the rank associated with a communicator.

Returns the user-ordered “rank” associated with the provided communicator.

**Parameters**

- **comm** – [in] Communicator to query
- **rank** – [out] Pointer to where the associated rank will be stored

**Returns**

Result code. See *Result Codes* for more details.

```
ncclResult_t ncclCommRegister(const ncclComm_t comm, void *buff, size_t size, void **handle)
```

*ncclResult\_t* **ncclCommDeregister**(const *ncclComm\_t* comm, void \*handle)

*ncclResult\_t* **ncclCommSuspend**(*ncclComm\_t* comm, int flags)

Suspend communicator operations to free resources.

Releases the resources selected by *flags*. The communicator cannot be used until *ncclCommResume* is called.

**Parameters**

- **comm** – [in] Communicator to suspend
- **flags** – [in] Bitmask of NCCL\_SUSPEND\_\* flags (e.g. NCCL\_SUSPEND\_MEM)

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommResume**(*ncclComm\_t* comm)

Resume a previously suspended communicator.

Reacquires every resource that was released by the matching *ncclCommSuspend* call.

**Parameters**

- **comm** – [in] Communicator to resume

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommMemStats**(*ncclComm\_t* comm, *ncclCommMemStat\_t* stat, uint64\_t \*value)

Query communicator memory statistics.

**Parameters**

- **comm** – [in] Communicator to query
- **stat** – [in] One of *ncclCommMemStat\_t* values
- **value** – [out] Pointer to receive the memory statistic value

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclCommWindowRegister**(*ncclComm\_t* comm, void \*buff, size\_t size, *ncclWindow\_t* \*win, int winFlags)

*ncclResult\_t* **ncclCommWindowDeregister**(*ncclComm\_t* comm, *ncclWindow\_t* win)

*ncclResult\_t* **ncclWinGetUserPtr**(*ncclComm\_t* comm, *ncclWindow\_t* win, void \*\*outUserPtr)

*ncclResult\_t* **pnccclWinGetUserPtr**(*ncclComm\_t* comm, *ncclWindow\_t* win, void \*\*outUserPtr)

*ncclResult\_t* **ncclRedOpCreatePreMulSum**(*ncclRedOp\_t* \*op, void \*scalar, *ncclDataType\_t* datatype, *ncclScalarResidence\_t* residence, *ncclComm\_t* comm)

Create a custom pre-multiplier reduction operator.

Creates a new reduction operator which pre-multiplies input values by a given scalar locally before reducing them with peer values via summation. For use only with collectives launched against *comm* and *datatype*. The *residence\** argument indicates how/when the memory pointed to by *scalar* will be dereferenced. Upon return, the newly created operator's handle is stored in *op*.

**Parameters**

- **op** – [out] Pointer to where newly created custom reduction operator is to be stored

- **scalar** – [in] Pointer to scalar value.
- **datatype** – [in] Scalar value datatype
- **residence** – [in] Memory type of the scalar value
- **comm** – [in] Communicator to associate with this custom reduction operator

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclRedOpDestroy**(*ncclRedOp\_t* op, *ncclComm\_t* comm)

Destroy custom reduction operator.

Destroys the reduction operator *op*. The operator must have been created by `ncclRedOpCreatePreMul` with the matching communicator *comm*. An operator may be destroyed as soon as the last RCCL function which is given that operator returns.

**Parameters**

- **op** – [in] Custom reduction operator is to be destroyed
- **comm** – [in] Communicator associated with this reduction operator

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclReduce**(const void \*sendbuff, void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, *ncclRedOp\_t* op, int root, *ncclComm\_t* comm, hipStream\_t stream)

Reduce.

Reduces data arrays of length *count* in *sendbuff* into *recvbuff* using *op* operation. *recvbuff\** may be NULL on all calls except for root device. *root\** is the rank (not the HIP device) where data will reside after the operation is complete. In-place operation will happen if *sendbuff* == *recvbuff*.

**Parameters**

- **sendbuff** – [in] Local device data buffer to be reduced
- **recvbuff** – [inout] Data buffer where result is stored (only for *root* rank). May be null for other ranks.
- **count** – [in] Number of elements in every send buffer
- **datatype** – [in] Data buffer element datatype
- **op** – [in] Reduction operator type
- **root** – [in] Rank where result data array will be stored
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclBcast**(void \*buff, size\_t count, *ncclDataType\_t* datatype, int root, *ncclComm\_t* comm, hipStream\_t stream)

(Deprecated) Broadcast (in-place)

Copies *count* values from *root* to all other devices. *root* is the rank (not the CUDA device) where data resides before the operation is started. This operation is implicitly in-place.

**Parameters**

- **buff** – [inout] Input array on *root* to be copied to other ranks. Output array for all ranks.
- **count** – [in] Number of elements in data buffer
- **datatype** – [in] Data buffer element datatype
- **root** – [in] Rank owning buffer to be copied to others
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclBroadcast**(const void \*sendbuff, void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, int root, *ncclComm\_t* comm, hipStream\_t stream)

Broadcast.

Copies *count* values from *sendbuff* on *root* to *recvbuff* on all devices. *root\** is the rank (not the HIP device) where data resides before the operation is started. *sendbuff\** may be NULL on ranks other than *root*. In-place operation will happen if *sendbuff* == *recvbuff*.

**Parameters**

- **sendbuff** – [in] Data array to copy (if *root*). May be NULL for other ranks
- **recvbuff** – [inout] Data array to store received array
- **count** – [in] Number of elements in data buffer
- **datatype** – [in] Data buffer element datatype
- **root** – [in] Rank of broadcast root
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclAllReduce**(const void \*sendbuff, void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, *ncclRedOp\_t* op, *ncclComm\_t* comm, hipStream\_t stream)

All-Reduce.

Reduces data arrays of length *count* in *sendbuff* using *op* operation, and leaves identical copies of result on each *recvbuff*. In-place operation will happen if *sendbuff* == *recvbuff*.

