



# NVIDIA<sup>®</sup>

## HPC Compilers

USER'S GUIDE

# TABLE OF CONTENTS

Preface.....	xii
Audience Description.....	xii
Compatibility and Conformance to Standards.....	xii
Organization.....	xiii
Hardware and Software Constraints.....	xiv
Conventions.....	xiv
Terms.....	xv
Related Publications.....	xvi
Chapter 1. Getting Started.....	1
1.1. Overview.....	1
1.2. Creating an Example.....	1
1.3. Invoking the Command-level NVIDIA HPC Compilers.....	2
1.3.1. Command-line Syntax.....	2
1.3.2. Command-line Options.....	3
1.4. Filename Conventions.....	3
1.4.1. Input Files.....	4
1.4.2. Output Files.....	5
1.5. Fortran, C++ and C Data Types.....	6
1.6. Platform-specific considerations.....	7
1.6.1. Using the NVIDIA HPC Compilers on Linux.....	7
1.7. Site-Specific Customization of the Compilers.....	7
1.7.1. Use siterc Files.....	7
1.7.2. Using User rc Files.....	7
1.8. Common Development Tasks.....	8
Chapter 2. Use Command-line Options.....	10
2.1. Command-line Option Overview.....	10
2.1.1. Command-line Options Syntax.....	10
2.1.2. Command-line Suboptions.....	11
2.1.3. Command-line Conflicting Options.....	11
2.2. Help with Command-line Options.....	11
2.3. Getting Started with Performance.....	12
2.3.1. Using -fast.....	12
2.3.2. Other Performance-Related Options.....	13
2.4. Frequently-used Options.....	13
2.5. Floating-point Subnormal.....	15
Chapter 3. Multicore CPU Optimization.....	17
3.1. Overview of Optimization.....	17
3.1.1. Local Optimization.....	18
3.1.2. Global Optimization.....	18
3.1.3. Loop Optimization: Unrolling, Vectorization and Parallelization.....	18

3.1.4. Interprocedural Analysis (IPA) and Optimization.....	18
3.1.5. Function Inlining.....	19
3.2. Getting Started with Optimization.....	19
3.2.1. -help.....	20
3.2.2. -Minfo.....	20
3.2.3. -Mneginfo.....	21
3.2.4. -dryrun.....	21
3.2.5. -v.....	21
3.3. Local and Global Optimization.....	21
3.3.1. -Msafepr.....	21
3.3.2. -O.....	22
3.4. Loop Unrolling using -Munroll.....	24
3.5. Vectorization using -Mvect.....	25
3.5.1. Vectorization Sub-options.....	25
3.5.2. Vectorization Example Using SIMD Instructions.....	27
3.6. Interprocedural Analysis and Optimization using -Mipa.....	29
3.6.1. Building a Program Without IPA – Single Step.....	30
3.6.2. Building a Program Without IPA – Several Steps.....	30
3.6.3. Building a Program Without IPA Using Make.....	30
3.6.4. Building a Program with IPA.....	31
3.6.5. Building a Program with IPA – Single Step.....	31
3.6.6. Building a Program with IPA – Several Steps.....	32
3.6.7. Building a Program with IPA Using Make.....	33
3.6.8. Questions about IPA.....	33
Chapter 4. Using Function Inlining.....	35
4.1. Automatic function inlining in C++ and C.....	35
4.2. Invoking Procedure Inlining.....	36
4.3. Using an Inline Library.....	37
4.4. Creating an Inline Library.....	37
4.4.1. Working with Inline Libraries.....	38
4.4.2. Dependencies.....	39
4.4.3. Updating Inline Libraries – Makefiles.....	39
4.5. Error Detection during Inlining.....	39
4.6. Examples.....	40
4.7. Restrictions on Inlining.....	40
Chapter 5. Using GPUs.....	42
5.1. Overview.....	42
5.2. Terminology.....	43
5.3. Execution Model.....	45
5.3.1. Host Functions.....	45
5.4. Memory Model.....	45
5.4.1. Separate Host and Accelerator Memory Considerations.....	47
5.4.1.1. Accelerator Memory.....	47

5.4.1.2. Staging Memory Buffer.....	47
5.4.1.3. Cache Management.....	48
5.4.1.4. Environment Variables Controlling Device Memory Management.....	48
5.4.2. Managed and Unified Memory Modes.....	50
5.4.2.1. Managed Memory Mode.....	50
5.4.2.2. Unified Memory Mode.....	52
5.4.3. Memory Pool Allocator.....	53
5.4.4. Interception of Deallocations.....	54
5.4.5. Command-line Options Selecting Compiler Memory Modes.....	54
5.5. Fortran pointers in device code.....	55
5.6. Calling routines in a compute kernel.....	56
5.7. Supported Processors and GPUs.....	56
5.8. CUDA Versions.....	57
5.9. Compute Capability.....	58
5.10. PTX JIT Compilation.....	58
Chapter 6. Using OpenACC.....	63
6.1. OpenACC Programming Model.....	63
6.1.1. Levels of Parallelism.....	63
6.1.2. Enable OpenACC Directives.....	64
6.1.3. OpenACC Support.....	64
6.1.4. OpenACC Extensions.....	64
6.2. Compiling an OpenACC Program.....	65
6.2.1. -[no]acc.....	65
6.2.2. -gpu.....	66
6.3. OpenACC for Multicore CPUs.....	69
6.4. OpenACC with CUDA Unified Memory.....	69
6.5. OpenACC Error Handling.....	73
6.6. OpenACC and CUDA Graphs.....	76
6.7. Environment Variables.....	82
6.8. Profiling Accelerator Kernels.....	83
6.9. OpenACC Runtime Libraries.....	84
6.9.1. Runtime Library Definitions.....	85
6.9.2. Runtime Library Routines.....	85
6.10. Supported Intrinsics.....	86
6.10.1. Supported Fortran Intrinsics Summary Table.....	86
6.10.2. Supported C Intrinsics Summary Table.....	88
Chapter 7. Using OpenMP.....	90
7.1. Environment Variables.....	91
7.2. Fallback Mode.....	93
7.3. Loop.....	94
7.4. OpenMP Subset.....	97
7.5. Using metadirective.....	106
7.6. Mapping target constructs to CUDA streams.....	108

7.7. Noncontiguous Array Sections.....	111
7.8. OpenMP with CUDA Unified Memory.....	111
7.9. Multiple Device Support.....	114
7.10. Interoperability with CUDA.....	114
7.11. Interoperability with Other OpenMP Compilers.....	115
7.12. GNU STL.....	115
Chapter 8. Using Stdpar.....	116
8.1. GPU Memory Modes.....	116
8.2. Stdpar C++.....	117
8.2.1. Introduction to Stdpar C++.....	117
8.2.2. NVC++ Compiler Parallel Algorithms Support.....	118
8.2.2.1. Enabling Parallel Algorithms with the -stdpar Option.....	118
8.2.3. Stdpar C++ Simple Example.....	119
8.2.4. Coding Guidelines for GPU-accelerating Parallel Algorithms.....	119
8.2.4.1. Parallel Algorithms and Device Function Annotations.....	120
8.2.4.2. Data Management in Parallel Algorithms.....	120
8.2.4.3. Parallel Algorithms and Function Pointers.....	124
8.2.4.4. Random Access Iterators.....	126
8.2.4.5. Interoperability with the C++ Standard Library.....	126
8.2.4.6. No Exceptions in GPU Code.....	126
8.2.5. NVC++ Experimental Features.....	126
8.2.5.1. Multi-dimensional Spans.....	127
8.2.5.2. Senders and Receivers.....	128
8.2.5.3. Linear Algebra.....	129
8.2.6. Stdpar C++ Larger Example: LULESH.....	131
8.2.7. Interoperability with OpenACC.....	132
8.2.7.1. Data Management Directives.....	132
8.2.7.2. External Device Function Annotations.....	138
8.2.8. Getting Started with Parallel Algorithms for GPUs.....	138
8.2.8.1. Supported NVIDIA GPUs.....	138
8.2.8.2. Supported CUDA Versions.....	139
8.3. Stdpar Fortran.....	139
8.3.1. Calling Routines in DO CONCURRENT on the GPU.....	139
8.3.2. GPU Data Management.....	140
8.3.3. Interoperability with OpenACC.....	141
8.3.4. Interoperability with CUDA Fortran.....	142
Chapter 9. PCAST.....	144
9.1. Overview.....	144
9.2. PCAST with a "Golden" File.....	145
9.3. PCAST with OpenACC.....	148
9.4. Limitations.....	153
9.5. Environment Variables.....	153
Chapter 10. Using MPI.....	155

10.1. Using Open MPI on Linux.....	155
10.2. Using MPI Compiler Wrappers.....	156
10.3. Testing and Benchmarking.....	156
Chapter 11. Creating and Using Libraries.....	157
11.1. Using builtin Math Functions in C++ and C.....	157
11.2. Using System Library Routines.....	158
11.3. Creating and Using Shared Object Files on Linux.....	158
11.3.1. Procedure to create a use a shared object file.....	158
11.3.2. Ldd Command.....	159
11.4. Using LIB3F.....	160
11.5. LAPACK, BLAS and FFTs.....	160
11.6. Linking with ScaLAPACK.....	160
11.7. The C++ Standard Template Library.....	160
11.8. NVIDIA Performance Libraries (NVPL).....	161
11.9. Linking with the nvmalloc Library.....	161
Chapter 12. Environment Variables.....	163
12.1. Setting Environment Variables.....	163
12.1.1. Setting Environment Variables on Linux.....	163
12.2. HPC Compiler Related Environment Variables.....	164
12.3. HPC Compilers Environment Variables.....	164
12.3.1. FORTRANOPT.....	165
12.3.2. FORT_FMT_RECL.....	165
12.3.3. GMON_OUT_PREFIX.....	165
12.3.4. LD_LIBRARY_PATH.....	166
12.3.5. MANPATH.....	166
12.3.6. NO_STOP_MESSAGE.....	166
12.3.7. PATH.....	166
12.3.8. NVCOMPILER_FPU_STATE.....	166
12.3.9. NVCOMPILER_TERM.....	168
12.3.10. NVCOMPILER_TERM_DEBUG.....	169
12.3.11. PWD.....	170
12.3.12. STATIC_RANDOM_SEED.....	170
12.3.13. TMP.....	170
12.3.14. TMPDIR.....	170
12.4. Using Environment Modules on Linux.....	170
12.5. Stack Traceback and JIT Debugging.....	171
Chapter 13. Distributing Files – Deployment.....	173
13.1. Deploying Applications on Linux.....	173
13.1.1. Runtime Library Considerations.....	173
13.1.2. 64-bit Linux Considerations.....	174
13.1.3. Linux Redistributable Files.....	174
13.1.4. Restrictions on Linux Portability.....	174
13.1.5. Licensing for Redistributable (REDIST) Files.....	175

Chapter 14. Inter-language Calling.....	176
14.1. Overview of Calling Conventions.....	176
14.2. Inter-language Calling Considerations.....	176
14.3. Functions and Subroutines.....	177
14.4. Upper and Lower Case Conventions, Underscores.....	178
14.5. Compatible Data Types.....	178
14.5.1. Fortran Named Common Blocks.....	179
14.6. Argument Passing and Return Values.....	180
14.6.1. Passing by Value (%VAL).....	180
14.6.2. Character Return Values.....	180
14.6.3. Complex Return Values.....	181
14.7. Array Indices.....	181
14.8. Examples.....	182
14.8.1. Example – Fortran Calling C.....	182
14.8.2. Example – C Calling Fortran.....	183
14.8.3. Example – C++ Calling C.....	184
14.8.4. Example – C Calling C ++.....	184
14.8.5. Example – Fortran Calling C++.....	185
14.8.6. Example – C++ Calling Fortran.....	186
Chapter 15. Programming Considerations for 64-Bit Environments.....	188
15.1. Data Types in the 64-Bit Environment.....	188
15.1.1. C++ and C Data Types.....	188
15.1.2. Fortran Data Types.....	188
15.2. Large Static Data in Linux.....	189
15.3. Large Dynamically Allocated Data.....	189
15.4. 64-Bit Array Indexing.....	189
15.5. Compiler Options for 64-bit Programming.....	190
15.6. Practical Limitations of Large Array Programming.....	191
15.7. Medium Memory Model and Large Array in C.....	191
15.8. Medium Memory Model and Large Array in Fortran.....	192
15.9. Large Array and Small Memory Model in Fortran.....	193
Chapter 16. C++ and C Inline Assembly and Intrinsics.....	195
16.1. Inline Assembly.....	195
16.2. Extended Inline Assembly.....	196
16.2.1. Output Operands.....	197
16.2.2. Input Operands.....	198
16.2.3. Clobber List.....	200
16.2.4. Additional Constraints.....	201
16.2.5. Simple Constraints.....	201
16.2.6. Machine Constraints.....	202
16.2.7. Multiple Alternative Constraints.....	204
16.2.8. Constraint Modifiers.....	205
16.3. Operand Aliases.....	206

16.4. Assembly String Modifiers.....	206
16.5. Extended Asm Macros.....	208
16.6. Intrinsics.....	208



## LIST OF FIGURES

Figure 1 Nsight Systems Report1 Timeline .....	81
Figure 2 Nsight Systems Report2 Timeline .....	82

## LIST OF TABLES

Table 1	NVIDIA HPC Compilers and Commands .....	xv
Table 2	Option Descriptions .....	6
Table 3	Examples of Using siterc and User rc Files .....	8
Table 4	Typical -fast Options .....	12
Table 5	Additional -fast Options .....	13
Table 6	Commonly Used Command-Line Options .....	14
Table 7	Default settings of -Mdaz and -Mflushz .....	15
Table 8	Typical -fast Options .....	19
Table 9	Additional -fast Options .....	20
Table 10	Example of Effect of Code Unrolling .....	24
Table 11	-Mvect Suboptions .....	26
Table 12	GPU Memory Modes .....	46
Table 13	Memory Management Environment Variables .....	48
Table 14	Pool Allocator Environment Variables .....	53
Table 15	Command-line Options Corresponding to Compiler Memory Modes .....	55
Table 16	Supported Environment Variables .....	82
Table 17	Accelerator Runtime Library Routines .....	85
Table 18	Supported Fortran Intrinsics .....	87
Table 19	Supported C Intrinsic Double Functions .....	88
Table 20	Supported C Intrinsic Float Functions .....	88
Table 21	Stdpar C++ Feature Differences for Memory Modes .....	122
Table 22	Experimental features information .....	126
Table 23	Supported Types for Tolerance Measurements .....	145
Table 24	PCAST_COMPARE Options .....	153

Table 25	NVIDIA HPC Compilers Environment Variable Summary .....	164
Table 26	Supported NVCOMPILER_FPU_STATE options .....	167
Table 27	Supported NVCOMPILER_TERM Values .....	168
Table 28	Fortran and C/C++ Data Type Compatibility .....	178
Table 29	Fortran and C/C++ Representation of the COMPLEX Type .....	179
Table 30	64-bit Compiler Options .....	190
Table 31	Effects of Options on Memory and Array Sizes .....	190
Table 32	64-Bit Limitations .....	191
Table 33	Simple Constraints .....	201
Table 34	x86_64 Machine Constraints .....	203
Table 35	Multiple Alternative Constraints .....	204
Table 36	Constraint Modifier Characters .....	205
Table 37	Assembly String Modifier Characters .....	207
Table 38	Intrinsic Header File Organization .....	209

# PREFACE

This guide is part of a set of manuals that describe how to use the NVIDIA HPC Fortran, C++ and C compilers. These compilers include the *NVFORTTRAN*, *NVC++* and *NVC* compilers. They work in conjunction with an assembler, linker, libraries and header files on your target system, and include a CUDA toolchain, libraries and header files for GPU computing. You can use the NVIDIA HPC compilers to develop, optimize and parallelize applications for NVIDIA GPUs and x86-64, OpenPOWER and Arm Server multicore CPUs.

The *NVIDIA HPC Compilers User's Guide* provides operating instructions for the NVIDIA HPC compilers command-level development environment. The *NVIDIA HPC Compilers Reference Manual* contains details concerning the NVIDIA compilers' interpretation of the Fortran, C++ and C language standards, implementation of language extensions, and command-level compilation. Users are expected to have previous experience with or knowledge of the Fortran, C++ and C programming languages. These guides do not teach the Fortran, C++ or C programming languages.

## Audience Description

This manual is intended for scientists and engineers using the NVIDIA HPC compilers. To use these compilers, you should be aware of the role of high-level languages, such as Fortran, C++ and C as well as parallel programming models such as CUDA, OpenACC and OpenMP in the software development process, and you should have some level of understanding of programming. The NVIDIA HPC compilers are available on a variety of NVIDIA GPUs and x86-64, OpenPOWER and Arm CPU-based platforms and operating systems. You need to be familiar with the basic commands available on your system.

## Compatibility and Conformance to Standards

Your system needs to be running a properly installed and configured version of the NVIDIA HPC compilers. For information on installing NVIDIA HPC compilers, refer to the Release Notes and Installation Guide included with your software.

For further information, refer to the following:

- ▶ *American National Standard Programming Language FORTRAN*, ANSI X3. -1978 (1978).

- ▶ *ISO/IEC 1539-1 : 1991, Information technology – Programming Languages – Fortran*, Geneva, 1991 (Fortran 90).
- ▶ *ISO/IEC 1539-1 : 1997, Information technology – Programming Languages – Fortran*, Geneva, 1997 (Fortran 95).
- ▶ *ISO/IEC 1539-1 : 2004, Information technology – Programming Languages – Fortran*, Geneva, 2004 (Fortran 2003).
- ▶ *ISO/IEC 1539-1 : 2010, Information technology – Programming Languages – Fortran*, Geneva, 2010 (Fortran 2008).
- ▶ *ISO/IEC 1539-1 : 2018, Information technology – Programming Languages – Fortran*, Geneva, 2018 (Fortran 2018).
- ▶ *Fortran 95 Handbook Complete ISO/ANSI Reference*, Adams et al, The MIT Press, Cambridge, Mass, 1997.
- ▶ *The Fortran 2003 Handbook*, Adams et al, Springer, 2009.
- ▶ *OpenACC Application Program Interface*, Version 2.7, November 2018, <http://www.openacc.org>.
- ▶ *OpenMP Application Program Interface*, Version 5.0, November 2018, <http://www.openmp.org>.
- ▶ *Programming in VAX Fortran*, Version 4.0, Digital Equipment Corporation (September, 1984).
- ▶ *IBM VS Fortran*, IBM Corporation, Rev. GC26-4119.
- ▶ *Military Standard, Fortran*, DOD Supplement to American National Standard Programming Language Fortran, ANSI x.3-1978, MIL-STD-1753 (November 9, 1978).
- ▶ *American National Standard Programming Language C*, ANSI X3.159-1989.
- ▶ *ISO/IEC 9899:1990, Information technology – Programming Languages – C*, Geneva, 1990 (C90).
- ▶ *ISO/IEC 9899:1999, Information technology – Programming Languages – C*, Geneva, 1999 (C99).
- ▶ *ISO/IEC 9899:2011, Information Technology – Programming Languages – C*, Geneva, 2011 (C11).
- ▶ *ISO/IEC 14882:2011, Information Technology – Programming Languages – C++*, Geneva, 2011 (C++11).
- ▶ *ISO/IEC 14882:2014, Information Technology – Programming Languages – C++*, Geneva, 2014 (C++14).
- ▶ *ISO/IEC 14882:2017, Information Technology – Programming Languages – C++*, Geneva, 2017 (C++17).

## Organization

This guide contains the essential information on how to use the NVIDIA HPC compilers and is divided into these sections:

**Getting Started** provides an introduction to the NVIDIA HPC compilers and describes their use and overall features.

**Use Command-line Options** provides an overview of the command-line options as well as task-related lists of options.

**Multicore CPU Optimization** describes multicore CPU optimizations and related compiler options.

**Using Function Inlining** describes how to use function inlining and shows how to create an inline library.

**Using OpenMP** describes how to use OpenMP for multicore CPU programming.

**Using OpenACC** describes how to use an NVIDIA GPU and gives an introduction to using OpenACC.

**Using Stdpar** describes how to use C++/Fortran Standard Language Parallelism for programming an NVIDIA GPU or multicore CPU.

**PCAST** describes how to use the Parallel Compiler Assisted Testing features of the HPC Compilers.

**Using MPI** describes how to use MPI with the NVIDIA HPC compilers.

**Creating and Using Libraries** discusses NVIDIA HPC compiler support libraries, shared object files, and environment variables that affect the behavior of the compilers.

**Environment Variables** describes the environment variables that affect the behavior of the NVIDIA HPC compilers.

**Distributing Files – Deployment** describes the deployment of your files once you have built, debugged and compiled them successfully.

**Inter-language Calling** provides examples showing how to place C language calls in a Fortran program and Fortran language calls in a C program.

**Programming Considerations for 64-Bit Environments** discusses issues of which programmers should be aware when targeting 64-bit processors.

**C++ and C Inline Assembly and Intrinsics** describes how to use inline assembly code in C++ and C programs, as well as how to use intrinsic functions that map directly to assembly machine instructions.

## Hardware and Software Constraints

This guide describes versions of the NVIDIA HPC compilers that target NVIDIA GPUs and x86-64, OpenPOWER and Arm CPUs. Details concerning environment-specific values and defaults and system-specific features or limitations are presented in the release notes delivered with the NVIDIA HPC compilers.

## Conventions

This guide uses the following conventions:

***italic***

is used for emphasis.

**Constant Width**

is used for filenames, directories, arguments, options, examples, and for language statements in the text, including assembly language statements.

**Bold**

is used for commands.

**[ item1 ]**

in general, square brackets indicate optional items. In this case item1 is optional. In the context of p/t-sets, square brackets are required to specify a p/t-set.

**{ item2 | item3 }**

braces indicate that a selection is required. In this case, you must select either item2 or item3.

**filename ...**

ellipsis indicate a repetition. Zero or more of the preceding item may occur. In this example, multiple filenames are allowed.

**FORTRAN**

Fortran language statements are shown in the text of this guide using a reduced fixed point size.

**C++ and C**

C++ and C language statements are shown in the text of this guide using a reduced fixed point size.

## Terms

A number of terms related to systems, processors, compilers and tools are used throughout this guide. For example:

accelerator	FMA	-mcmmodel=medium	shared library
AVX	host	-mcmmodel=small	SIMD
CUDA	hyperthreading (HT)	MPI	SSE
device	large arrays	MPICH	static linking
driver	linux86-64	NUMA	x86-64
DWARF	LLVM	OpenPOWER	Arm
dynamic library	multicore	ppc64le	Aarch64

The following table lists the NVIDIA HPC compilers and their corresponding commands:

Table 1 NVIDIA HPC Compilers and Commands

Compiler or Tool	Language or Function	Command
NVFORTTRAN	ISO/ANSI Fortran 2003	nvfortran
NVC++	ISO/ANSI C++17 with GNU compatibility	nvc++
NVC	ISO/ANSI C11	nvc

In general, the designation *NVFORTRAN* is used to refer to the NVIDIA Fortran compiler, and *nvfortran* is used to refer to the command that invokes the compiler. A similar convention is used for each of the NVIDIA HPC compilers.

For simplicity, examples of command-line invocation of the compilers generally reference the *nvfortran* command, and most source code examples are written in Fortran. Use of *NVC++* and *NVC* is consistent with *NVFORTRAN*, though there are command-line options and features of these compilers that do not apply to *NVFORTRAN*, and vice versa.

There are a wide variety of x86-64 CPUs in use. Most of these CPUs are forward-compatible, but not backward-compatible, meaning that code compiled to target a given processor will not necessarily execute correctly on a previous-generation processor.

A table listing the processor options that NVIDIA HPC compilers support is available in the Release Notes. The table also includes the features utilized by the compilers that distinguish them from a compatibility standpoint.

In this manual, the convention is to use "x86-64" to specify the group of CPUs that are x86-compatible, 64-bit enabled, and run a 64-bit operating system. x86-64 processors can differ in terms of their support for various prefetch, SSE and AVX instructions. Where such distinctions are important with respect to a given compiler option or feature, it is explicitly noted in this manual.

## Related Publications

The following documents contain additional information related to the NVIDIA HPC compilers.

- ▶ *System V Application Binary Interface Processor Supplement* by AT&T UNIX System Laboratories, Inc. (Prentice Hall, Inc.).
- ▶ *System V Application Binary Interface X86-64 Architecture Processor Supplement*.
- ▶ *OpenPOWER ABI for Linux Supplement, Power Architecture 64-Bit ELF V2 ABI Specification*, [http://openpowerfoundation.org/wp-content/uploads/2016/03/ABI64BitOpenPOWERv1.1\\_16July2015\\_pub4.pdf](http://openpowerfoundation.org/wp-content/uploads/2016/03/ABI64BitOpenPOWERv1.1_16July2015_pub4.pdf).
- ▶ *Fortran 95 Handbook Complete ISO/ANSI Reference*, Adams et al, The MIT Press, Cambridge, Mass, 1997.
- ▶ *Programming in VAX Fortran, Version 4.0*, Digital Equipment Corporation (September, 1984).
- ▶ *IBM VS Fortran*, IBM Corporation, Rev. GC26-4119.
- ▶ *The C Programming Language* by Kernighan and Ritchie (Prentice Hall).
- ▶ *C: A Reference Manual* by Samuel P. Harbison and Guy L. Steele Jr. (Prentice Hall, 1987).
- ▶ *The Annotated C++ Reference Manual* by Margaret Ellis and Bjarne Stroustrup, AT&T Bell Laboratories, Inc. (Addison-Wesley Publishing Co., 1990).



# Chapter 1.

## GETTING STARTED

This section describes how to use the NVIDIA HPC compilers.

### 1.1. Overview

The command used to invoke a compiler, such as the `nvfortran` command, is called a *compiler driver*. The compiler driver controls the following phases of compilation: preprocessing, compiling, assembling, and linking. Once a file is compiled and an executable file is produced, you can execute, debug, or profile the program on your system.

In general, using an NVIDIA HPC compiler involves three steps:

1. Produce program source code in a file containing a `.f` extension or another appropriate extension, as described in [Input Files](#). This program may be one that you have written or one that you are modifying.
2. Compile the program using the appropriate compiler command.
3. Execute, debug, or profile the executable file on your system.

You might also want to deploy your application, though this is not a required step.

The NVIDIA HPC compilers allow many variations on these general program development steps. These variations include the following:

- ▶ Stop the compilation after preprocessing, compiling or assembling to save and examine intermediate results.
- ▶ Provide options to the driver that control compiler optimization or that specify various features or limitations.
- ▶ Include as input intermediate files such as preprocessor output, compiler output, or assembler output.

