

NVIDIA Collective Communication Library (NCCL)

Release Notes

Table of Contents

Chapter 1. NCCL Overview	1
Chapter 2. NCCL Release 2.28.3	3
Chapter 3. NCCL Release 2.27.7	6
Chapter 4. NCCL Release 2.27.6	8
Chapter 5. NCCL Release 2.27.5	10
Chapter 6. NCCL Release 2.27.3	12
Chapter 7. NCCL Release 2.26.5	15
Chapter 8. NCCL Release 2.26.2	17
Chapter 9. NCCL Release 2.25.1	20
Chapter 10. NCCL Release 2.24.3	21
Chapter 11. NCCL Release 2.23.4	23
Chapter 12. NCCL Release 2.22.3	25
Chapter 13. NCCL Release 2.21.5	27
Chapter 14. NCCL Release 2.20.5	29
Chapter 15. NCCL Release 2.20.3	30
Chapter 16. NCCL Release 2.19.3	32
Chapter 17. NCCL Release 2.18.5	34
Chapter 18. NCCL Release 2.18.3	36
Chapter 19. NCCL Release 2.18.1	38
Chapter 20. NCCL Release 2.17.1	40
Chapter 21. NCCL Release 2.16.5	42
Chapter 22. NCCL Release 2.16.2	44
Chapter 23. NCCL Release 2.15.5	46
Chapter 24. NCCL Release 2.15.1	47
Chapter 25. NCCL Release 2.14.3	48
Chapter 26. NCCL Release 2.13.4	50
Chapter 27. NCCL Release 2.12.12	52
Chapter 28. NCCL Release 2.12.10	54
Chapter 29. NCCL Release 2.12.7	55

Chapter 30. NCCL Release 2.11.4	56
Chapter 31. NCCL Release 2.10.3	57
Chapter 32. NCCL Release 2.9.9	58
Chapter 33. NCCL Release 2.9.8	59
Chapter 34. NCCL Release 2.9.6	60
Chapter 35. NCCL Release 2.8.4	61
Chapter 36. NCCL Release 2.8.3	62
Chapter 37. NCCL Release 2.7.8	64
Chapter 38. NCCL Release 2.7.6	66
Chapter 39. NCCL Release 2.7.5	67
Chapter 40. NCCL Release 2.7.3	69
Chapter 41. NCCL Release 2.6.4	71
Chapter 42. NCCL Release 2.5.6	72
Chapter 43. NCCL Release 2.4.8	73
Chapter 44. NCCL Release 2.4.7	75
Chapter 45. NCCL Release 2.4.2	77
Chapter 46. NCCL Release 2.3.7	78
Chapter 47. NCCL Release 2.3.5	79
Chapter 48. NCCL Release 2.3.4	80
Chapter 49. NCCL Release 2.2.13	81
Chapter 50. NCCL Release 2.2.12	82
Chapter 51. NCCL Release 2.1.15	84
Chapter 52. NCCL Release 2.1.4	86
Chapter 53. NCCL Release 2.1.2	87
Chapter 54. NCCL Release 2.0.5	88
Chapter 55. NCCL Release 2.0.4	89
Chapter 56. NCCL Release 2.0.2	90

Chapter 1. NCCL Overview

The NVIDIA® Collective Communications Library ™ (NCCL) (pronounced "Nickel") is a library of multi-GPU collective communication primitives that are topology-aware and can be easily integrated into applications.

Collective communication algorithms employ many processors working in concert to aggregate data. NCCL is not a full-blown parallel programming framework; rather, it is a library focused on accelerating collective communication primitives. The following collective operations are currently supported:

- AllReduce
- Broadcast
- Reduce
- AllGather
- ReduceScatter

Tight synchronization between communicating processors is a key aspect of collective communication. CUDA® based collectives would traditionally be realized through a combination of CUDA memory copy operations and CUDA kernels for local reductions. NCCL, on the other hand, implements each collective in a single kernel handling both communication and computation operations. This allows for fast synchronization and minimizes the resources needed to reach peak bandwidth.

NCCL conveniently removes the need for developers to optimize their applications for specific machines. NCCL provides fast collectives over multiple GPUs both within and across nodes. It supports a variety of interconnect technologies including PCIe, NVLink™ , InfiniBand Verbs, and IP sockets. NCCL also automatically patterns its communication strategy to match the system's underlying GPU interconnect topology.

Next to performance, ease of programming was the primary consideration in the design of NCCL. NCCL uses a simple C API, which can be easily accessed from a variety of programming languages. NCCL closely follows the popular collectives API defined by MPI (Message Passing Interface). Anyone familiar with MPI will thus find NCCL API very natural to use. In a minor departure from MPI, NCCL collectives take a "stream" argument which provides direct integration with the CUDA programming model. Finally, NCCL is compatible with virtually any multi-GPU parallelization model, for example:

- single-threaded
- multi-threaded, for example, using one thread per GPU

multi-process, for example, MPI combined with multi-threaded operation on GPUs

NCCL has found great application in deep learning frameworks, where the AllReduce collective is heavily used for neural network training. Efficient scaling of neural network training is possible with the multi-GPU and multi node communication provided by NCCL.

Chapter 2. NCCL Release 2.28.3

This is the NCCL 2.28.3 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.28.3 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 12.x and CUDA 13.x.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added Device API, an experimental device-side API to integrate NCCL communication directly into application kernels. The signatures and functionality may evolve in future releases. No ABI compatibility is currently guaranteed for this API. Applications must be recompiled with each new NCCL release.
- Symmetric kernels were reimplemented using the device API and have support for aggregating symmetric operations using ncclGroupStart/End.
- Added new host collective APIs: ncclAlltoAll, ncclScatter, and ncclGather.
- Added experimental support for CMake as an alternative to building with Makefiles. Known issues to be fixed in an upcoming release: pkg.build and Device API are not supported.
- Reduced SM utilization within a single (MN)NVL domain by using CE (Copy Engine) Collectives. Freed up more SM capacity for the application when using ncclAlltoAll, ncclGather, or ncclScatter. To enable the feature, register buffers into symmetric windows and use the NCCL CTA POLICY ZERO flag in the communicator config t.
- Decreased max CTA count from 32 to 16 on Blackwell. This improvement decreased SM overhead by 50%, but it may cause a performance drop on Blackwell. It can be overridden by setting NCCL MIN CTAS=32 and NCCL MAX CTAS=32, or setting the same in the communicator config. Based on community feedback, future versions may consider different trade-offs between performance and SM overhead.

- Network Plugin added version 11 which supports per-communicator init/finalize. Now passes information about communication operations to be executed on the network end point. Added multi-request net API to help anticipate and optimize multiple send/ recv requests.
- Profiler Plugin added support for API events (group, collective, and p2p) and for tracking kernel launches. Added Inspector Profiler example plugin for always-on performance monitoring and a hook to Google's CoMMA profiler on github.
- ► Tuner Plugin exposed NCCL tuning constants with ncclTunerConstants v5 t and added NVL Domain Information API.
- Plugin system added support for multiple plugin types from a single shared object.
- ▶ Added new option NCCL MNNVL CLIQUE ID=-2 which will use rack serial number to partition the MNNVL clique. This will limit NVLink domains to GPUs within a single rack.
- Added NCCL NETDEVS POLICY to control how NET devices are assigned to GPUs. The default (AUTO) is the policy used in previous versions.
- ▶ Added NCCL SINGLE PROC MEM REG ENABLE control variable to enable NVLS UB registration in the "one process, multiple ranks" case as opt in.
- ▶ Moved nChannelsPerNetPeer into ncclConfig. NCCL NCHANNELS PER NET PEER can override the value in ncclConfig.
- Enabled PXN over C2C by default, which can improve performance for Grace-Blackwell platforms by allowing NCCL to leverage the NIC attached to a peer GPU over NVLINK, C2C, and PCIe. It can be disabled by setting NCCL PXN C2C=0.

Fixed Issues

The following issues have been resolved in NCCL 2.28.3:

- Allowed FP8 support for non-reductive operations on pre-sm90 devices.
- ► Fixed NVLS+CollNet and temporarily disabled COLLNET CHAIN for >8 GPUs.
- Fixed considering non-running interfaces for socket traffic. Now NCCL will not attempt to use interfaces that do not have the IFF RUNNING bit.
- Improved response to RoCE link flaps. Instead of reporting an "unknown event" it will now report "GID table changed."
- Moved libvirt bridge interface to the end of possible interfaces so that it is considered last. These interfaces are usually virtual bridges to relay traffic to containers running on the host and cannot be used for traffic to a remote node and are therefore unsuitable.

Updating the GPG Repository Key

To best ensure the security and reliability of our RPM and Debian package repositories, NVIDIA is updating and rotating the signing keys used by apt, dnf/yum, and zypper package managers beginning on April 27, 2022. Failure to update your repository signing

keys will result in package management errors when attempting to access or install NCCL packages. To ensure continued access to the latest NCCL release, please follow the updated NCCL installation guide.