**Parameters**

- **sendbuff** – [in] Input data array to reduce
- **recvbuff** – [inout] Data array to store reduced result array
- **count** – [in] Number of elements in data buffer
- **datatype** – [in] Data buffer element datatype
- **op** – [in] Reduction operator
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclAllReduceWithBias**(const void \*sendbuff, void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, *ncclRedOp\_t* op, *ncclComm\_t* comm, hipStream\_t stream, const void \*acc)

All-Reduce-with-Bias.

Reduces data arrays of length *count* in *sendbuff* using *op* operation, and leaves identical copies of result on each *recvbuff*. In-place operation will happen if *sendbuff* == *recvbuff*.

#### Parameters

- **sendbuff** – [in] Input data array to reduce
- **recvbuff** – [inout] Data array to store reduced result array
- **count** – [in] Number of elements in data buffer
- **datatype** – [in] Data buffer element datatype
- **op** – [in] Reduction operator
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on
- **acc** – [in] Bias data array to reduce

#### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclReduceScatter**(const void \*sendbuff, void \*recvbuff, size\_t recvcount, *ncclDataType\_t* datatype, *ncclRedOp\_t* op, *ncclComm\_t* comm, hipStream\_t stream)

Reduce-Scatter.

Reduces data in *sendbuff* using *op* operation and leaves reduced result scattered over the devices so that *recvbuff* on rank *i* will contain the *i*-th block of the result. Assumes *sendcount* is equal to *n ranks*\**recvcount*, which means that *sendbuff* should have a size of at least *n ranks*\**recvcount* elements. In-place operations will happen if *recvbuff* == *sendbuff* + rank \* *recvcount*.

#### Parameters

- **sendbuff** – [in] Input data array to reduce
- **recvbuff** – [inout] Data array to store reduced result subarray
- **recvcount** – [in] Number of elements each rank receives
- **datatype** – [in] Data buffer element datatype
- **op** – [in] Reduction operator
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

#### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclAllGather**(const void \*sendbuff, void \*recvbuff, size\_t sendcount, *ncclDataType\_t* datatype, *ncclComm\_t* comm, hipStream\_t stream)

All-Gather.

Each device gathers *sendcount* values from other GPUs into *recvbuff*, receiving data from rank *i* at offset *i*\**sendcount*. Assumes *recvcount* is equal to *n ranks*\**sendcount*, which means that *recvbuff* should have a size of at least *n ranks*\**sendcount* elements. In-place operations will happen if *sendbuff* == *recvbuff* + rank \* *sendcount*.

**Parameters**

- **sendbuff** – [in] Input data array to send
- **recvbuff** – [inout] Data array to store the gathered result
- **sendcount** – [in] Number of elements each rank sends
- **datatype** – [in] Data buffer element datatype
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclAlltoAll**(const void \*sendbuff, void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, *ncclComm\_t* comm, hipStream\_t stream)

All-to-All.

Each device sends count values to all other devices and receives count values from all other devices. Data to send to destination rank *j* is taken from `sendbuff+j*count` and data received from source rank *i* is placed at `recvbuff+i*count`.

**Parameters**

- **sendbuff** – [in] Data array to send (contains blocks for each other rank)
- **recvbuff** – [inout] Data array to receive (contains blocks from each other rank)
- **count** – [in] Number of elements to send between each pair of ranks
- **datatype** – [in] Data buffer element datatype
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclAlltoAllv**(const void \*sendbuff, const size\_t sendcounts[], const size\_t sdispls[], void \*recvbuff, const size\_t recvcounts[], const size\_t rdispls[], *ncclDataType\_t* datatype, *ncclComm\_t* comm, hipStream\_t stream)

All-To-Allv.

Device (*i*) sends `sendcounts[j]` of data from offset `sdispls[j]` to device (*j*). At the same time, device (*i*) receives `recvcounts[j]` of data from device (*j*) to be placed at `rdispls[j]`. `sendcounts`, `sdispls`, `recvcounts` and `rdispls` are all measured in the units of `datatype`, not bytes. In-place operation will happen if `sendbuff == recvbuff`.

**Parameters**

- **sendbuff** – [in] Data array to send (contains blocks for each other rank)
- **sendcounts** – [in] Array containing number of elements to send to each participating rank
- **sdispls** – [in] Array of offsets into `sendbuff` for each participating rank
- **recvbuff** – [inout] Data array to receive (contains blocks from each other rank)
- **recvcounts** – [in] Array containing number of elements to receive from each participating rank
- **rdispls** – [in] Array of offsets into `recvbuff` for each participating rank

- **datatype** – [in] Data buffer element datatype
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

#### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclGather**(const void \*sendbuff, void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, int root, *ncclComm\_t* comm, hipStream\_t stream)

Gather.

Each rank sends count elements from sendbuff to the root rank. On the root rank, data from rank i is placed at recvbuff + i\*count. On non-root ranks, recvbuff is not used. root is the rank where data will be gathered.

In-place operations will happen if sendbuff == recvbuff + root \* count.

#### Parameters

- **sendbuff** – [in] Data array to send
- **recvbuff** – [inout] Data array to recv
- **count** – [in] Number of elements
- **datatype** – [in] Data buffer element datatype
- **root** – [in] Rank of gather root
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

#### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclScatter**(const void \*sendbuff, void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, int root, *ncclComm\_t* comm, hipStream\_t stream)

Scatter.

On the root rank, count elements from sendbuff+i\*count are sent to rank i. On non-root ranks, sendbuff is not used. Each rank receives count elements into recvbuff. root is the rank that will distribute the data.

In-place operations will happen if recvbuff == sendbuff + root \* count.

#### Parameters

- **sendbuff** – [in] Data array to send
- **recvbuff** – [inout] Data array to recv
- **count** – [in] Number of elements
- **datatype** – [in] Data buffer element datatype
- **root** – [in] Rank of scatter root
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

#### Returns

Result code. See *Result Codes* for more details.

```
ncclResult_t ncclAllToAll(const void *sendbuff, void *recvbuff, size_t count, ncclDataType_t datatype,
                          ncclComm_t comm, hipStream_t stream)
    __attribute__((deprecated("ncclAllToAll is replaced with ncclAlltoAll and will
                              be removed in the future")))
```

All-To-All.