### 1.2. Creating an Example

Let's look at a simple example of using the NVIDIA Fortran compiler to create, compile, and execute a program that prints:

hello

### 1. Create your program.

For this example, suppose you enter the following simple Fortran program in the file `hello.f`:

```
print *, "hello"
end
```

### 2. Compile the program.

When you created your program, you called it `hello.f`. In this example, we compile it from a shell command prompt using the default `nvfortran` driver option. Use the following syntax:

```
$ nvfortran hello.f
```

By default, the executable output is placed in the file `a.out`. However, you can specify an output file name by using the `-o` option.

To place the executable output in the file `hello`, use this command:

```
$ nvfortran -o hello hello.f
```

### 3. Execute the program.

To execute the resulting hello program, simply type the filename at the command prompt and press the **Return** or **Enter** key on your keyboard:

```
$ hello
```

Below is the expected output:

```
hello
```

## 1.3. Invoking the Command-level NVIDIA HPC Compilers

To translate and link a Fortran, C, or C++ program, the `nvfortran`, `nvc` and `nvc++` commands do the following:

1. Preprocess the source text file.
2. Check the syntax of the source text.
3. Generate an assembly language file.
4. Pass control to the subsequent assembly and linking steps.

### 1.3.1. Command-line Syntax

The compiler command-line syntax, using `nvfortran` as an example, is:

```
nvfortran [options] [path]filename [...]
```

Where:

#### **options**

is one or more command-line options, all of which are described in detail in [Use Command-line Options](#).

**path**

is the pathname to the directory containing the file named by filename. If you do not specify the path for a filename, the compiler uses the current directory. You must specify the path separately for each filename not in the current directory.

**filename**

is the name of a source file, preprocessed source file, assembly-language file, object file, or library to be processed by the compilation system. You can specify more than one [path]filename.

## 1.3.2. Command-line Options

The command-line options control various aspects of the compilation process. For a complete alphabetical listing and a description of all the command-line options, refer to [Use Command-Line Options](#).

The following list provides important information about proper use of command-line options.

- ▶ Command-line options and their arguments are case sensitive.
- ▶ The compiler drivers recognize characters preceded by a hyphen (-) as command-line options. For example, the `-Mlist` option specifies that the compiler creates a listing file.



The convention for the text of this manual is to show command-line options using a dash instead of a hyphen; for example, you see `-Mlist`.

- ▶ The order of options and the filename is flexible. That is, you can place options before and after the filename argument on the command line. However, the placement of some options is significant, such as the `-l` option, in which the order of the filenames determines the search order.



If two or more options contradict each other, the last one in the command line takes precedence.

- ▶ You may write linker options into a text file prefixed with the '@' symbol, e.g. `@file`, and pass that file to the compiler as an option. The contents of `@file` are passed to the linker.

```
$ echo "foo.o bar.o" > ./option_file.rsp
$ nvcc @./option_files.rsp
```

The above will pass "foo.o bar.o" to the compiler as linker arguments.

## 1.4. Filename Conventions

The NVIDIA HPC compilers use the filenames that you specify on the command line to find and to create input and output files. This section describes the input and output filename conventions for the phases of the compilation process.

## 1.4.1. Input Files

You can specify assembly-language files, preprocessed source files, Fortran/C/C++ source files, object files, and libraries as inputs on the command line. The compiler driver determines the type of each input file by examining the filename extensions.

The drivers use the following conventions:

- filename.f**  
indicates a Fortran source file.
- filename.F**  
indicates a Fortran source file that can contain macros and preprocessor directives (to be preprocessed).
- filename.FOR**  
indicates a Fortran source file that can contain macros and preprocessor directives (to be preprocessed).
- filename.F90**  
indicates a Fortran 90/95 source file that can contain macros and preprocessor directives (to be preprocessed).
- filename.F95**  
indicates a Fortran 90/95 source file that can contain macros and preprocessor directives (to be preprocessed).
- filename.f90**  
indicates a Fortran 90/95 source file that is in freeform format.
- filename.f95**  
indicates a Fortran 90/95 source file that is in freeform format.
- filename.cuf**  
indicates a Fortran 90/95 source file in free format with CUDA Fortran extensions.
- filename.CUF**  
indicates a Fortran 90/95 source file in free format with CUDA Fortran extensions and that can contain macros and preprocessor directives (to be preprocessed).
- filename.c**  
indicates a C source file that can contain macros and preprocessor directives (to be preprocessed).
- filename.C**  
indicates a C++ source file that can contain macros and preprocessor directives (to be preprocessed).
- filename.i**  
indicates a preprocessed C or C++ source file.
- filename.cc**  
indicates a C++ source file that can contain macros and preprocessor directives (to be preprocessed).
- filename.cpp**  
indicates a C++ source file that can contain macros and preprocessor directives (to be preprocessed).
- filename.s**  
indicates an assembly-language file.
- filename.o**  
(Linux) indicates an object file.

**filename.a**

(Linux) indicates a library of object files.

**filename.so**

(Linux only) indicates a library of shared object files.

The driver passes files with .s extensions to the assembler and files with .o, .so and .a extensions to the linker. Input files with unrecognized extensions, or no extension, are also passed to the linker.

Files with a .F (Capital F) or .FOR suffix are first preprocessed by the Fortran compilers and the output is passed to the compilation phase. The Fortran preprocessor functions like cpp for C programs, but is built in to the Fortran compilers rather than implemented through an invocation of cpp. This design ensures consistency in the preprocessing step regardless of the type or revision of operating system under which you are compiling.

Any input files not needed for a particular phase of processing are not processed. For example, if on the command line you specify an assembly-language file (filename.s) and the -S option to stop before the assembly phase, the compiler takes no action on the assembly language file. Processing stops after compilation and the assembler does not run. In this scenario, the compilation must have been completed in a previous pass which created the .s file. For a complete description of the -S option, refer to [Output Files](#).

In addition to specifying primary input files on the command line, code within other files can be compiled as part of include files using the INCLUDE statement in a Fortran source file or the preprocessor #include directive in Fortran source files that use a .F extension or C++ and C source files.

When linking a program with a library, the linker extracts only those library components that the program needs. The compiler drivers link in several libraries by default. For more information about libraries, refer to [Create and Use Libraries](#).

## 1.4.2. Output Files

By default, an executable output file produced by one of the NVIDIA HPC compilers is placed in the file a.out. As the [Hello example](#) shows, you can use the -o option to specify the output file name.

If you use option -F (Fortran only), -P (C/C++ only), -S or -c, the compiler produces a file containing the output of the last completed phase for each input file, as specified by the option supplied.

The output file is a preprocessed source file, an assembly-language file, or an unlinked object file respectively. Similarly, the -E option does not produce a file, but displays the preprocessed source file on the standard output. Using any of these options, the -o option is valid only if you specify a single input file. If no errors occur during processing, you can use the files created by these options as input to a future invocation of any of the NVIDIA compiler drivers.

The following table lists the stop-after options and the output files that the compilers create when you use these options. It also indicates the accepted input files.

Table 2 Option Descriptions

Option	Stop After	Input	Output
-E	preprocessing	Source files	preprocessed file to standard out
-F	preprocessing	Source files. This option is not valid for <code>nvc</code> or <code>nvc++</code> .	preprocessed file ( <code>.f</code> )
-P	preprocessing	Source files. This option is not valid for <code>nvfortran</code> .	preprocessed file ( <code>.i</code> )
-S	compilation	Source files or preprocessed files	assembly-language file ( <code>.s</code> )
-c	assembly	Source files, or preprocessed files, or assembly-language files	unlinked object file ( <code>.o</code> or <code>.obj</code> )
none	linking	Source files, or preprocessed files, assembly-language files, object files, or libraries	executable file ( <code>a.out</code> )

If you specify multiple input files or do not specify an object filename, the compiler uses the input filenames to derive corresponding default output filenames of the following form, where *filename* is the input filename without its extension:

**filename.f**

indicates a preprocessed file, if you compiled a Fortran file using the -F option.

**filename.i**

indicates a preprocessed file, if you compiled using the -P option.

**filename.lst**

indicates a listing file from the -Mlist option.

**filename.o or filename.obj**

indicates a object file from the -c option.

**filename.s**

indicates an assembly-language file from the -S option.



Unless you specify otherwise, the destination directory for any output file is the current working directory. If the file exists in the destination directory, the compiler overwrites it.

The following example demonstrates the use of output filename extensions.

```
$ nvfortran -c proto.f proto1.F
```

This produces the output files `proto.o` and `proto1.o`, which are binary object files. Prior to compilation, the file `proto1.F` is preprocessed because it has a `.F` filename extension.

## 1.5. Fortran, C++ and C Data Types

The NVIDIA Fortran, C++ and C compilers recognize scalar and aggregate data types. A scalar data type holds a single value, such as the integer value 42 or the real value 112.6. An aggregate data type consists of one or more scalar data type objects, such as an array of integer values.

## 1.6. Platform-specific considerations

The NVIDIA HPC Compilers are supported on x86-64, OpenPOWER and 64-bit Arm multicore CPUs running Linux.

### 1.6.1. Using the NVIDIA HPC Compilers on Linux

#### Linux Header Files

The Linux system header files contain many GNU gcc extensions. The NVIDIA HPC C++ and C compilers support many of these extensions and can compile most programs that the GNU compilers can compile. A few header files not interoperable with the NVIDIA compilers have been rewritten.

If you are using the NVIDIA HPC C++ or C compilers, please make sure that the supplied versions of these include files are found before the system versions. This hierarchy happens by default unless you explicitly add a `-I` option that references one of the system `include` directories.

## 1.7. Site-Specific Customization of the Compilers

If you are using the NVIDIA HPC Compilers and want all your users to have access to specific libraries or other files, there are special files that allow you to customize the compilers for your site.

### 1.7.1. Use `siterc` Files

The NVIDIA HPC Compiler command-level drivers utilize a file named `siterc` to enable site-specific customization of the behavior of the NVIDIA compilers. The `siterc` file is located in the `bin` subdirectory of the NVIDIA HPC Compilers installation directory. Using `siterc`, you can control how the compiler drivers invoke the various components in the compilation tool chain.

### 1.7.2. Using User `rc` Files

In addition to the `siterc` file, user `rc` files can reside in a given user's home directory, as specified by the user's `HOME` environment variable. You can use these files to control the respective NVIDIA HPC Compilers. All of these files are optional.

On Linux, these files are named `.mynvfortranrc`, `.mynvccrc`, and `.mynvcc++rc`.

The following examples show how you can use these `rc` files to tailor a given installation for a particular purpose on `Linux_x86_64` targets. The process is similar with obvious substitutions for `ppc64le` and `aarch64` targets.

Table 3 Examples of Using siterc and User rc Files

To do this...	Add the line shown to the indicated file(s)
Make available to all linux compilations the libraries found in /opt/newlibs/64	<b>set SITELIB=/opt/newlibs/64;</b> to /opt/nv/Linux_x86_64/24.9/compilers/bin/siterc
Add to all linux compilations a new library path: /opt/local/fast	<b>append SITELIB=/opt/local/fast;</b> to /opt/nv/Linux_x86_64/24.9/compilers/bin/siterc
With linux compilations, change -Mmpi to link in /opt/mympi/64/libmpix.a	<b>set MPILIBDIR=/opt/mympi/64; set MPILIBNAME=mpix;</b> to /opt/nv/Linux_x86_64/24.9/compilers/bin/siterc
Build a Fortran executable for linux that resolves shared objects in the relative directory ./REDIST	<b>set RPATH=./REDIST;</b> to ~/.mynvfortranrc

## 1.8. Common Development Tasks

Now that you have a brief introduction to the compiler, let's look at some common development tasks that you might wish to perform.

- ▶ When you compile code you can specify a number of options on the command line that define specific characteristics related to how the program is compiled and linked, typically enhancing or overriding the default behavior of the compiler. For a list of the most common command line options and information on all the command line options, refer to [Use Command-line Options](#).
- ▶ Code optimization for multicore CPUs allows the compiler to organize your code for efficient execution. While possibly increasing compilation time and making the code more difficult to debug, these techniques typically produce code that runs significantly faster than code that does not use them. For more information on optimization refer to [Multicore CPU Optimization](#).
- ▶ Function inlining, a special type of optimization, replaces a call to a function or a subroutine with the body of the function or subroutine. This process can speed up execution by eliminating parameter passing and the function or subroutine call and return overhead. In addition, function inlining allows the compiler to optimize the function with the rest of the code. However, function inlining may also result in much larger code size with no increase in execution speed. For more information on function inlining, refer to [Using Function Inlining](#).
- ▶ A library is a collection of functions or subprograms used to develop software. Libraries contain "helper" code and data, which provide services to independent programs, allowing code and data to be shared and changed in a modular fashion. The functions and programs in a library are grouped for ease of use and linking.



When creating your programs, it is often useful to incorporate standard libraries or proprietary ones. For more information on this topic, refer to [Creating and Using Libraries](#).

- ▶ Environment variables define a set of dynamic values that can affect the way running processes behave on a computer. It is often useful to use these variables to set and pass information that alters the default behavior of the NVIDIA HPC Compilers and the executables which they generate. For more information on these variables, refer to [Environment Variables](#).
- ▶ Deployment, though possibly an infrequent task, can present some unique issues related to concerns of porting the code to other systems. Deployment, in this context, involves distribution of a specific file or set of files that are already compiled and configured. The distribution must occur in such a way that the application executes accurately on another system which may not be configured exactly the same as the system on which the code was created. For more information on what you might need to know to successfully deploy your code, refer to [Distributing Files – Deployment](#).
- ▶ An intrinsic is a function available in a given language whose implementation is handled specially by the compiler. Intrinsics make using processor-specific enhancements easier because they provide a C++ and C language interface to assembly instructions. In doing so, the compiler manages details that the user would normally have to be concerned with, such as register names, register allocations, and memory locations of data.

# Chapter 2.

## USE COMMAND-LINE OPTIONS

A command line option allows you to control specific behavior when a program is compiled and linked. This section describes the syntax for properly using command-line options and provides a brief overview of a few of the more common options.

### 2.1. Command-line Option Overview

Before looking at all the command-line options, first become familiar with the syntax for these options. There are a large number of options available to you, yet most users only use a few of them. So, start simple and progress into using the more advanced options.

By default, the NVIDIA HPC Compilers generate code that is optimized for the type of processor on which compilation is performed, the compilation host. Before adding options to your command-line, review [Help with Command-line Options](#) and [Frequently-used Options](#).

#### 2.1.1. Command-line Options Syntax

On a command-line, options need to be preceded by a hyphen (-). If the compiler does not recognize an option, you get an unknown switch error. The error can be downgraded to a warning by adding the `-noswitcherror` option.

This document uses the following notation when describing options:

**[item]**

Square brackets indicate that the enclosed item is optional.

**{item | item}**

Braces indicate that you must select one and only one of the enclosed items. A vertical bar (|) separates the choices.

...

Horizontal ellipses indicate that zero or more instances of the preceding item are valid.

## 2.1.2. Command-line Suboptions

Some options accept several suboptions. You can specify these suboptions either by using the full option statement multiple times or by using a comma-separated list for the suboptions.

The following two command lines are equivalent:

```
nvfortran -Mvect=simd -Mvect=noaltcode
```

```
nvfortran -Mvect=simd,noaltcode
```

## 2.1.3. Command-line Conflicting Options

Some options have an opposite or negated counterpart. For example, both `-Mvect` and `-Mnovect` are available. `-Mvect` enables vectorization and `-Mnovect` disables it. If you used both of these commands on a command line, they would conflict.



When you use conflicting options on a command line, the last encountered option takes precedence over any previous one.

The conflicting options rule is important for a number of reasons.

- Some options, such as `-fast`, include other options. Therefore, it is possible for you to be unaware that you have conflicting options.
- You can use this rule to create makefiles that apply specific flags to a set of files, as shown in the following example.

Example: Makefiles with Options

In this makefile fragment, `CCFLAGS` uses vectorization. `CCNOVECTFLAGS` uses the flags defined for `CCFLAGS` but disables vectorization.

```
CCFLAGS=c -Mvect=simd
CCNOVECTFLAGS=$(CCFLAGS) -Mnovect
```

## 2.2. Help with Command-line Options

If you are just getting started with the NVIDIA HPC Compilers, it is helpful to know which options are available, when to use them, and which options most users find effective.

Using `-help`

The `-help` option is useful because it provides information about all options supported by a given compiler.

You can use `-help` in one of three ways:

- Use `-help` with no parameters to obtain a list of all the available options with a brief one-line description of each.

- Add a parameter to `-help` to restrict the output to information about a specific option. The syntax for this usage is:

```
-help <command line option>
```

Suppose you use the following command to restrict the output to information about the `-fast` option:

```
$ nvfortran -help -fast
```

The output you see is similar to:

```
-fast Common optimizations; includes -O2 -Munroll=c:1 -Mnoframe -Mlre
```

In the following example, we add the `-help` parameter to restrict the output to information about the help command. The usage information for `-help` shows how groups of options can be listed or examined according to function.

```
$ nvfortran -help -help
      -help[=groups|asm|debug|language|linker|opt|other|overall|phase|
prepro|
      suffix|switch|target|variable]
```

- Add a parameter to `-help` to restrict the output to a specific set of options or to a building process. The syntax for this usage is this:

```
-help=<subgroup>
```

## 2.3. Getting Started with Performance

This section provides a quick overview of a few of the command-line options that are useful in improving multicore CPU performance.

### 2.3.1. Using `-fast`

The NVIDIA HPC Compilers implement a wide range of options that allow users a fine degree of control on each optimization phase. When it comes to optimization of code, the quickest way to start is to use the option `-fast`. These options create a generally optimal set of flags. They incorporate optimization options to enable use of vector streaming SIMD instructions for 64-bit targets. They enable vectorization with SIMD instructions, cache alignment, and flush to zero mode.



The contents of the `-fast` option are host-dependent. Further, you should use these options on both compile and link command lines.

The following table shows the typical `-fast` options.

Table 4 Typical `-fast` Options

Use this option...	To do this...
<code>-O2</code>	Specifies a code optimization level of 2.
<code>-Munroll=c:1</code>	Unrolls loops, executing multiple instances of the original loop during each iteration.

Use this option...	To do this...
-Mnoframe	Do not generate code to set up a stack frame. Note: With this option, a stack trace does not work.
-Mlre	Enable loop-carried redundancy elimination.
-Mpre	Enable partial redundancy elimination

On most modern CPUs the `-fast` also includes the options shown in this table:

Table 5 Additional `-fast` Options

Use this option...	To do this...
-Mvect=simd	Generates packed SIMD instructions.
-Mcache_align	Aligns long objects on cache-line boundaries.
-Mflushz	Sets flush-to-zero mode.
-M[no]vect	Controls automatic vector pipelining.

To see the specific behavior of `-fast` for your target, use the following command:

```
$ nvfortran -help -fast
```

## 2.3.2. Other Performance-Related Options

While `-fast` is designed to be the quickest route to best performance, it is limited to routine boundaries. Depending on the nature and writing style of the source code, the compiler often can perform further optimization by knowing the global context of usage of a given routine. For instance, determining the possible value range of actual parameters of a routine could enable a loop to be vectorized; similarly, determining static occurrence of calls helps to decide which routine is beneficial to inline.

These types of global optimizations are under control of Interprocedural Analysis (IPA) in NVIDIA HPC Compilers. Option `-Mipa` enables Interprocedural Analysis. `-Mipa=fast` is the recommended option to get best performances for global optimization. You can also add the suboption `inline` to enable automatic global inlining across files. You might consider using `-Mipa=fast,inline`. This option for interprocedural analysis and global optimization can improve performance.

For more information on optimization, refer to [Multicore CPU Optimization](#). For specific information about these options, refer to the 'Optimization Controls' section of the [HPC Compilers Reference User Guide, docs.nvidia.com/hpc-sdk/compilers/pdf/hpc249ref.pdf](https://docs.nvidia.com/hpc-sdk/compilers/pdf/hpc249ref.pdf).

## 2.4. Frequently-used Options

In addition to overall performance, there are a number of other options that many users find useful when getting started. The following table provides a brief summary of these options.

Table 6 Commonly Used Command-Line Options

Use this option...	To do this...
<b>-acc</b>	Enable parallelization using OpenACC directives. By default the compilers will parallelize and offload OpenACC regions to an NVIDIA GPU. Use <b>-acc=multicore</b> to parallelize OpenACC regions for execution on all the cores of a multicore CPU.
<b>-fast</b>	This option creates a generally optimal set of flags for targets that support SIMD capability. It incorporates optimization options to enable use of vector streaming SIMD instructions, cache alignment and flushz.
<b>-g</b>	Instructs the compiler to include symbolic debugging information in the object module; sets the optimization level to zero unless a <b>-O</b> option is present on the command line. Conversely, to prevent the generation of DWARF information, use the <b>-Mnodwarf</b> option.
<b>-gopt</b>	Instructs the compiler to include symbolic debugging information in the object file, and to generate optimized code identical to that generated when <b>-g</b> is not specified.
<b>-gpu</b>	Control the type of GPU for which code is generated, the version of CUDA to be targeted, and several other aspects of GPU code generation.
<b>-help</b>	Provides information about available options.
<b>-mcmmodel=medium</b>	Enables medium=medium code generation for 64-bit targets, which is useful when the data space of the program exceeds 4GB.
<b>-mp</b>	Enable parallelization using OpenMP directives. By default the compilers will parallelize OpenMP regions for execution on all the cores of a multicore CPU. Use <b>-mp=gpu</b> to parallelize OpenMP regions for offload to an NVIDIA GPU.
<b>-Mconcur</b>	Instructs the compiler to enable auto-concurrentization of loops. If specified, the compiler uses multiple CPU cores to execute loops that it determines to be parallelizable; thus, loop iterations are split to execute optimally in a multithreaded execution context.
<b>-Minfo</b>	Instructs the compiler to produce information on standard error.
<b>-Minline</b>	Enables function inlining.
<b>-Mipa=fast,inline</b>	Enables interprocedural analysis and optimization. Also enables automatic procedure inlining.
<b>-Mkeepasm</b>	Keeps the generated assembly files.
<b>-Munroll</b>	Invokes the loop unroller to unroll loops, executing multiple instances of the loop during each iteration. This also sets the optimization level to 2 if the level is set to less than 2, or if no <b>-O</b> or <b>-g</b> options are supplied.
<b>-M[no]vect</b>	Enables [Disables] the code vectorizer.
<b>--[no_]exceptions</b>	Removes exception handling from user code. For C++, declares that the functions in this file generate no C++ exceptions, allowing more optimal code generation.
<b>-o</b>	Names the output file.
<b>-O &lt;level&gt;</b>	Specifies code optimization level where <level> is 0, 1, 2, 3, or 4.

Use this option...	To do this...
<b>-stdpar</b>	Enable parallelization and offloading of Standard C++ and Fortran parallel constructs to NVIDIA GPUs; default is <code>-stdpar=gpu</code> .
<b>-tp &lt;target&gt;</b>	Specify a CPU target other than the compilation host CPU.
<b>-Wl, &lt;option&gt;</b>	Compiler driver passes the specified options to the linker.

## 2.5. Floating-point Subnormal

Starting with the 22.7 release of the NV HPC SDK the default setting of how floating-point denormal (IEEE 754 terminology "subnormal") values are processed at runtime across both x86\_64 and aarch64 processors has been changed to be more consistent.

Denormal values can be both operands to, and results of, floating-point operations. The x86\_64 ISA differentiate between the two categories, operands and results, and use the terminology "daz" denormals are zeros for operands, and "flushz" flush to zero for results. The Arm V8 ISA as defined can differentiate between the two categories, but currently the processors that NV HPC SDK support only have a single setting for both operands and results and is defined as "fz" in the floating-point status and control register.

The NV HPC SDK C, C++, and Fortran compilers have command line switches `-M[no]daz` and `-M[no]flushz`, which when specified for the C/C++ main function or the Fortran main program affect how denormals are handled by the processor at runtime. The values of these two command line switches are passed to the runtime library to configure the floating-point status and control register at program startup.

NV HPC SDK supports x86\_64 processors from both Intel and AMD, and ArmV8.1 and later processors. The following table summarizes the default settings of the `-Mdaz` and `-Mflushz` command line switches pre and post the 22.7 release.

Table 7 Default settings of `-Mdaz` and `-Mflushz`

	Pre 22.7 defaults	22.7 defaults
Intel	<code>-Mdaz</code> <code>-Mnoflushz</code>	<code>-Mdaz</code> <code>-Mflushz</code>
AMD	<code>-Mnodaz</code> <code>-Mnoflushz</code>	<code>-Mdaz</code> <code>-Mflushz</code>
Arm processors	<code>-Mnodaz</code>	<code>-Mdaz</code>

With the NV HPC SDK 22.7 release, the default handling of denormals operands and results is to treat them as zero, as if the main function/program were compiled with `-Mdaz -Mflushz`. Consequently, these changes can potentially affect applications that are dependent on subnormal values being non-zero.

Along with the change to the default treatment of denormal values, users now have the ability to configure the floating-point status and control register through the

NVCOMPILER\_FPU\_STATE environment variable - effectively overriding how the program was originally compiled. For further information, see the description of the [NVCOMPILER\\_FPU\\_STATE](#) environment variable.



# Chapter 3.

## MULTICORE CPU OPTIMIZATION

Source code that is readable, maintainable, and produces correct results is not always organized for efficient execution. Normally, the first step in the program development process involves producing code that executes and produces the correct results. This first step usually involves compiling without much worry about optimization. After code is compiled and debugged, code optimization and parallelization become an issue.

Invoking one of the NVIDIA HPC Compiler commands with certain options instructs the compiler to generate optimized code. Optimization is not always performed since it increases compilation time and may make debugging difficult. However, optimization produces more efficient code that usually runs significantly faster than code that is not optimized.

The compilers optimize code according to the specified optimization level. You can use a number of options to specify the optimization levels, including `-O`, `-Mvect`, `-Mipa` and `-Mconcur`. In addition, you can use several of the `-M<nvflag>` switches to control specific types of optimization.

This chapter describes the overall effect of the optimization options supported by the NVIDIA HPC Compilers, and basic usage of several options.

### 3.1. Overview of Optimization

In general, optimization involves using transformations and replacements that generate more efficient code. This is done by the compiler and involves replacements that are independent of the particular target processor's architecture as well as replacements that take advantage of the x86-64/OpenPOWER architecture, instruction set and registers.

For discussion purposes, we categorize optimization:

- Local Optimization

- Global Optimization

- Loop Optimization

- Interprocedural Analysis (IPA) and Optimization

- Optimization Through Function Inlining

### 3.1.1. Local Optimization

A basic block is a sequence of statements in which the flow of control enters at the beginning and leaves at the end without the possibility of branching, except at the end. Local optimization is performed on a block-by-block basis within a program's basic blocks.