Chapter 3. NCCL Release 2.27.7

This is the NCCL 2.27.7 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.27.7 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- This NCCL release supports CUDA 12.2, CUDA 12.4, and CUDA 12.9. The provided prebuilt binaries should work with other CUDA 12.x versions as well. CUDA 13.0 is supported as well.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

Added support for CUDA 13.0.

Fixed Issues

The following issues have been resolved in NCCL 2.27.7:

Fixed initialization failures in certain configurations when attempting to load fp8specific symmetric multicast kernels on GPUs older than Blackwell.

Known Issues

A note for users on MNNVL systems: please ensure an adequate stack size for NCCL threads. While the default Linux stack size limit of 8192 KB is known to be sufficient, we've seen crashes if the limit is changed to "unlimited," as it causes the glibc library to unexpectedly decrease the stack size of NCCL's background threads to just 2048 KB. Use ulimit -s in bash to print the current limit; if needed, reset it to 8192 KB using ulimit -s 8192. (You also need to ensure that the new setting is propagated to other nodes when launching a multi-node NCCL job).

Chapter 4. NCCL Release 2.27.6

This is the NCCL 2.27.6 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.27.6 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- This NCCL release supports CUDA 12.2, CUDA 12.4, and CUDA 12.9. The provided prebuilt binaries should work with other CUDA 12.x versions as well.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Improved support for DirectNIC (CX8): XDR speed detection and the reporting of the RDMA interfaces only when DirectNIC is enabled.
- Completed P2C (PXN over C2C) support, which is now preferred over regular PXN and extends to send/receive operations as well. This feature is currently preliminary and is disabled by default; use NCCL PXN C2C=1 to enable.
- Added support for compilation with GCC 14.

Fixed Issues

The following issues have been resolved in NCCL 2.27.6:

- Fixed the unloading of network plugins that also provide tuner capability.
- ▶ Fixed the change of the current device across the calls to ncclCommDestroy() and ncclCommAbort().

Known Issues

A note for users on MNNVL systems: please ensure an adequate stack size for NCCL threads. While the default Linux stack size limit of 8192 KB is known to be sufficient, we've seen crashes if the limit is changed to "unlimited," as it causes the glibc library

to unexpectedly decrease the stack size of NCCL's background threads to just 2048 KB. Use ulimit -s in bash to print the current limit; if needed, reset it to 8192 KB using ulimit -s 8192. (You also need to ensure that the new setting is propagated to other nodes when launching a multi-node NCCL job).

Updating the GPG Repository Key

Chapter 5. NCCL Release 2.27.5

This is the NCCL 2.27.5 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.27.5 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- This NCCL release supports CUDA 12.2, CUDA 12.4, and CUDA 12.9. The provided prebuilt binaries should work with other CUDA 12.x versions as well.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Optimized the network performance on GB200 systems by alternating the direction of the rings and the NIC to GPU assignment across communicators to limit unnecessary sharing.
- Optimized the network performance on DGX B200 systems by adjusting the bandwidths provided to the graph search algorithm.
- Added an example tuner plugin with CSV-based overrides.

Fixed Issues

The following issues have been resolved in NCCL 2.27.5:

- Fixed the detection of C2C links in case GPU Direct RDMA is disabled between a GPU and a NIC.
- Fixed PXN support on MNNVL systems, where NCCL would try (and fail) to share regular host memory across multiple nodes.
- Further reduced the overheads of CUDA graph capturing, which increased in NCCL 2.26.2 for large graphs.
- Enabled fp8 reductions in symmetric kernels on Blackwell with CUDA 12.8.

- Restored the plugin name handling logic to make it possible to specify a path to the plugin.
- ▶ Restored the ability to change NCCL COLLNET ENABLE during execution
- Removed an x86 dependency from the example profiler.

Chapter 6. NCCL Release 2.27.3

This is the NCCL 2.27.3 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.27.3 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- This NCCL release supports CUDA 12.2, CUDA 12.4, and CUDA 12.9. The provided prebuilt binaries should work with other CUDA 12.x versions as well.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- ► Added symmetric memory API (ncclCommWindowRegister with the NCCL WIN COLL SYMMETRIC flag, ncclCommWindowDeregister). This capability is on by default (NCCL WIN ENABLE=0 can be used to disable) and depends on CUMEM being operational.
- Implemented specialized kernels taking advantage of symmetrically registered memory. Initial support includes P2P and NVLS connectivity, floating point types up to 32 bits, sum as the reduction operator, and one collective operation per group.
- Added support for DGX Spark.
- Added support for DirectNIC (CX8) to the internal IB plugin.
- Added support for communicator shrinking (ncclCommShrink).
- ▶ Added support for loading multiple network plugins (NCCL NET PLUGIN now accepts a comma-separated list of plugins to load).
- Added NVLS+IB SHARP support for AllGather and ReduceScatter with user buffer registration.
- Decreased the NVLS channel count to 24 on Blackwell systems with multiple NVLink domains per communicator.

- Enabled fine-tuning of NCCL behavior per communicator using new ncclconfig t members collnetEnable, CTAPolicy, and nvlsCTAs (the environment variable NCCL CTA POLICY can also be used to force a particular CTA policy).
- Implemented several enhancements to the profiler in terms of expanding the amount of information being provided at initialization, providing GPU-generated timestamps for the GPU kernel events, adding support for network-defined event updates, and more.
- Optimized the performance of collective operations on GB200 systems by increasing to 16 the number of NICs used to communicate between MNNVL domains.
- Added support for more complex MNNVL topologies.
- Disabled implicit fallback to a potentially much slower network transport if the MNNVL fabric initialization was unsuccessful. Such failures typically indicate a misconfigured IMEX support on the system; NCCL MNNVL ENABLE=0 can be used to continue without MNNVL.
- Disabled the creation of fused NICs for physical devices that haven't been merged.
- Improved support for platforms with C2C connectivity by enabling GPUDirect RDMA for the NICs by default (NCCL NET GDR C2C=0 can be used to disable) and by adding support for P2C (PXN over C2C) and the LL128 protocol.
- Extended NCCL fault tolerance to support more complex multithreaded scenarios during initialization.
- Enabled ncclimplicitOrderLaunch for CUDA 12.9+.
- Improved the netSocket transport latency and provided finer control over, e.g., send/ receive buffer size.
- Improved the readability of the CPU affinity in the debug output.

Fixed Issues

The following issues have been resolved in NCCL 2.27.3:

- ▶ Enabled graceful fallback if NVLS initialization fails (NCCL NVLS ENABLE=1 can be used to preserve the old behavior of returning an error).
- Fixed several issues in the profiler support: reporting the correct number of channels used, keeping track of an event identifier, improving backward compatibility, and more.
- Fixed multiple issues related to MNNVL support: a hang in alltoall-like communication patterns at a scale of over 80 ranks, NCCL P2P DISABLE=1 leaving MNNVL enabled (it now implies NCCL MNNVL ENABLE=0), an initialization failure when NCCL TOPO FILE was being used, failure to exclude non-local NICs in the graph search, and an incompatibility of the SHM transport with MNNVL.
- Fixed a potential race condition when mixing graph and non-graph execution.
- Improved PXN connection code to avoid duplicate and unused connections.

- Fixed several issues in RAS related to race conditions and memory corruption.
- ▶ Fixed a potential memory corruption in ncclCommSplit with resource sharing.
- Fixed a small memory leak and over-sychronization in asynchronous graph upload.
- ▶ Added a check for out-of-memory conditions in ncclMemAlloc.
- Cleaned up the NCCL socket code: made it more willing to retry in case of errors during connection establishment, added error checking in a few instances, and improved the debug output.
- ▶ Switched NCCL DEBUG FILE to line buffering.
- Fixed other minor issues, primarily in the graph search code and the internal IB plugin

Chapter 7. NCCL Release 2.26.5

This is the NCCL 2.26.5 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.26.5 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 12.2, CUDA 12.4, and CUDA 12.9.

Fixed Issues

The following issues have been resolved in NCCL 2.26.5:

- Minimized the performance impact of the device kernel profiling support when the profiler plugin is not loaded.
- Reduced the overheads of CUDA graph capturing, which increased in NCCL 2.26.2 for large graphs.
- Fixed the exchange of enhanced connection establishment (ECE) options to address potential slowdowns on networks utilizing RoCE.
- Added testing if cuMem host allocations work and if not, disabling them. Enabled by default since NCCL 2.24 if the CUDA driver version is at least 12.6, such allocations rely on NUMA support, which is by default not available under Docker. We recommend invoking Docker with --cap-add SYS NICE to enable it.
- Worked around a potential hang in alltoall-like communication patterns on MNNVL systems at a scale of over 80 ranks.
- ► Fixed an initialization error when running with NCCL NET GDR C2C=1 on multiple MNNVL domains with non-uniform network configurations across nodes.
- Fixed the printing of sub-seconds in the debug log when using a custom NCCL DEBUG TIMESTAMP FORMAT setting.