Device (i) send (j)th block of data to device (j) and be placed as (i)th block. Each block for sending/receiving has *count* elements, which means that *recvbuff* and *sendbuff* should have a size of *nrank*\**count* elements. In-place operation is NOT supported. It is the user's responsibility to ensure that *sendbuff* and *recvbuff* are distinct.

*Deprecated:*

*ncclAllToAll* is replaced with *ncclAlltoAll* and will be removed in the future.

#### Parameters

- **sendbuff** – [in] Data array to send (contains blocks for each other rank)
- **recvbuff** – [inout] Data array to receive (contains blocks from each other rank)
- **count** – [in] Number of elements to send between each pair of ranks
- **datatype** – [in] Data buffer element datatype
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

#### Returns

Result code. See *Result Codes* for more details.

```
ncclResult_t ncclAllToAllv(const void *sendbuff, const size_t sendcounts[], const size_t sdispls[], void
                             *recvbuff, const size_t recvcounts[], const size_t rdispls[], ncclDataType_t
                             datatype, ncclComm_t comm, hipStream_t stream)
    __attribute__((deprecated("ncclAllToAllv is replaced with ncclAlltoAllv and
                              will be removed in the future")))
```

All-To-Allv.

Device (i) sends *sendcounts*[j] of data from offset *sdispls*[j] to device (j). At the same time, device (i) receives *recvcounts*[j] of data from device (j) to be placed at *rdispls*[j]. *sendcounts*, *sdispls*, *recvcounts* and *rdispls* are all measured in the units of *datatype*, not bytes. In-place operation will happen if *sendbuff* == *recvbuff*.

*Deprecated:*

*ncclAllToAllv* is replaced with *ncclAlltoAllv* and will be removed in the future.

#### Parameters

- **sendbuff** – [in] Data array to send (contains blocks for each other rank)
- **sendcounts** – [in] Array containing number of elements to send to each participating rank
- **sdispls** – [in] Array of offsets into *sendbuff* for each participating rank
- **recvbuff** – [inout] Data array to receive (contains blocks from each other rank)
- **recvcounts** – [in] Array containing number of elements to receive from each participating rank
- **rdispls** – [in] Array of offsets into *recvbuff* for each participating rank
- **datatype** – [in] Data buffer element datatype

- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclSend**(const void \*sendbuff, size\_t count, *ncclDataType\_t* datatype, int peer, *ncclComm\_t* comm, hipStream\_t stream)

Send.

Send data from *sendbuff* to rank *peer*. Rank *peer* needs to call `ncclRecv` with the same *datatype* and the same *count* as this rank. This operation is blocking for the GPU. If multiple `ncclSend` and `ncclRecv` operations need to progress concurrently to complete, they must be fused within a `ncclGroupStart` / `ncclGroupEnd` section.

**Parameters**

- **sendbuff** – [in] Data array to send
- **count** – [in] Number of elements to send
- **datatype** – [in] Data buffer element datatype
- **peer** – [in] Peer rank to send to
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclRecv**(void \*recvbuff, size\_t count, *ncclDataType\_t* datatype, int peer, *ncclComm\_t* comm, hipStream\_t stream)

Receive.

Receive data from rank *peer* into *recvbuff*. Rank *peer* needs to call `ncclSend` with the same *datatype* and the same *count* as this rank. This operation is blocking for the GPU. If multiple `ncclSend` and `ncclRecv` operations need to progress concurrently to complete, they must be fused within a `ncclGroupStart` / `ncclGroupEnd` section.

**Parameters**

- **recvbuff** – [inout] Data array to receive
- **count** – [in] Number of elements to receive
- **datatype** – [in] Data buffer element datatype
- **peer** – [in] Peer rank to send to
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclPutSignal**(const void \*localbuff, size\_t count, *ncclDataType\_t* datatype, int peer, *ncclWindow\_t* peerWin, size\_t peerWinOffset, int sigIdx, int ctx, unsigned int flags, *ncclComm\_t* comm, hipStream\_t stream)

*ncclResult\_t* **pnccIPutSignal**(const void \*localbuff, size\_t count, *ncclDataType\_t* datatype, int peer, *ncclWindow\_t* peerWin, size\_t peerWinOffset, int sigIdx, int ctx, unsigned int flags, *ncclComm\_t* comm, hipStream\_t stream)

*ncclResult\_t* **ncclSignal**(int peer, int sigIdx, int ctx, unsigned int flags, *ncclComm\_t* comm, hipStream\_t stream)

*ncclResult\_t* **pnccISignal**(int peer, int sigIdx, int ctx, unsigned int flags, *ncclComm\_t* comm, hipStream\_t stream)

*ncclResult\_t* **ncclWaitSignal**(int nDesc, *ncclWaitSignalDesc\_t* \*signalDescs, *ncclComm\_t* comm, hipStream\_t stream)

*ncclResult\_t* **pnccIWaitSignal**(int nDesc, *ncclWaitSignalDesc\_t* \*signalDescs, *ncclComm\_t* comm, hipStream\_t stream)

*ncclResult\_t* **mscclLoadAlgo**(const char \*mscclAlgoFilePath, *mscclAlgoHandle\_t* \*mscclAlgoHandle, int rank)

MSCCL Load Algorithm.

*Deprecated:*

This function has been removed from the public API.

Load MSCCL algorithm file specified in *mscclAlgoFilePath* and return its handle via *mscclAlgoHandle*. This API is expected to be called by MSCCL scheduler instead of end users.

#### Parameters

- **mscclAlgoFilePath** – [in] Path to MSCCL algorithm file
- **mscclAlgoHandle** – [out] Returned handle to MSCCL algorithm
- **rank** – [in] Current rank

#### Returns

Result code. See *Result Codes* for more details.

*ncclResult\_t* **mscclRunAlgo**(const void \*sendBuff, const size\_t sendCounts[], const size\_t sDisPls[], void \*recvBuff, const size\_t recvCounts[], const size\_t rDisPls[], size\_t count, *ncclDataType\_t* dataType, int root, int peer, *ncclRedOp\_t* op, *mscclAlgoHandle\_t* mscclAlgoHandle, *ncclComm\_t* comm, hipStream\_t stream)

MSCCL Run Algorithm.

*Deprecated:*

This function has been removed from the public API.