The NVIDIA HPC Compilers perform many types of local optimization including: algebraic identity removal, constant folding, common sub-expression elimination, redundant load and store elimination, scheduling, strength reduction, and peephole optimizations.

### 3.1.2. Global Optimization

This optimization is performed on a subprogram/function over all its basic blocks. The optimizer performs control-flow and data-flow analysis for an entire program unit. All loops, including those formed by ad hoc branches such as IFs or GOTOs, are detected and optimized.

Global optimization includes: constant propagation, copy propagation, dead store elimination, global register allocation, invariant code motion, and induction variable elimination.

### 3.1.3. Loop Optimization: Unrolling, Vectorization and Parallelization

The performance of certain classes of loops may be improved through vectorization or unrolling options. Vectorization transforms loops to improve memory access performance and make use of packed SSEvector instructions which perform the same operation on multiple data items concurrently. Unrolling replicates the body of loops to reduce loop branching overhead and provide better opportunities for local optimization, vectorization and scheduling of instructions. Performance for loops on systems with multiple processors may also improve using the parallelization features of the NVIDIA HPC Compilers.

### 3.1.4. Interprocedural Analysis (IPA) and Optimization

Interprocedural analysis (IPA) allows use of information across function call boundaries to perform optimizations that would otherwise be unavailable. For example, if the actual argument to a function is in fact a constant in the caller, it may be possible to propagate that constant into the callee and perform optimizations that are not valid if the dummy argument is treated as a variable. A wide range of optimizations are enabled or improved by using IPA, including but not limited to data alignment optimizations, argument removal, constant propagation, pointer disambiguation, pure function detection, F90/F95 array shape propagation, data placement, empty function removal, automatic function inlining, inlining of functions from pre-compiled libraries, and interprocedural optimization of functions from pre-compiled libraries.

### 3.1.5. Function Inlining

This optimization allows a call to a function to be replaced by a copy of the body of that function. This optimization will sometimes speed up execution by eliminating the function call and return overhead. Function inlining may also create opportunities for other types of optimization. Function inlining is not always beneficial. When used improperly it may increase code size and generate less efficient code.

## 3.2. Getting Started with Optimization

The first concern should be getting the program to execute and produce correct results. To get the program running, start by compiling and linking without optimization. Add `-O0` to the compile line to select no optimization; or add `-g` to debug the program easily and isolate any coding errors exposed during porting.

To get started quickly with optimization, a good set of options to use with any of the NVIDIA HPC compilers is `-fast`. For example:

```
$ nvfortran -fast -Mipa=fast,inline prog.f
```

For all of the NVIDIA HPC Fortran, C++ and C compilers, the `-fast -Mipa=fast,inline` options generally produce code that is well-optimized without the possibility of significant slowdowns due to pathological cases.

- ▶ The `-fast` option is an aggregate option that includes a number of individual NVIDIA compiler options; which compiler options are included depends on the target for which compilation is performed.
- ▶ The `-Mipa=fast,inline` option invokes interprocedural analysis (IPA), including several IPA suboptions. The `inline` suboption enables automatic inlining with IPA. If you do not wish to use automatic inlining, you can compile with `-Mipa=fast` and use several IPA suboptions without inlining.

These aggregate options incorporate a generally optimal set of flags for targets that support SIMD capability, including vectorization with SIMD instructions, cache alignment, and flushz.

The following table shows the typical `-fast` options.

Table 8 Typical `-fast` Options

Use this option...	To do this...
<code>-O2</code>	Specifies a code optimization level of 2 and <code>-Mvect=SIMD</code> .
<code>-Munroll=c:1</code>	Unrolls loops, executing multiple instances of the original loop during each iteration.
<code>-Mnoframe</code>	Indicates to not generate code to set up a stack frame. Note With this option, a stack trace does not work.
<code>-Mlre</code>	Indicates loop-carried redundancy elimination.
<code>-Mautoinline</code>	Enables automatic function inlining in C & C++.

Use this option...	To do this...
-Mpre	Indicates partial redundancy elimination

On modern multicore CPUs the `-fast` also typically includes the options shown in the following table:

Table 9 Additional `-fast` Options

Use this option...	To do this...
-Mvect=simd	Generates packed SSE and AVX instructions.
-Mcache_align	Aligns long objects on cache-line boundaries.
-Mflushz	Sets flush-to-zero mode.

By experimenting with individual compiler options on a file-by-file basis, further significant performance gains can sometimes be realized. However, depending on the coding style, individual optimizations can sometimes cause slowdowns, and must be used carefully to ensure performance improvements.

There are other useful command line options related to optimization and parallelization, such as `-help`, `-Minfo`, `-Mneginfo`, `-dryrun`, and `-v`.

### 3.2.1. `-help`

As described in [Help with Command-Line Options](#), you can see a specification of any command-line option by invoking any of the NVIDIA HPC Compilers with `-help` in combination with the option in question, without specifying any input files.

For example, you might want information on `-O`:

```
$ nvfortran -help -O
```

The resulting output is similar to this:

```
-O Set opt level. All -O1 optimizations plus traditional scheduling and
  global scalar optimizations performed
```

Or you can see the full functionality of `-help` itself, which can return information on either an individual option or groups of options:

```
$ nvfortran -help -help
```

The resulting output is similar to this:

```
-help[=groups|asm|debug|language|linker|opt|other|overall|
phase|prepro|suffix|switch|target|variable]
Show compiler switches
```

### 3.2.2. `-Minfo`

You can use the `-Minfo` option to display compile-time optimization listings. When this option is used, the NVIDIA HPC Compilers issue informational messages to standard error (stderr) as compilation proceeds. From these messages, you can determine which loops are optimized using unrolling, SIMD vectorization, parallelization, GPU

offloading, interprocedural optimizations and various miscellaneous optimizations. You can also see where and whether functions are inlined.

### 3.2.3. -Mneginfo

You can use the `-Mneginfo` option to display informational messages to standard error (stderr) that explain why certain optimizations are inhibited.

### 3.2.4. -dryrun

The `-dryrun` option can be useful as a diagnostic tool if you need to see the steps used by the compiler driver to preprocess, compile, assemble and link in the presence of a given set of command line inputs. When you specify the `-dryrun` option, these steps are printed to standard error (stderr) but are not actually performed. For example, you can use this option to inspect the default and user-specified libraries that are searched during the link phase, and the order in which they are searched by the linker.

### 3.2.5. -v

The `-v` option is similar to `-dryrun`, except each compilation step is performed and not simply printed.

## 3.3. Local and Global Optimization

This section describes local and global optimization.

### 3.3.1. -Msafepttr

The `-Msafepttr` option can significantly improve performance of C++ and C programs in which there is known to be no pointer aliasing. For obvious reasons, this command-line option must be used carefully. There are a number of suboptions for `-Msafepttr`:

- ▶ `-Msafepttr=all` – All pointers are safe. Equivalent to the default setting: `-Msafepttr`.
- ▶ `-Msafepttr=arg` – Function formal argument pointers are safe. Equivalent to `-Msafepttr=dummy`.
- ▶ `-Msafepttr=global` – Global pointers are safe.
- ▶ `-Msafepttr=local` – Local pointers are safe. Equivalent to `-Msafepttr=auto`.
- ▶ `-Msafepttr=static` – Static local pointers are safe.

If your C++ or C program has pointer aliasing and you also want automating inlining, then compiling with `-Mipa=fast` or `-Mipa=fast,inline` includes pointer aliasing optimizations. IPA may be able to optimize some of the alias references in your program and leave intact those that cannot be safely optimized.

### 3.3.2. -O

Using the NVIDIA HPC Compiler commands with the `-O<level>` option (the capital O is for Optimize), you can specify any integer level from 0 to 4.

#### -O0

Level zero specifies no optimization. A basic block is generated for each language statement. At this level, the compiler generates a basic block for each statement.

Performance will almost always be slowest using this optimization level. This level is useful for the initial execution of a program. It is also useful for debugging, since there is a direct correlation between the program text and the code generated. To enable debugging, include `-g` on your compile line.

#### -O1

Level one specifies local optimization. Scheduling of basic blocks is performed. Register allocation is performed.

Local optimization is a good choice when the code is very irregular, such as code that contains many short statements containing IF statements and does not contain loops (DO or DO WHILE statements). Although this case rarely occurs, for certain types of code, this optimization level may perform better than level-two (`-O2`).

#### -O

When no level is specified, level two global optimizations are performed, including traditional scalar optimizations, induction recognition, and loop invariant motion. No SIMD vectorization is enabled.

#### -O2

Level two specifies global optimization. This level performs all level-one local optimization as well as level two global optimization described in `-O`. In addition, more advanced optimizations such as SIMD code generation, cache alignment, and partial redundancy elimination are enabled.

#### -O3

Level three specifies aggressive global optimization. This level performs all level-one and level-two optimizations and enables more aggressive hoisting and scalar replacement optimizations that may or may not be profitable.

-O4

Level four performs all level-one, level-two, and level-three optimizations and enables hoisting of guarded invariant floating point expressions.

### Types of Optimizations

The NVIDIA HPC Compilers perform many different types of local optimizations, including but not limited to:

- Algebraic identity removal
- Constant folding
- Common subexpression elimination
- Local register optimization
- Peephole optimizations
- Redundant load and store elimination
- Strength reductions

Level-two optimization (`-O2` or `-O`) specifies global optimization. The `-fast` option generally specifies global optimization; however, the `-fast` switch varies from release to release, depending on a reasonable selection of switches for any one particular release. The `-O` or `-O2` level performs all level-one local optimizations as well as global optimizations. Control flow analysis is applied and global registers are allocated for all functions and subroutines. Loop regions are given special consideration. This optimization level is a good choice when the program contains loops, the loops are short, and the structure of the code is regular.

The NVIDIA HPC Compilers perform many different types of global optimizations, including but not limited to:

- Branch to branch elimination
- Constant propagation
- Copy propagation
- Dead store elimination
- Global register allocation
- Induction variable elimination
- Invariant code motion

You can explicitly select the optimization level on the command line. For example, the following command line specifies level-two optimization which results in global optimization:

```
$ nvfortran -O2 prog.f
```

The default optimization level changes depending on which options you select on the command line. For example, when you select the `-g` debugging option, the default optimization level is set to level-zero (`-O0`). However, if you need to debug optimized code, you can use the `-gopt` option to generate debug information without perturbing

optimization. For a description of the default levels, refer to Default Optimization Levels.

The `-fast` option includes `-O2` on all targets. If you want to override the default for `-fast` with `-O3` while maintaining all other elements of `-fast`, simply compile as follows:

```
$ nvfortran -fast -O3 prog.f
```

### 3.4. Loop Unrolling using `-Munroll`

This optimization unrolls loops, which reduces branch overhead, and can improve execution speed by creating better opportunities for instruction scheduling. A loop with a constant count may be completely unrolled or partially unrolled. A loop with a non-constant count may also be unrolled. A candidate loop must be an innermost loop containing one to four blocks of code.

The following example shows the use of the `-Munroll` option:


```
$ nvfortran -Munroll prog.f
```

The `-Munroll` option is included as part of `-fast` on all targets. The loop unroller expands the contents of a loop and reduces the number of times a loop is executed. Branching overhead is reduced when a loop is unrolled two or more times, since each iteration of the unrolled loop corresponds to two or more iterations of the original loop; the number of branch instructions executed is proportionately reduced. When a loop is unrolled completely, the loop's branch overhead is eliminated altogether.

Loop unrolling may be beneficial for the instruction scheduler. When a loop is completely unrolled or unrolled two or more times, opportunities for improved scheduling may be presented. The code generator can take advantage of more possibilities for instruction grouping or filling instruction delays found within the loop.

#### Examples Showing Effect of Unrolling

The following side-by-side examples show the effect of code unrolling on a segment that computes a dot product.



This example is only meant to represent how the compiler can transform the loop; it is not meant to imply that the programmer needs to manually change code. In fact, manually unrolling your code can sometimes inhibit the compiler's analysis and optimization.

Table 10 Example of Effect of Code Unrolling

Dot Product Code	Unrolled Dot Product Code
<pre>REAL*4 A(100), B(100), Z INTEGER I DO I=1, 100   Z = Z + A(i) * B(i)</pre>	<pre>REAL*4 A(100), B(100), Z INTEGER I DO I=1, 100, 2   Z = Z + A(i) * B(i)</pre>



Dot Product Code	Unrolled Dot Product Code
<pre>END DO END</pre>	<pre>      Z = Z + A(i+1) * B(i+1) END DO END</pre>

Using the `-Minfo` option, the compiler informs you when a loop is being unrolled. For example, a message similar to the following, indicating the line number, and the number of times the code is unrolled, displays when a loop is unrolled:

```
dot:
  5, Loop unrolled 5 times
```

Using the `c:<m>` and `n:<m>` sub-options to `-Munroll`, or using `-Mnounroll`, you can control whether and how loops are unrolled on a file-by-file basis. For more information on `-Munroll`, refer to [Use Command-line Options](#).

## 3.5. Vectorization using `-Mvect`

The `-Mvect` option is included as part of `-fast` on all multicore CPU targets. If your program contains computationally-intensive loops, the `-Mvect` option may be helpful. If in addition you specify `-Minfo`, and your code contains loops that can be vectorized, the compiler reports relevant information on the optimizations applied.

When an NVIDIA HPC Compiler command is invoked with the `-Mvect` option, the vectorizer scans code searching for loops that are candidates for high-level transformations such as loop distribution, loop interchange, cache tiling, and idiom recognition (replacement of a recognizable code sequence, such as a reduction loop, with optimized code sequences or function calls). When the vectorizer finds vectorization opportunities, it internally rearranges or replaces sections of loops (the vectorizer changes the code generated; your source code's loops are not altered). In addition to performing these loop transformations, the vectorizer produces extensive data dependence information for use by other phases of compilation and detects opportunities to use vector or packed SIMD instructions on processors where these are supported.

The `-Mvect` option can speed up code which contains well-behaved countable loops which operate on large floating point arrays in Fortran and their C++ and C counterparts. However, it is possible that some codes will show a decrease in performance when compiled with the `-Mvect` option due to the generation of conditionally executed code segments, inability to determine data alignment, and other code generation factors. For this reason, it is recommended that you check carefully whether particular program units or loops show improved performance when compiled with this option enabled.

### 3.5.1. Vectorization Sub-options

The vectorizer performs high-level loop transformations on countable loops. A loop is countable if the number of iterations is set only before loop execution and cannot be modified during loop execution. Some of the vectorizer transformations can be

controlled by arguments to the `-Mvect` command line option. The following sections describe the arguments that affect the operation of the vectorizer. In addition, some of these vectorizer operations can be controlled from within code using directives and pragmas.

The vectorizer performs the following operations:

- ▶ Loop interchange
- ▶ Loop splitting
- ▶ Loop fusion
- ▶ Generation of SIMD instructions on CPUs where these are supported
- ▶ Generation of prefetch instructions on processors where these are supported
- ▶ Loop iteration peeling to maximize vector alignment
- ▶ Alternate code generation

The following table lists and briefly describes some of the `-Mvect` suboptions.

Table 11 -Mvect Suboptions

Use this option ...	To instruct the vectorizer to do this ...
<code>-Mvect=altcode</code>	Generate appropriate code for vectorized loops.
<code>-Mvect=[no]assoc</code>	Perform[disable] associativity conversions that can change the results of a computation due to a round-off error. For example, a typical optimization is to change one arithmetic operation to another arithmetic operation that is mathematically correct, but can be computationally different and generate faster code. This option is provided to enable or disable this transformation, since a round-off error for such associativity conversions may produce unacceptable results.
<code>-Mvect=fuse</code>	Enable loop fusion.
<code>-Mvect=gather</code>	Enable vectorization of indirect array references.
<code>-Mvect=idiom</code>	Enable idiom recognition.
<code>-Mvect=levels:&lt;n&gt;</code>	Set the maximum next level of loops to optimize.
<code>-Mvect=nocond</code>	Disable vectorization of loops with conditions.
<code>-Mvect=partial</code>	Enable partial loop vectorization via inner loop distribution.
<code>-Mvect=prefetch</code>	Automatically generate prefetch instructions when vectorizable loops are encountered, even in cases where SSE/SIMD instructions are not generated.
<code>-Mvect=short</code>	Enable short vector operations.
<code>-Mvect=simd</code>	Automatically generate packed SSE (Streaming SIMD Extensions)/SIMD, and prefetch instructions when vectorizable loops are encountered. SIMD instructions, first introduced on Pentium III and AthlonXP processors, operate on single-precision floating-point data.

Use this option ...	To instruct the vectorizer to do this ...
-Mvect=sizelimit:n	Limit the size of vectorized loops.
-Mvect=sse	Equivalent to -Mvect=simd.
-Mvect=uniform	Perform consistent optimizations in both vectorized and residual loops. Be aware that this may affect the performance of the residual loop.



Inserting **no** in front of an option disables the option. For example, to disable the generation of SIMD instructions, compile with -Mvect=nosimd.

### 3.5.2. Vectorization Example Using SIMD Instructions

One of the most important vectorization options is -Mvect=simd. When you use this option, the compiler automatically generates SIMD vector instructions, where possible, when targeting processors on which these instructions are supported. This process can improve performance by several factors compared with the equivalent scalar code. All of the NVIDIA HPC Fortran, C++ and C compilers support this capability.

In the program in [Vector operation using SIMD instructions](#), the vectorizer recognizes the vector operation in subroutine 'loop' when either compiler switch -Mvect=simd or -fast is used. This example shows the compilation, informational messages, and runtime results using SIMD instructions on an Intel Core i7 7800X Skylake system, along with issues that affect SIMD performance.

Loops vectorized using SIMD instructions operate much more efficiently when processing vectors that are aligned to a cache-line boundary. You can cause unconstrained data objects of size 16 bytes or greater to be cache-aligned by compiling with the -Mcache\_align switch. An unconstrained data object is a data object that is not a common block member and not a member of an aggregate data structure.



For stack-based local variables to be properly aligned, the main program or function must be compiled with -Mcache\_align.

The -Mcache\_align switch has no effect on the alignment of Fortran allocatable or automatic arrays. If you have arrays that are constrained, such as vectors that are members of Fortran common blocks, you must specifically pad your data structures to ensure proper cache alignment. You can use -Mcache\_align for only the beginning address of each common block to be cache-aligned.

The following examples show the results of compiling the sample code in [Vector operation using SIMD instructions](#) both with and without the option -Mvect=simd.

#### Vector operation using SIMD instructions

```

program vector_op
  parameter (N = 9999)
  real*4 x(N), y(N), z(N), w(N)
  do i = 1, n
    y(i) = i
    z(i) = 2*i
    w(i) = 4*i
  end do
end program

```

```

enddo
do j = 1, 200000
  call loop(x,y,z,w,1.0e0,N)
enddo
print *, x(1),x(771),x(3618),x(6498),x(9999)
end

```

```

subroutine loop(a,b,c,d,s,n)
  integer i, n
  real*4 a(n), b(n), c(n), d(n),s
  do i = 1, n
    a(i) = b(i) + c(i) - s * d(i)
  enddo
end

```

Assume the preceding program is compiled as follows, where `-Mvect=nosimd` disables SIMD vectorization:

```

% nvfortran -fast -Mvect=nosimd -Minfo vadd.f -Mfree -o vadd
vector_op:
  4, Loop unrolled 16 times
    Generated 1 prefetches in scalar loop
  9, Loop not vectorized/parallelized: contains call
loop:
  18, Loop unrolled 8 times
    FMA (fused multiply-add) instruction(s) generated

```

The following output shows a sample result if the generated executable is run and timed on an Intel Core i7 7800X Skylake system:

```

$ /bin/time vadd
-1.000000      -771.0000      -3618.000      -6498.000
-9999.000
0.99user 0.01system 0:01.18elapsed 84%CPU (0avgtext+0avgdata 3120maxresident)k
7736inputs+0outputs (4major+834minor)pagefaults 0swaps

$ /bin/time vadd
-1.000000      -771.0000      -3618.000      -6498.000
-9999.000
2.31user 0.00system 0:02.57elapsed 89%CPU (0avgtext+0avgdata 6976maxresident)k
8192inputs+0outputs (4major+149minor)pagefaults 0swaps

```

Now, recompile with vectorization enabled, and you see results similar to these:

```

% nvfortran -fast -Minfo vadd.f -Mfree -o vadd
vector_op:
  4, Loop not vectorized: may not be beneficial
    Unrolled inner loop 8 times
    Residual loop unrolled 7 times (completely unrolled)
    Generated 1 prefetches in scalar loop
  9, Loop not vectorized/parallelized: contains call
loop:
  18, Generated 2 alternate versions of the loop
    Generated vector simd code for the loop
    Generated 3 prefetch instructions for the loop
    Generated vector simd code for the loop
    Generated 3 prefetch instructions for the loop
    Generated vector simd code for the loop
    Generated 3 prefetch instructions for the loop
    FMA (fused multiply-add) instruction(s) generated

```

Notice the informational messages for the loop at line 18. The first line of the message indicates that two alternate versions of the loop were generated. The loop count and alignments of the arrays determine which of these versions is executed. The next several

lines indicate the loop was vectorized and that prefetch instructions have been generated for three loads to minimize latency of data transfers from main memory.

Executing again, you should see results similar to the following:

```
$ /bin/time vadd-simd
-1.000000      -771.0000      -3618.000      -6498.000
-9999.000
0.27user 0.00system 0:00.29elapsed 93%CPU (0avgtext+0avgdata 3124maxresident)k
0inputs+0outputs (0major+838minor)pagefaults 0swaps

$ /bin/time vadd-simd
-1.000000      -771.0000      -3618.000      -6498.000
-9999.000
0.62user 0.00system 0:00.65elapsed 95%CPU (0avgtext+0avgdata 6976maxresident)k
0inputs+0outputs (0major+151minor)pagefaults 0swaps
```

The SIMD result is 3.7 times faster than the equivalent non-SIMD version of the program.

Speed-up realized by a given loop or program can vary widely based on a number of factors:

- ▶ When the vectors of data are resident in the data cache, performance improvement using SIMD instructions is most effective.
- ▶ If data is aligned properly, performance will be better in general than when using SIMD operations on unaligned data.
- ▶ If the compiler can guarantee that data is aligned properly, even more efficient sequences of SIMD instructions can be generated.
- ▶ The efficiency of loops that operate on single-precision data can be higher. SIMD instructions can operate on four single-precision elements concurrently, but only two double-precision elements.



Compiling with `-Mvect=simd` can result in numerical differences from the executables generated with less optimization. Certain vectorizable operations, for example dot products, are sensitive to order of operations and the associative transformations necessary to enable vectorization (or parallelization).

## 3.6. Interprocedural Analysis and Optimization using `-Mipa`

The NVIDIA HPC Fortran, C++ and C compilers use interprocedural analysis (IPA) that results in minimal changes to makefiles and the standard edit-build-run application development cycle. Other than adding `-Mipa` to the command line, no other changes are required. For reference and background, the process of building a program without IPA is described later in this section, followed by the minor modifications required to use IPA with the NVIDIA compilers. While the NVC compiler is used here to show how IPA works, similar capabilities apply to each of the NVIDIA HPC Fortran, C++ and C compilers.

### 3.6.1. Building a Program Without IPA – Single Step

Using the `nvc` command-level compiler driver, multiple source files can be compiled and linked into a single executable with one command. The following example compiles and links three source files:

```
% nvc -o a.out file1.c file2.c file3.c
```

In actuality, the `nvc` driver executes several steps to produce the assembly code and object files corresponding to each source file, and subsequently to link the object files together into a single executable file. This command is roughly equivalent to the following commands performed individually:

```
% nvc -S -o file1.s file1.c
% as -o file1.o file1.s
% nvc -S -o file2.s file2.c
% as -o file2.o file2.s
% nvc -S -o file3.s file3.c
% as -o file3.o file3.s
% nvc -o a.out file1.o file2.o file3.o
```

If any of the three source files is edited, the executable can be rebuilt with the same command line:

```
% nvc -o a.out file1.c file2.c file3.c
```



This always works as intended, but has the side-effect of recompiling all of the source files, even if only one has changed. For applications with a large number of source files, this can be time-consuming and inefficient.

### 3.6.2. Building a Program Without IPA – Several Steps

It is also possible to use individual `nvc` commands to compile each source file into a corresponding object file, and one to link the resulting object files into an executable:

```
% nvc -c file1.c
% nvc -c file2.c
% nvc -c file3.c
% nvc -o a.out file1.o file2.o file3.o
```

The `nvc` driver invokes the compiler and assembler as required to process each source file, and invokes the linker for the final link command. If you modify one of the source files, the executable can be rebuilt by compiling just that file and then relinking:

```
% nvc -c file1.c
% nvc -o a.out file1.o file2.o file3.o
```

### 3.6.3. Building a Program Without IPA Using Make

The program compilation and linking process can be simplified greatly using the `make` utility on systems where it is supported. Suppose you create a `makefile` containing the following lines:

```
a.out: file1.o file2.o file3.o
    nvc $(OPT) -o a.out file1.o file2.o file3.o
file1.o: file1.c
    nvc $(OPT) -c file1.c
file2.o: file2.c
    nvc $(OPT) -c file2.c
```

```
file3.o: file3.c
nvc $(OPT) -c file3.c
```

It is then possible to type a single make command:

```
% make
```

The make utility determines which object files are out of date with respect to their corresponding source files, and invokes the compiler to recompile only those source files and to relink the executable. If you subsequently edit one or more source files, the executable can be rebuilt with the minimum number of recompilations using the same single make command.

### 3.6.4. Building a Program with IPA

Interprocedural analysis and optimization (IPA) by the NVIDIA HPC Compilers alters the standard and make utility command-level interfaces as little as possible. IPA occurs in three phases:

- ▶ **Collection:** Create a summary of each function or procedure, collecting the useful information for interprocedural optimizations. This is done during the compile step if the `-Mipa` switch is present on the command line; summary information is collected and stored in the object file.
- ▶ **Propagation:** Process all the object files to propagate the interprocedural summary information across function and file boundaries. This is done during the link step, when all the object files are combined, if the `-Mipa` switch is present on the link command line.
- ▶ **Recompile/Optimization:** Recompile each of the object files with the propagated interprocedural information, producing a specialized object file. This process is also performed during the link step when the `-Mipa` switch is present on the link command line.

When linking with `-Mipa`, the NVIDIA HPC Compilers automatically regenerate IPA-optimized versions of each object file, essentially recompiling each file. If there are IPA-optimized objects from a previous build, the compilers will minimize the recompile time by reusing those objects if they are still valid. They will still be valid if the IPA-optimized object is newer than the original object file, and the propagated IPA information for that file has not changed since it was optimized.

After each object file has been recompiled, the regular linker is invoked to build the application with the IPA-optimized object files. The IPA-optimized object files are saved in the same directory as the original object files, for use in subsequent program builds.