Chapter 8. NCCL Release 2.26.2

This is the NCCL 2.26.2 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.26.2 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 12.2, CUDA 12.4, and CUDA 12.8.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added support for the profiling of CUDA kernel start and ends and of the network plugin. Improved profiling granularity and added support for graph capturing.
- Added implicit launch order, which allows to prevent deadlocks when using multiple NCCL communicators per device by implicitly ordering NCCL operations using the host program order. Disabled by default, set NCCL LAUNCH ORDER IMPLICIT=1 to enable.
- Added a complementary mechanism to detect host threads racing to launch to the same device. Enabled by default, set NCCL LAUNCH RACE FATAL=0 to disable.
- Significantly accelerated the PAT algorithm by separating the computation and execution of PAT steps on different warps.
- Added support for setting QoS per communicator, using a new traffic class config to allow the application to select a particular traffic class for a given communicator. For the IB/RoCE plugin, existing config variables such as NCCL IB SL and NCCL IB TC take precedence.
- Added support for enabling GPU Direct RDMA specifically on C2C platforms. Disabled by default, set NCCL NET GDR C2C=1 to enable.
- Keep user buffer registration enabled as much as possible, disabling it only when a communicator has more than one rank per node on any node.

- Report operation counts in RAS separately for each collective type and provide details about missing communicator ranks.
- Added support for timestamps to NCCL diagnostic messages. This is enabled by default for WARN messages, use NCCL DEBUG TIMESTAMP LEVELS to specify the levels that should include a timestamp. Use NCCL DEBUG TIMESTAMP FORMAT to adjust the format of the timestamps.
- Reduced the memory usage with NVLink SHARP (NVLS).
- Improved algorithm/protocol selection on recent Intel CPUs such as Emerald Rapids and Sapphire Rapids.
- Improved channel scheduling when mixing LL and Simple operations.
- Added support for comment lines (starting with #) in the nccl.conf file.
- Make user buffer registration problems print an INFO instead of a WARN.

Fixed Issues

The following issues have been resolved in NCCL 2.26.2:

- Fixed a potential hang during connection setup when multiple communicators share resources.
- ► Fixed a performance regression when using NCCL CROSS NIC=1.
- Made the GID index detection code more resilient to avoid issues with containers.
- Fixed a potential crash when creating a non-blocking communicator after a nonblocking collective operation on another communicator.
- Fixed shared memory usage on recent Blackwell GPUs.
- Fixed an issue where reloading the IB SHARP plugin could result in an error.
- Made the failures to auto-merge NICs non-fatal, such as when mixing IB and RoCE devices.
- Fixed a potential hang in ncclCommAbort and reduced the abort time by up to two orders of magnitude.
- ▶ Fixed a crash when libnccl.so was dynamically unloaded.
- Fixed a hang if the network plugin returned an error.
- Fixed a hang on heterogeneous architectures by harmonizing tuning choices.
- Fixed a potential crash in case of a failed communicator initialization or termination.
- Fixed a potential bug during a group launch of multiple communicators.
- Fixed a bug where, under rare circumstances, certain variables specified in the config file could be ignored.

Updating the GPG Repository Key

To best ensure the security and reliability of our RPM and Debian package repositories, NVIDIA is updating and rotating the signing keys used by apt, dnf/yum, and zypper

package managers beginning on April 27, 2022. Failure to update your repository signing keys will result in package management errors when attempting to access or install NCCL packages. To ensure continued access to the latest NCCL release, please follow the updated NCCL installation guide.

Chapter 9. NCCL Release 2.25.1

This is the NCCL 2.25.1 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.25.1 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 12.2, CUDA 12.4, and CUDA 12.8.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added support for Blackwell.
- Improved NVTX instrumentation.

Updating the GPG Repository Key

Chapter 10. NCCL Release 2.24.3

This is the NCCL 2.24.3 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.24.3 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 12.2, CUDA 12.4, and CUDA 12.6.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Leverage user buffer registration for inter-node copy operations using the ring or tree algorithm.
- Added RAS subsystem and ncclras tool to allow for a report of the NCCL state in case of a hang.
- Added support for 8bit floating point data types (e5m2 and e4m3).
- Added NIC fusion: fuse multiple NICs together as a single, larger one.
- Retry in case of socket connection failure (unreachable host).
- Retry in case of IB QP connection failure.
- Improved support for external network plugins: allow plugins to force a flush, indicate when completion is not needed, allow for full offload of allgather operations when using one GPU per node.
- ► Extended NCCL ALGO/NCCL PROTO syntax and strictly enforce their values.
- Made host memory allocation use cumem functions by default.

Fixed Issues

The following issues have been resolved in NCCL 2.24.3:

Return ncclinvalidusage when NCCL SOCKET IFNAME is set to an incorrect value.

- Fixed PAT tuning.
- Fixed hangs when running with different CPU architectures.
- Fixed FD leak in UDS.
- Fixed crash when mixing buffer registration and graph buffer registration.
- ▶ Made ncclSend/ncclRecv communication with buffer registration functional on network plugins relying on dmabuf for buffer registration.
- Fixed crash in IB code caused by uninitialized fields.
- ► Fixed case where ncclSend/ncclRecv would return ncclSuccess in non-blocking mode even though the operation was not enqueued onto the stream (Github issue 1495).

Chapter 11. NCCL Release 2.23.4

This is the NCCL 2.23.4 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.23.4 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 12.2, CUDA 12.4, and CUDA 12.6.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Add a new scalable initialization APL
- Accelerate bootstrap operations.
- Add new PAT (Parallel Aggregated Trees) algorithm for allgather and reduce scatter operations using one GPU per node.
- Accelerate intra-node communication when buffers are registered.
- Add profiler plugin API.
- Make CUDA calls within graph allocation asynchronous.
- Use fatal RDMA asynchronous events to stop network operations early.
- Disable GPU Direct P2P on AMD CPUs when using more than 2 GPUs.
- Add parameter to set the location of the user configuration file.
- Increase default IB timeout from 18 to 20.

Fixed Issues

The following issues have been resolved in NCCL 2.23.4:

- Fixed GPU Direct RDMA check on linux kernels 6.6+.
- Fixed performance regression when mixing small and large operations.
- ► Fixed crash in topology detection when devices have a NUMA ID of -1.

- ▶ Fixed Tree graph search when NCCL_CROSS_NIC is set to 1...
- Fixes for IB operation on multi-node NVLink systems.

Chapter 12. NCCL Release 2.22.3

This is the NCCL 2.22.3 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.22.3 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 12.2, CUDA 12.4, and CUDA 12.5.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Lazy connection establishment for collective operations to speed up init and lower memory usage when not using all algorithms.
- Accelerated Intra-node NVLink detection for faster init.
- Init profiling, reporting time spent in different init phases.
- Support for PCI p2p on split PCI switches.
- Cost estimation API.
- Net/IB: separate traffic class for fifo messages.
- Net/IB: support for IB router.
- Optimizations and fixes for device network offload (unpack).
- Support ncclGroupStart/End for ncclCommAbort/Destroy.
- Improved Tuner API.
- Allow net and tuner plugins to be statically linked.

Fixed Issues

The following issues have been resolved in NCCL 2.22.3:

- Disabled port fusion in heterogeneous systems.
- Fixed NVLS graphs search for dual NICs.

- Fixed crash with collnetDirect.
- Fixed hang when nodes have different CPU types.
- Fixed performance of registered send/recv operations.
- ▶ Ensured ncclCommDestroy is non-blocking if ncclCommFinalize has been called.
- Improved error reporting during SHM segment creation.
- Improved support of various compilers.

Chapter 13. NCCL Release 2.21.5

This is the NCCL 2.21.5 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.21.5 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- This NCCL release supports CUDA 11.0, CUDA 12.2, CUDA 12.4, and CUDA 12.5.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added support for user buffer registrations on IB SHARP operations with one GPU per node.
- Improved support for Multi-node NVLink systems, add NVLink SHARP support and multi-click support.
- Added support for dynamic GID detection on RoCE.
- Reduced memory usage when NVLink SHARP is enabled.
- Improved tuner plugin loading.
- Added signature to communicator objects to help detect the corruption of communicator objects or invalid communicator pointers.

Fixed Issues

The following issues have been resolved in NCCL 2.21.5:

- Fixed IB SHARP rail mapping when using split communicators.
- Fixed mismatch crash during bootstrap due to TCP packet reordering.
- Fixed hang with heterogeneous systems, causing the crossNic value to be different between nodes.
- Fixed minCompCap/maxCompCap computation.

Chapter 14. NCCL Release 2.20.5

This is the NCCL 2.20.5 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.20.5 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 11.0, CUDA 12.2, and CUDA 12.4.

Fixed Issues

The following issues have been resolved in NCCL 2.20.5:

▶ Fixed crash in UDS connection when using ncclCommSplit.