Run MSCCL algorithm specified by *mscclAlgoHandle*. The parameter list merges all possible parameters required by different operations as this is a general-purposed API. This API is expected to be called by MSCCL scheduler instead of end users.

#### Parameters

- **sendBuff** – [in] Data array to send
- **sendCounts** – [in] Array containing number of elements to send to each participating rank
- **sDisPls** – [in] Array of offsets into *sendbuff* for each participating rank

- **recvBuff** – [inout] Data array to receive
- **recvCounts** – [in] Array containing number of elements to receive from each participating rank
- **rDisPls** – [in] Array of offsets into *recvbuff* for each participating rank
- **count** – [in] Number of elements
- **dataType** – [in] Data buffer element datatype
- **root** – [in] Root rank index
- **peer** – [in] Peer rank index
- **op** – [in] Reduction operator
- **mscclAlgoHandle** – [in] Handle to MSCCL algorithm
- **comm** – [in] Communicator group object to execute on
- **stream** – [in] HIP stream to execute collective on

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **mscclUnloadAlgo**(*mscclAlgoHandle\_t* mscclAlgoHandle)

MSCCL Unload Algorithm.

*Deprecated:*

This function has been removed from the public API.

Unload MSCCL algorithm previous loaded using its handle. This API is expected to be called by MSCCL scheduler instead of end users.

**Parameters**

**mscclAlgoHandle** – [in] Handle to MSCCL algorithm to unload

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclGroupStart**()

Group Start.

Start a group call. All calls to RCCL until `ncclGroupEnd` will be fused into a single RCCL operation. Nothing will be started on the HIP stream until `ncclGroupEnd`.

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclGroupEnd**()

Group End.

End a group call. Start a fused RCCL operation consisting of all calls since `ncclGroupStart`. Operations on the HIP stream depending on the RCCL operations need to be called after `ncclGroupEnd`.

**Returns**

Result code. See *Result Codes* for more details.

*ncclResult\_t* **ncclGroupSimulateEnd**(*ncclSimInfo\_t* \*simInfo)

*ncclResult\_t* **pnccclGroupSimulateEnd**(*ncclSimInfo\_t* \*simInfo)

*group* **rccl\_result\_code**

The various result codes that RCCL API calls may return

**Enums**

enum **ncclResult\_t**

Result type.

Return codes aside from `ncclSuccess` indicate that a call has failed

*Values:*

enumerator **ncclSuccess**

No error

enumerator **ncclUnhandledCudaError**

Unhandled HIP error

enumerator **ncclSystemError**

Unhandled system error

enumerator **ncclInternalError**

Internal Error - Please report to RCCL developers

enumerator **ncclInvalidArgument**

Invalid argument

enumerator **ncclInvalidUsage**

Invalid usage

enumerator **ncclRemoteError**

Remote process exited or there was a network error

enumerator **ncclInProgress**

RCCL operation in progress

enumerator **ncclNumResults**

Number of result types

*group* **rccl\_config\_type**

Structure that allows for customizing Communicator behavior via `ncclCommInitRankConfig`

## Defines

**NCCL\_CONFIG\_INITIALIZER**

*group* **rccl\_api\_version**

API call that returns RCCL version

*group* **rccl\_api\_communicator**

API calls that operate on communicators. Communicators objects are used to launch collective communication operations. Unique ranks between 0 and N-1 must be assigned to each HIP device participating in the same Communicator. Using the same HIP device for multiple ranks of the same Communicator is not supported at this time.

*group* **rccl\_api\_errcheck**

API calls that check for errors

*group* **rccl\_api\_comminfo**

API calls that query communicator information

*group* **rccl\_api\_suspend\_resume**

APIs to suspend and resume communicator resources, and to query per-communicator memory statistics.

## Defines

**NCCL\_SUSPEND\_MEM**

Communicator suspend flag: release dynamic GPU memory allocations.

## Enums

enum **ncclCommMemStat\_t**

Communicator memory statistic selector.

Identifier passed to *ncclCommMemStats* to choose which memory counter to read.

*Values:*

enumerator **ncclStatGpuMemSuspend**

Allocated GPU memory that can be suspended (bytes)

enumerator **ncclStatGpuMemSuspended**

GPU memory suspended? (0=active, 1=suspended)

enumerator **ncclStatGpuMemPersist**

Allocated GPU memory that cannot be suspended (bytes)

enumerator **ncclStatGpuMemTotal**

Total allocated GPU memory tracked by NCCL (bytes)

group **rccl\_api\_enumerations**

Enumerations used by collective communication calls

## Enums

enum **ncclRedOp\_dummy\_t**

Dummy reduction enumeration.

Dummy reduction enumeration used to determine value for `ncclMaxRedOp`

*Values:*

enumerator **ncclNumOps\_dummy**

enum **ncclRedOp\_t**

Reduction operation selector.

Enumeration used to specify the various reduction operations `ncclNumOps` is the number of built-in `ncclRedOp_t` values and serves as the least possible value for dynamic `ncclRedOp_t` values constructed by `ncclRedOpCreate` functions.

`ncclMaxRedOp` is the largest valid value for `ncclRedOp_t` and is defined to be the largest signed value (since compilers are permitted to use signed enums) that won't grow `sizeof(ncclRedOp_t)` when compared to previous RCCL versions to maintain ABI compatibility.

*Values:*

enumerator **ncclSum**

Sum

enumerator **ncclProd**

Product

enumerator **ncclMax**

Max

enumerator **ncclMin**

Min

enumerator **ncclAvg**

Average

enumerator **ncclNumOps**

Number of built-in reduction ops

enumerator **ncclMaxRedOp**

Largest value for ncclRedOp\_t

enum **ncclDataType\_t**

Data types.

Enumeration of the various supported datatype

*Values:*

enumerator **ncclInt8**

enumerator **ncclChar**

enumerator **ncclUInt8**

enumerator **ncclInt32**

enumerator **ncclInt**

enumerator **ncclUInt32**

enumerator **ncclInt64**

enumerator **ncclUInt64**

enumerator **ncclFloat16**

enumerator **ncclHalf**

enumerator **ncclFloat32**

enumerator **ncclFloat**

enumerator **ncclFloat64**

enumerator **ncclDouble**

enumerator **ncclBfloat16**

enumerator **ncclFloat8e4m3**

enumerator **ncclFloat8e5m2**

enumerator **ncclNumTypes**

*group* **rccl\_api\_custom\_redop**

API calls relating to creation/destroying custom reduction operator that pre-multiplies local source arrays prior to reduction

## Enums

enum **ncclScalarResidence\_t**

Location and dereferencing logic for scalar arguments.