### 3.6.5. Building a Program with IPA – Single Step

By adding the `-Mipa` command line switch, several source files can be compiled and linked with interprocedural optimizations with one command:

```
% nvc -Mipa=fast -o a.out file1.c file2.c file3.c
```

Just like compiling without `-Mipa`, the driver executes several steps to produce the assembly and object files to create the executable:

```
% nvc -Mipa=fast -S -o file1.s file1.c
% as -o file1.o file1.s
```

```
% nvc -Mipa=fast -S -o file2.s file2.c
% as -o file2.o file2.s
% nvc -Mipa=fast -S -o file3.s file3.c
% as -o file3.o file3.s
% nvc -Mipa=fast -o a.out file1.o file2.o file3.o
```

In the last step, an IPA linker is invoked to read all the IPA summary information and perform the interprocedural propagation. The IPA linker reinvokes the compiler on each of the object files to recompile them with interprocedural information. This creates three new objects with mangled names:

```
file1_ipa5_a.out.oo.o, file2_ipa5_a.out.oo.o, file3_ipa5_a.out.oo.o
```

The system linker is then invoked to link these IPA-optimized objects into the final executable. Later, if one of the three source files is edited, the executable can be rebuilt with the same command line:

```
% nvc -Mipa=fast -o a.out file1.c file2.c file3.c
```

This works, but again has the side-effect of compiling each source file, and recompiling each object file at link time.

### 3.6.6. Building a Program with IPA – Several Steps

Just by adding the `-Mipa` command-line switch, it is possible to use individual `nvc` commands to compile each source file, followed by a command to link the resulting object files into an executable:

```
% nvc -Mipa=fast -c file1.c
% nvc -Mipa=fast -c file2.c
% nvc -Mipa=fast -c file3.c
% nvc -Mipa=fast -o a.out file1.o file2.o file3.o
```

The `nvc` driver invokes the compiler and assembler as required to process each source file, and invokes the IPA linker for the final link command. If you modify one of the source files, the executable can be rebuilt by compiling just that file and then relinking:

```
% nvc -Mipa=fast -c file1.c
% nvc -Mipa=fast -o a.out file1.o file2.o file3.o
```

When the IPA linker is invoked, it will determine that the IPA-optimized object for `file1.o` (`file1_ipa5_a.out.oo.o`) is stale, since it is older than the object `file1.o`; and hence it needs to be rebuilt, and reinvokes the compiler to generate it. In addition, depending on the nature of the changes to the source file `file1.c`, the interprocedural optimizations previously performed for `file2` and `file3` may now be inaccurate. For instance, IPA may have propagated a constant argument value in a call from a function in `file1.c` to a function in `file2.c`; if the value of the argument has changed, any optimizations based on that constant value are invalid. The IPA linker determines which, if any, of the previously created IPA-optimized objects need to be regenerated; and, as appropriate, reinvokes the compiler to regenerate them. Only those objects that are stale or which have new or different IPA information are regenerated. This approach saves compile time.



### 3.6.7. Building a Program with IPA Using Make

As shown earlier, programs can be built with IPA using the make utility. Just add the command-line switch `-Mipa`, as shown here:

```
OPT=-Mipa=fast
a.out: file1.o file2.o file3.o
    nvc $(OPT) -o a.out file1.o file2.o file3.o
file1.o: file1.c
    nvc $(OPT) -c file1.c
file2.o: file2.c
    nvc $(OPT) -c file2.c
file3.o: file3.c
    nvc $(OPT) -c file3.c
```

Using the single `make` command invokes the compiler to generate any of the object files that are out-of-date, then invokes `nvc` to link the objects into the executable. At link time, `nvc` calls the IPA linker to regenerate any stale or invalid IPA-optimized objects.

```
% make
```

### 3.6.8. Questions about IPA

**Question:** Why is the object file so large?

**Answer:** An object file created with `-Mipa` contains several additional sections. One is the summary information used to drive the interprocedural analysis. In addition, the object file contains the compiler internal representation of the source file, so the file can be recompiled at link time with interprocedural optimizations. There may be additional information when inlining is enabled. The total size of the object file may be 5-10 times its original size. The extra sections are not added to the final executable.

**Question:** What if I compile with `-Mipa` and link without `-Mipa`?

**Answer:** The NVIDIA HPC Compilers generate a legal object file, even when the source file is compiled with `-Mipa`. If you compile with `-Mipa` and link without `-Mipa`, the linker is invoked on the original object files. A legal executable is generated. While this executable does not have the benefit of interprocedural optimizations, any other optimizations do apply.

**Question:** What if I compile without `-Mipa` and link with `-Mipa`?

**Answer:** At link time, the IPA linker must have summary information about all the functions or routines used in the program. This information is created only when a file is compiled with `-Mipa`. If you compile a file without `-Mipa` and then try to get interprocedural optimizations by linking with `-Mipa`, the IPA linker will issue a message that some routines have no IPA summary information, and will proceed to run the system linker using the original object files. If some files were compiled with `-Mipa` and others were not, it will determine the safest approximation of the IPA summary information for those files not compiled with `-Mipa`, and use that to recompile the other files using interprocedural optimizations.

**Question:** Can I build multiple applications in the same directory with `-Mipa`?

**Answer:** Yes. Suppose you have three source files: `main1.c`, `main2.c`, and `sub.c`, where `sub.c` is shared between the two applications. Suppose you build the first application with `-Mipa`, using this command:

```
% nvc -Mipa=fast -o app1 main1.c sub.c
```

The IPA linker creates two IPA-optimized object files and uses them to build the first application.

```
main1_ipa4_app1.o sub_ipa4_app1.o
```

Now suppose you build the second application using this command:

```
% nvc -Mipa=fast -o app2 main2.c sub.c
```

The IPA linker creates two more IPA-optimized object files:

```
main2_ipa4_app2.o sub_ipa4_app2.o
```



There are now three object files for `sub.c`: the original `sub.o`, and two IPA-optimized objects, one for each application in which it appears.

**Question:** How is the mangled name for the IPA-optimized object files generated?

**Answer:** The mangled name has `'_ipa'` appended, followed by the decimal number of the length of the executable file name, followed by an underscore and the executable file name itself. The suffix is changed to `.oo` so that linking `*.o` does not pull in the IPA-optimized objects. If the IPA linker determines that the file would not benefit from any interprocedural optimizations, it does not have to recompile the file at link time, and uses the original object.

**Question:** Can I use parallel make environments (e.g., `pmake`) with IPA?

**Answer:** No. IPA is not compatible with parallel make environments.

# Chapter 4.

## USING FUNCTION INLINING

Function inlining replaces a call to a function or a subroutine with the body of the function or subroutine. This can speed up execution by eliminating parameter passing and function/subroutine call and return overhead. It also allows the compiler to optimize the function with the rest of the code. Note that using function inlining indiscriminately can result in much larger code size and no increase in execution speed.

The NVIDIA HPC compilers provide two categories of inlining:

- ▶ **Automatic function inlining** – In C++ and C, you can inline static functions with the `inline` keyword by using the `-Mautoinline` option, which is included with `-fast`.
- ▶ **Function inlining** – You can inline functions which were extracted to the inline libraries in Fortran, C++ and C. There are two ways of enabling function inlining: with and without the `lib` suboption. For the latter, you create inline libraries, for example using the `nvfortran` compiler driver and the `-o` and `-Mextract` options.

There are important restrictions on inlining. Inlining only applies to certain types of functions. Refer to [Restrictions on Inlining](#) for more details on function inlining limitations.

This section describes how to use the following options related to function inlining:

```
-Mautoinline  
-Mextract  
-Minline  
-Mnoinline  
-Mrecursive
```

### 4.1. Automatic function inlining in C++ and C

To enable automatic function inlining in C++ and C for static functions with the `inline` keyword, use the `-Mautoinline` option (included in `-fast`). Use `-Mnoautoinline` to disable it.

These `-Mautoinline` suboptions let you determine the selection criteria, where `n` loosely corresponds to the number of lines in the procedure:

**maxsize:n**

Automatically inline functions size `n` and less

**totalsize:n**

Limit automatic inlining to total size of `n`

## 4.2. Invoking Procedure Inlining

To invoke the procedure inliner, use the `-Minline` option. If you do not specify an inline library, the compiler performs a special prepass on all source files named on the compiler command line before it compiles any of them. This pass extracts procedures that meet the requirements for inlining and puts them in a temporary inline library for use by the compilation pass.

Several `-Minline` suboptions let you determine the selection criteria for procedures to be inlined. These suboptions include:

**except:func**

Inlines all eligible procedures except `func`, a procedure in the source text. You can use a comma-separated list to specify multiple procedure.

**[name:]func**

Inlines all procedures in the source text whose name matches `func`. You can use a comma-separated list to specify multiple procedures.

**[maxsize:]n**

A numeric option is assumed to be a size. Procedures of size `n` or less are inlined, where `n` loosely corresponds to the number of lines in the procedure. If both `n` and `func` are specified, then procedures matching the given name(s) or meeting the size requirements are inlined.

**reshape**

Fortran subprograms with array arguments are not inlined by default if the array shape does not match the shape in the caller. Use this option to override the default.

**smallsize:n**

Always inline procedures of size smaller than `n` regardless of other size limits.

**totalsize:n**

Stop inlining in a procedure when the procedure's total size with inlining reaches the `n` specified.

**[lib:]file.ext**

Instructs the inliner to inline the procedures within the library file `file.ext`. If no inline library is specified, procedures are extracted from a temporary library created during an extract prepass.



Tip Create the library file using the `-Mextract` option.

If you specify both a procedure name and a `maxsize n`, the compiler inlines procedures that match the procedure name *or* have `n` or fewer statements.

If a name is used without a keyword, then a name with a period is assumed to be an inline library and a name without a period is assumed to be a procedure name. If a number is used without a keyword, the number is assumed to be a size.

Inlining can be disabled with `-Mnoinline`.

In the following example, the compiler inlines procedures with fewer than approximately 100 statements in the source file `myprog.f` and writes the executable code in the default output file `a.out`.

```
$ nvfortran -Minline=maxsize:100 myprog.f
```

## 4.3. Using an Inline Library

If you specify one or more inline libraries on the command line with the `-Minline` option, the compiler does not perform an initial extract pass. The compiler selects functions to inline from the specified inline library. If you also specify a size or function name, all functions in the inline library meeting the selection criteria are selected for inline expansion at points in the source text where they are called.

If you do not specify a function name or a size limitation for the `-Minline` option, the compiler tries to inline every function in the inline library that matches a function in the source text.

In the following example, the compiler inlines the function `proc` from the inline library `lib.il` and writes the executable code in the default output file `a.out`.

```
$ nvfortran -Minline=name:proc,lib:lib.il myprog.f
```

The following command line is equivalent to the preceding line, with the exception that in the following example does not use the keywords `name:` and `lib:`. You typically use keywords to avoid name conflicts when you use an inline library name that does not contain a period. Otherwise, without the keywords, a period informs the compiler that the file on the command line is an inline library.

```
$ nvfortran -Minline=proc,lib.il myprog.f
```

## 4.4. Creating an Inline Library

You can create or update an inline library using the `-Mextract` command-line option. If you do not specify selection criteria with the `-Mextract` option, the compiler attempts to extract all procedures.

Several `-Mextract` options let you determine the selection criteria for creating or updating an inline library. These selection criteria include:

### **func**

Extracts the procedure `func`. you can use a comma-separated list to specify multiple procedures.

### **[name:]func**

Extracts the procedure whose name matches `func`, a procedure in the source text.

**[size:]n**

Limits the size of the extracted procedures to those with a statement count less than or equal to `n`, the specified size.



The size `n` may not exactly equal the number of statements in a selected procedure; the size parameter is merely a rough gauge.

**[lib:]ext.lib**

Stores the extracted information in the library directory `ext.lib`.

If no inline library is specified, procedures are extracted to a temporary library created during an extract prepass for use during the compilation stage.

When you use the `-Mextract` option, only the extract phase is performed; the compile and link phases are not performed. The output of an extract pass is a library of procedures available for inlining. This output is placed in the inline library file specified on the command line with the `-o` filename specification. If the library file exists, new information is appended to it. If the file does not exist, it is created. You can use a command similar to the following:

```
$ nvfortran -Mextract=lib:lib.il myfunc.f
```

You can use the `-Minline` option with the `-Mextract` option. In this case, the extracted library of procedures can have other procedures inlined into the library. Using both options enables you to obtain more than one level of inlining. In this situation, if you do not specify a library with the `-Minline` option, the inline process consists of two extract passes. The first pass is a hidden pass implied by the `-Minline` option, during which the compiler extracts procedures and places them into a temporary library. The second pass uses the results of the first pass but puts its results into the library that you specify with the `-o` option.

### 4.4.1. Working with Inline Libraries

An inline library is implemented as a directory with each inline function in the library stored as a file using an encoded form of the inlinable function.

A special file named `TOC` in the inline library directory serves as a table of contents for the inline library. This is a printable, ASCII file which you can examine to locate information about the library contents, such as names and sizes of functions, the source file from which they were extracted, the version number of the extractor which created the entry, and so on.

Libraries and their elements can be manipulated using ordinary system commands.

- ▶ Inline libraries can be copied or renamed.
- ▶ Elements of libraries can be deleted or copied from one library to another.
- ▶ The `ls` or `dir` command can be used to determine the last-change date of a library entry.

## 4.4.2. Dependencies

When a library is created or updated using one of the NVIDIA HPC compilers, the last-change date of the library directory is updated. This allows a library to be listed as a dependence in a makefile and ensures that the necessary compilations are performed when a library is changed.

## 4.4.3. Updating Inline Libraries – Makefiles

If you use inline libraries you must be certain that they remain up-to-date with the source files into which they are inlined. One way to assure inline libraries are updated is to include them in a makefile.

The makefile fragment in the following example assumes the file `utils.f` contains a number of small functions used in the files `parser.f` and `alloc.f`.

This portion of the makefile:

- ▶ Maintains the inline library `utils.il`.
- ▶ Updates the library whenever you change `utils.f` or one of the include files it uses.
- ▶ Compiles `parser.f` and `alloc.f` whenever you update the library.

### Sample Makefile

```
SRC = mydir
FC = nvfortran
FFLAGS = -O2
main.o: $(SRC)/main.f $(SRC)/global.h
    $(FC) $(FFLAGS) -c $(SRC)/main.f
utils.o: $(SRC)/utils.f $(SRC)/global.h $(SRC)/utils.h
    $(FC) $(FFLAGS) -c $(SRC)/utils.f
utils.il: $(SRC)/utils.f $(SRC)/global.h $(SRC)/utils.h
    $(FC) $(FFLAGS) -Mextract=15 -o utils.il $(SRC)/utils.f
parser.o: $(SRC)/parser.f $(SRC)/global.h utils.il
    $(FC) $(FFLAGS) -Minline=utils.il -c $(SRC)/parser.f
alloc.o: $(SRC)/alloc.f $(SRC)/global.h utils.il
    $(FC) $(FFLAGS) -Minline=utils.il -c $(SRC)/alloc.f
myprog: main.o utils.o parser.o alloc.o
    $(FC) -o myprog main.o utils.o parser.o alloc.o
```

## 4.5. Error Detection during Inlining

You can specify the `-Minfo=inline` option to request inlining information from the compiler when you invoke the inliner. For example:

```
$ nvfortran -Minline=mylib.il -Minfo=inline myext.f
```

## 4.6. Examples

Assume the program `dhry` consists of a single source file `dhry.f`. The following command line builds an executable file for `dhry` in which `proc7` is inlined wherever it is called:

```
$ nvfortran dhry.f -Minline=proc7
```

The following command lines build an executable file for `dhry` in which `proc7` plus any functions of approximately 10 or fewer statements are inlined (one level only).



The specified functions are inlined only if they are previously placed in the inline library, `temp.il`, during the extract phase.

```
$ nvfortran dhry.f -Mextract=lib:temp.il
$ nvfortran dhry.f -Minline=10,proc7,temp.il
```

Using the same source file `dhry.f`, the following example builds an executable for `dhry` in which all functions of roughly ten or fewer statements are inlined. Two levels of inlining are performed. This means that if function A calls function B, and B calls C, and both B and C are inlinable, then the version of B which is inlined into A will have had C inlined into it.

```
$ nvfortran dhry.f -Minline=maxsize:10
```

## 4.7. Restrictions on Inlining

The following Fortran subprograms cannot be extracted:

- ▶ Main or BLOCK DATA programs.
- ▶ Subprograms containing alternate return, assigned GO TO, DATA, SAVE, or EQUIVALENCE statements.
- ▶ Subprograms containing FORMAT statements.
- ▶ Subprograms containing multiple entries.

A Fortran subprogram is not inlined if any of the following applies:

- ▶ It is referenced in a statement function.
- ▶ A common block mismatch exists; in other words, the caller must contain all common blocks specified in the callee, and elements of the common blocks must agree in name, order, and type (except that the caller's common block can have additional members appended to the end of the common block).
- ▶ An argument mismatch exists; in other words, the number and type (size) of actual and formal parameters must be equal.
- ▶ A name clash exists, such as a call to subroutine `xyz` in the extracted subprogram and a variable named `xyz` in the caller.

The following types of C and C++ functions cannot be inlined:

- ▶ Functions which accept a variable number of arguments



Certain C/C++ functions can only be inlined into the file that contains their definition:

- ▶ Static functions
- ▶ Functions which call a static function
- ▶ Functions which reference a static variable

# Chapter 5.

## USING GPUS

An NVIDIA GPU can be used as an accelerator to which a CPU can offload data and executable kernels to perform compute-intensive calculations. This section gives an overview of options for programming NVIDIA GPUs with NVIDIA's HPC Compilers and covers topics that affect GPU programming when using one or more of the GPU programming models.

### 5.1. Overview

With the NVIDIA HPC Compilers you can program NVIDIA GPUs using certain standard language constructs, OpenACC directives, OpenMP directives, or CUDA Fortran language extensions. GPU programming with standard language constructs or directives allows you to create high-level GPU-accelerated programs without the need to explicitly initialize the GPU, manage data or program transfers between the host and GPU, or initiate GPU startup and shutdown. Rather, all of these details are implicit in the programming model and are managed by the NVIDIA HPC SDK Fortran, C++ and C compilers. GPU programming with CUDA extensions gives you access to all NVIDIA GPU features and full control over data management and offloading of compute-intensive loops and kernels.

The NVCC++ compiler supports automatic offload of C++17 Parallel Algorithms invocations to NVIDIA GPUs under control of the `-stdpar` compiler option. See the Blog post *Accelerating Standard C++ with GPUs* for details on using this feature. The NVFORTRAN compiler supports automatic offload to NVIDIA GPUs of certain Fortran array intrinsics and patterns of array syntax, including use of Volta and Ampere architecture Tensor Cores for appropriate intrinsics. See the Blog post *Bringing Tensor Cores to Standard Fortran* for details on using this feature.

The NVFORTRAN compiler supports CUDA programming in Fortran. See the *NVIDIA CUDA Fortran Programming Guide* for complete details on how to use CUDA Fortran. The NVCC compiler supports CUDA programming in C and C++ in combination with a host C++ compiler on your system. See the *CUDA C++ Programming Guide* for an introduction and overview of how to use NVCC and CUDA C++.

The NVFORTRAN, NVCC++ and NVCC compilers all support directive-based programming of NVIDIA GPUs using OpenACC. OpenACC is an accelerator

programming model that is portable across operating systems and various host CPUs and types of accelerators, including both NVIDIA GPUs and multicore CPUs. OpenACC directives allow a programmer to migrate applications incrementally to accelerator targets using standards-compliant Fortran, C++ or C that remains completely portable to other compilers and systems. It allows the programmer to augment information available to the compilers, including specification of data local to an accelerator region, guidance on mapping of loops onto an accelerator, and similar performance-related details.

The NVFORTRAN, NVC++, and NVC compilers support a subset of the OpenMP Application Program Interface for CPUs and GPUs. OpenMP applications properly structured for GPUs, meaning they expose massive parallelism and have relatively little or no synchronization in GPU-side code segments, should compile and execute with performance on par with or close to equivalent OpenACC. Codes that are not well-structured for GPUs may perform poorly but should execute correctly.

In user-directed accelerator programming the user specifies the regions of a host program to be targeted for offloading to an accelerator. The bulk of a user's program, as well as regions containing constructs that are not supported on the targeted accelerator, are executed on the host.

## 5.2. Terminology

Clear and consistent terminology is important in describing any programming model. This section provides definitions of the terms required for you to effectively use this section and the associated programming model.

### **Accelerator**

a parallel processor, such as a GPU or a CPU running in multicore mode, to which a CPU can offload data and executable kernels to perform compute-intensive calculations.

### **Compute intensity**

for a given loop, region, or program unit, the ratio of the number of arithmetic operations performed on computed data divided by the number of memory transfers required to move that data between two levels of a memory hierarchy.

### **Compute region**

a structured block defined by a compute construct. A *compute construct* is a structured block containing loops which are compiled for the accelerator. A compute region may require device memory to be allocated and data to be copied from host to device upon region entry, and data to be copied from device to host memory and device memory deallocated upon exit. The dynamic range of a compute construct, including any code in procedures called from within the construct, is the compute region. In this release, compute regions may not contain other compute regions or data regions.

### **Construct**

a structured block identified by the programmer or implicitly defined by the language. Certain actions may occur when program execution reaches the start and end of a construct, such as device memory allocation or data movement between the host and device memory. Loops in a compute construct are targeted for execution on

the accelerator. The dynamic range of a construct including any code in procedures called from within the construct, is called a *region*.

### **CUDA**

stands for Compute Unified Device Architecture; CUDA C++ and Fortran language extensions and API calls can be used to explicitly control and program an NVIDIA GPU.

### **Data region**

a region defined by a data construct, or an implicit data region for a function or subroutine containing directives. Data regions typically require device memory to be allocated and data to be copied from host to device memory upon entry, and data to be copied from device to host memory and device memory deallocated upon exit. Data regions may contain other data regions and compute regions.

### **Device**

a general reference to any type of accelerator.

### **Device memory**

memory attached to an accelerator which is physically separate from the host memory.

### **Directive**

in C, a `#pragma`, or in Fortran, a specially formatted comment statement that is interpreted by a compiler to augment information about or specify the behavior of the program.

### **DMA**

Direct Memory Access, a method to move data between physically separate memories; this is typically performed by a DMA engine, separate from the host CPU, that can access the host physical memory as well as an IO device or GPU physical memory.

### **GPU**

a Graphics Processing Unit; one type of accelerator device.

### **Host**

the main CPU that in this context has an attached accelerator device. The host CPU controls the program regions and data loaded into and executed on the device.

### **Loop trip count**

the number of times a particular loop executes.

### **Private data**

with respect to an iterative loop, data which is used only during a particular loop iteration. With respect to a more general region of code, data which is used within the region but is not initialized prior to the region and is re-initialized prior to any use after the region.

### **Region**

the dynamic range of a construct, including any procedures invoked from within the construct.

### **Structured block**

in C++ or C, an executable statement, possibly compound, with a single entry at the top and a single exit at the bottom. In Fortran, a block of executable statements with a single entry at the top and a single exit at the bottom.

### **Vector operation**

a single operation or sequence of operations applied uniformly to each element of an array.

**Visible device copy**

a copy of a variable, array, or subarray allocated in device memory, that is visible to the program unit being compiled.

## 5.3. Execution Model

The execution model targeted by the NVIDIA HPC Compilers is host-directed execution with an attached accelerator device, such as a GPU. The bulk of a user application executes on the host. Compute intensive regions are offloaded to the accelerator device under control of the host. The accelerator device executes kernels, which may be as simple as a tightly-nested loop, or as complex as a subroutine, depending on the accelerator hardware.

### 5.3.1. Host Functions

Even in accelerator-targeted regions, the host must orchestrate the execution; it

- ▶ allocates memory on the accelerator device
- ▶ initiates data transfer
- ▶ sends the kernel code to the accelerator
- ▶ passes kernel arguments
- ▶ queues the kernel
- ▶ waits for completion
- ▶ transfers results back to the host
- ▶ deallocates memory



In most cases, the host can queue a sequence of kernels to be executed on the device, one after the other.

## 5.4. Memory Model

The most significant difference between a *host-only* program and a *host+accelerator* program is that the memory on the accelerator can be completely separate from host memory, which is the case on many GPUs. For example:

- ▶ The host cannot read or write accelerator memory directly because it is not mapped into the virtual memory space of the host.
- ▶ All data movement between host memory and accelerator memory must be performed by the host through runtime library calls that explicitly move data between the separate memories.
- ▶ In general it is not valid for the compiler to assume the accelerator can read or write host memory directly. This is well-defined starting with the OpenACC 2.7 and OpenMP 5.0 specifications.

The systems with the latest GPUs provide a unified single address space between CPU and GPU for some or all memory regions, as detailed in the [Managed and Unified](#)

**Memory Modes** subsection below. In these systems data can be accessed from host and accelerator subprograms without the need for explicit data movement.

The NVIDIA HPC Compilers support the following system memory modes:

Table 12 GPU Memory Modes

Memory Mode	Description	Compiler flags
Separate	All data accessed in host and accelerator programs are in separate (CPU and GPU) memories. Data in the application need to be physically moved between CPU and GPU memory either by adding explicit annotations or by relying on a compiler to detect and migrate the data.	<b>-gpu=mem:separate</b>
Managed	Dynamically allocated host data are placed in CUDA Managed Memory which is a unified single address space between host and accelerator programs and can therefore be accessed on device without explicit data movement. All other data (host, stack, or global data) remain in separate memory.	<b>-gpu=mem:managed</b>
Unified	All host data are placed in a unified single address space between the host and accelerator subprograms; no explicit data movements are required. This mode is intended for targets with full CUDA Unified Memory capability and it may utilize CUDA Managed Memory for dynamic allocations.	<b>-gpu=mem:unified</b>

If the memory mode is not selected explicitly by passing one of the above **-gpu=mem:\*** options, the compiler selects a default memory mode. The default memory mode for Stdpar is explained in [Using Stdpar](#). When Stdpar is not enabled, the default memory mode is Separate Memory. Memory modes may have specific semantics in each programming language and the compilers can sometimes implicitly determine the data movement that's required. More details can be found in the subsections of each programming model.

The following options **-gpu=[no]managed**, **-gpu=[no]unified** and **-gpu=pinned** are deprecated but still accepted. Refer to [Command-line Options Selecting Compiler Memory Modes](#) for compatibility between the current and deprecated memory specific flags.