Updating the GPG Repository Key

Chapter 15. NCCL Release 2.20.3

This is the NCCL 2.20.3 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.20.3 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 11.0, CUDA 12.2, and CUDA 12.3.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Improved ring algorithm (alternating rings), to remove bottleneck on systems where each GPU only has one local NIC, like the DGX H100.
- Added support for user buffer registration for network send/recv operations.
- Optimized aggregated operations to better utilize all channels.
- Added support for large Broadcom PCI gen5 switches, flattening the reported twolevel topology.
- Added support for inter-node NVLink communication.
- Add support for port fusion in NET/IB.
- Added support for ReduceScatter and AllGather using IB SHARP.

Fixed Issues

The following issues have been resolved in NCCL 2.20.3:

Fixed hang during A2A connection.

Updating the GPG Repository Key

To best ensure the security and reliability of our RPM and Debian package repositories, NVIDIA is updating and rotating the signing keys used by apt, dnf/yum, and zypper package managers beginning on April 27, 2022. Failure to update your repository signing

keys will result in package management errors when attempting to access or install NCCL packages. To ensure continued access to the latest NCCL release, please follow the updated NCCL installation guide.

Chapter 16. NCCL Release 2.19.3

This is the NCCL 2.19.3 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.19.3 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ► This NCCL release supports CUDA 11.0, CUDA 12.0, CUDA 12.2, and CUDA 12.3.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added local user buffer registration for NVLink SHARP.
- Improved performance on Hopper GPUs (H800/H100).
- Added tuning plugin support.
- Added NET API v7 to allow for device side packet reordering; removed v4 plugin support.
- Added support for RoCE ECE.
- Added support for C2C links.
- Disabled network flush by default on H100.

Fixed Issues

The following issues have been resolved in NCCL 2.19.3:

- Better detection of shared memory allocation failures and report them instead of crashing later with a "Bus error."
- Fixed missing thread unlock in bootstrap code.

Known Issues

- Alltoall performance at scale may see degradation for medium-size operations. Setting NCCL NCHANNELS PER NET PEER=1 should workaround the problem.
- ► A hang may occur during alltoall connection establishment between multiple nodes (more likely at scale). This hang can be avoided by setting NCCL CUMEM ENABLE=0.

Updating the GPG Repository Key

Chapter 17. NCCL Release 2.18.5

This is the NCCL 2.18.5 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.18.5 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 11.0, CUDA 12.0, and CUDA 12.2.

Fixed Issues

The following issues have been resolved in NCCL 2.18.5:

- Fixed NVLS search issues.
- Increased Max IB network interfaces to 32.
- Fixed inconsistent network device ordering when creating communicators with only one GPU per node.
- Try to have different GPUs use all network interfaces on systems with more than one network interface per GPU.

Known Issues

Send/receive communication using CUDA VISIBLE DEVICES and PXN only works if the GPU mappings to local ranks is the same across nodes. Disabing PXN for Send/ Receive communication can workaround the issue (NCCL P2P PXN LEVEL=0).

Updating the GPG Repository Key

Chapter 18. NCCL Release 2.18.3

This is the NCCL 2.18.3 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.18.3 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- This NCCL release supports CUDA 11.0, CUDA 12.0, CUDA 12.1, and CUDA 12.2.

Fixed Issues

The following issues have been resolved in NCCL 2.18.3:

- Fixed data corruption on DGX/HGX H100 systems when using LL128 protocol.
- Fixed hang with IB SHARP and bfloat 16 on systems with less than one NIC per GPU.
- Fixed regression in initialization time.
- Fixed data corruption with IB SHARP on H100 platforms when combining multiple GPUs per process and multiple processes per node.
- Fixed crash when shared memory creation fails.
- Fixed Avg operation with IB SHARP when using Collnet/Chain algorithm.
- Fixed performance for all-to-all operations at large scale on systems with more than one NIC per GPU.
- Fixed performance on DGX H800.
- Fixed race condition in connection progress that caused a crash.
- Fixed network flush with IB SHARP.
- ▶ Fixed PXN operation when CUDA VISIBLE DEVICES is set.
- Fixed performance of aggregated reduceScatter/allGather operations.

Known Issues

▶ Send/receive communication using CUDA VISIBLE DEVICES and PXN only works if the GPU mappings to local ranks is the same across nodes. Disabing PXN for Send/ Receive communication can workaround the issue (NCCL P2P PXN LEVEL=0).

Updating the GPG Repository Key

Chapter 19. NCCL Release 2.18.1

This is the NCCL 2.18.1 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.18.1 has been tested with the following:

- Deep learning framework containers. Refer to the Support Matrix for the supported container version.
- ▶ This NCCL release supports CUDA 11.0, CUDA 12.0, and CUDA 12.1.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Add inter-node algorithms for NVLink SHARP: NVLink SHARP + IB SHARP (NVLS), NVLink SHARP + Tree (NVLSTREE).
- Add ncclCommSplit primitive, with optional resource sharing.
- Add support for memory management using cuMem functions (disabled by default).
- Add option to use multiple QPs in round-robin mode.

Known Issues

On systems where a NIC shares a PCI switch with only one GPU (like on HGX H100), the Tree algorithm will make data transit through the CPU, making the LL128 protocol unsafe. This could result in data corruption. You can workaround this issue by setting the following:

NCCL IB PCI RELAXED ORDERING=0

Another solution is to disable the LL128 protocol with the following:

NCCL PROTO=^LL128

- On systems with less than 1 NIC per GPU, running bfloat16 reductions with IB SHARP will cause a hang.
- Running allreduce on H100 platforms with IB SHARP with multiple GPUs per process can result in data corruption if all GPUs on the node are not part of the same process, or cannot access other GPUs buffers directly.

Fixed Issues

The following issues have been resolved in NCCL 2.18.1:

- Fixed hangs with irregular send/receive patterns (e.g., alltoallv).
- Use all NICs for Send/Receive operations on systems with more than one NIC per GPU.
- Increased number of channels on H100 for network communication when bandwidth is not limited by NVLink bandwidth.
- Improved error reporting in case of IB Verbs errors.
- Fixed context creation for progress thread.
- Fixed hang in commReclaim.
- Fixed performance issue when NVB was disabled.

Updating the GPG Repository Key

Chapter 20. NCCL Release 2.17.1

This is the NCCL 2.17.1 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.17.1 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 11.0, CUDA 12.0, and CUDA 12.1.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Add support for NVLink SHARP Reduction / Broadcast to accelerate intra-node allreduce operations.
- Add new fields in the communicator configuration structure: Cooperative Group Array cluster size, minimum and maximum number of CTAs, network plugin name to use.
- Update NVTX3 includes.

Known Issues

On systems where a NIC shares a PCI switch with only one GPU (like on HGX H100), the Tree algorithm will make data transit through the CPU, making the LL128 protocol unsafe. This could result in data corruption. You can workaround this issue by setting the following:

```
NCCL_IB_PCI_RELAXED_ORDERING=0
```

Another solution is to disable the LL128 protocol with the following:

NCCL PROTO=^LL128

Performance is sub-optimal for NCCL communicators using 1 NIC per node on DGX H100/HGX H100 + NDR. It can be worked around by setting the following parameter:

```
NCCL MIN NCHANNELS=4
```

Fixed Issues

The following issues have been resolved in NCCL 2.17.1:

- Fix crash when one CollNet (SHARP) rail fails to initialize.
- Re-enable the LL128 protocol on H100 when we use PXN to close rings.

Updating the GPG Repository Key

Chapter 21. NCCL Release 2.16.5

This is the NCCL 2.16.5 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.16.5 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 11.0, CUDA 11.8, and CUDA 12.0..

Known Issues

Performance is sub-optimal for NCCL communicators using 1 NIC per node on DGX H100/HGX H100 + NDR. It can be worked around by setting the following parameter: NCCL MIN NCHANNELS=4

Fixed Issues

The following issues have been resolved in NCCL 2.16.5:

- Fix speed of IB NDR links
- Fix handling of EINTR in socket polling
- Improve proxy progress scheduling
- Fix resource cleanup
- Fix double free on case of init failure
- Fix crash in case of communicator abort
- Revert performance optimization causing significant regressions on some AMD platforms

Updating the GPG Repository Key

To best ensure the security and reliability of our RPM and Debian package repositories, NVIDIA is updating and rotating the signing keys used by apt, dnf/yum, and zypper package managers beginning on April 27, 2022. Failure to update your repository signing

keys will result in package management errors when attempting to access or install NCCL packages. To ensure continued access to the latest NCCL release, please follow the updated NCCL installation guide.

Chapter 22. NCCL Release 2.16.2

This is the NCCL 2.16.2 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.16.2 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 11.0, CUDA 11.8, and CUDA 12.0..