Enumeration specifying memory location of the scalar argument. Based on where the value is stored, the argument will be dereferenced either while the collective is running (if in device memory), or before the `ncclRedOpCreate()` function returns (if in host memory).

*Values:*

enumerator **ncclScalarDevice**

Scalar is in device-visible memory

enumerator **ncclScalarHostImmediate**

Scalar is in host-visible memory

*group* **rccl\_collective\_api**

Collective communication operations must be called separately for each communicator in a communicator clique.

They return when operations have been enqueued on the HIP stream. Since they may perform inter-CPU synchronization, each call has to be done from a different thread or process, or need to use Group Semantics (see below).

*group* **msccl\_api**

API calls relating to the optional MSCCL algorithm datapath

## Typedefs

typedef int **mscclAlgoHandle\_t**

Opaque handle to MSCCL algorithm.

*group* **rccl\_group\_api**

When managing multiple GPUs from a single thread, and since RCCL collective calls may perform inter-CPU synchronization, we need to “group” calls for different ranks/devices into a single call.

Grouping RCCL calls as being part of the same collective operation is done using `ncclGroupStart` and `ncclGroupEnd`. `ncclGroupStart` will enqueue all collective calls until the `ncclGroupEnd` call, which will wait for all calls to be complete. Note that for collective communication, `ncclGroupEnd` only guarantees that the operations are enqueued on the streams, not that the operation is effectively done.

Both collective communication and `ncclCommInitRank` can be used in conjunction of `ncclGroupStart/ncclGroupEnd`, but not together.

Group semantics also allow to fuse multiple operations on the same device to improve performance (for aggregated collective calls), or to permit concurrent progress of multiple send/receive operations.

*page deprecated*

Global *mscclLoadAlgo* (const char \*mscclAlgoFilePath, mscclAlgoHandle\_t \*mscclAlgoHandle, int rank)

This function has been removed from the public API.

Global *mscclRunAlgo* (const void \*sendBuff, const size\_t sendCounts[], const size\_t sDisPls[], void \*recvBuff, const size\_t recvCounts[], const size\_t rDisPls[], size\_t count, ncclDataType\_t dataType, int root, int peer, ncclRedOp\_t op, mscclAlgoHandle\_t mscclAlgoHandle, ncclComm\_t comm, hipStream\_t stream)

This function has been removed from the public API.

Global *mscclUnloadAlgo* (mscclAlgoHandle\_t mscclAlgoHandle)

This function has been removed from the public API.

Global *ncclAllToAll* (const void \*sendbuff, void \*recvbuff, size\_t count, ncclDataType\_t datatype, ncclComm\_t comm, hipStream\_t stream) \_\_attribute\_\_((deprecated("ncclAllToAll is replaced with ncclAlltoAll and will be removed in the future")))

ncclAllToAll is replaced with ncclAlltoAll and will be removed in the future.

Global *ncclAllToAllv* (const void \*sendbuff, const size\_t sendcounts[], const size\_t sdispls[], void \*recvbuff, const size\_t recvcounts[], const size\_t rdispls[], ncclDataType\_t datatype, ncclComm\_t comm, hipStream\_t stream) \_\_attribute\_\_((deprecated("ncclAllToAllv is replaced with ncclAlltoAllv and will be removed in the future")))

ncclAllToAllv is replaced with ncclAlltoAllv and will be removed in the future.

*dir src*



## RCCL ENVIRONMENT VARIABLES

This section describes the most important RCCL environment variables, which are grouped by functionality.

### 11.1 Configuration and setup

The configuration and setup environment variables for RCCL are collected in the following table.

Environment variable	Values
<b>NCCL_CONF_FILE</b> Specifies the path to the RCCL configuration file.	String path to configuration file Default: <code>~/.rccl.conf</code> or <code>/etc/rccl.conf</code>
<b>NCCL_HOSTID</b> Sets the host identifier for multi-node communication.	String value for host identification Used for host hash generation
<b>NCCL_BOOTSTRAP_BIDIR_ALLGATHER</b> Enables the bidirectional ring AllGather (N/2 steps) on the socket OOB path during bootstrap. The unidirectional ring (N-1 steps) is kept as a fallback. Has no effect when net OOB is in use.	0: Force unidirectional ring. 1: Force bidirectional ring (default).

## 11.2 Logging and debugging

The logging and debugging environment variables for RCCL are collected in the following table.

Environment variable	Values
<p><b>NCCL_DEBUG</b> Controls debug logging in RCCL for troubleshooting and monitoring collective communication operations.</p>	<p>These are the logging levels in RCCL set via NCCL_DEBUG. Each logging level contains all logging for levels below it. The default logging level is ERROR.</p> <p>NONE: No logging is printed.  ERROR: These messages report when a fatal condition has occurred in RCCL and the operation can't continue.  VERSION: <code>librccl</code> version info is printed during the initialization phase.  WARN: Prints warnings about unusual conditions that could lead to unexpected results.  INFO: Prints standard logging messages about status and operations performed.  ABORT: Unused.  TRACE: Prints trace-level logging of function calls and parameters. Only active when <code>librccl</code> is built using <code>ENABLE_TRACE</code>.</p>
<p><b>NCCL_DEBUG_SUBSYS</b> Controls which subsystems generate debug output.</p>	<p>These are the logging subsystems set via NCCL_DEBUG_SUBSYS. These can be set as a comma-separated list, and can be inverted using the <code>^</code> prefix. The default subsystem set is INIT, BOOTSTRAP, and ENV.</p> <p>INIT: Prints during the initialization phase.  COLL: Prints during execution of collectives.  P2P: Prints logs related to peer-to-peer setup or communication.  SHM: Prints logs related to shared memory.  NET: Prints logs related to network setup or communication.  GRAPH: Prints logs related to parsing the topology of the network.  TUNING: Prints logs related to the tuner plugin.  ENV: Prints logs related to environment variables.  ALLOC: Prints logs related to memory allocation.  CALL: Prints logs for function calls (TRACE only).  PROXY: Prints logs related to the proxy thread.  NVLS: Not valid for AMD/RCCL.  BOOTSTRAP: Prints logs related to the bootstrapping phase of initialization.  REG: Prints logs related to registration and deregistration of transport initialization.  PROFILE: Prints logs related to the profiling/timing info.  RAS: Prints logs related to RAS.  VERBS: Prints logs related to IB/Verbs.  ALL: Activates all logging subsystems.</p>
<p><b>NCCL_WARN_ENABLE_DEBUG_INFO</b> Converts all WARN level logs to INFO level logs.</p>	<p>0: Default value. Variable is not enabled.  1: Enable the variable.</p>