The compiler implicitly defines the following macros corresponding to the memory mode it compiles for:

- ▶ When the code is compiled for Separate Memory Mode, the compiler defines **\_\_NVCOMPILER\_GPU\_SEPARATE\_MEM** macro.
- ▶ When the code is compiled for Managed Memory Mode, the compiler defines **\_\_NVCOMPILER\_GPU\_MANAGED\_MEM** macro.
- ▶ When the code is compiled for Unified Memory Mode, the compiler defines **\_\_NVCOMPILER\_GPU\_UNIFIED\_MEM** macro. If CUDA Managed Memory is utilised, the compiler defines additionally **\_\_NVCOMPILER\_GPU\_MANAGED\_MEM**.

When a binary is compiled for one memory mode it may need to be run on a system with specific memory capabilities as follows:

- ▶ Applications compiled for Separate Memory Mode can run on any CUDA platforms.
- ▶ Applications compiled for Managed Memory Mode must be run on platforms with CUDA Managed Memory or full CUDA Unified Memory capabilities.
- ▶ Applications compiled for Unified Memory Mode must be run on platforms with full CUDA Unified Memory.



Memory allocated in the accelerator subprogram can't be accessed or deallocated from the host.

### 5.4.1. Separate Host and Accelerator Memory Considerations

The programmer must be aware of the potentially separate memories for many reasons, including but not limited to:

- ▶ Memory bandwidth between host memory and accelerator memory determines the compute intensity required to effectively accelerate a given region of code.
- ▶ Limited size of accelerator memory may prohibit offloading of regions of code that operate on very large amounts of data.

#### 5.4.1.1. Accelerator Memory

On the accelerator side, current GPUs implement a weak memory model. In particular, they do not support memory coherence between threads unless those threads are parallel only at the synchronous level and the memory operations are separated by an explicit barrier. Otherwise, if one thread updates a memory location and another reads the same location, or two threads store a value to the same location, the hardware does not guarantee the results. While the results of running such a program might be inconsistent, it is not accurate to say that the results are incorrect. By definition, such programs are defined as being in error. While a compiler can detect some potential errors of this nature, it is nonetheless possible to write an accelerator region that produces inconsistent numerical results.

Stack data in accelerator subprograms are allocated per thread. Stack data from one thread are not accessible by the other threads.

#### 5.4.1.2. Staging Memory Buffer

Memory transfers between the accelerator and host may not always be asynchronous with respect to the host, even if the chosen programming model (for instance, OpenACC) declares that. This limitation may be due to the specific GPU and host memory architectures.

In order to help the host program proceed while a memory transfer to or from the accelerator is underway, the NVIDIA HPC Compilers Runtime maintains a designated staging memory area, also known as a pinned buffer. This memory area is registered with the CUDA API, which makes it suitable for asynchronous memory transfers between the GPU and the host. When an asynchronous memory transfer is started,

the data being transferred is staged through the pinned buffer. Multiple asynchronous operations on the same data can be issued - in that case, the runtime system will operate on the data staged in the pinned buffer, not on the original host memory. When the host program issues an explicit or implicit synchronization request, the data is moved from the pinned buffer to its destination transparently to the application.

The runtime has the discretion to enable or disable the pinned buffer depending on the host and GPU memory architecture. Also, the size of the pinned buffer is determined by the runtime system as appropriate. The user can control some of these decisions using environment variables at the start of the application. Please refer to [Environment Variables Controlling Device Memory Management](#) to learn more.

### 5.4.1.3. Cache Management

Some current GPUs have a software-managed cache, some have hardware-managed caches, and most have hardware caches that can be used only in certain situations and are limited to read-only data. In low-level programming models such as CUDA, it is up to the programmer to manage these caches. The OpenACC programming model provides directives the programmer can use as hints to the compiler for cache management.

### 5.4.1.4. Environment Variables Controlling Device Memory Management

This section summarizes the environment variables that NVIDIA HPC Compilers use to control device memory management.

The following table contains the environment variables that are currently supported and provides a brief description of each.

Table 13 Memory Management Environment Variables

Environment Variable	Use
NVCOMPILER_ACC_BUFFERSIZE	For NVIDIA CUDA devices, this defines the size of the pinned buffer used to transfer data between host and device.
NVCOMPILER_ACC_CUDA_CTX_SCHED	For NVIDIA CUDA devices, sets flags to be used when creating a new CUDA context. By default, the <b>CU_CTX_SCHED_YIELD</b> flag is used. Please refer to the <a href="#">CUDA Toolkit Documentation</a> for the detailed description of the <b>cuCtxCreate</b> function and the possible flag values.
NVCOMPILER_ACC_CUDA_HEAPSIZE	For NVIDIA CUDA devices, sets the heap size limit for <b>malloc()</b> when called on device.
NVCOMPILER_ACC_CUDA_MAX_L2_FETCH_GRANUL	For NVIDIA CUDA devices, sets the maximum L2 cache fetch granularity size in bytes. A correct value is an integer between 0 and 128.
NVCOMPILER_ACC_CUDA_MEMALLOCSYNC	For NVIDIA CUDA devices, when set to a non-zero integer value, enables CUDA asynchronous memory allocations from the default CUDA memory pool as described in the <a href="#">CUDA Toolkit Documentation</a> . By



Environment Variable	Use
	default, an internal NVIDIA HPC Runtime memory pool is used instead.
NVCOMPILER_ACC_CUDA_MEMALLOCASYNC_POOL	For NVIDIA CUDA devices, sets the size of the default CUDA memory pool for asynchronous allocations if the <b>NVCOMPILER_ACC_CUDA_MEMALLOCASYNC</b> environment variable is also set to a non-zero integer value.
NVCOMPILER_ACC_CUDA_NOCOPY	Disables the use of the pinned buffer when transferring user data between host and NVIDIA CUDA devices. When this variable is set to a non-zero integer value, user data will be transferred directly bypassing the pinned buffer. Asynchronous execution of such data transfers can be limited when this setting is in effect.
NVCOMPILER_ACC_CUDA_PIN	For NVIDIA CUDA devices, enables host memory pinning at data directives. When host memory is pinned, data transfers to and from the device can be asynchronous, which can potentially improve program performance. A non-zero integer value enables this mechanism. A value of <b>2</b> or greater additionally disallows unpinning the host data after it is pinned. A value of <b>3</b> or greater also enables pinning the whole array referenced in a data directive (provided that the size of the array is known), rather than its subarray specified in the data directive. By default, host data referenced at data directives is not pinned unless directed by the compiler at compile-time; refer to <a href="#">Command-line Options Selecting Compiler Memory Modes</a> for more information about the compile-time memory modes.
NVCOMPILER_ACC_CUDA_PINSIZE	For NVIDIA CUDA devices, sets the host memory pinning granularity. If host memory pinning is enabled with the <b>NVCOMPILER_ACC_CUDA_PIN</b> environment variable, the runtime will attempt to use this setting to pin larger regions of memory at once, thus potentially lowering the cost of pinning memory when the program needs to pin multiple data regions separately. The maximum allowed value is 1 MB. By default, single byte pinning granularity is used.
NVCOMPILER_ACC_CUDA_PRINTFIFOSIZE	For NVIDIA CUDA devices, sets the buffer size for formatted output calls on device. In particular, it controls the buffer size for the <b>printf</b> C function.
NVCOMPILER_ACC_CUDA_STACKSIZE	For NVIDIA CUDA devices, sets the stack size limit for device threads.
NVCOMPILER_ACC_DEV_MEMORY	For NVIDIA CUDA devices, when set to a valid non-zero size value, enables the use of a device memory pool and sets its size. By default, the device memory pool is not used.
NVCOMPILER_ACC_MEM_MANAGE	For NVIDIA CUDA devices, when set to the integer value 0, disables the use of an internal

Environment Variable	Use
	device memory manager. By default, the device memory manager is enabled. It maintains a list of deallocated chunks of device memory in an attempt to efficiently reuse them for future allocations.

## 5.4.2. Managed and Unified Memory Modes

The NVIDIA HPC Compilers support interoperability with [CUDA Unified Memory](#).

This feature is available with the x86-64, OpenPOWER and Arm Server compilers.

Unified memory provides a single address space for CPU and GPU; data movement between CPU and GPU memories is implicitly handled by the NVIDIA CUDA driver.

Whenever data is accessed on the CPU or the GPU, it could trigger a data transfer if the last time it was accessed was not on the same device. In some cases, page thrashing may occur and impact performance. An introduction to CUDA Unified Memory is available on [Parallel Forall](#).

### 5.4.2.1. Managed Memory Mode

In Managed Memory Mode, all Fortran, C++ and C explicit allocation statements (e.g. **allocate**, **new**, and **malloc**, respectively) in a program unit are replaced by equivalent CUDA managed data allocation calls that place the data in CUDA Managed Memory. The result is that OpenACC and OpenMP data clauses and directives are not needed to manage data movement. They are essentially ignored and can be omitted. For Stdpar this is the minimal required memory mode since there are no specific annotations for data used in the parallel region.

To enable Managed Memory Mode, add the option **-gpu=mem:managed** to the compiler and linker command lines.

When a program allocates managed memory, it allocates host pinned memory as well as device memory thus making allocate and free operations somewhat more expensive and data transfers somewhat faster. A memory pool allocator is used to mitigate the overhead of the allocate and deallocate operations. More details can be found in [Memory Pool Allocator](#).

Managed Memory Mode has the following limitations:

- ▶ Use of managed memory applies only to dynamically-allocated data.
- ▶ Given an allocatable aggregate with a member that points to local, global, or static data, compiling with **-gpu=mem:managed** and attempting to access memory through that pointer from the compute kernel will cause a failure at runtime.
- ▶ C++ virtual functions are not supported.
- ▶ The **-gpu=mem:managed** compiler option must be used to compile the files in which variables (accessed from GPU) are allocated, even if there is no code to accelerate on the GPU in the source file.

- ▶ When linking multiple translation units, the application must ensure that all data are deallocated using the scheme corresponding to their allocation. For example if the data are allocated in managed memory the deallocation must be performed using CUDA API calls for managed memory. More details and extra compiler support is detailed in [Interception of Deallocations](#).

Managed Memory Mode has the following additional limitations when used with NVIDIA Kepler GPUs:

- ▶ Data motion on Kepler GPUs is achieved through fast pinned asynchronous data transfers; from the program's perspective, however, the transfers are synchronous.
- ▶ The NVIDIA HPC Compiler Runtime enforces synchronous execution of kernels when **-gpu=mem:managed** is used on a system with a Kepler GPU. This situation may result in slower performance because of the extra synchronizations and decreased overlap between CPU and GPU.
- ▶ The total amount of managed memory is limited to the amount of available device memory on Kepler GPUs.

#### Memory Allocations/Deallocations Automatically Changed to Managed Memory

When the compiler utilizes CUDA Managed Memory capability either with **-gpu=mem:managed** or **-gpu=mem:unified**, the following explicit allocations/deallocations are automatically changed into **cudaMallocManaged/cudaFree**-type allocations/deallocations:

- ▶ For C++:
  - ▶ All calls to global **operator new** and **operator delete** that allocate or deallocate memory, such as:
 

```
operator new(std::size_t size)
operator new(std::size_t size, const std::nothrow_t &nothrow_value)
operator new(std::size_t size, std::align_val_t align)
operator new(std::size_t size, std::align_val_t align, const
  std::nothrow_t &nothrow_value)
operator delete(void *p)
operator delete(void *p, std::size_t size)
operator delete(void *p, std::align_val_t align)
operator delete(void *p, std::size_t size, std::align_val_t align)
operator delete(void *p, const std::nothrow_t &nothrow_value)
operator delete(void *p, std::align_val_t align, const std::nothrow_t
  &nothrow_value)
```
  - ▶ All the array forms of the above overloads.
  - ▶ All calls to **malloc/free** functions.
- ▶ For C: all calls to **malloc/free** functions.
- ▶ For Fortran:
  - ▶ All allocations of automatic arrays.
  - ▶ all **allocate/deallocate** statements with allocatable arrays or pointer variables.

### 5.4.2.2. Unified Memory Mode

In Unified Memory Mode, the requirements for the program are further relaxed compared to Managed Memory Mode. Specifically, not only is dynamically allocated system memory accessible on the GPU, but global and local memory are also accessible.

To enable this feature, add the option `-gpu=mem:unified` to the compiler and linker command lines.

Programs compiled with `-gpu=mem:unified` must be run on systems that support full CUDA Unified Memory capability. At this time, full CUDA Unified Memory is supported on NVIDIA Grace Hopper Superchip systems and Linux x86-64 systems running with the Heterogeneous Memory Management (HMM) feature enabled in the Linux kernel. Details about these platforms are available in the following blog posts on the NVIDIA website: [Simplifying GPU Programming for HPC with NVIDIA Grace Hopper Superchip](#) and [Simplifying GPU Application Development with Heterogeneous Memory Management](#).

In Unified Memory Mode, the compiler assumes that any system memory is accessible on the GPU. Even so, the compiler may generate managed memory allocations for explicit data allocations when it considers them beneficial for program performance. If you would like to enforce or prohibit the use of managed memory for dynamic allocations pass `-gpu=mem:unified:[no]managedalloc` to compilation and linking.

Unified Memory Mode has the following limitations:

- ▶ Unified memory support for OpenACC, OpenMP and Stdpar Fortran is not mix-and-match; all object files containing OpenACC/OpenMP directives or Fortran **DO CONCURRENT** constructs must be compiled and linked with `-gpu=mem:unified` to ensure correct execution.
- ▶ C++ virtual functions are not supported.

#### Transitioning to Unified Memory Mode

Applications transitioning to architectures that support Unified Memory Mode can be recompiled with `-gpu=mem:unified` without any code modifications.

The programmer should be aware that in Unified Memory Mode, the whole program state becomes essentially shared between the CPU and the GPU. By implication, modifications to program variables made on the GPU are visible on the CPU. That is, the GPU does not operate on a copy of the data even if the program contains respective directives, but instead the GPU operates directly on the data in system memory. To understand the importance of this idea, consider the following OpenACC C program:

```
int x[N];
void foo() {
    #pragma acc enter data create(x[0:N])
    #pragma acc parallel loop
    for (int i = 0; i < N; i++) {
```

```

    x[i] = i;
  }
}

```

When compiled in Separate Memory Mode, in the `foo()` function a copy of the array `x` is created in GPU memory and initialized as written in the `loop` construct. When `-gpu=mem:unified` is added, however, the compiler ignores the `acc enter data` construct, and the `loop` construct initializes the array `x` in system memory.

Another implication of which to be aware, *asynchronous* code execution on the GPU can introduce race conditions over access to program data. More details about code patterns to avoid when writing application sources for Unified Memory Mode can be found in the sections about specific programming models of this guide e.g. OpenACC, OpenMP, or CUDA Fortran.

### 5.4.3. Memory Pool Allocator

Dynamic memory allocations may be made using `cudaMallocManaged()`, a routine which has higher overhead than allocating non-managed memory using `cudaMalloc()`. The more calls to `cudaMallocManaged()`, the more significant the impact on performance.

To mitigate the overhead of `cudaMallocManaged()` or other CUDA allocation API calls, there is a pool allocator enabled by default in the presence of the `-gpu=mem:managed`, `-gpu=mem:separate:pinnedalloc`, or `-gpu=mem:unified` compiler options. It can be disabled, or its behavior modified, using these environment variables:

Table 14 Pool Allocator Environment Variables

Environment Variable	Use
NVCOMPILER_ACC_POOL_ALLOC	Disable the pool allocator. The pool allocator is enabled by default; to disable it, set NVCOMPILER_ACC_POOL_ALLOC to 0.
NVCOMPILER_ACC_POOL_SIZE	Set the of the pool. The default size is 1GB but other sizes (i.e., 2GB, 100MB, 500KB, etc.) can be used. The actual pool size is set such that the size is the nearest, smaller number in the Fibonacci series compared to the provided or default size. If necessary, the pool allocator will add more pools but only up to the NVCOMPILER_ACC_POOL_THRESHOLD value.
NVCOMPILER_ACC_POOL_ALLOC_MAXSIZE	Set the maximum size for allocations. The default maximum size for allocations is 500MB but another size (i.e., 100KB, 10MB, 250MB, etc.) can be used as long as it is greater than or equal to 16B.

Environment Variable	Use
NVCOMPILER_ACC_POOL_ALLOC_MINSIZE	Set the minimum size for allocation blocks. The default size is 128B but other sizes can be used. The size must be greater than or equal to 16B.
NVCOMPILER_ACC_POOL_THRESHOLD	Set the percentage of total device memory that the pool allocator can occupy. Values from 0 to 100 are accepted. The default value is 50, corresponding to 50% of device memory.



Note that where the size is specified if the unit suffix (B, KB, MB or GB) is omitted, the value is set by default in bytes.

### 5.4.4. Interception of Deallocations

While NVIDIA HPC Compilers facilitate the use of managed or pinned memory automatically, the application must ensure that memory is deallocated using the API which "matches" the API used to allocate said memory. For example, if **cudaMallocManaged** is used to allocate, then **cudaFree** must be used to deallocate; if **cudaMallocHost** is used for allocations, **cudaFreeHost** must be used for deallocations. Understanding this requirement is particularly important when third party or standard libraries are used; these libraries may have been compiled without any memory mode settings which sets up a situation where the deallocation routines in the libraries may not match the allocations made. When data is deallocated with an unmatching API call, the application may exhibit undefined behavior including crashing. To mitigate this issue, the compiler supports an interception mode in which calls to the standard deallocation function (e.g. free in C, delete in C++, or deallocate in Fortran) are inspected by the runtime and, if the memory is not detected as being system-allocated, the runtime replaces the standard deallocation function with the deallocation API corresponding to the allocation scheme in use. To activate this interception mode, use the **-gpu=interceptdeallocations** compiler flag. The interception is enabled by default for Stdpar in the presence of managed memory allocations. To deactivate the interception use the **-gpu=nointerceptdeallocations** compiler switch. This interception can incur extra runtime overhead.

### 5.4.5. Command-line Options Selecting Compiler Memory Modes

The following table maps the new memory model flags to their deprecated equivalents.

Table 15 Command-line Options Corresponding to Compiler Memory Modes

Current Flags	Deprecated Flags	Brief Description
<code>-gpu=mem:managed</code>	<code>-gpu=managed</code>	Managed Memory Mode
<code>-gpu=mem:managed -stdpar</code>	<code>-gpu=nounified -stdpar</code>	Managed Memory Mode
<code>-gpu=mem:unified</code>	<code>-gpu=unified</code>	Unified Memory Mode
<code>-gpu=mem:unified:managedalloc</code>	<code>-gpu=unified,managed</code>	Unified Memory Mode, all dynamically allocated data are implicitly in CUDA Managed Memory.
<code>-gpu=mem:unified:nomanagedalloc</code>	<code>-gpu=unified,nomanaged</code>	Unified Memory Mode, CUDA Managed Memory is not used implicitly.
<code>-gpu=mem:separate</code>	<code>-gpu=nomanaged</code>	Separate Memory Mode
<code>-gpu=mem:separate</code>	<code>-gpu=nounified</code>	Separate Memory Mode
<code>-gpu=mem:separate</code>	<code>-gpu=nomanaged,nounified</code>	Separate Memory Mode
<code>-gpu=mem:separate:pinnedalloc</code>	<code>-gpu=pinned</code>	Separate Memory Mode, dynamically allocated data are in CPU pinned memory implicitly.

## 5.5. Fortran pointers in device code

A Fortran pointer variable is implemented with a pointer and a descriptor, where the descriptor (often called a "dope vector") holds the array bounds and strides for each dimension, among other information, such as the size for each element and whether the pointer is associated. A Fortran scalar pointer has no bounds information, but does have a minimal descriptor. In Fortran, referring to the pointer variable always refers to the pointer target. There is no syntax to explicitly refer to the pointer and descriptor that implement the pointer variable.

Fortran allocatable arrays and variables are implemented much the same way as pointer arrays and variables. Much of the discussion below applies both to allocatables and pointers.

In OpenACC and OpenMP, when a pointer variable reference appears in a data clause, it's the pointer target that gets allocated or moved to device memory. The pointer and descriptor are neither allocated nor moved.

When a pointer variable is declared in a module declaration section and appears in an `!$acc declare create()` or `!$omp declare target to()` directive, then the

pointer and descriptor are statically allocated in device memory. When the pointer variable appears in a data clause, the pointer target is allocated or copied to the device, and the pointer and descriptor are 'attached' to the device copy of the data. If the pointer target is already present in device memory, no new memory is allocated or copied, but the pointer and descriptor are still 'attached', making the pointer valid in device memory. An important side effect of adding **declare create** in the module declaration section is that when the program executes an 'allocate' statement for the pointer (or allocatable), memory is allocated in both CPU and device memory. This means the newly allocated data is already present in device memory. To get values from CPU to device memory or back, you'll have to use **update** directives.

When a pointer variable is used in an OpenACC or OpenMP compute construct, the compiler creates a private copy of the pointer and descriptor for each thread, unless the pointer variable was in a module as described above. The private pointer and descriptor will contain information about the device copy of the pointer target. In the compute construct, the pointer variables may be used pretty much as they can in host code outside a compute construct. However, there are some limitations. The program can do a pointer assignment to the pointer, changing the pointer, but that will only change the private pointer for that thread. The modified pointer in the compute construct will not change the corresponding pointer and descriptor in host memory.

## 5.6. Calling routines in a compute kernel

Using explicit interfaces is a common occurrence when writing Fortran applications. Here are some cases where doing so is required for GPU programming.

- ▶ Explicit interfaces are required when using OpenACC **routine bind** or OpenMP **declare variant**.
- ▶ Fortran **do concurrent** requires routines to be **pure** which creates the need for an explicit interface.

## 5.7. Supported Processors and GPUs

This NVIDIA HPC Compilers release supports x86-64, OpenPOWER and Arm Server CPUs. Cross-compilation across the different families of CPUs is not supported, but you can use the `-tp=<target>` flag as documented in the man pages to specify a target processor within a family.

To direct the compilers to generate code for NVIDIA GPUs, use the `-acc` flag to enable OpenACC directives, the `-mp=gpu` flag to enable OpenMP directives, the `-stdpar` flag for standard language parallelism, and the `-cuda` flag for CUDA Fortran. Use the `-gpu` flag to select specific options for GPU code generation. You can then use the generated code on any supported system with CUDA installed that has a CUDA-enabled GeForce, Quadro, or Tesla card.

For more information on these flags as they relate to accelerator technology, refer to [Compiling an OpenACC Program](#).



For a complete list of supported CUDA GPUs, refer to the NVIDIA website at: [http://www.nvidia.com/object/cuda\\_learn\\_products.html](http://www.nvidia.com/object/cuda_learn_products.html)

## 5.8. CUDA Versions

The NVIDIA HPC compilers use components from NVIDIA's CUDA Toolkit to build programs for execution on an NVIDIA GPU. The NVIDIA HPC SDK puts the CUDA Toolkit components into an HPC SDK installation sub-directory; the HPC SDK currently bundles two versions of recently-released Toolkits.

You can compile a program for an NVIDIA GPU on any system supported by the HPC compilers. You will be able to run that program only on a system with an NVIDIA GPU and an installed NVIDIA CUDA driver. NVIDIA HPC SDK products do not contain CUDA device drivers. You must download and install the appropriate [CUDA Driver from NVIDIA](#).

The NVIDIA HPC SDK utility `nvaccelinfo` prints the driver version as its first line of output. You can use it to find out which version of the CUDA Driver is installed on your system.

The NVIDIA HPC SDK 24.9 includes components from the following versions of the CUDA Toolkit:

- ▶ CUDA 11.8
- ▶ CUDA 12.4

If you are compiling a program for GPU execution on a system *without* an installed CUDA driver, the compiler selects the version of the CUDA Toolkit to use based on the value of the `DEFCUDAVERSION` variable contained in a file called `localrc` which is created during installation of the HPC SDK.

If you are compiling a program for GPU execution on a system *with* an installed CUDA driver, the compiler detects the version of the CUDA driver and selects the appropriate CUDA Toolkit version to use from those bundled with the HPC SDK.

The compilers look for a CUDA Toolkit version in the `/opt/nvidia/hpc_sdk/target/24.9/cuda` directory that matches the version of the CUDA Driver installed on the system. If an exact match is not found, the compiler searches for the closest match. For CUDA Driver versions 11.2 through 11.8, the compiler will use the CUDA 11.8 Toolkit. For CUDA Driver versions 12.0 and later, the compiler will use the newest CUDA 12.x Toolkit.

You can change the compiler's default selection of CUDA Toolkit version using a compiler option. Add the `cudaX.Y` sub-option to `-gpu` where `X.Y` denotes the CUDA version. Using a compiler option changes the CUDA Toolkit version for one invocation of the compiler. For example, to compile an OpenACC C file with the CUDA 11.8 Toolkit you would use:

```
nvc -acc -gpu=cuda11.8
```

## 5.9. Compute Capability

The compilers can generate code for NVIDIA GPU compute capabilities 3.5 through 8.6. The compilers construct a default list of compute capabilities that matches the compute capabilities supported by the GPUs found on the system used in compilation. If there are no GPUs detected, the compilers generate code for every supported compute capability.

You can override the default by specifying one or more compute capabilities using either command-line options or an `rcfile`.

To change the default with a command-line option, provide a comma-separated list of compute capabilities to the `-gpu` option.

To change the default with an `rcfile`, set the **DEFCOMPUTECAP** value to a blank-separated list of compute capabilities in the `siterc` file located in your installation's `bin` directory:

```
set DEFCOMPUTECAP=60 70;
```

Alternatively, if you don't have permissions to change the `siterc` file, you can add the **DEFCOMPUTECAP** definition to a separate `.mynvrc` file in your home directory.

The generation of device code can be time consuming, so you may notice an increase in compile time as the number of compute capabilities increases.

## 5.10. PTX JIT Compilation

As of HPC SDK 22.9, support for PTX JIT compilation is enabled in all compilers for relocatable device code mode. This means that applications built with `-gpu=rdc` (that is, with relocatable device code enabled, which is the default mode) are forward-compatible with newer GPUs thanks to the embedded PTX code. The embedded PTX code is dynamically compiled when the application runs on a GPU architecture newer than the architecture specified at compile time.

The support for PTX JIT compilation is enabled automatically, which means that you do not need to change the compiler invocation command lines for your existing projects.

Use scenarios

- ▶ As an example, you can compile your application targeting the Ampere GPU without having to worry about the Hopper GPU architecture. Once the application runs on a Hopper GPU, it will seamlessly use the embedded PTX code.
- ▶ In CUDA Fortran, or with the CUDA Interoperability mode enabled, you can mix in object files compiled with the CUDA NVCC compiler containing PTX code. This PTX code from NVCC will be handled by the JIT compiler alongside the PTX code contained in object files produced by the HPC SDK compilers. When using the CUDA NVCC compiler, the relocatable device code generation must be enabled explicitly using the NVCC `--relocatable-device-code true`

switch, as explained in the [CUDA Compiler Driver guide](#). For information about CUDA Interoperability, please refer to <https://docs.nvidia.com/hpc-sdk/compiler/hpc-compilers-user-guide/index.html#openmp-interop-cuda>. The CUDA Fortran Programming Guide is available here: <https://docs.nvidia.com/hpc-sdk/compiler/cuda-fortran-prog-guide>.