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Add support for CUDA 12.0
- Make socket support more resistant to network scanners
- ▶ Improve performance on large CUDA graphs, reducing dependencies
- Compile with profiling API by default
- Extend NVTX instrumentation with call arguments

Fixed Issues

The following issues have been resolved in NCCL 2.16.2:

- ▶ Various fixes to ncclCommAbort
- Make service thread polling resistant to EINTR
- Adjust inter-socket AMD bandwidth model to favor faster paths

Updating the GPG Repository Key

To best ensure the security and reliability of our RPM and Debian package repositories, NVIDIA is updating and rotating the signing keys used by apt, dnf/yum, and zypper package managers beginning on April 27, 2022. Failure to update your repository signing keys will result in package management errors when attempting to access or install

NCCL packages. To ensure continued access to the latest NCCL release, please follow the updated NCCL installation guide.

Chapter 23. NCCL Release 2.15.5

This is the NCCL 2.15.5 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.15.5 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 10.2, CUDA 11.0, and CUDA 11.8.

Fixed Issues

The following issues have been resolved in NCCL 2.15.5:

- Fix crash with CollnetChain on some node topologies
- Fix hang when interleaving the capture of different graphs
- Fix hang during init in multi-threaded mode
- Fix potential data corruption with LL128 protocol on unaligned buffers
- Fix CPU usage during preconnect
- Fixes double-free in the error path for ncclCommInitAll
- Workaround hang on H100 with Ring/LL128 on 2 GPUs

Updating the GPG Repository Key

Chapter 24. NCCL Release 2.15.1

This is the NCCL 2.15.1 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.15.1 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 10.2, CUDA 11.0, and CUDA 11.8.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

Support for H100 (sm90).

Fixed Issues

The following issues have been resolved in NCCL 2.15.1:

Make sure NCCL kernel honors user stream priorities.

Updating the GPG Repository Key

Chapter 25. NCCL Release 2.14.3

This is the NCCL 2.14.3 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.14.3 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 10.2, CUDA 11.0, and CUDA 11.7.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Add support for improved fault tolerance: non-blocking mode, new init function with configuration, new finalize function.
- Reintroduce the collnet+chain algorithm.
- Add LL protocol for intra-node send/recv communication, and inter-node (disabled by default).
- Communicate through the network within a node instead of shared memory if performance would be better that way.

Fixed Issues

The following issues have been resolved in NCCL 2.14.3:

- Wait for CUDA graph destruction before freeing communicator.
- Remove aggressive polling during enqueue.
- Fix DMABUF fallback on MOFED 5.4 and earlier.
- ► Fix NCCL DEBUG FILE functionality.

Updating the GPG Repository Key

To best ensure the security and reliability of our RPM and Debian package repositories, NVIDIA is updating and rotating the signing keys used by apt, dnf/yum, and zypper

package managers beginning on April 27, 2022. Failure to update your repository signing keys will result in package management errors when attempting to access or install NCCL packages. To ensure continued access to the latest NCCL release, please follow the updated NCCL installation guide.

Chapter 26. NCCL Release 2.13.4

This is the NCCL 2.13.4 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.13.4 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 10.2, CUDA 11.0, and CUDA 11.7.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Optimize CUDA graph launch; avoid launching a CPU callback for intra-node operations.
- Simplify kernel common code to improve the latency of send/recv operations.
- Strengthen CUDA streams semantics.
- Change NET API to v6, to add dmabuf support.
- ▶ Add ncclGetLastError() function.
- ▶ Add ncclRemoteError code and use it for remote network errors.
- ▶ Support the use of a different NCCL NET parameter per communicator.
- Add support for SHM and P2P transfers using cudaMemcpy.

Fixed Issues

The following issues have been resolved in NCCL 2.13.4:

- Fix multi-receive size encoding which could cause flush to be skipped in corner cases mixing zero-bytes send/receive operations and non-zero-bytes send/receive operations.
- Replace busy polling in the bootstrap thread waiting for ranks to check in by a blocking accept.

Updating the GPG Repository Key

Chapter 27. NCCL Release 2.12.12

This is the NCCL 2.12.12 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.12.12 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- This NCCL release supports CUDA 10.2, CUDA 11.0, CUDA 11.6, and CUDA 11.7.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Improved allreduce performance when you have more than one network interface per GPU and you need to use PXN to close rings.
- ▶ Added P2P DIRECT DISABLE parameter to disable direct access to pointers within a process.
- Added support for PCI Gen5 on 5.4 kernels.

Fixed Issues

The following issues have been resolved in NCCL 2.12.12:

- Fixed hang on cubemesh topologies.
- Fixed random crash in init due to uninitialized struct.
- ► Fixed crash when setting NCCL SET THREAD NAME.

Updating the GPG Repository Key

To best ensure the security and reliability of our RPM and Debian package repositories, NVIDIA is updating and rotating the signing keys used by apt, dnf/yum, and zypper package managers beginning on April 27, 2022. Failure to update your repository signing keys will result in package management errors when attempting to access or install

NCCL packages. To ensure continued access to the latest NCCL release, please follow the updated NCCL installation guide.

Chapter 28. NCCL Release 2.12.10

This is the NCCL 2.12.10 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.12.10 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ► This NCCL release supports <u>CUDA 10.2</u>, <u>CUDA 11.0</u>, and <u>CUDA 11.6</u>.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

Improved error messages for network errors, with hostname instead of IP address.

Fixed Issues

The following issues have been resolved in NCCL 2.12.10:

- Fixed bug with collNet.
- Fixed bug with zero-byte send/recv operations.
- ► Fixed NCCL PARAM implementation taking a lock every time a parameter value was queried.
- ► Fixed bug when setting NCCL IB QPS PER CONNECTION to more than one.

Chapter 29. NCCL Release 2.12.7

This is the NCCL 2.12.7 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.12.7 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 10.2, CUDA 11.0, and CUDA 11.6.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added NVLink-optimized network communication to keep traffic rail-local (PXN).
- Improved alltoall latency by aggregating messages within a node to a given destination.
- Added new v5 plugin API with grouped receives and tags, keeping compatibility for v4 plugins.
- Added naming of NCCL threads to help debugging.
- Added support for Relaxed Ordering for IB.
- Added profiling and timing infrastructure.

Fixed Issues

The following issues have been resolved in NCCL 2.12.7:

Fixed NVLink detection and avoid data corruption when some NVLinks are down.

Chapter 30. NCCL Release 2.11.4

This is the NCCL 2.11.4 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.11.4 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ► This NCCL release supports CUDA 10.2, CUDA 11.0, CUDA 11.4, CUDA 11.5, and CUDA 11.6.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added new API for creating a reduction operation which multiplies the input by a rank-specific scalar before doing an inter-rank summation (see: ncclRedOpCreatePreMulSum).
- Improved CollNet (SHARP) performance of ncclAllReduce when captured in a CUDA Graph via user buffer registration.
- Added env NCCL NET PLUGIN="<suffix>" to allow the user a way to choose among multiple NCCL net plugins by substituting into libnccl-net-<suffix>.so.

Fixed Issues

The following issues have been resolved in NCCL 2.11.4:

- Fixed memory leak of NVB connections.
- Fixed crash of ncclGroup() containing mixed datatypes/operations (GitHub issue #560, introduced in NCCL 2.10.3).
- Fixed topology detection of IB Virtual Functions (SR-IOV).

Chapter 31. NCCL Release 2.10.3

This is the NCCL 2.10.3 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.10.3 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 10.2, CUDA 11.0, and CUDA 11.4.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added ncclAvg operation
- Added support for bfloat 16 type
- Added support for multiple IB queue pairs
- Added WSL2 support (single GPU only)
- Improved performance for aggregation
- Improved performance for medium sizes
- Improved network error reporting
- Added NCCL_NET environment variable to use a specific network
- Added auto-load of XML topology from default location

Fixed Issues

The following issues have been resolved in NCCL 2.10.3:

- Fixed graph search on cubemesh topologies.
- Fixed all-to-all affinity to improve latency.
- Fixed hang in cubemesh NVB connections during initialization.

Chapter 32. NCCL Release 2.9.9

This is the NCCL 2.9.9 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.9.9 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ► This NCCL release supports CUDA 10.1, CUDA 10.2, CUDA 11.0, and CUDA 11.3.

Fixed Issues

The following issues have been resolved in NCCL 2.9.9:

- Fixed crash when setting NCCL MAX P2P NCHANNELS.
- Fixed hang during sendrecv dynamic connection on cubemesh topologies.
- ► Fixed compilation with TRACE=1.

Chapter 33. NCCL Release 2.9.8

This is the NCCL 2.9.8 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.9.8 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ► This NCCL release supports CUDA 10.1, CUDA 10.2, CUDA 11.0, and CUDA 11.3.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

Added support for nvidia-peermem module.

Fixed Issues

The following issues have been resolved in NCCL 2.9.8:

- Fixed error when injecting a topology file without vendor/device information.
- Fixed crash in bootstrap error case.
- Fixed memory leaks.
- Fixed collnet clean-up issue.

Chapter 34. NCCL Release 2.9.6

This is the NCCL 2.9.6 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.9.6 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ► This NCCL release supports CUDA 10.1, CUDA 10.2, CUDA 11.0, and CUDA 11.3.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added support for CUDA graphs
- Improved performance for CollNet (SHARP)
- ► Fuse PCI Gen4 switches showing a two-level hierarchy into a single level.
- Improve NIC balancing for communicators using a single GPU per node.