## 11.2. Logging and debugging

79

**NCCL\_DEBUG\_TIMESTAMP\_LEVELS**  
The timestamp levels for

A set of NCCL\_DEBUG levels can have a timestamp prepended set as a comma-separated list which can be inverted using the `^` prefix. The

## 11.3 Algorithm and protocol control

The algorithm and protocol control environment variables for RCCL are collected in the following table.

Environment variable	Values
NCCL_ALGO Forces specific algorithm selection for collectives.	Algorithm name string Used to override automatic algorithm selection
NCCL_PROTO Forces specific protocol selection for communication.	Protocol name string Used to override automatic protocol selection

## 11.4 Network and topology

The network and topology environment variables for RCCL are collected in the following table.

Environment variable	Values
<p>NCCL_IB_HCA Specifies InfiniBand device:port to use.</p>	<p>Device specification string Prefix with ^ for exclusion, = for exact match</p>
<p>NCCL_IB_GID_INDEX Defines the Global ID index used in RoCE mode.</p>	<p>Integer value (default: -1) See InfiniBand show_gids command for valid values</p>
<p>NCCL_SOCKET_IFNAME Specifies which IP interfaces to use for communication.</p>	<p>Interface prefix string or list Multiple prefixes separated by , Prefix with ^ for exclusion, = for exact match Example: eth (all eth interfaces), =eth0 (exact match)</p>
<p>NCCL_SOCKET_FAMILY Forces IPv4/IPv6 interface selection.</p>	<p>AF_INET: Force IPv4 AF_INET6: Force IPv6 Unset: Use first available</p>
<p>NCCL_NET_MERGE_LEVEL Controls network device merging behavior.</p>	<p>Integer value specifying merge level Default: PATH_PORT</p>
<p>NCCL_NET_FORCE_MERGE Forces merging of network devices.</p>	<p>String specifying forced merge configuration</p>
<p>RCCL_IB_SPLIT_DATA_THRESHOLD Minimum message size (in bytes) before the payload is split across multiple NICs/QPs. Smaller messages use one QP for data to reduce latency. This variable can be leveraged when NIC Fusion (NCCL_NET_MERGE_LEVEL) and/or data splitting on QPs (NCCL_IB_SPLIT_DATA_ON_QPS) is enabled.</p>	<p>Integer value in bytes (default: 128) N: Split only when message size &gt;= N bytes</p>
<p>NCCL_RINGS Defines custom ring topology.</p>	<p>Ring topology specification string Overrides automatic topology detection</p>
<p>RCCL_TREES <b>11.4. Network and topology</b> Defines custom tree topology.</p>	<p>Tree topology specification string Alternative to ring topology</p>

## 11.5 Development and testing (advanced)

The development and testing environment variables for RCCL are collected in the following table. These variables are primarily intended for debugging and development purposes.

Environment variable	Values
CUDA_LAUNCH_BLOCKING Controls CUDA kernel launch blocking behavior.	0: Non-blocking launches 1 or non-zero: Blocking launches
NCCL_COMM_ID Enables multi-process mode in test applications.	Any non-empty value enables multi-process mode Used with test executables for distributed testing

## 11.6 Multi-communicator ordering

When an application uses multiple RCCL communicators on the same device, collective operations may execute in an unpredictable order unless the application adds explicit synchronization between streams.

Environment variable	Values
NCCL_LAUNCH_ORDER_IMPLICIT Serializes RCCL operations across different communicators on the same device according to their host-side launch sequence. This provides deterministic execution order for multi-communicator workloads such as chained collectives where one operation's output feeds into the next.	0: Disabled (default). 1: Enabled. Operations execute in host launch order.

## LICENSE

### Attributions

Contains contributions from NVIDIA.

Copyright (c) 2015-2020, NVIDIA CORPORATION. All rights reserved. Modifications Copyright (c) 2019-2026 Advanced Micro Devices, Inc. All rights reserved. Modifications Copyright (c) Microsoft Corporation. Licensed under the MIT License.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- Neither the name of NVIDIA CORPORATION, Lawrence Berkeley National Laboratory, the U.S. Department of Energy, nor the names of their contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS “AS IS” AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

The U.S. Department of Energy funded the development of this software under subcontract 7078610 with Lawrence Berkeley National Laboratory.

This code also includes files from the NVIDIA Tools Extension SDK project.

See:

<https://github.com/NVIDIA/NVTX>

for more information and license details.



## ATTRIBUTIONS

Contains contributions from NVIDIA.

Copyright (c) 2015-2020, NVIDIA CORPORATION. All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- Neither the name of NVIDIA CORPORATION, Lawrence Berkeley National Laboratory, the U.S. Department of Energy, nor the names of their contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS “AS IS” AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

The U.S. Department of Energy funded the development of this software under subcontract 7078610 with Lawrence Berkeley National Laboratory.

This code also includes files from the NVIDIA Tools Extension SDK project.