By default, the compiler will choose the compute capability that matches the GPU on the system where the code is being compiled. For code that is going to run on the system where it is compiled, we recommend letting the compiler set the compute capability.

When the default won't work, we recommend compiling applications for a range of compute capabilities that the application is expected to run against, for example, using the `-gpu=ccll` compiler option. When running the application on a system that supports one of those compute capabilities, the CUDA driver minor version is allowed to be less than the version of the CUDA toolkit used at compile time, as covered in section [CUDA Versions](#).

#### Performance considerations

PTX JIT compilation, when it occurs, can have a start-up overhead for the application. The JIT compiler keeps a cached copy of the produced device code, which reduces the overhead on subsequent runs. Please refer to the [CUDA Programming Guide](#) for detailed information about how the JIT compiler works.

#### Known limitations

In general, in order for PTX JIT compilation to work, the CUDA driver installed on the deployment system must be at least of the version that matches the CUDA toolkit used to compile the application. This requirement is stricter than those explained in section [CUDA Versions](#).

For example, as explained in that section, the compilers will use the CUDA 11.8 toolkit that is shipped as part of the HPC SDK toolkit when the CUDA driver installed in the system is at least 11.2. However, while the CUDA 11.2 driver is commonly sufficient to run the application, it will not be able to compile the PTX code produced by the CUDA 11.8 toolkit. This means that any deployment system where the PTX JIT compilation is expected to be used must have at least the CUDA 11.8 driver installed. Please refer to the [CUDA Compatibility](#) guide for further information about the CUDA Driver compatibility with CUDA Toolkits.

When the application is expected to run on a newer GPU architecture than specified at compile time, we recommend having a CUDA driver installed on the deployment system matching the CUDA toolkit used to build the application. One way to achieve that is to use the `NVHPC_CUDA_HOME` environment variable at compile time to provide a specific CUDA toolkit.

Below are a few examples of how the PTX version incompatibility can be diagnosed and fixed. As a general rule, if the CUDA driver is unable to run the application due to incompatible PTX, the application will terminate with an error message indicating the cause. OpenACC and OpenMP applications will in most cases suggest compiler flags to target the current CUDA installation.

### OpenACC

Consider this program that we will compile for Volta GPU and attempt to run on an Ampere GPU, on a system that has CUDA 11.5 installed:

```
#include <stdio.h>
#define N 1000
int array[N];
int main() {
#pragma acc parallel loop copy(array[0:N])
    for(int i = 0; i < N; i++) {
        array[i] = 3.0;
    }
    printf("Success!\n");
}
```

When we build the program, HPC SDK will choose the CUDA 11.8 toolkit that is included as the default. When we attempt to run it, it fails because code generated with 11.8 does not work with the 11.5 driver:

```
$ nvc -acc -gpu=cc70 app.c
$ ./a.out
Accelerator Fatal Error: This file was compiled: -acc=gpu -gpu=cc70
Rebuild this file with -gpu=cc80 to use NVIDIA Tesla GPU 0
File: /tmp/app.c
Function: main:3
Line: 3
```

From the error message it follows that the system is unable to execute the Volta GPU instructions on the current system. The embedded Volta PTX could not be compiled, which implies a CUDA driver incompatibility. A way to fix this is to use the installed CUDA 11.5 toolkit at compile time:

```
$ export NVHPC_CUDA_HOME=/usr/local/cuda-11.5
$ nvc -acc -gpu=cc70 app.c
$ ./a.out
Success!
```

### OpenMP

Likewise, an OpenMP program will compile but not run:

```
#include <stdio.h>
#define N 1000
int array[N];
int main() {
#pragma omp target loop
    for(int i = 0; i < N; i++) {
```

```

    array[i] = 0;

}
printf("Success!\n");
}

```

```

$ nvc -mp=gpu -gpu=cc70 app.c
$ ./a.out
Accelerator Fatal Error: Failed to find device function 'nvkernel_main_F1L3_2'!
File was compiled with: -gpu=cc70
Rebuild this file with -gpu=cc80 to use NVIDIA Tesla GPU 0
File: /tmp/app.c
Function: main:3
Line: 3

```

We can also fix it by having **NVHPC\_CUDA\_HOME** point at the matching CUDA toolkit location:

```

$ export NVHPC_CUDA_HOME=/usr/local/cuda-11.5
$ nvc -acc -gpu=cc70 app.c
$ ./a.out
Success!

```

## C++

In contrast to OpenACC and OpenMP applications that simply terminate when PTX JIT encounters an insufficient CUDA driver version, C++ applications throw a system exception when there is a PTX incompatibility:

```

#include <vector>
#include <algorithm>
#include <execution>
#include <iostream>
#include <assert.h>
int main() {
    std::vector<int> x(1000, 0);
    x[1] = -20;
    auto result = std::count(std::execution::par, x.begin(), x.end(), -20);
    assert(result == 1);
    std::cout << "Success!" << std::endl;
}

```

```

$ nvcc++ -stdpar -gpu=cc70 app.cpp
$ ./a.out
terminate called after throwing an instance of 'thrust::system::system_error'
what():  after reduction step 1: cudaErrorUnsupportedPtxVersion: the provided
PTX was compiled with an unsupported toolchain.
Aborted (core dumped)

```

The exception message contains a direct reference to an incompatible PTX, which in turn implies a mismatch between the CUDA toolkit and the CUDA driver version.

We can fix it similarly by setting **NVHPC\_CUDA\_HOME**:

```

$ export NVHPC_CUDA_HOME=/usr/local/cuda-11.5

```

```
$ nvc++ -stdpar -gpu=cc70 app.cpp  
$ ./a.out  
Success!
```

# Chapter 6.

## USING OPENACC

This chapter gives an overview of directive-based OpenACC programming in which compiler directives are used to specify regions of code in Fortran, C and C++ programs to be offloaded from a *host* CPU to an NVIDIA GPU. For complete details on using OpenACC with NVIDIA GPUs, see the *OpenACC Getting Started Guide*.

### 6.1. OpenACC Programming Model

With the emergence of GPU architectures in high performance computing, programmers want the ability to program using a familiar, high level programming model that provides both high performance and portability to a wide range of computing architectures. OpenACC emerged in 2011 as a programming model that uses high-level compiler directives to expose parallelism in the code and parallelizing compilers to build the code for a variety of parallel accelerators.

This chapter will not attempt to describe OpenACC itself. For that, please refer to the OpenACC specification on the OpenACC [www.openacc.org](http://www.openacc.org) website. Here, we will discuss differences between the OpenACC specification and its implementation by the NVIDIA HPC Compilers.

Other resources to help you with your parallel programming including video tutorials, course materials, code samples, a best practices guide and more are available on the OpenACC website.

#### 6.1.1. Levels of Parallelism

OpenACC supports three levels of parallelism:

- ▶ an outer *doall* (fully parallel) loop level
- ▶ a *workgroup* or *threadblock* (worker parallel) loop level
- ▶ an inner *synchronous* (SIMD or vector) loop level

Each level can be multidimensional with 2 or 3 dimensions, but the domain must be strictly rectangular. The *synchronous* level may not be fully implemented with SIMD or vector operations, so explicit synchronization is supported and required across this level. No synchronization is supported between parallel threads across the *doall* level.

The OpenACC execution model on the device side exposes these levels of parallelism and the programmer is required to understand the difference between, for example, a fully parallel loop and a loop that is vectorizable but requires synchronization across iterations. All fully parallel loops can be scheduled for any of *doall*, *workgroup* or *synchronous* parallel execution, but by definition SIMD vector loops that require synchronization can only be scheduled for synchronous parallel execution.

### 6.1.2. Enable OpenACC Directives

NVIDIA HPC compilers enable OpenACC directives with the `-acc` and `-gpu` command line options. For more information on these options refer to [Compiling an OpenACC Program](#).

`_OPENACC` macro

The `_OPENACC` macro name is defined to have a value `yyyymm` where `yyyy` is the year and `mm` is the month designation of the version of the OpenACC directives supported by the implementation. For example, the version for November, 2017 is 201711. All OpenACC compilers define this macro when OpenACC directives are enabled.

### 6.1.3. OpenACC Support

The NVIDIA HPC Compilers implement most features of OpenACC 2.7 as defined in *The OpenACC Application Programming Interface*, Version 2.7, November 2018, <http://www.openacc.org>, with the exception that the following OpenACC 2.7 features are not supported:

- ▶ nested parallelism
- ▶ declare link
- ▶ enforcement of the **cache** clause restriction that all references to listed variables must lie within the region being cached
- ▶ Subarrays and composite variables in **reduction** clauses
- ▶ The **self** clause
- ▶ The **default** clause on data constructs

### 6.1.4. OpenACC Extensions

The NVIDIA Fortran compiler supports an extension to the `collapse` clause on the loop construct. The OpenACC specification defines `collapse`:

```
collapse (n)
```

NVIDIA Fortran supports the use of the identifier `force` within `collapse`:

```
collapse (force:n)
```

Using `collapse (force:n)` instructs the compiler to enforce collapsing parallel loops that are not perfectly nested.



## 6.2. Compiling an OpenACC Program

Several compiler options are applicable specifically when working with OpenACC. These options include `-acc`, `-gpu`, and `-Minfo`.

### 6.2.1. `-[no]acc`

Enable [disable] OpenACC directives. The following suboptions may be used following an equals sign ("="), with multiple sub-options separated by commas:

**gpu**

OpenACC directives are compiled for GPU execution only.

**host**

Compile for serial execution on the host CPU.

**multicore**

Compile for parallel execution on the host CPU.

**legacy**

Suppress warnings about deprecated NVIDIA accelerator directives.

**[no]autopar**

Enable [disable] loop autoparallelization within `acc parallel`. The default is to autoparallelize, that is, to enable loop autoparallelization.

**[no]routineseq**

Compile every routine for the device. The default behavior is to not treat every routine as a `seq` directive.

**strict**

Instructs the compiler to issue warnings for non-OpenACC accelerator directives.

**sync**

Ignore `async` clauses

**verystrict**

Instructs the compiler to fail with an error for any non-OpenACC accelerator directive.

**[no]wait**

Wait for each device kernel to finish. Kernel launching is blocked by default unless the `async` clause is used.

Default

By default OpenACC directives are compiled for GPU and sequential CPU host execution (i.e. equivalent to explicitly setting `-acc=gpu,host`).

## Usage

The following command-line requests that OpenACC directives be enabled and that an error be issued for any non-OpenACC accelerator directive.

```
$ nvfortran -acc=verystrict prog.f
```

## Predefined Macros

The following macros corresponding to the target compiled for are added implicitly:

- ▶ `__NVCOMPILER_OPENACC_GPU` when the OpenACC directives are compiled for GPU.
- ▶ `__NVCOMPILER_OPENACC_MULTICORE` when the OpenACC directives are compiled for multicore CPU.
- ▶ `__NVCOMPILER_OPENACC_HOST` when the OpenACC directives are compiled for serial execution on CPU.

## 6.2.2. -gpu

Used in combination with the `-acc`, `-cuda`, `-mp`, and `-stdpar` flags to specify options for GPU code generation. The following sub-options may be used following an equals sign ("`=`"), with multiple sub-options separated by commas:

### **autocompare**

Automatically compare CPU vs GPU results at execution time: implies redundant **ccXY**

Generate code for a device with compute capability X.Y. Multiple compute capabilities can be specified, and one version will be generated for each. By default, the compiler will detect the compute capability for each installed GPU. Use `-help -gpu` to see the valid compute capabilities for your installation.

### **ccall**

Generate code for all compute capabilities supported by this platform and by the selected or default CUDA Toolkit.

### **ccall-major**

Compile for all major supported compute capabilities.

### **ccnative**

Detects the visible GPUs on the system and generates codes for them. If no device is available, the compute capability matching NVCC's default will be used.

### **cudaX.Y**

Use CUDA X.Y Toolkit compatibility, where installed

### **[no]debug**

Enable [disable] debug information generation in device code

### **deepcopy**

Enable full deep copy of aggregate data structures in OpenACC; Fortran only

### **fastmath**

Use routines from the fast math library

**[no]flushz**

Enable [disable] flush-to-zero mode for floating point computations on the GPU

**[no]fma**

Generate [do not generate] fused multiply-add instructions; default at `-O3`

**[no]implicitsections**

Change [Do not change] array element references in a data clause into an array section. In C++, the `implicitsections` option will change `update device(a[n])` to `update device(a[0:n])`. In Fortran, it will change `enter data copyin(a(n))` to `enter data copyin(a(:n))`. The default behavior, `noimplicitsections`, can also be changed using rcfiles; for example, one could add `set IMPLICITSECTIONS=0`; to `siterc` or another rcfile.

**[no]interceptdeallocations**

Intercept [Do not intercept] calls to standard library memory deallocations (e.g. `free`) and call the corresponding CUDA memory deallocation version if address is in pinned or managed memory, regular version otherwise.

**keep**

Keep the kernel files (.cubin, .ptx, source)

**[no]lineinfo**

Enable [disable] GPU line information generation

**loadcache:{L1|L2}**

Choose what hardware level cache to use for global memory loads; options include the default, `L1`, or `L2`

**[no]managed**

Allocate [do not allocate] any dynamically allocated data in CUDA Managed memory. Use `-gpu=nomanaged` with `-stdpar` to prevent that flag's implicit use of `-gpu=managed` when CUDA Managed memory capability is detected. This option is deprecated.

**maxregcount:n**

Specify the maximum number of registers to use on the GPU; leaving this blank indicates no limit

**mem:{separate|managed|unified}**

Select GPU memory mode for the generated binary. This controls CUDA memory capability to be utilised such as separate GPU memory only (`separate`), GPU Managed Memory for the dynamically allocated data (`managed`), or system memory aka full CUDA Unified Memory (`unified`). Use of Managed or Unified Memory facilitates simpler programming by eliminating the need to detect all data to be copied into and outside of the code region executing on the GPU.

**pinned**

Use CUDA Pinned Memory. This option is deprecated.

**ptxinfo**

Print PTX info

**[no]rdc**

Generate [do not generate] relocatable device code.

**redundant**

Redundant CPU/GPU execution

**safecache**

Allow variable-sized array sections in cache directives; compiler assumes they fit into CUDA shared memory

**sm\_XY**

Generate code for a device with compute capability X.Y. Multiple compute capabilities can be specified, and one version will be generated for each. By default, the compiler will detect the compute capability for each installed GPU. Use `-help-gpu` to see the valid compute capabilities for your installation.

**stacklimit:<l>nostacklimit**

Sets the limit (l) of stack variables in a procedure or kernel, in KB. This option is deprecated.

**[no]unified**

Compile [do not compile] for CUDA Unified memory capability, where system memory is accessible from the GPU. This mode utilizes system and managed memory for dynamically allocated data unless explicit behavior is set through `-gpu=[no]managed`. Use `-gpu=nounified` with `-stdpar` to prevent that flag's implicit use of `-gpu=unified` when CUDA Unified memory capability is detected. This option must appear in both the compile and link lines. This option is deprecated.

**[no]unroll**

Enable [disable] automatic inner loop unrolling; default at `-O3`

**zeroinit**

Initialize allocated device memory with zero

## Usage

In the following example, the compiler generates code for NVIDIA GPUs with compute capabilities 6.0 and 7.0.

```
$ nvfortran -acc -gpu=cc60,cc70 myprog.f
```

The compiler automatically invokes the necessary software tools to create the kernel code and embeds the kernels in the object file.

To link in the appropriate GPU libraries, you must link an OpenACC program with the `-acc` flag, and similarly for `-cuda`, `-mp`, or `-stdpar`.

## DWARF Debugging Formats

Use the `-g` option to enable generation of full DWARF information on both the host and device; in the absence of other optimization flags, `-g` sets the optimization level to zero. If a `-O` option raises the optimization level to one or higher, only GPU line information is generated in device code even when `-g` is specified. To enforce full DWARF generation for device code at optimization levels above zero, use the `debug` sub-option to `-gpu`. Conversely, to prevent the generation of dwarf information for device code, use the

**nodebug** sub-option to **-gpu**. Both **debug** and **nodebug** can be used independently of **-g**.

## 6.3. OpenACC for Multicore CPUs

The NVIDIA OpenACC compilers support the option `-acc=multicore`, to set the target accelerator for OpenACC programs to the host multicore CPU. This will compile OpenACC compute regions for parallel execution across the cores of the host processor or processors. The host multicore CPU will be treated as a shared-memory accelerator, so the data clauses (**copy**, **copyin**, **copyout**, **create**) will be ignored and no data copies will be executed.

By default, `-acc=multicore` will generate code that will use all the available cores of the processor. If the compute region specifies a value in the **num\_gangs** clause, the minimum of the **num\_gangs** value and the number of available cores will be used. At runtime, the number of cores can be limited by setting the environment variable **ACC\_NUM\_CORES** to a constant integer value. The number of cores can also be set with the `void acc_set_num_cores(int numcores)` runtime call. If an OpenACC compute construct appears lexically within an OpenMP parallel construct, the OpenACC compute region will generate sequential code. If an OpenACC compute region appears dynamically within an OpenMP region or another OpenACC compute region, the program may generate many more threads than there are cores, and may produce poor performance.

The `-acc=multicore` option differs from the `-acc=host` option in that `-acc=host` generates sequential host CPU code for the OpenACC compute regions.

## 6.4. OpenACC with CUDA Unified Memory

When developing OpenACC source for a target supporting CUDA Unified Memory, you can take advantage of a simplified approach to programming because there is no need for data clauses and directives, either in full or in part, depending on the exact memory capability the target supports and the compiler options used.

The discussion in this section assumes you have become familiar with the Separate, Managed, and Unified Memory Modes covered in the [Memory Model](#) and [Managed and Unified Memory Modes](#) sections.

In Managed Memory Mode, only dynamically-allocated data are implicitly managed by the CUDA runtime; OpenACC data clauses and directives are therefore not needed for movement of this "managed" data. Data clauses and directives are still required to handle static data (C static and extern variables, Fortran module, common block and save variables) and function local data.

In Unified Memory Mode, all data is managed by the CUDA runtime. Explicit data clauses and directives are no longer required to indicate which data should reside in GPU memory. All variables are accessible from the OpenACC compute regions

executing on the GPU. The NVHPC compiler implementation closely adheres to the shared memory mode detailed in the OpenACC specification, meaning that **copy**, **copyin**, **copyout**, and **create** clauses will not result in any device allocation or data transfer. The **device\_resident** clause is still honored as in discrete memory mode and results in an allocation of data only accessible from device code. Device memory can also be allocated or deallocated in OpenACC programs in Unified Memory Mode by using the **acc\_malloc** or **acc\_free** API calls.

### Understanding Data Movement

In the absence of visible data clauses or directives, when the compiler encounters a compute construct it attempts to determine what data is required for correct execution of the region on the GPU. When the compiler is unable to determine the size and shape of data needing to be accessible on the device, it behaves as follows:

- ▶ In Separate Memory Mode, the compiler emits an error requesting an explicit data clause be added to specify size/shape of the data to be copied.
- ▶ In Managed Memory Mode (**-gpu=mem:managed**), the compiler assumes the data is allocated in managed memory and thus is accessible from the device; if this assumption is wrong, if the data was defined globally or is located on the CPU stack, the program may fail at runtime.
- ▶ In Unified Memory Mode (**-gpu=mem:unified**), all data is accessible from the device making information about size and shape unnecessary.

Take the following example in C:

```
void set(int* ptr, int i, int j, int dim){
    int idx = i * dim + j;
    return ptr[idx] = someval(i, j);
}

void fill2d(int* ptr, int dim){
#pragma acc parallel loop
    for (int i = 0; i < dim; i++)
        for (int j = 0; j < dim; j++)
            set(ptr, i, j, dim);
}
```

In Separate Memory Mode, the only way to guarantee correctness for this example is to change the line with the **acc** directive as follows:

```
#pragma acc parallel loop create(ptr[0:dim*dim]) copyout(ptr[0:dim*dim])
```

This change explicitly instructs the OpenACC implementation about the precise data segment used within the parallel loop.

In Unified Memory Mode, that is, by compiling with **-acc -gpu=mem:unified** and executing on a platform with unified memory capability, the **create** and **copyout** clauses are not required.

The next example, in Fortran, illustrates how a global variable can be accessed in an OpenACC routine without requiring any explicit annotation.

```
module m
integer :: globmin = 1234
contains
subroutine findmin(a)
!$acc routine seq
  integer, intent(in) :: a(:)
  integer :: i
  do i = 1, size(a)
    if (a(i) .lt. globmin) then
      globmin = a(i)
    endif
  end do
end subroutine
end module m
```

Compile the example above for Unified Memory Mode:

```
nvfortran -acc -gpu=mem:unified example.f90
```

The source does not need any OpenACC directives to access module variable **globmin**, to either read or update its value, in the routine invoked from CPU and GPU. Moreover, any access to **globmin** will be made to the same exact instance of the variable from CPU and GPU; its value is synchronized automatically. In Separate or Managed Memory Modes, such behavior can only be achieved with a combination of OpenACC **declare** and **update** directives in the source code.

In most cases, migrating existing OpenACC applications written for Separate Memory Mode should be a seamless process requiring no source changes. Some data access patterns, however, may lead to different results produced during application execution in Unified Memory Mode.

Applications which rely on having separate data copies in GPU memory to conduct temporary computations on the GPU -- without maintaining data synchronization with the CPU -- pose a challenge for migration to Unified Memory.

For the following Fortran example, the value of variable **c** after the last loop will differ depending on whether the example is compiled with or without **-gpu=mem:unified**.

```
b(:) = ...
c = 0

!$acc kernels copyin(b) copyout(a)
!$acc loop
do i = 1, N
  b(i) = b(i) * i
end do
!$acc loop
do i = 1, N
  a(i) = b(i) + i
end do
!$acc end kernels

do i = 1, N
```

```
c = c + a(i) + b(i)
end do
```

Without Unified Memory, array **b** is copied into the GPU memory at the beginning of the OpenACC **kernels** region. It is then updated in the GPU memory and used to compute elements of array **a**. As instructed by the data clause **copyin(b)**, **b** is not copied back to the CPU memory at the end of the **kernels** region and therefore its initial value is used in the computation of **c**. With **-acc -gpu=mem:unified**, the updated value of **b** in the first loop is automatically visible in the last loop leading to a different value of **c** at its end.

### Implications of Asynchronous Execution

Additional complexities can arise when dealing with asynchronous execution, particularly when CPU-GPU shared data is accessed within **async** compute regions instead of using an independent data copy on GPU. The programmer should be especially careful about accessing local variables in asynchronous GPU code. Unless the GPU code execution is explicitly synchronized before the end of the scope in which local variables are defined, the GPU can access stale data thus resulting in undefined behavior. Consider the following OpenACC C example, where a local array is used to hold temporary data on the GPU:

```
void bar() {
    int x[N];
    #pragma acc enter data create(x[0:N]) async
    #pragma acc parallel loop async
    for (int i = 0; i < N; i++)
        x[i] = i;
    ...
    #pragma acc exit data delete(x[0:N]) async
}
```

When compiled for Separate Memory Mode, the **bar()** function creates a copy of the array **x** in GPU memory and initializes it as written in the **loop** construct. That copy is eventually deleted. In Unified Memory Mode, however, the compiler ignores the **acc enter data** and **acc exit data** directives, so the **loop** construct executed on the GPU accesses the array **x** in local CPU memory. Moreover, since all constructs in this example are made asynchronous, the access to **x** on the GPU leads to undefined behavior of the program because the variable **x** goes out of scope once the **bar()** function finishes.

### Performance Considerations

In Unified Memory Mode, the OpenACC runtime may leverage data action information such as **create/delete** or **copyin/copyout** to communicate preferable data placement to the CUDA runtime by means of memory hint APIs as elaborated in the following blog post on the NVIDIA website: [Simplifying GPU Application Development with Heterogeneous Memory Management](#). Such actions originate either from explicit data



clauses in the source code or via implicit data movement generated by the compiler. This approach can minimize the amount of automatic data migration and may let a developer fine-tune application performance. For the C example above, while adding the data clauses `create(ptr[0:dim*dim])` and `copyout(ptr[0:dim*dim])` becomes optional with `-gpu=mem:unified`, their uses in the OpenACC `parallel loop` directive may improve performance.

## 6.5. OpenACC Error Handling

The OpenACC specification provides a mechanism to allow you to intercept errors triggered during execution on a GPU and execute a specific routine in response before the program exits. For example, if an MPI process fails while allocating memory on the GPU, the application may want to call `MPI_Abort` to shut down all the other processes before the program exits. This section explains how to take advantage of this feature.

To intercept errors the application must give a callback routine to the OpenACC runtime. To provide the callback, the application calls `acc_set_error_routine` with a pointer to the callback routine.

The interface is the following, where `err_msg` contains a description of the error:

```
typedef void (*exitroutinetype)(char *err_msg);
extern void acc_set_error_routine(exitroutinetype callback_routine);
```

When the OpenACC runtime detects a runtime error, it will invoke the `callback_routine`.



This feature is not the same as error recovery. If the callback routine returns to the application, the behavior is decidedly undefined.

Let's look at this feature in more depth using an example.

Take the MPI program below and run it with two processes. Process 0 tries to allocate a large array on the GPU, then sends a message to the second process to acknowledge the success of the operation. Process 1 waits for the acknowledgment and terminates upon receiving it.

```
#include <stdio.h>
#include <stdlib.h>
#include "mpi.h"

#define N 2147483648

int main(int argc, char **argv)
{
    int rank, size;

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    int ack;
    if(rank == 0) {
        float *a = (float*) malloc(sizeof(float) * N);
```

```

#pragma acc enter data create(a[0:N])
#pragma acc parallel loop independent
    for(int i = 0; i < N; i++) {
        a[i] = i *0.5;
    }
#pragma acc exit data copyout(a[0:N])
    printf("I am process %d, I have initialized a vector of size %ld bytes on
the GPU. Sending acknowledgment to process 1.", rank, N);
    ack = 1;
    MPI_Send(&ack, 1, MPI_INT, 1, 0, MPI_COMM_WORLD);
} else if(rank == 1) {
    MPI_Recv(&ack, 1, MPI_INT, 0, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
    printf("I am process %d, I have received the acknowledgment from process 0
that data in the GPU has been initialized.\n", rank, N);
    fflush(stdout);
}

// do some more work

MPI_Finalize();

return 0;
}

```

We compile the program with:

```
$ mpicc -acc -o error_handling_mpi error_handling_mpi.c
```

If we run this program with two MPI processes, the output will look like the following:

```

$ mpirun -n 2 ./error_handling_mpi
Out of memory allocating -8589934592 bytes of device memory
total/free CUDA memory: 11995578368/11919294464
Present table dump for device[1]:
NVIDIA Tesla GPU 0, compute capability 3.7, threadid=1
...empty...
call to cuMemAlloc returned error 2: Out of memory

-----
Primary job terminated normally, but 1 process returned
a non-zero exit code.. Per user-direction, the job has been aborted.
-----
-----
mpirun detected that one or more processes exited with non-zero status,
thus causing the job to be terminated.