Fixed Issues

The following issues have been resolved in NCCL 2.9.6:

- Fixed bootstrap hang in case of reordered packets causing connections to be inverted.
- Fix locking issue causing NCCL calls to block until previous operations were complete.

Chapter 35. NCCL Release 2.8.4

This is the NCCL 2.8.4 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.8.4 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ► This NCCL release supports CUDA 10.1, CUDA 10.2, CUDA 11.0, CUDA 11.1, and CUDA 11.2.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

Added support for Zhaoxin CPUs

Known Issues

Send/receive operations have a number of limitations:

Using send/receive operations in combination to launch work on multiple GPUs from a single process can fail or hang if the GPUs process different amounts of data. Setting NCCL LAUNCH MODE=PARALLEL can work around the issue, but can also cause other problems. For more information, see the NCCL User Guide section Troubleshooting > Known Issues > Concurrency Between NCCL and CUDA calls.

Fixed Issues

The following issues have been resolved in NCCL 2.8.4:

Fixed hang for some imbalanced send/recv operation (alltoally).

Chapter 36. NCCL Release 2.8.3

This is the NCCL 2.8.3 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.8.3 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- This NCCL release supports CUDA 10.1, CUDA 10.2, CUDA 11.0, CUDA 11.1, and CUDA 11.2.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Optimized Tree performance on A100
- Improved performance for aggregated operations
- Improved performance for all-to-all operations at scale
- Reduced memory usage for all-to-all operations at scale
- Optimized all-to-all performance on DGX-1

Known Issues

Send/receive operations have a number of limitations:

Using send/receive operations in combination to launch work on multiple GPUs from a single process can fail or hang if the GPUs process different amounts of data. Setting NCCL LAUNCH MODE=PARALLEL can work around the issue, but can also cause other problems. For more information, see the NCCL User Guide section Troubleshooting > Known Issues > Concurrency Between NCCL and CUDA calls.

Fixed Issues

The following issues have been resolved in NCCL 2.8.3:

Hang in LL128 protocol after 2^31 steps.

- ▶ Topology injection error when using fewer GPUs than described. (github issue #379)
- Protocol mismatch causing hangs or crashes when using one GPU per node. (github issue #394)

Chapter 37. NCCL Release 2.7.8

This is the NCCL 2.7.8 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.7.8 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- This NCCL release supports CUDA 10.1, CUDA 10.2, CUDA 11.0, amd CUDA 11.1.

Known Issues

Send/receive operations have a number of limitations:

- Using send/receive operations in combination to launch work on multiple GPUs from a single process can fail or hang if the GPUs process different amounts of data. Setting NCCL LAUNCH MODE=PARALLEL can work around the issue, but can also cause other problems. For more information, see the NCCL User Guide section Troubleshooting > Known Issues > Concurrency Between NCCL and CUDA calls.
- Aggregation is not supported on a given source-destination pair, meaning that there can only be one send/receive per source-destination pair.
- Each source-destination pair allocates a dedicated fifo of 4 MB (see NCCL BUFFSIZE variable) which can use a significant amount of GPU memory at scale.
- When using GPU Direct RDMA, each point-to-point connection will also use resources on the GPU PCI address space, which we might run out of on some models. If that happens, consider disabling GPU Direct RDMA (NCCL NET GDR LEVEL=0) or reduce the per-peer buffer size (NCCL BUFFSIZE).
- Send/receive operations are not yet optimized for the DGX-1 NVLink topology.

Fixed Issues

The following issues have been resolved in NCCL 2.7.8:

operations.

Chapter 38. NCCL Release 2.7.6

This is the NCCL 2.7.6 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.7.6 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ▶ This NCCL release supports CUDA 10.1, CUDA 10.2, and CUDA 11.0.

Known Issues

Send/receive operations have a number of limitations:

- Using send/receive operations in combination to launch work on multiple GPUs from a single process can fail or hang if the GPUs process different amounts of data. Setting NCCL LAUNCH MODE=PARALLEL can work around the issue, but can also cause other problems. For more information, see the NCCL User Guide section Troubleshooting > Known Issues > Concurrency Between NCCL and CUDA calls.
- Aggregation is not supported on a given source-destination pair, meaning that there can only be one send/receive per source-destination pair.
- ► Each source-destination pair allocates a dedicated fifo of 4 MB (see NCCL BUFFSIZE variable) which can use a significant amount of GPU memory at scale.
- When using GPU Direct RDMA, each point-to-point connection will also use resources on the GPU PCI address space, which we might run out of on some models. If that happens, consider disabling GPU Direct RDMA (NCCL NET GDR LEVEL=0) or reduce the per-peer buffer size (NCCL BUFFSIZE).
- Send/receive operations are not yet optimized for the DGX-1 NVLink topology.

Fixed Issues

The following issues have been resolved in NCCL 2.7.6:

Fixed crash when NVswitch is not visible inside a virtual machine.

Chapter 39. NCCL Release 2.7.5

This is the NCCL 2.7.5 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Compatibility

NCCL 2.7.5 has been tested with the following:

- Deep learning framework containers. Refer to the Support Matrix for the supported container version.
- ▶ This NCCL release supports CUDA 10.1, CUDA 10.2, and CUDA 11.0.

Known Issues

Send/receive operations have a number of limitations:

- Using send/receive operations in combination to launch work on multiple GPUs from a single process can fail or hang if the GPUs process different amounts of data. Setting NCCL LAUNCH MODE=PARALLEL can work around the issue, but can also cause other problems. For more information, see the NCCL User Guide section Troubleshooting > Known Issues > Concurrency Between NCCL and CUDA calls.
- Aggregation is not supported on a given source-destination pair, meaning that there can only be one send/receive per source-destination pair.
- ► Each source-destination pair allocates a dedicated fifo of 4 MB (see NCCL BUFFSIZE variable) which can use a significant amount of GPU memory at scale.
- When using GPU Direct RDMA, each point-to-point connection will also use resources on the GPU PCI address space, which we might run out of on some models. If that happens, consider disabling GPU Direct RDMA (NCCL NET GDR LEVEL=0) or reduce the per-peer buffer size (NCCL BUFFSIZE).
- Send/receive operations are not yet optimized for the DGX-1 NVLink topology.
- Running inside a virtual machine on an NVswitch platform can cause a crash if NVswitch is not visible inside the virtual machine.

Fixed Issues

The following issues have been resolved in NCCL 2.7.5:

- ▶ Minor fixes for A100 platforms.
- ▶ Add proper message for invalid GroupEnd call.

Chapter 40. NCCL Release 2.7.3

This is the NCCL 2.7.3 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added support for A100 GPU and related platforms
- Added support for CUDA 11
- Added support for send/receive operations (beta)

Compatibility

NCCL 2.7.3 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- ► This NCCL release supports <u>CUDA 10.1</u>, <u>CUDA 10.2</u>, and <u>CUDA 11.0</u>.

Known Issues

Send/receive operations have a number of limitations:

- Using send/receive operations in combination to launch work on multiple GPUs from a single process can fail or hang if the GPUs process different amounts of data. Setting NCCL LAUNCH MODE=PARALLEL can work around the issue, but can also cause other problems. For more information, see the NCCL User Guide section Troubleshooting > Known Issues > Concurrency Between NCCL and CUDA calls.
- Aggregation is not supported on a given source-destination pair, meaning that there can only be one send/receive per source-destination pair.
- Each source-destination pair allocates a dedicated fifo of 4 MB (see NCCL BUFFSIZE variable) which can use a significant amount of GPU memory at scale.
- When using GPU Direct RDMA, each point-to-point connection will also use resources on the GPU PCI address space, which we might run out of on some models. If that

happens, consider disabling GPU Direct RDMA (NCCL_NET_GDR_LEVEL=0) or reduce the per-peer buffer size (NCCL BUFFSIZE).

▶ Send/receive operations are not yet optimized for the DGX-1 NVLink topology.

Fixed Issues

The following issues have been resolved in NCCL 2.7.3:

Fixed crash when only a subset of GPUs are visible within a container (#326).

Chapter 41. NCCL Release 2.6.4

This is the NCCL 2.6.4 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added support for in-network collectives (beta)
- Added support for adaptive routing on Infiniband
- Added support for topology dump/injection through XML files
- Added speed detection for PCI, Infiniband and Ethernet cards
- Added CPU detection for AMD CPUs and ARM

Chapter 42. NCCL Release 2.5.6

This is the NCCL 2.5.6 release notes. For previous NCCL release notes, refer to the NCCL Archives.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added new protocol to improve performance on small operations.
- Improved topology detection and tree/ring creation (#179, #262).
- Improved multi-node tree performance by sending/receiving from different GPUs.
- Added model-based tuning to switch between the different algorithms and protocols.
- Reworked P2P/SHM detection in containers (#155, #248).
- Added detection for duplicate CUDA devices and return an error (#231).
- Added tuning for Google Cloud's gVNIC platform.