For more information and license details, see <https://github.com/NVIDIA/NVTX>



## M

mscclAlgoHandle\_t (C++ type), 51, 74  
 mscclLoadAlgo (C++ function), 68  
 mscclRunAlgo (C++ function), 68  
 mscclUnloadAlgo (C++ function), 69

## N

NCCL\_API\_MAGIC (C macro), 50  
 NCCL\_COMM\_NULL (C macro), 50  
 NCCL\_CONFIG\_INITIALIZER (C macro), 51, 71  
 NCCL\_CONFIG\_UNDEF\_INT (C macro), 50  
 NCCL\_CONFIG\_UNDEF\_PTR (C macro), 50  
 NCCL\_CTA\_POLICY\_DEFAULT (C macro), 50  
 NCCL\_CTA\_POLICY\_EFFICIENCY (C macro), 51  
 NCCL\_CTA\_POLICY\_ZERO (C macro), 51  
 NCCL\_H\_ (C macro), 49  
 NCCL\_MAJOR (C macro), 49  
 NCCL\_MINOR (C macro), 49  
 NCCL\_PATCH (C macro), 49  
 NCCL\_REVOKE\_DEFAULT (C macro), 51  
 NCCL\_SHRINK\_ABORT (C macro), 51  
 NCCL\_SHRINK\_DEFAULT (C macro), 51  
 NCCL\_SIM\_INFO\_INITIALIZER (C macro), 51  
 NCCL\_SPLIT\_NOCOLOR (C macro), 50  
 NCCL\_SUFFIX (C macro), 50  
 NCCL\_SUSPEND\_MEM (C macro), 51, 71  
 NCCL\_UNDEF\_FLOAT (C macro), 50  
 NCCL\_UNIQUE\_ID\_BYTES (C macro), 50  
 NCCL\_VERSION (C macro), 50  
 NCCL\_VERSION\_CODE (C macro), 50  
 NCCL\_WIN\_COLL\_SYMMETRIC (C macro), 50  
 NCCL\_WIN\_DEFAULT (C macro), 50  
 NCCL\_WIN\_REQUIRED\_ALIGNMENT (C macro), 50  
 NCCL\_WIN\_STRICT\_ORDERING (C macro), 50  
 ncclAllGather (C++ function), 39, 63  
 ncclAllReduce (C++ function), 38, 62  
 ncclAllReduceWithBias (C++ function), 62  
 ncclAllToAll (C++ function), 41, 65  
 ncclAlltoAll (C++ function), 64  
 ncclAllToAllv (C++ function), 66  
 ncclAlltoAllv (C++ function), 64  
 ncclBcast (C++ function), 37, 61  
 ncclBroadcast (C++ function), 38, 62  
 ncclComm\_t (C++ type), 42, 51  
 ncclCommAbort (C++ function), 36, 56  
 ncclCommCount (C++ function), 36, 59  
 ncclCommCuDevice (C++ function), 36, 59  
 ncclCommDeregister (C++ function), 59  
 ncclCommDestroy (C++ function), 36, 56  
 ncclCommFinalize (C++ function), 56  
 ncclCommGetAsyncError (C++ function), 59  
 ncclCommGetUniqueId (C++ function), 58  
 ncclCommGrow (C++ function), 58  
 ncclCommInitAll (C++ function), 35, 56  
 ncclCommInitRank (C++ function), 35, 55  
 ncclCommInitRankConfig (C++ function), 55  
 ncclCommInitRankScalable (C++ function), 58  
 ncclCommMemStat\_t (C++ enum), 52, 71  
 ncclCommMemStat\_t::ncclStatGpuMemPersist  
   (C++ enumerator), 52, 71  
 ncclCommMemStat\_t::ncclStatGpuMemSuspend  
   (C++ enumerator), 52, 71  
 ncclCommMemStat\_t::ncclStatGpuMemSuspended  
   (C++ enumerator), 52, 71  
 ncclCommMemStat\_t::ncclStatGpuMemTotal (C++  
   enumerator), 52, 71  
 ncclCommMemStats (C++ function), 60  
 ncclCommRegister (C++ function), 59  
 ncclCommResume (C++ function), 60  
 ncclCommRevoke (C++ function), 57  
 ncclCommShrink (C++ function), 57  
 ncclCommSplit (C++ function), 57  
 ncclCommSuspend (C++ function), 60  
 ncclCommUserRank (C++ function), 36, 59  
 ncclCommWindowDeregister (C++ function), 60  
 ncclCommWindowRegister (C++ function), 60  
 ncclConfig\_t (C++ struct), 47  
 ncclConfig\_t::blocking (C++ member), 47  
 ncclConfig\_t::cgaClusterSize (C++ member), 47  
 ncclConfig\_t::collnetEnable (C++ member), 48  
 ncclConfig\_t::commName (C++ member), 48  
 ncclConfig\_t::CTAPolicy (C++ member), 48  
 ncclConfig\_t::graphUsageMode (C++ member), 48  
 ncclConfig\_t::magic (C++ member), 47