```

Process 0 failed while allocating memory on the GPU and terminated unexpectedly with an error. In this case **mpirun** was able to identify that one of the processes failed, so it shut down the remaining process and terminated the application. A simple two-process program like this is straightforward to debug. In a real world application though, with hundreds or thousands of processes, having a process exit prematurely may cause the application to hang indefinitely. Therefore it would be ideal to catch the failure of a process, control the termination of the other processes, and provide a useful error message.

We can use the OpenACC error handling feature to improve the previous program and correctly terminate the application in case of failure of an MPI process.

In the following sample code, we have added an error handling callback routine that will shut down the other processes if a process encounters an error while executing on the GPU. Process 0 tries to allocate a large array into the GPU and, if the operation is successful, process 0 will send an acknowledgment to process 1. Process 0 calls the

OpenACC function **acc\_set\_error\_routine** to set the function **handle\_gpu\_errors** as an error handling callback routine. This routine prints a message and calls **MPI\_Abort** to shut down all the MPI processes. If process 0 successfully allocates the array on the GPU, process 1 will receive the acknowledgment. Otherwise, if process 0 fails, it will terminate itself and trigger the call to **handle\_gpu\_errors**. Process 1 is then terminated by the code executed in the callback routine.

```
#include <stdio.h>
#include <stdlib.h>
#include "mpi.h"

#define N 2147483648

typedef void (*exitroutinetype)(char *err_msg);
extern void acc_set_error_routine(exitroutinetype callback_routine);

void handle_gpu_errors(char *err_msg) {
    printf("GPU Error: %s", err_msg);
    printf("Exiting...\n\n");
    MPI_Abort(MPI_COMM_WORLD, 1);
    exit(-1);
}

int main(int argc, char **argv)
{
    int rank, size;

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    int ack;
    if(rank == 0) {
        float *a = (float*) malloc(sizeof(float) * N);

        acc_set_error_routine(&handle_gpu_errors);

#pragma acc enter data create(a[0:N])
#pragma acc parallel loop independent
        for(int i = 0; i < N; i++) {
            a[i] = i * 0.5;
        }
#pragma acc exit data copyout(a[0:N])
        printf("I am process %d, I have initialized a vector of size %ld bytes on
the GPU. Sending acknowledgment to process 1.", rank, N);
        fflush(stdout);
        ack = 1;
        MPI_Send(&ack, 1, MPI_INT, 1, 0, MPI_COMM_WORLD);
    } else if(rank == 1) {
        MPI_Recv(&ack, 1, MPI_INT, 0, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
        printf("I am process %d, I have received the acknowledgment from process 0
that data in the GPU has been initialized.\n", rank, N);
        fflush(stdout);
    }

    // more work

    MPI_Finalize();
}
```

```
    return 0;
}
```

Again, we compile the program with:

```
$ mpicc -acc -o error_handling_mpi error_handling_mpi.c
```

We run the program with two MPI processes and obtain the output below:

```
$ mpirun -n 2 ./error_handling_mpi
Out of memory allocating -8589934592 bytes of device memory
total/free CUDA memory: 11995578368/11919294464
Present table dump for device[1]:
NVIDIA Tesla GPU 0, compute capability 3.7, threadid=1
...empty...
GPU Error: call to cuMemAlloc returned error 2: Out of memory
Exiting...
```

```
-----
MPI_ABORT was invoked on rank 0 in communicator MPI_COMM_WORLD
with errorcode 1.
```

This time the error on the GPU was intercepted by the application which managed it with the error handling callback routine. In this case the routine printed some information about the problem and called **MPI\_Abort** to terminate the remaining processes and avoid any unexpected behavior from the application.

## 6.6. OpenACC and CUDA Graphs

NVIDIA provides an optimized model for work submission onto GPUs called CUDA Graphs. A graph is a series of operations, such as kernel launches and other stream-oriented tasks, connected by their dependencies. A graph can be defined once, "captured", then launched repeatedly. This has potential benefits in reducing launch latencies and other overheads associated with kernel setup.

A complete write-up explaining CUDA Graphs and the CUDA API for graph definition, instantiation, and execution can be found in Chapter 3 of the CUDA C Programming Guide. In OpenACC, we currently expose just the minimal set of operations to allow capture and replay of a graph containing OpenACC compute regions and data directives. The code executed between a "begin capture" call, **accx\_begin\_capture\_async()**, and the "end capture" call, **accx\_end\_capture\_async()**, is called the capture region.

The CUDA graph API captures (or records) all the device work between **accx\_begin\_capture\_async** and **accx\_end\_capture\_async**. The host code in the capture region will be executed once normally, with the exception that no device work is actually executed on the device. Instead, a graph object is created that can be used to replay the captured work multiple times.



Graph capture is similar to a closure concept in many programming languages, like lambda-functions in C++. In lambda-function terms, CUDA graphs capture all the variables by value. That means that all the **FIRSTPRIVATE** scalars, array shapes, and those derived types, arrays and scalar addresses for data resident on the GPU, are baked into the graph object and cannot be altered. The device data behind the pointers,

of course, can be updated by the graph execution normally, and updated by the host between replays.

It is important to understand both what can and cannot be captured within a CUDA Graph capture region:

- ▶ Asynchronous data clauses including data create can be captured. The OpenACC runtime will use the stream-ordered `cudaMallocAsync()` call in the capture region for variables which need allocation in data clauses, an API call allowed in CUDA Graphs.
- ▶ Asynchronous compute regions, preferably ACC parallel regions, can be captured. For ACC kernels regions, verify that no work is performed on the host. Host compute sections cannot be captured.
- ▶ Asynchronous ACC update host (self) and update device directives can be captured. The host and device addresses which are captured must be valid during the graph replay/execution.
- ▶ Since only the device work is captured and replayed, any data dependencies between the host and device inside the capture region are erroneous. For example, downloading data from the device, processing it on host and uploading it back to the device within the capture region is invalid.
- ▶ Host code, even host code containing conditionals, can occur within a capture region. Note though that the path taken through the host code will be the path captured by the graph, i.e. the conditionals must likely be consistent during the replay for correct results. Host code which updates host variables, such as `i=i+1` will not be captured in the graph, which might affect proper indexing into device-side arrays or other kernel arguments.
- ▶ Similarly, device work initiated in host code loops can be captured in the CUDA Graph. The graph will not contain a notion of looping, just the sequence of device operations submitted to the device during the loop.
- ▶ Subroutine and function calls within a capture region, which contain further compute regions or other work which runs on the device, are captured. Care must be taken that the device data addresses passed to the kernels are valid throughout graph execution, and don't come and go based on stack addresses or something similar.
- ▶ Codes which double-buffer, or ping-pong between source and destination arrays that are input on odd iterations, and output on even iterations, can be accommodated by capturing two graphs: one per even iteration, one per odd iteration.
- ▶ Many CUDA library calls, like cublas, etc. can occur in a captured region. Setup for the library calls, such as creating handles, and computing and allocating workspace requirements, should be done before the capture region.
- ▶ Graph capturing is thread-safe with respect to each async queue. Host threads can independently capture graphs using different async queues.

The OpenACC API follows the basic portion of the CUDA Graph API fairly closely. The major difference is OpenACC includes the `cudaGraphInstantiate()` call as part of the end capture function.

From Fortran, the graph type is defined in the OpenACC module:

```
type, bind(c) :: acc_graph_t
type(c_ptr) :: graph
```

```

    type(c_ptr) :: graph_exec
end type acc_graph_t

```

These subroutines are available in the OpenACC runtime. Here, pGraph is type(acc\_graph\_t) and async is just the asynchronous queue value:

```

subroutine accx_async_begin_capture( async )
subroutine accx_async_end_capture( async, pGraph )
subroutine accx_graph_launch( pGraph, async )
subroutine accx_graph_delete( pGraph )
type(c_ptr) function accx_get_graph( pGraph )
type(c_ptr) function accx_get_graph_exec( pGraph )

```

From C, the graph type is defined in OpenACC.h:

```

typedef struct { void *graph; void *graph_exec; } acc_graph_t;

```

These void functions are available in the OpenACC runtime:

```

extern void accx_async_begin_capture(long async);
extern void accx_async_end_capture(long async, acc_graph_t *pgraph);
extern void accx_graph_launch(acc_graph_t *pgraph, long async);
extern void accx_graph_delete(acc_graph_t *pgraph);
extern void *accx_get_graph(acc_graph_t *pgraph);
extern void *accx_get_graph_exec(acc_graph_t *pgraph);

```

We will use a simple Fortran example code which demonstrates some of the modifications needed to use CUDA Graphs from OpenACC. The original serial code for a conjugate gradient iterative solver:

```

subroutine RunCG(N, A, b, x, tol, max_iter)
  implicit none
  integer, intent(in) :: N, max_iter
  real(WP), intent(in) :: A(N, N), b(N), tol
  real(WP), intent(inout) :: x(N)

  real(WP) :: alpha, rr0, rr
  real(WP), allocatable :: Ax(:), r(:), p(:)
  integer :: it, i

  allocate(Ax(N), r(N), p(N))

  call symmatvec(N, N, A, x, Ax)
  do i = 1, N
    r(i) = b(i) - Ax(i)
    p(i) = r(i)
  enddo
  rr0 = dot(N, r, r)

  do it = 1, max_iter
    call symmatvec(N, N, A, p, Ax)
    alpha = rr0 / dot(N, p, Ax)

    do i = 1, N
      x(i) = x(i) + alpha * p(i)
      r(i) = r(i) - alpha * Ax(i)
    enddo

    rr = dot(N, r, r)

    print*, "Iteration ", it, " residual: ", sqrt(rr)
    if (sqrt(rr) <= tol) then
      deallocate(Ax, r, p)
      return
    endif
    do i = 1, N
      p(i) = r(i) + (rr / rr0) * p(i)
    enddo
  enddo

```

```

    rr0 = rr
enddo

deallocate(Ax, r, p)

end subroutine RunCG

```

For this exercise we wish to put the **do it = 1,max\_iter** work for each iteration into a CUDA graph. Step one is to port the code to OpenACC, keeping in mind that we want to use asynchronous queues. We annotate the dot function with OpenACC directives like this:

```

function dot(N, x, y) result(r)
  integer, intent(in) :: N
  real(WP), intent(in) :: x(N), y(N)
  integer :: i
  real(WP) :: r

  r = 0.d0
  !$acc parallel loop present(x, y) reduction(+:r) async(1)
  do i = 1, N
    r = r + x(i) * y(i)
  enddo
  !$acc wait(1)
end function dot

```

We write the symmetric matrix multiply like this:

```

subroutine symmatvec(M, N, AT, x, Ax)
  implicit none
  integer, intent(in) :: M, N
  real(WP), intent(in) :: AT(N, M), x(N)
  real(WP), intent(out) :: Ax(M)

  integer :: i, j
  real(WP) :: s

  ! Note: Since A is symmetric, we can use the "transpose"
  ! for better memory access here
  !$acc parallel loop gang present(AT, x, Ax) async(1)
  do i = 1, M
    s = 0.d0
    !$acc loop vector reduction(+:s)
    do j = 1, N
      s = s + AT(j,i) * x(j)
    end do
    Ax(i) = s
  end do
end subroutine

```

And now our main loop of the conjugate gradient solver looks like this:

```

do it = 1, max_iter
  call symmatvec(N, N, A, p, Ax)
  alpha = rr0 / dot(N, p, Ax)

  !$acc parallel loop gang vector async(1)
  do i = 1, N
    x(i) = x(i) + alpha * p(i)
    r(i) = r(i) - alpha * Ax(i)
  enddo

  rr = dot(N, r, r)

  print*, "Iteration ", it, " residual: ", sqrt(rr)
  if (sqrt(rr) <= tol) exit

  !$acc parallel loop gang vector async(1)

```

```

do i = 1, N
  p(i) = r(i) + (rr / rr0) * p(i)
enddo
rr0 = rr
enddo

```

Step 2 is to prepare the code for running under CUDA Graphs. There is a lot of host code executing in the main loop. While the `dot()` function runs on the GPU, the rest of the statement `alpha = rr0 / dot(...)` runs on the host. Similarly, the 2nd `dot()` call returns its value to the host. The print statement occurs on the host, as does the residual check. Finally, this iteration's value for `rr` is moved to `rr0` in the last statement of the loop, on the host.

The dot product is tricky. We wish to compute the dot product on the GPU, and leave the result on the GPU, so the reduction variable must be present on the GPU. Here, we change the function call to a subroutine, and remove the initialization which is outside of the parallel region:

```

subroutine dot(N, x, y, r)
  implicit none
  integer, intent(in) :: N
  real(WP), intent(in) :: x(N), y(N)
  integer :: i
  real(WP) :: r

  !$acc parallel loop present(x, y, r) reduction(+:r) async(1)
  do i = 1, N
    r = r + x(i) * y(i)
  enddo
end subroutine dot

```

We add one serial kernel to do some of the swapping between `rr0` and `rr`, as well as zeroing out the scalar that will hold the dot product reduction, and move the print and check outside of the GPU capture region, replaced by a update host operation. The finished loop, complete with graph control, looks like this:

```

do it = 1, max_iter
  if (it .eq. 1) then ! First time capture
    call accx_async_begin_capture(1)

    call symmatvec(N, N, A, p, Ax)
    call dot(N, p, Ax, rden)

    !$acc serial async(1)
    rr0 = rr
    alpha = rr0 / rden
    rden = 0.0d0
    rr = 0.0d0
    !$acc end serial

    !$acc parallel loop gang vector async(1)
    do i = 1, N
      x(i) = x(i) + alpha * p(i)
      r(i) = r(i) - alpha * Ax(i)
    enddo

    call dot(N, r, r, rr)

    !$acc update host(rr) async(1)

    !$acc parallel loop gang vector async(1)
    do i = 1, N
      p(i) = r(i) + (rr / rr0) * p(i)
    enddo
  end if
end do

```



```

    call accx_async_end_capture(1, graph)
endif
! Always launch, then wait
call accx_graph_launch(graph, 1)
!$acc wait(1)

rra(it) = rr
if (sqrt(rr) <= tol) exit
enddo

```

Step 3 is to compile, run, and profile the result. No special compiler options are needed besides `-acc=gpu`. When running, you may be advised to set the **NVCOMPILER\_ACC\_USE\_GRAPH** environment variable. This is currently necessary to properly set the OpenACC runtime for graph capture. Failure to abide by the guidelines above may result in wrong answers, which can be hard to debug. See the following sections on how to use environment variables to help. A common issue is that the pointers passed to the device kernels during graph playback will be the same every time. Make sure that is the case between iterations in the code without graph capture.

The Nsight Systems tool has very good support for profiling CUDA graphs. The timeline view will provide information on whether you have reduced the launch overhead gaps between the GPU kernels. [Figure 1](#) shows a timeline of the iterations of the original OpenACC loop:

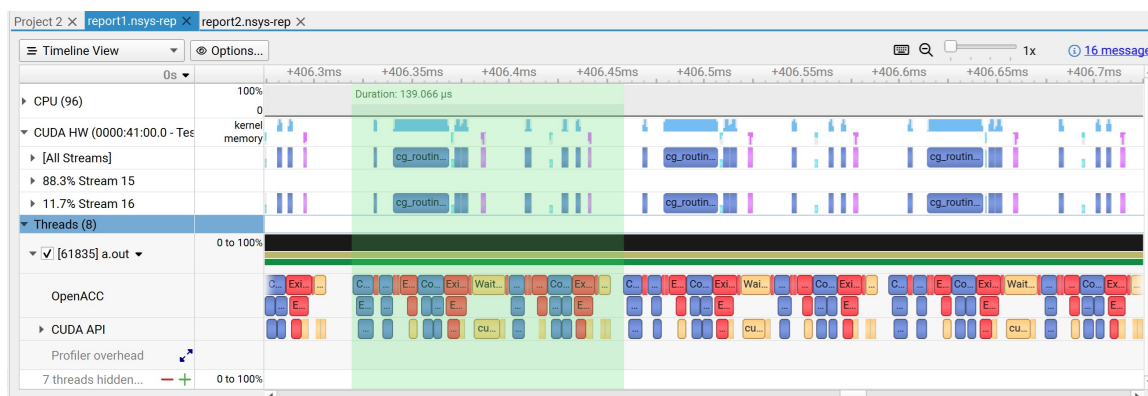


Figure 1 Nsight Systems Report1 Timeline

[Figure 2](#) shows a timeline of the iterations when using CUDA Graphs. When the size  $N$  is less than a few thousand, launch latency becomes a major contributor to the overall time and here we can see about a 2x speedup:

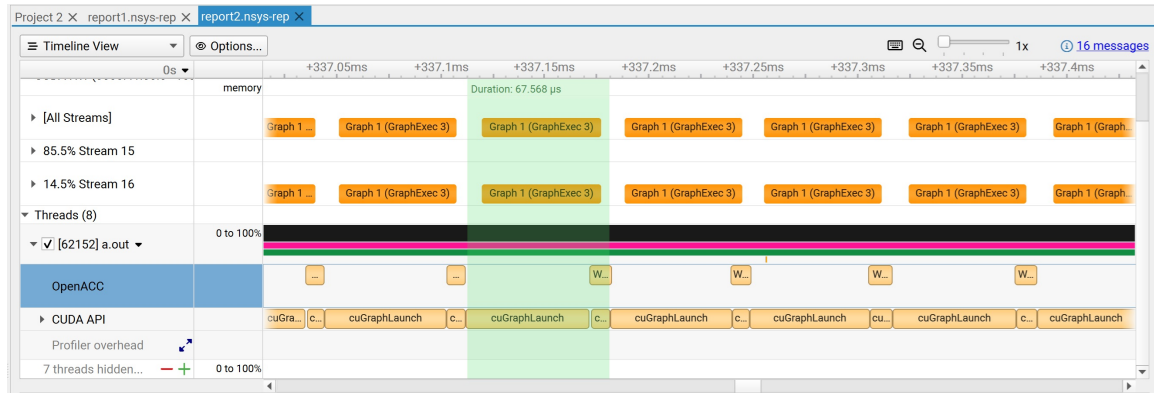


Figure 2 Nsight Systems Report2 Timeline

You can see a more-detailed trace of the CUDA Graph components by adding the `--cuda-graph-trace=node` option to the nsys profile command.

The above loop demonstrates several of the guidelines outlined at the top of this section, namely, capturing compute regions, whether at the top level or in subprogram units, capturing data movement, and restructuring code regions to minimize or eliminate the host code within a capture region. And the minimal API to begin capture, end capture, then launch the captured graph.

## 6.7. Environment Variables

This section summarizes the environment variables that NVIDIA OpenACC supports. These environment variables are user-setable environment variables that control behavior of accelerator-enabled programs at execution. These environment variables must comply with these rules:

- ▶ The names of the environment variables must be upper case.
- ▶ The values of environment variables are case insensitive and may have leading and trailing white space.
- ▶ The behavior is implementation-defined if the values of the environment variables change after the program has started, even if the program itself modifies the values.

The following table contains the environment variables that are currently supported and provides a brief description of each.

Table 16 Supported Environment Variables

Use this environment variable...	To do this...
NVCOMPILER_ACC_CUDA_PROFS	Set to 1 (or any positive value) to tell the runtime environment to insert an 'atexit(cuProfilerStop)' call upon exit. This behavior may be desired in the case where a profile is incomplete or where a message is issued to call <code>cudaProfilerStop()</code> .
NVCOMPILER_ACC_DEVICE_NUM	Sets the default device number to use. NVCOMPILER_ACC_DEVICE_NUM. Specifies the default device

Use this environment variable...	To do this...
	number to use when executing accelerator regions. The value of this environment variable must be a nonnegative integer between zero and the number of devices attached to the host.
ACC_DEVICE_NUM	Legacy name. Superseded by NVCOMPILER_ACC_DEVICE_NUM.
NVCOMPILER_ACC_DEVICE_TYPE	Sets the default device type to use for OpenACC regions. NVCOMPILER_ACC_DEVICE_TYPE. Specifies which accelerator device to use when executing accelerator regions when the program has been compiled to use more than one different type of device. The value of this environment variable is implementation-defined, and in the NVIDIA OpenACC implementation may be the strings NVIDIA, MULTICORE or HOST
ACC_DEVICE_TYPE	Legacy name. Superseded by NVCOMPILER_ACC_DEVICE_TYPE.
NVCOMPILER_ACC_GANGLIMIT	For NVIDIA CUDA devices, this defines the maximum number of gangs (CUDA thread blocks) that will be launched by a kernel.
NVCOMPILER_ACC_NOTIFY	With no argument, a debug message will be written to stderr for each kernel launch and/or data transfer. When set to an integer value, the value is used as a bit mask to print information about:  1: kernel launches 2: data transfers 4: region entry/exit 8: wait operations or synchronizations with the device 16: device memory allocates and deallocates
NVCOMPILER_ACC_PROFLIB	Enables 3rd party tools interface using the new profiler dynamic library interface.
NVCOMPILER_ACC_SYNCHRONOUS	Disables asynchronous launches and data movement.
NVCOMPILER_ACC_TIME	Enables a lightweight profiler to measure data movement and accelerator kernel execution time and print a summary at the end of program execution.

## 6.8. Profiling Accelerator Kernels

### Support for Profiler/Trace Tool Interface

The NVIDIA HPC Compilers support the OpenACC Profiler/Trace Tools Interface. This is the interface used by the NVIDIA profilers to collect performance measurements of OpenACC programs.

## Using NVCOMPILER\_ACC\_TIME

Setting the environment variable `NVCOMPILER_ACC_TIME` to a nonzero value enables collection and printing of simple timing information about the accelerator regions and generated kernels.



Turn off all CUDA Profilers (NVIDIA's Visual Profiler, NVPROF, CUDA\_PROFILE, etc) when enabling `NVCOMPILER_ACC_TIME`, they use the same library to gather performance data and cannot be used concurrently.

## Accelerator Kernel Timing Data

```
bb04.f90
  s1
15: region entered 1 times
    time(us): total=1490738
              init=1489138 region=1600
              kernels=155 data=1445
    w/o init: total=1600 max=1600
              min=1600 avg=1600
18: kernel launched 1 times
    time(us): total=155 max=155 min=155 avg=155
```

In this example, a number of things are occurring:

- ▶ For each accelerator region, the file name `bb04.f90` and subroutine or function name `s1` is printed, with the line number of the accelerator region, which in the example is 15.
- ▶ The library counts how many times the region is entered (1 in the example) and the microseconds spent in the region (in this example 1490738), which is split into initialization time (in this example 1489138) and execution time (in this example 1600).
- ▶ The execution time is then divided into kernel execution time and data transfer time between the host and GPU.
- ▶ For each kernel, the line number is given, (18 in the example), along with a count of kernel launches, and the total, maximum, minimum, and average time spent in the kernel, all of which are 155 in this example.

## 6.9. OpenACC Runtime Libraries

This section provides an overview of the user-callable functions and library routines that are available for use by programmers to query the accelerator features and to control behavior of accelerator-enabled programs at runtime.



In Fortran, none of the OpenACC runtime library routines may be called from a PURE or ELEMENTAL procedure.

## 6.9.1. Runtime Library Definitions

There are separate runtime library files for Fortran, and for C++ and C.

### C++ and C Runtime Library Files

In C++ and C, prototypes for the runtime library routines are available in a header file named `accel.h`. All the library routines are `extern` functions with 'C' linkage. This file defines:

- ▶ The prototypes of all routines in this section.
- ▶ Any data types used in those prototypes, including an enumeration type to describe types of accelerators.

### Fortran Runtime Library Files

In Fortran, interface declarations are provided in a Fortran include file named `accel_lib.h` and in a Fortran module named `accel_lib`. These files define:

- ▶ Interfaces for all routines in this section.
- ▶ Integer parameters to define integer kinds for arguments to those routines.
- ▶ Integer parameters to describe types of accelerators.

## 6.9.2. Runtime Library Routines

**Table 17** lists and briefly describes the runtime library routines supported by the NVIDIA HPC Compilers in addition to the standard OpenACC runtime API routines.

Table 17 Accelerator Runtime Library Routines

This Runtime Library Routine...	Does this...
<code>acc_allocs</code>	Returns the number of arrays allocated in data or compute regions.
<code>acc_bytesalloc</code>	Returns the total bytes allocated by data or compute regions.
<code>acc_bytesin</code>	Returns the total bytes copied in to the accelerator by data or compute regions.
<code>acc_bytesout</code>	Returns the total bytes copied out from the accelerator by data or compute regions.
<code>acc_copyins</code>	Returns the number of arrays copied in to the accelerator by data or compute regions.
<code>acc_copyouts</code>	Returns the number of arrays copied out from the accelerator by data or compute regions.
<code>acc_disable_time</code>	Tells the runtime to stop profiling accelerator regions and kernels.
<code>acc_enable_time</code>	Tells the runtime to start profiling accelerator regions and kernels, if it is not already doing so.

This Runtime Library Routine...	Does this...
acc_exec_time	Returns the number of microseconds spent on the accelerator executing kernels.
acc_frees	Returns the number of arrays freed or deallocated in data or compute regions.
acc_get_device	Returns the type of accelerator device used to run the next accelerator region, if one is selected.
acc_get_device_num	Returns the number of the device being used to execute an accelerator region.
acc_get_free_memory	Returns the total available free memory on the attached accelerator device.
acc_get_memory	Returns the total memory on the attached accelerator device.
acc_get_num_devices	Returns the number of accelerator devices of the given type attached to the host.
acc_kernels	Returns the number of accelerator kernels launched since the start of the program.
acc_present_dump	Summarizes all data present on the current device.
acc_present_dump_all	Summarizes all data present on all devices.
acc_regions	Returns the number of accelerator regions entered since the start of the program.
acc_total_time	Returns the number of microseconds spent in accelerator compute regions and in moving data for accelerator data regions.

## 6.10. Supported Intrinsics

An intrinsic is a function available in a given language whose implementation is handled specifically by the compiler. Typically, an intrinsic substitutes a sequence of automatically-generated instructions for the original function call. Since the compiler has an intimate knowledge of the intrinsic function, it can better integrate it and optimize it for the situation.

Intrinsics make the use of processor-specific enhancements easier because they provide a language interface to assembly instructions. In doing so, the compiler manages things that the user would normally have to be concerned with, such as register names, register allocations, and memory locations of data.