Compatibility

NCCL 2.5.6 has been tested with the following:

- ▶ Deep learning framework containers. Refer to the <u>Support Matrix</u> for the supported container version.
- This NCCL release supports <u>CUDA 9.0</u>, <u>CUDA 10.0</u>, <u>CUDA 10.1</u>, and <u>CUDA 10.2</u>.

Fixed Issues

The following issues have been resolved in NCCL 2.5.6:

- Sporadic NCCL error "ring 0 does not loop back to start" (#179).
- NCCL doesn't form proper rings on GCP V100s (#262).

Chapter 43. NCCL Release 2.4.8

This is the NCCL 2.4.8 release notes. This release includes fixes from the previous NCCL 2.4.x releases as well as the following additional changes. For previous NCCL release notes, see the archived NCCL Release Notes.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

Improved socket transport performance by splitting transfer over multiple sockets.



Note: This feature adds two new environment variables NCCL SOCKET NTHREADS and NCCL NSOCKS PERTHREAD for users to tune NCCL performance on socket-based networks. See the NCCL documentation for more details.

Compatibility

NCCL 2.4.8 has been tested with the following:

- Deep learning framework 19.05 containers
- This NCCL release supports; CUDA 9.0, CUDA 9.2, CUDA 10.0, and CUDA 10.1.

Fixed Issues

The following issues have been resolved in NCCL 2.4.8:

Suboptimal performance with TCP over high bandwidth networks. (GitHub issue #209)

- On single node Power systems with 4 GPUs, some performance regressions have been observed compared to NCCL 2.4.2. These will be addressed in future NCCL releases.
- By default, NCCL does not enable direct P2P communication through different PCIe root ports on Intel Skylake CPU and later. This is due to a known performance issue when using P2P on these CPU versions. There is now a new BIOS and performance

tuning option available (PCIe Peer-to-Peer Serialization) from Intel and their OEM vendors that resolves this P2P bandwidth issue. If the BIOS performance tuning option has been enabled, then NCCL direct P2P connections can be re-enabled by setting NCCL_P2P_LEVEL=5.

Chapter 44. NCCL Release 2.4.7

This is the NCCL 2.4.7 release notes. This release includes fixes from the previous NCCL 2.4.x releases as well as the following additional changes. For previous NCCL release notes, see the archived NCCL Release Notes.

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Improved bootstrap socket connection reliability at scale.
- Added detection of IBM/Power NVLink bridge device.
- Added NUMA support to PCI-E distance calculations on x86 architectures.



Note: This adds a new level (5) for the NCCL P2P LEVEL and NCCL NET GDR LEVEL environment variables. See the NCCL documentation for more details.

Added the NCCL IGNORE CPU AFFINITY environment variable.

Compatibility

NCCL 2.4.7 has been tested with the following:

- Deep learning framework 19.04 containers
- This NCCL release supports; CUDA 9.0, CUDA 9.2, CUDA 10.0, and CUDA 10.1.

Fixed Issues

The following issues have been resolved in NCCL 2.4.7:

- Fixed hostname hashing issue. (GitHub issue #187)
- Fixed memory leaks. (GitHub issue #180)
- Fixed compiler warning. (GitHub issue #178)
- Replaced non-standard variable length arrays. (GitHub issue #171)
- ► Fixed Tree and Shared Memory crash. (GitHub PR #185)
- Fixed hangs during long running jobs.
- ► Fixed the NCCL RINGS environment variable handling.

▶ Added extra checks to catch duplicate calls to ncclCommDestroy(). (GitHub issue #191)

- On single node Power systems with 4 GPUs, some performance regressions have been observed compared to NCCL 2.4.2. These will be addressed in future NCCL releases.
- ▶ By default, NCCL does not enable direct P2P communication through different PCle root ports on Intel Skylake CPU and later. This is due to a known performance issue when using P2P on these CPU versions. There is now a new BIOS and performance tuning option available (PCIe Peer-to-Peer Serialization) from Intel and their OEM vendors that resolves this P2P bandwidth issue. If the BIOS performance tuning option has been enabled, then NCCL direct P2P connections can be re-enabled by setting NCCL P2P LEVEL=5.

Chapter 45. NCCL Release 2.4.2

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Implemented tree-based algorithms for better All Reduce performance at scale and with small and medium size messages.
- Support for external network plugins (e.g., libfabric).
- ► Add ncclCommGetAsyncError() function to report errors happening during collective operations.
- Add ncclCommAbort () function to destroy a communicator, aborting any outstanding operations.
- ▶ Support different ranks having a different CUDA VISIBLE DEVICES.
- Add a best-effort mechanism to check for size mismatch among collective calls.

- Support communication between Mesos containers (Github issue #155).
- ► Fix case where posix fallocate() returns EINTR (Github issue #137).
- NCCL threads no longer escape the CPU affinity set by the user or job scheduler.

Chapter 46. NCCL Release 2.3.7

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Minor tuning of the LL threshold for multi-node jobs.
- Improve performance of initialization on large number of ranks.

Fixed Issues

Fixed issue causing "WARN: Message truncated" errors.

Chapter 47. NCCL Release 2.3.5

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

▶ This release is open-sourced with no new features or fixes.

Chapter 48. NCCL Release 2.3.4

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Improve performance tuning on large number of ranks.
- ▶ Add NCCL P2P LEVEL and NCCL IB GDR LEVEL knobs to finely control when to use GPU Direct P2P and GPU Direct RDMA.
- Reduce setup time for large scale jobs.
- Increased maximum number of rings supported to 16.
- Added a runtime NCCL version API: ncclGetVersion().
- ▶ Added NCCL_DEBUG_SUBSYS to allow filtering of NCCL DEBUG=INFO logging from different subsystems.
- Support for Turing based systems.

- Fix hang on Power platforms.
- Fix low inter-node bandwidth issue on multi-DGX2 systems.
- Fix crash when used with PID isolator.

Chapter 49. NCCL Release 2.2.13

Key Features and Enhancements

There were no new features and enhancements in this release.

Using NCCL 2.2.13

Ensure you are familiar with the following notes when using this release.

- ▶ If NCCL returns an error code, set the environment variable NCCL DEBUG to WARN to receive an explicit error message.
- ▶ The NCCL 2.x API is different from NCCL 1.x. Some porting may be needed for NCCL 1.x applications to work correctly. Refer to the migration documentation in the NCCL Developer Guide.
- Starting in 2.2, NCCL supports collective aggregation. You can put multiple NCCL collective operations in between ncclGroupStart() and ncclGroupEnd() to enable this feature.

Known Issues

- Using multiple processes in conjunction with multiple threads to manage the different GPUs may in some cases cause ncclCommInitRank to fail while establishing IPCs (cudaIpcOpenMemHandle). This problem does not appear when using only processes or only threads. This issue is fixed in recent driver versions, therefore, consider updating to the latest driver.
- NCCL uses CUDA 9 cooperative group launch by default, which may induce increased latencies in multi-threaded programs. See the NCCL LAUNCH MODE knob to restore the original behavior.
- Driver version 390 can cause data corruption when used together with GPU Direct RDMA. Disabling GPU Direct RDMA by setting NCCL IB CUDA SUPPORT=0, or upgrading to 396.26 or newer driver should resolve the issue.

Fixed Issues

► Fix crash in child processes after calling ncclCommDestroy.

Chapter 50. NCCL Release 2.2.12

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added support for collective operations aggregation.
- Added ncclBroadcast function.

Using NCCL 2.2.12

Ensure you are familiar with the following notes when using this release.

- ▶ If NCCL returns an error code, set the environment variable NCCL DEBUG to WARN to receive an explicit error message.
- ▶ The NCCL 2.x API is different from NCCL 1.x. Some porting may be needed for NCCL 1.x applications to work correctly. Refer to the migration documentation in the NCCL Developer Guide.
- Starting in 2.2, NCCL supports collective aggregation. You can put multiple NCCL collective operations in between ncclGroupStart() and ncclGroupEnd() to enable this feature.

- Using multiple processes in conjunction with multiple threads to manage the different GPUs may in some cases cause ncclCommInitRank to fail while establishing IPCs (cudaIpcOpenMemHandle). This problem does not appear when using only processes or only threads. This issue is fixed in recent driver versions, therefore, consider updating to the latest driver.
- NCCL uses CUDA 9 cooperative group launch by default, which may induce increased latencies in multi-threaded programs. See the NCCL LAUNCH MODE knob to restore the original behavior.
- Driver version 390 and later can cause data corruption when used together with GPU Direct RDMA. Disabling GPU Direct RDMA by setting NCCL IB CUDA SUPPORT=0 or reverting to driver 387 should resolve the issue.

- ▶ No longer clear the CPU affinity during initialization functions
- Fix various large scale issues
- Reduce the size of the library
- Fix crash or hang with PyTorch relative to the usage of calls to fork

Chapter 51. NCCL Release 2.1.15

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

Added a variable to control InfiniBand Traffic Class and Retry Count.