ncclConfig\_t::maxCTAs (C++ member), 48  
 ncclConfig\_t::minCTAs (C++ member), 47  
 ncclConfig\_t::nChannelsPerNetPeer (C++ member), 48  
 ncclConfig\_t::netName (C++ member), 48  
 ncclConfig\_t::numRmaCtx (C++ member), 48  
 ncclConfig\_t::nvlinkCentricSched (C++ member), 48  
 ncclConfig\_t::nvlsCTAs (C++ member), 48  
 ncclConfig\_t::shrinkShare (C++ member), 48  
 ncclConfig\_t::size (C++ member), 47  
 ncclConfig\_t::splitShare (C++ member), 48  
 ncclConfig\_t::trafficClass (C++ member), 48  
 ncclConfig\_t::version (C++ member), 47  
 ncclDataType\_t (C++ enum), 44, 53, 73  
 ncclDataType\_t::ncclBfloat16 (C++ enumerator), 45, 54, 73  
 ncclDataType\_t::ncclChar (C++ enumerator), 44, 53, 73  
 ncclDataType\_t::ncclDouble (C++ enumerator), 45, 54, 73  
 ncclDataType\_t::ncclFloat (C++ enumerator), 45, 54, 73  
 ncclDataType\_t::ncclFloat16 (C++ enumerator), 44, 54, 73  
 ncclDataType\_t::ncclFloat32 (C++ enumerator), 45, 54, 73  
 ncclDataType\_t::ncclFloat64 (C++ enumerator), 45, 54, 73  
 ncclDataType\_t::ncclFloat8e4m3 (C++ enumerator), 45, 54, 73  
 ncclDataType\_t::ncclFloat8e5m2 (C++ enumerator), 45, 54, 73  
 ncclDataType\_t::ncclHalf (C++ enumerator), 44, 54, 73  
 ncclDataType\_t::ncclInt (C++ enumerator), 44, 53, 73  
 ncclDataType\_t::ncclInt32 (C++ enumerator), 44, 53, 73  
 ncclDataType\_t::ncclInt64 (C++ enumerator), 44, 54, 73  
 ncclDataType\_t::ncclInt8 (C++ enumerator), 44, 53, 73  
 ncclDataType\_t::ncclNumTypes (C++ enumerator), 45, 54, 73  
 ncclDataType\_t::ncclUInt32 (C++ enumerator), 44, 54, 73  
 ncclDataType\_t::ncclUInt64 (C++ enumerator), 44, 54, 73  
 ncclDataType\_t::ncclUInt8 (C++ enumerator), 44, 53, 73  
 ncclGather (C++ function), 40, 65  
 ncclGetErrorString (C++ function), 42, 58  
 ncclGetLastError (C++ function), 59  
 ncclGetUniqueId (C++ function), 35, 55  
 ncclGetVersion (C++ function), 42, 55  
 ncclGroupEnd (C++ function), 42, 69  
 ncclGroupSimulateEnd (C++ function), 69  
 ncclGroupStart (C++ function), 41, 69  
 ncclMemAlloc (C++ function), 55  
 ncclMemFree (C++ function), 55  
 ncclPutSignal (C++ function), 67  
 ncclRecv (C++ function), 40, 67  
 ncclRedOp\_dummy\_t (C++ enum), 52, 72  
 ncclRedOp\_dummy\_t::ncclNumOps\_dummy (C++ enumerator), 52, 72  
 ncclRedOp\_t (C++ enum), 43, 52, 72  
 ncclRedOp\_t::ncclAvg (C++ enumerator), 44, 53, 72  
 ncclRedOp\_t::ncclMax (C++ enumerator), 44, 53, 72  
 ncclRedOp\_t::ncclMaxRedOp (C++ enumerator), 44, 53, 72  
 ncclRedOp\_t::ncclMin (C++ enumerator), 44, 53, 72  
 ncclRedOp\_t::ncclNumOps (C++ enumerator), 44, 53, 72  
 ncclRedOp\_t::ncclProd (C++ enumerator), 44, 53, 72  
 ncclRedOp\_t::ncclSum (C++ enumerator), 43, 53, 72  
 ncclRedOpCreatePreMulSum (C++ function), 60  
 ncclRedOpDestroy (C++ function), 61  
 ncclReduce (C++ function), 37, 61  
 ncclReduceScatter (C++ function), 38, 63  
 ncclResetDebugInit (C++ function), 59  
 ncclResult\_t (C++ enum), 43, 51, 70  
 ncclResult\_t::ncclInProgress (C++ enumerator), 43, 52, 70  
 ncclResult\_t::ncclInternalError (C++ enumerator), 43, 52, 70  
 ncclResult\_t::ncclInvalidArgument (C++ enumerator), 43, 52, 70  
 ncclResult\_t::ncclInvalidUsage (C++ enumerator), 43, 52, 70  
 ncclResult\_t::ncclNumResults (C++ enumerator), 43, 52, 70  
 ncclResult\_t::ncclRemoteError (C++ enumerator), 43, 52, 70  
 ncclResult\_t::ncclSuccess (C++ enumerator), 43, 51, 70  
 ncclResult\_t::ncclSystemError (C++ enumerator), 43, 51, 70  
 ncclResult\_t::ncclUnhandledCudaError (C++ enumerator), 43, 51, 70  
 ncclScalarResidence\_t (C++ enum), 54, 74  
 ncclScalarResidence\_t::ncclScalarDevice (C++ enumerator), 54, 74  
 ncclScalarResidence\_t::ncclScalarHostImmediate (C++ enumerator), 54, 74  
 ncclScatter (C++ function), 40, 65  
 ncclSend (C++ function), 39, 67

ncclSignal (C++ *function*), 68  
 ncclSimInfo\_t (C++ *struct*), 48  
 ncclSimInfo\_t::estimatedTime (C++ *member*), 49  
 ncclSimInfo\_t::magic (C++ *member*), 49  
 ncclSimInfo\_t::size (C++ *member*), 49  
 ncclSimInfo\_t::version (C++ *member*), 49  
 ncclUniqueId (C++ *struct*), 42, 49  
 ncclUniqueId::internal (C++ *member*), 49  
 ncclWaitSignal (C++ *function*), 68  
 ncclWaitSignalDesc\_t (C++ *struct*), 49  
 ncclWaitSignalDesc\_t::ctx (C++ *member*), 49  
 ncclWaitSignalDesc\_t::opCnt (C++ *member*), 49  
 ncclWaitSignalDesc\_t::peer (C++ *member*), 49  
 ncclWaitSignalDesc\_t::sigIdx (C++ *member*), 49  
 ncclWindow\_t (C++ *type*), 51  
 ncclWinGetUserPtr (C++ *function*), 60

## P

pnccclCommGetUniqueId (C++ *function*), 58  
 pnccclCommGrow (C++ *function*), 58  
 pnccclCommShrink (C++ *function*), 58  
 pnccclGroupSimulateEnd (C++ *function*), 69  
 pnccclMemAlloc (C++ *function*), 55  
 pnccclMemFree (C++ *function*), 55  
 pnccclPutSignal (C++ *function*), 67  
 pnccclSignal (C++ *function*), 68  
 pnccclWaitSignal (C++ *function*), 68  
 pnccclWinGetUserPtr (C++ *function*), 60

## R

RCCL\_ALLREDUCE\_WITH\_BIAS (C *macro*), 50  
 RCCL\_ALLTOALLV (C *macro*), 50  
 RCCL\_BFLOAT16 (C *macro*), 50  
 RCCL\_FLOAT8 (C *macro*), 50  
 RCCL\_GATHER\_SCATTER (C *macro*), 50  
 RCCL\_SUSPEND\_RESUME (C *macro*), 50