Using NCCL 2.1.15

Ensure you are familiar with the following notes when using this release.

- ▶ The NCCL 2.x API is different from NCCL 1.x. Some porting may be needed for NCCL 1.x applications to work correctly. Refer to the migration documentation in the NCCL Developer Guide.
- Starting in 2.2, NCCL supports collective aggregation. You can put multiple collectives in between ncclGroupStart() and ncclGroupEnd() to enable this feature.

Known Issues

- ▶ If NCCL returns an error code, set the environment variable NCCL DEBUG to WARN to receive an explicit error message.
- Using multiple processes in conjunction with multiple threads to manage the different GPUs may in some cases cause ncclCommInitRank to fail while establishing IPCs (cudaIpcOpenMemHandle). This problem does not appear when using only processes or only threads.
- NCCL uses CUDA 9 cooperative group launch by default, which may induce increased latencies in multi-threaded programs. See the NCCL LAUNCH MODE knob to restore the original behavior.

- Fixed CPU usage and scheduling of NCCL network threads.
- Fixed CUDA launch crash when mixing different types of GPUs in a node.
- Fixed a performance problem on Skylake CPUs.
- Fixed hanging issues with cudaFree and inter-node communication.

- ▶ Restored library installation path to /usr/lib/x86_64-linux-gnu in debian packages.
- Fixed RoCEv2 failure when using a non-zero GID.
- ▶ No longer link to stdc++ library statically as this can cause issues with C++ applications.
- Fixed PyTorch hanging issues when using multiple rings and many back-to-back broadcast operations.

Chapter 52. NCCL Release 2.1.4

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- Added support for InfiniBand GID selection, enabling the use of RoCE v2.
- Added support for InfiniBand Service Level (SL) selection.

Using NCCL 2.1.4

Ensure you are familiar with the following notes when using this release.

The NCCL 2.x API is different from NCCL 1.x. Some porting may be needed for NCCL 1.x applications to work correctly. Refer to the migration documentation in the NCCL Developer Guide.

Known Issues

- ▶ If NCCL returns an error code, set the environment variable NCCL DEBUG to WARN to receive an explicit error message.
- Using multiple processes in conjunction with multiple threads to manage the different GPUs may in some cases cause ncclCommInitRank to fail while establishing IPCs (cudalpcopenMemHandle). This problem does not appear when using only processes or only threads.
- ▶ NCCL uses CUDA® 9 cooperative group launch by default, which may induce increased latencies in multi-threaded programs. See the NCCL LAUNCH MODE knob in the NCCL Developer Guide to restore the original behavior.
- ▶ NCCL 2.1.4-1 embeds libstdc++ and exports its symbols. This can break C++ applications.

- Fixed bug causing CUDA IPC to fail in some situations.
- Fixed bug causing a crash when p2p mappings are exhausted instead of returning an error.

Chapter 53. NCCL Release 2.1.2

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- New algorithms for improved latency communication
- RoCE support

Using NCCL 2.1.2

Ensure you are familiar with the following notes when using this release.

The NCCL 2.x API is different from NCCL 1.x. Some porting may be needed for NCCL 1.x applications to work correctly. Refer to the migration documentation in the NCCL Developer Guide.

Known Issues

- ▶ If NCCL returns an error code, set the environment variable NCCL DEBUG to WARN to receive an explicit error message.
- Using multiple processes in conjunction with multiple threads to manage the different GPUs may in some cases cause ncclCommInitRank to fail while establishing IPCs (cudaIpcOpenMemHandle). This problem does not appear when using only processes or only threads.
- NCCL uses CUDA 9 cooperative group launch by default, which may induce increased latencies in multi-threaded programs. See the NCCL LAUNCH MODE knob in the NCCL <u>Developer Guide</u> to restore the original behavior.

Fixed Issues

 NCCL now uses CUDA 9 cooperative groups to launch the CUDA kernels, fixing a long-standing issue with CUDA and NCCL operations being interleaved and potentially causing hangs.

Chapter 54. NCCL Release 2.0.5

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- NCCL 2.0.5 provides support for intra-node and inter-node communication.
- NCCL optimizes intra-node communication using NVLink, PCI express, and shared memory.
- Between nodes, NCCL implements fast transfers over sockets or InfiniBand verbs.
- GPU-to-GPU and GPU-to-Network direct transfers, using the GPU Direct technology, are extensively used when the hardware topology permits it.

Using NCCL 2.0.5

Ensure you are familiar with the following notes when using this release.

- ▶ The NCCL 2.0 API is different from NCCL 1.x. Some porting may be needed for NCCL 1.x applications to work correctly. Refer to the migration documentation in the NCCL Developer Guide.
- NCCL 2.0.5 has the new configuration file support. The NCCL environment variables can now be set in ~/.nccl.conf and /etc/nccl.conf.
- ▶ Values defined in ~/.nccl.conf take precedence over ones in /etc/nccl.conf.
- The syntax for each line of the NCCL configuration file is <NCCL VAR NAME>=<VALUE>.

- ▶ If NCCL returns any error code, set the environment variable NCCL DEBUG to WARN to receive an explicit error message.
- RoCE support is experimental.
- Using multiple processes in conjunction with multiple threads to manage the different GPUs may in some cases cause ncclCommInitRank to fail while establishing IPCs (cudaIpcOpenMemHandle). This problem does not appear when using only processes or only threads.

Chapter 55. NCCL Release 2.0.4

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- NCCL 2.0.4 provides support for intra-node and inter-node communication.
- NCCL optimizes intra-node communication using NVLink, PCI express, and shared memory.
- Between nodes, NCCL implements fast transfers over sockets or InfiniBand verbs.
- GPU-to-GPU and GPU-to-Network direct transfers, using the GPU Direct technology, are extensively used when the hardware topology permits it.

Using NCCL 2.0.4

Ensure you are familiar with the following notes when using this release.

- ▶ The NCCL 2.0 API is different from NCCL 1.x. Some porting may be needed for NCCL 1.x applications to work correctly. Refer to the migration documentation in the NCCL Developer Guide.
- NCCL 2.0.4 has the new configuration file support. The NCCL environment variables can now be set in ~/.nccl.conf and /etc/nccl.conf.
- ▶ Values defined in ~/.nccl.conf take precedence over ones in /etc/nccl.conf.
- The syntax for each line of the NCCL configuration file is <NCCL VAR NAME>=<VALUE>.

- ▶ If NCCL returns any error code, set the environment variable NCCL DEBUG to WARN to receive an explicit error message.
- RoCE is not supported.
- Using multiple processes in conjunction with multiple threads to manage the different GPUs may in some cases cause ncclCommInitRank to fail while establishing IPCs (cudaIpcOpenMemHandle). This problem does not appear when using only processes or only threads.

Chapter 56. NCCL Release 2.0.2

Key Features and Enhancements

This NCCL release includes the following key features and enhancements.

- NCCL 2.0.2 provides support for intra-node and inter-node communication.
- NCCL optimizes intra-node communication using NVLink, PCI express, and shared memory.
- Between nodes, NCCL implements fast transfers over sockets or InfiniBand verbs.
- GPU-to-GPU and GPU-to-Network direct transfers, using the GPU Direct technology, is extensively used when the hardware topology permits it.

Using NCCL 2.0.2

Ensure you are familiar with the following notes when using this release.

▶ The NCCL 2.0 API is different from NCCL 1.x. Some porting may be needed for NCCL 1.x applications to work correctly. Refer to the migration documentation in the NCCL Developer Guide.

- NCCL 2.0.2 is known to not work with CUDA driver 384.40 and later.
- ▶ If NCCL returns any error code, set the environment variable NCCL DEBUG to WARN to receive an explicit error message.
- NCCL 2.0.2 does not support RoCE, that is, InfiniBand cards using Ethernet as link layer. The presence of an RoCE card on a node will make NCCL fail even when run within the node.
- Using multiple processes in conjunction with multiple threads to manage the different GPUs may in some cases cause ncclCommInitRank to fail while establishing IPCs (cudaIpcOpenMemHandle). This problem does not appear when using only processes or only threads.

Notice

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

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

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

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

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

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

Trademarks

NVIDIA, the NVIDIA logo, and cuBLAS, CUDA, CUDA Toolkit, cuDNN, DALI, DIGITS, DGX, DGX-1, DGX-2, DGX Station, DLProf, GPU, Jetson, Kepler, Maxwell, NCCL, Nsight Compute, Nsight Systems, NVCaffe, NVIDIA Deep Learning SDK, NVIDIA Developer Program, NVIDIA GPU Cloud, NVLink, NVSHMEM, PerfWorks, Pascal, SDK Manager, Tegra, TensorRT, TensorRT Inference Server, Tesla, TF-TRT, Triton Inference Server, Turing, and Volta are trademarks and/or registered trademarks of NVIDIA Corporation in the United States and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

Copyright

 $^{\hbox{\scriptsize @}}$ 2025 NVIDIA CORPORATION and affiliates. All rights reserved.