This section contains an overview of the Fortran and C intrinsics that the accelerator supports.

### 6.10.1. Supported Fortran Intrinsics Summary Table

**Table 18** is an alphabetical summary of the supported Fortran intrinsics that the accelerator supports. These functions are specific to Fortran 90/95 unless otherwise specified.

In most cases support is provided for all the data types for which the intrinsic is valid. When support is available for only certain data types, the middle column of the table specifies which ones, using the following codes:

I for integer

S for single precision real

C for single precision complex

D for double precision real

Z for double precision complex

Table 18 Supported Fortran Intrinsics

This intrinsic		Return value
ABS	I,S,D	absolute value of the argument.
ACOS		arccosine of the specified argument.
AINT		truncation of the argument to a whole number.
ANINT		nearest whole number of the real argument.
ASIN		arcsine of the argument.
ATAN		arctangent of the argument.
ATAN2		angle in radians of the complex value first-argument + i*second-argument.
COS	S,D,C,Z	cosine of the argument.
COSH		hyperbolic cosine of the argument.
DBLE	S,D	conversion of the argument to double precision real.
DPROD		double precision product of two single precision arguments.
EXP	S,D,C,Z	natural exponential value of the argument.
IAND		result of logical AND of the two integer arguments.
IEOR		result of the boolean exclusive OR of the two integer arguments.
INT	I,S,D	conversion of the argument to integer type.
IOR		result of the boolean inclusive OR of the two integer arguments.
LOG	S,D,C,Z	base-e (natural logarithm) of the argument.
LOG10		base-10 logarithm of the argument.
MAX		maximum value of the arguments.
MIN		minimum value of the arguments.
MOD	I	remainder of the first argument divided by the second argument.
NINT		nearest integer of the real argument.
NOT		logical complement of the integer argument.
REAL	I,S,D	conversion of the argument to real.
SIGN		absolute value of first argument times the sign of second argument.
SIN	S,D,C,Z	sine of the argument.
SINH		hyperbolic sine of the argument.
SQRT	S,D,C,Z	square root of the argument.

This intrinsic		Return value
TAN		tangent of the argument.
TANH		hyperbolic tangent of the argument.

## 6.10.2. Supported C Intrinsic Summary Table

This section contains two alphabetical summaries – one for double functions and a second for float functions. These lists contain only those C intrinsics that the accelerator supports.

Table 19 Supported C Intrinsic Double Functions

This intrinsic	Return value
acos	arccosine of the argument.
asin	arcsine of the argument.
atan	arctangent of the argument.
atan2	arctangent of $y/x$ , where $y$ is the first argument, $x$ the second.
cos	cosine of the argument.
cosh	hyperbolic cosine of the argument.
exp	exponential value of the argument.
fabs	absolute value of the argument.
fmax	maximum value of the two arguments
fmin	minimum value of the two arguments
log	natural logarithm of the argument.
log10	base-10 logarithm of the argument.
pow	value of the first argument raised to the power of the second argument.
sin	value of the sine of the argument.
sinh	hyperbolic sine of the argument.
sqrt	square root of the argument.
tan	tangent of the argument.
tanh	hyperbolic tangent of the argument.

Table 20 Supported C Intrinsic Float Functions

This intrinsic	Return value
acosf	arccosine of the argument.
asinf	arcsine of the argument.
atanf	arctangent of the argument.
atan2f	arctangent of $y/x$ , where $y$ is the first argument, $x$ the second.



This intrinsic	Return value
cosf	cosine of the argument.
coshf	hyperbolic cosine of the argument.
expf	exponential value of the argument.
fabsf	absolute value of the argument.
logf	natural logarithm of the argument.
log10f	base-10 logarithm of the argument.
powf	value of the first argument raised to the power of the second argument.
sinf	value of the sine of the argument.
sinhf	hyperbolic sine of the argument.
sqrtf	square root of the argument.
tanf	tangent of the argument.
tanhf	hyperbolic tangent of the argument.

# Chapter 7.

## USING OPENMP

OpenMP is a specification for a set of compiler directives, an applications programming interface (API), and a set of environment variables that can be used to specify parallel execution in Fortran, C++, and C programs. For general information about using OpenMP and to obtain a copy of the OpenMP specification, refer to the [OpenMP organization's website](#).

The NVFORTRAN, NVC++, and NVC compilers support a subset of the OpenMP Application Program Interface for CPUs and GPUs. In defining this subset, we have focused on OpenMP 5.0 features that will enable CPU and GPU targeting for OpenMP applications with a goal of encouraging programming practices that are portable and scalable. For features that are to be avoided, wherever possible, the directives and API calls related to those features are parsed and ignored to maximize portability. Where ignoring such features is not possible, or could result in ambiguous or incorrect execution, the compilers emit appropriate error messages at compile- or run-time.

OpenMP applications properly structured for GPUs, meaning they expose massive parallelism and have relatively little or no synchronization in GPU-side code segments, should compile and execute with performance on par with or close to equivalent OpenACC. Codes that are not well-structured for GPUs may perform poorly but should execute correctly.

Use the **-mp** compiler switch to enable processing of OpenMP directives and pragmas. The most important sub-options to **-mp** are the following:

- ▶ **gpu**: OpenMP directives are compiled for GPU execution plus multicore CPU fallback; this feature is supported on NVIDIA V100 or later GPUs.
- ▶ **multicore**: OpenMP directives are compiled for multicore CPU execution only; this sub-option is the default.

### Predefined Macros

The following macros corresponding to the offload target compiled for are added implicitly:

- ▶ **\_\_NVCOMPILER\_OPENMP\_GPU** when OpenMP target directives are compiled for GPU.

- `__NVCOMPILER_OPENMP_MULTICORE` when OpenMP target directives are compiled for multicore CPU.

## 7.1. Environment Variables

The OpenMP specification includes many environment variables related to program execution.

### Thread affinity

One important environment variable is `OMP_PROC_BIND`. It controls the OpenMP CPU thread affinity policy. When thread affinity is disabled, the operating system is free to move threads between the available CPU cores. When thread affinity is enabled, each thread is bound to a subset of the available CPU cores. The environment variable `OMP_PLACES` can be used to specify how a subset of the available CPU cores is determined for each thread. When set to a valid value, this environment variable will enable thread affinity and override the default thread affinity policy.

Binding threads to certain CPU cores is often beneficial for application performance, because that can improve the CPU cache hit rate and limit memory transactions between different NUMA nodes. Therefore, it is important to consider enabling thread affinity for your application.

The default value of `OMP_PROC_BIND` is `false`. Thus, thread affinity is disabled by default. This is a conservative setting that allows certain classes of applications (such as OpenMP + MPI) to create multiple processes without taking special care of the thread affinity policy to avoid binding threads in different processes to the same CPU cores.

The following table explains the simplest possible values of `OMP_PROC_BIND`. For the comprehensive explanation of `OMP_PROC_BIND` and `OMP_PLACES`, please refer to the OpenMP specification.

Value	Behavior
<code>OMP_PROC_BIND=false</code>	Thread affinity is disabled unless <code>OMP_PLACES</code> is set to a valid value. When thread affinity is disabled, the operating system is free to assign threads to any available CPU core at any time of the application execution. This is the default value.
<code>OMP_PROC_BIND=true</code>	Thread affinity is enabled. Unless <code>OMP_PLACES</code> is set, the implementation attempts to assign threads optimally to CPU cores to maximize the cache hit rate and minimize the number of memory transactions between NUMA nodes.

Device offload

Another important environment variable to understand is **OMP\_TARGET\_OFFLOAD**. Use this environment variable to affect the behavior of execution on host and device including host fallback. The following table explains the behavior determined by each of the values to which you can set this environment variable.

Value	Behavior
<b>OMP_TARGET_OFFLOAD=DEFAULT</b>	Try to execute on a GPU; if a supported GPU is not available, fallback to the host
<b>OMP_TARGET_OFFLOAD=DISABLED</b>	Do not execute on the GPU even if one is available; execute on the host
<b>OMP_TARGET_OFFLOAD=MANDATORY</b>	Execute on a GPU or terminate the program

Number of teams on device

When an application offloads an **omp target teams** construct to the GPU, the number of teams is calculated automatically unless the construct has a **num\_teams** clause. The automatic setting of the number of teams can be limited to a maximum value provided by the **OMP\_NUM\_TEAMS** environment variable. The same maximum value can also be set by the application at run time with the function **omp\_set\_num\_teams**.

Value	Behavior
<b>OMP_NUM_TEAMS=&lt;positive_integer&gt;</b>	Maximum number of teams on device

For the comprehensive explanation of **OMP\_NUM\_TEAMS**, please refer to the OpenMP specification.

Number of threads in teams

An **omp target teams** construct offloaded to the GPU creates a league of teams each consisting of a certain number of threads. The number of threads is the same for all teams in the league, and is calculated automatically unless the construct has a **thread\_limit** clause.

The environment variable **OMP\_TEAMS\_THREAD\_LIMIT** can be used to limit the maximum number of threads in teams. The same maximum value can be set by the application with the runtime function **omp\_set\_teams\_thread\_limit**.

For NVIDIA GPUs, we recommend using values that are multiples of 32 (which is the size of the GPU thread warp). That equally applies to the **OMP\_TEAMS\_THREAD\_LIMIT** environment variable, the **omp\_set\_teams\_thread\_limit** function and the **thread\_limit** clause. For any other value, the actual limit on the number of threads per team will likely be rounded down to the nearest multiple of 32. The same guidance applies to the **num\_threads** clause as well.

Value	Behavior
<code>OMP_TEAMS_THREAD_LIMIT=&lt;positive_integer&gt;</code>	Maximum number of threads in teams

For the comprehensive explanation of `OMP_TEAMS_THREAD_LIMIT`, please refer to the OpenMP specification.

Forcing the number of device teams and threads

In certain situations, for instance for debugging or performance tuning, it may be desirable to specify an exact number of teams and threads on the GPU. While OpenMP offers a number of convenient ways to control that, e.g. the `num_teams` and `thread_limit` clauses, as well as the environment variables described above, they do not guarantee an exact teams and threads configuration.

The NVIDIA HPC OpenMP Runtime supports the `NVCOMPILER_OMP_CUDA_GRID` environment variable. When set, it requests the runtime to use the exact number of teams and threads per team when running OpenMP compute constructs on the GPU. Essentially, its effect is to use a specific CUDA grid configuration for any kernel, bypassing runtime and compiler guidance.

Value	Behavior
<code>NVCOMPILER_OMP_CUDA_GRID=&lt;num_blocks&gt;,&lt;num_th</code>	The <code>&lt;num_blocks&gt;</code> and <code>&lt;num_threads&gt;</code> must be positive integers. They are used to form a CUDA grid when running GPU kernels associated with <code>omp target</code> compute constructs.

However, even with an exact CUDA grid specified, the runtime may still use a corrected configuration if that is necessary for a successful kernel launch.

Please refer to the [CUDA C++ Programming Guide](#) for the detailed explanation of how the CUDA kernel execution configurations work.

## 7.2. Fallback Mode

The HPC compilers support host fallback of OpenMP `target` regions when no GPU is present or `OMP_TARGET_OFFLOAD` is set to `DISABLED`. Execution should always be correct but the performance of the target region may not always be optimal when run on the host. OpenMP target regions prescriptively structured for optimal execution on GPUs may not perform well when run on the dissimilar architecture of the CPU. To provide performance portability between host and device, we recommend use of the `loop` construct.

firstprivates with `nowait` not supported for host execution

There is currently a limitation on the use of the `nowait` clause on target regions intended for execution on the host (`-mp` or `-mp=cpu` with `OMP_TARGET_OFFLOAD=DISABLED`). If the target region references variables having the `firstprivate` data-sharing attribute, their concurrent updates are not guaranteed to be safe. To work around this limitation, when running on the host, we recommend avoiding the `nowait` clause on such target regions or equivalently using the `taskwait` construct immediately following the region.

## 7.3. Loop

The HPC compilers support the `loop` construct with an extension to the default binding thread set mechanism specified by OpenMP in order to allow the compilers the freedom to analyze loops and dependencies to generate highly parallel code for CPU and GPU targets. In other words, the compilers map `loop` to either teams or to threads, as the compiler chooses, unless the user explicitly specifies otherwise. The mapping selected is specific to each target architecture even within the same executable (i.e., GPU offload and host fallback) thereby facilitating performance portability.

The shape of the parallelism offered by NVIDIA's GPUs, consisting of thread blocks and three dimensions of threads therein, differs from the multi-threaded vector parallelism of modern CPUs. The following table summarizes the OpenMP mapping to NVIDIA GPUs and multicore CPUs:

Construct	CPU	GPU
<code>!\$omp target</code>		starts offload
<code>!\$omp teams</code>	single team	CUDA thread blocks in grid
<code>!\$omp parallel</code>	CPU threads	CUDA threads within thread block
<code>!\$omp simd</code>	hint for vector instructions	<code>simdlen(1)</code>

HPC programs need to leverage all available parallelism to achieve performance. The programmer can attempt to become an expert in the intricacies of each target architecture and use that knowledge to structure programs accordingly. This prescriptive model can be successful but tends to increase source code complexity and often requires restructuring for each new target architecture. Here's an example where a programmer explicitly requests the steps the compiler should take to map parallelism to two targets:

```
#ifdef TARGET_GPU
    #pragma omp target teams distribute reduction(max:error)
#else
    #pragma omp parallel for reduction(max:error)
#endif
for( int j = 1; j < n-1; j++) {
#ifdef TARGET_GPU
    #pragma omp parallel for reduction(max:error)
#endif
    for( int i = 1; i < m-1; i++ ) {
        Anew[j][i] = 0.25f * ( A[j][i+1] + A[j][i-1]
```

```

        + A[j-1][i] + A[j+1][i]);
    error = fmaxf( error, fabsf(Anew[j][i]-A[j][i]));
}

```

An alternative is for the programmer to focus on exposing parallelism in a program and allowing a compiler to do the mapping onto the target architectures. The HPC compilers' implementation of **loop** supports this descriptive model. In this example, the programmer specifies the loop regions to be parallelized by the compiler and the compilers parallelize **loop** across teams and threads:

```

#pragma omp target teams loop reduction(max:error)
for( int j = 1; j < n-1; j++) {
    #pragma omp loop reduction(max:error)
    for( int i = 1; i < m-1; i++ ) {
        Anew[j][i] = 0.25f * ( A[j][i+1] + A[j][i-1]
                             + A[j-1][i] + A[j+1][i]);
        error = fmaxf( error, fabsf(Anew[j][i]-A[j][i]));
    }
}

```

The programmer's tuning tool with **loop** is the **bind** clause. The following table extends the previous mapping example:

Construct	CPU	GPU
<b>!\$omp loop bind(teams)</b>	threads	CUDA thread blocks and threads
<b>!\$omp loop bind(parallel)</b>	threads	CUDA threads
<b>!\$omp loop bind(thread)</b>	single thread (useful for vector instructions)	single thread

Orphaned **loop** constructs within a single file are supported; a binding region of either **parallel** or **thread** must be specified with such loops via the **bind** clause. The compilers support **loop** regions containing procedure calls as long as the callee does not contain OpenMP directives.

Here are a few additional examples using **loop**. We also show examples of the type of information the compiler would provide when using the **-Minfo** compiler option.

Use of **loop** in Fortran:

```

!$omp target teams loop
do nllloc_blk = 1, nllloc_blksize
  do igp = 1, ngpown
    do ig_blk = 1, ig_blksize
      do ig = ig_blk, ncouls, ig_blksize
        do n1_loc = nllloc_blk, ntband_dist, nllloc_blksize
          !expensive computation codes
        enddo
      enddo
    enddo
  enddo
enddo

$ nvfortran test.f90 -mp=gpu -Minfo=mp
42, !$omp target teams loop
42, Generating "nvkernel MAIN_F1L42_1" GPU kernel
   Generating Tesla code
43, Loop parallelized across teams ! blockidx%x

```

```

44, Loop run sequentially
45, Loop run sequentially
46, Loop run sequentially
47, Loop parallelized across threads(128) ! threadidx%x
42, Generating Multicore code
43, Loop parallelized across threads

```

### Use of **loop**, **collapse**, and **bind**:

```

!$omp target teams loop collapse(3)
do nllloc_blk = 1, nllloc_blksize
  do igp = 1, ngpown
    do ig_blk = 1, ig_blksize
      !$omp loop bind(parallel) collapse(2)
      do ig = ig_blk, ncouls, ig_blksize
        do n1_loc = nllloc_blk, ntband_dist, nllloc_blksize
          !expensive computation codes
        enddo
      enddo
    enddo
  enddo
enddo

$ nvfortran test.f90 -mp=gpu -Minfo=mp

42, !$omp target teams loop
42, Generating "nvkernel_MAIN__F1L42_1" GPU kernel
   Generating Tesla code
43, Loop parallelized across teams collapse(3) ! blockidx%x
44,   ! blockidx%x collapsed
45,   ! blockidx%x collapsed
47, Loop parallelized across threads(128) collapse(2) ! threadidx%x
48,   ! threadidx%x collapsed
42, Generating Multicore code
43, Loop parallelized across threads

```

### Use of **loop**, **collapse**, and **bind(thread)**:

```

!$omp target teams loop collapse(3)
do nllloc_blk = 1, nllloc_blksize
  do igp = 1, ngpown
    do ig_blk = 1, ig_blksize
      !$omp loop bind(thread) collapse(2)
      do ig = ig_blk, ncouls, ig_blksize
        do n1_loc = nllloc_blk, ntband_dist, nllloc_blksize
          ! expensive computation codes
        enddo
      enddo
    enddo
  enddo
enddo

$ nvfortran test.f90 -mp=gpu -Minfo=mp

42, !$omp target teams loop
42, Generating "nvkernel_MAIN__F1L42_1" GPU kernel
   Generating Tesla code
43, Loop parallelized across teams, threads(128) collapse(3) ! blockidx%x
   threadidx%x
44,   ! blockidx%x threadidx%x collapsed
45,   ! blockidx%x threadidx%x collapsed
47, Loop run sequentially
48,   collapsed
42, Generating Multicore code
43, Loop parallelized across threads

```



## 7.4. OpenMP Subset

This section contains the subset of OpenMP 5.0 features that the HPC compilers support. We have attempted to define this subset of features to be those that enable, where possible, OpenMP-for-GPU application performance that closely mirrors the success NVIDIA has seen with OpenACC. Almost every feature supported on NVIDIA GPUs is also supported on multicore CPUs, although the reverse is not true. Most constructs from OpenMP 3.1 and OpenMP 4.5 that apply to multicore CPUs are supported for CPU targets, and some features from OpenMP 5.0 are supported as well.

OpenMP target offload to NVIDIA GPUs is supported on NVIDIA V100 or later GPUs.

The section numbers below correspond to the section numbers in the OpenMP Application Programming Interface Version 5.0 November 2018 document.

### 2. Directives

#### 2.3 Variant Directives

##### 2.3.4 Metadirectives

The **target\_device/device** context selector is supported with the **kind(host|nohost|cpu|gpu)** and **arch(nvtx|nvtx64)** trait selectors. The **arch** trait property **nvtx** is an alias for **nvtx64**; any other **arch** trait properties are treated as not matching or are ignored. The **isa** selector is treated as not matching or is ignored; no support is provided to select a context based on NVIDIA GPU compute capability.

The **implementation** context selector is supported with the **vendor(nvidia)** trait selector.

The **user** context selector is supported with the **condition(expression)** trait selector including dynamic **user** traits.

The syntax **begin/end metadirective** is not supported.

##### 2.3.5 Declare Variant Directive

The **device** context selector is supported with the **kind(host|nohost|cpu|gpu)** and **arch(nvtx|nvtx64)** trait selectors. The **arch** trait property **nvtx** is an alias for **nvtx64**; any other **arch** trait properties are treated as not matching or are ignored. The **isa** selector is also treated as not matching or is ignored; no support is provided to select a context based on NVIDIA GPU compute capability.

The **implementation** context selector is supported with the **vendor(nvidia)** trait selector; all other implementation trait selectors are treated as not matching.

The syntax **begin/end declare variant** is supported for C/C++.

### 2.4 Requires Directive

The **requires** directive has limited support. The requirement clauses **unified\_address** and **unified\_shared\_memory** are accepted but have no effect. To activate OpenMP unified shared memory programming a command-line option needs to be passed in (refer to [OpenMP with CUDA Unified Memory](#) for more details).

## 2.5 Internal Control Variables

ICV support is as follows.

- ▶ **dyn-var**, **nthread-var**, **thread-limit-var**, **max-active-levels-var**, **active-levels-var**, **levels-var**, **run-sched-var**, **dyn-sched-var**, and **stacksize-var** are supported
- ▶ **place-partition-var**, **bind-var**, **wait-policy-var**, **display-affinity-var**, **default-device-var**, and **target-offload-var** are supported only on the CPU
- ▶ **affinity-format-var** is supported only on the CPU; its value is immutable
- ▶ **max-task-priority-var**, **def-allocator-var** are not supported
- ▶ **cancel-var** is not supported; it always returns false

## 2.6 Parallel Construct

Support for **parallel** construct clauses is as follows.

- ▶ The **num\_threads**, **default**, **private**, **firstprivate**, and **shared** clauses are supported
- ▶ The **reduction** clause is supported as described in 2.19.5
- ▶ The **if** and **copyin** clauses are supported only for CPU targets; the compiler emits an error for GPU targets
- ▶ The **proc\_bind** clause is supported only for CPU targets; it is ignored for GPU targets
- ▶ The **allocate** clause is ignored

## 2.7 Teams Construct

The **teams** construct is supported only when nested within a **target** construct that contains no statements, declarations, or directives outside the **teams** construct, or as a combined **target teams** construct. The **teams** construct is supported for GPU targets. If the **target** construct falls back to CPU mode, the number of teams is one. Support for **teams** construct clauses is as follows.

- ▶ The **num\_teams**, **thread\_limit**, **default**, **private**, and **firstprivate** clauses are supported
- ▶ The **reduction** clause is supported as described in 2.19.5
- ▶ The **shared** clause is supported for CPU targets and is supported for GPU targets in unified-memory mode
- ▶ The **allocate** clause is ignored

## 2.8 Worksharing Constructs

### 2.8.1 Sections Construct

The **sections** construct is supported only for CPU targets; the compiler emits an error for GPU targets. Support for **sections** construct clauses is as follows.

- ▶ The **private** and **firstprivate** clauses are supported
- ▶ The **reduction** clause is supported as described in 2.19.5
- ▶ The **lastprivate** clause is supported; the optional **lastprivate** modifier is not supported
- ▶ The **allocate** clause is ignored

### 2.8.2 Single Construct

Support for **single** construct clauses is as follows.

- ▶ The **private**, **firstprivate**, and **nowait** clauses are supported
- ▶ The **copyprivate** clause is supported only for CPU targets; the compiler emits an error for GPU targets
- ▶ The **allocate** clause is ignored

### 2.8.3 Workshare Construct

The **workshare** construct is supported in Fortran only for CPU targets; the compiler emits an error for GPU targets.

## 2.9 Loop-Related Constructs

### 2.9.2 Worksharing-Loop Construct (for/do)

Support for worksharing **for** and **do** construct clauses is as follows.

- ▶ The **private**, **firstprivate**, and **collapse** clauses are supported
- ▶ The **reduction** clause is supported as described in 2.19.5
- ▶ The **schedule** clause is supported; the optional modifiers are not supported
- ▶ The **lastprivate** clause is supported; the optional **lastprivate** modifier is not supported
- ▶ The **ordered** clause is supported only for CPU targets; **ordered(n)** clause is not supported
- ▶ The **linear** clause is not supported
- ▶ The **order(concurrent)** clause is ignored
- ▶ The **allocate** clause is ignored

### 2.9.3 SIMD Directives

The **simd** construct can be used to provide tuning hints for CPU targets; the **simd** construct is ignored for GPU targets. Support for **simd** construct clauses is as follows.

- ▶ The **reduction** clause is supported as described in 2.19.5
- ▶ The **lastprivate** clause is supported; the optional **lastprivate** modifier is not supported

- ▶ The **if**, **simdlen**, and **linear** clauses are not supported
- ▶ The **safelen**, **aligned**, **nontemporal**, and **order (concurrent)** clauses are ignored

The composite **for simd** and **do simd** constructs are supported for CPU targets; they are treated as **for** and **do** directives for GPU targets. Supported **simd** clauses are supported on the composite constructs for the CPU. Any **simd** clauses are ignored for GPU targets.

The **declare simd** directive is ignored.

## 2.9.4 Distribute Directives

The **distribute** construct is supported within a **teams** construct. Support for **distribute** construct clauses is as follows:

- ▶ The **private**, **firstprivate**, **collapse**, and **dist\_schedule (static [ , chunksize] )** clauses are supported
- ▶ The **lastprivate** clause is not supported
- ▶ The **allocate** clause is ignored

The **distribute simd** construct is treated as a **distribute** construct and is supported for GPU targets; valid supported **distribute** clauses are accepted; **simd** clauses are ignored. The **distribute simd** construct is not supported for CPU targets.

The **distribute parallel for** or **distribute parallel do** constructs are supported for GPU targets. Valid supported **distribute** and **parallel** and **for** or **do** clauses are accepted. The **distribute parallel for** or **distribute parallel do** constructs are not supported for CPU targets.

The **distribute parallel for simd** or **distribute parallel do simd** constructs are treated as **distribute parallel for** or **distribute parallel do** constructs and are supported for GPU targets. These are not supported for CPU targets.

## 2.9.5 Loop Construct

Support for **loop** construct clauses is as follows.

- ▶ The **private**, **bind**, and **collapse** clauses are supported
- ▶ The **reduction** clause is supported as described in 2.19.5
- ▶ The **order (concurrent)** clause is assumed
- ▶ The **lastprivate** clause is not supported

## 2.10 Tasking Constructs

### 2.10.1 Task Construct

The **task** construct is supported for CPU targets. The compiler emits an error when it encounters **task** within a **target** construct. Support for **task** construct clauses is as follows:

- ▶ The **if**, **final**, **default**, **private**, **firstprivate**, and **shared** clauses are supported
- ▶ The **depend**(**[dependmodifier]**, **dependtype** : **list**) clause is supported as described in 2.17.11

#### 2.10.4 Taskyield Construct

The **taskyield** construct is supported for CPU targets; it is ignored for GPU targets.

### 2.11 Memory Management Directives

The memory management allocators, memory management API routines, and memory management directives are not supported

### 2.12 Device Directives

#### 2.12.1 Device Initialization

Depending on how the program is compiled and linked, device initialization may occur at the first **target** construct or API routine call, or may occur implicitly at program startup.

#### 2.12.2 Target Data Construct

The **target data** construct is supported for GPU targets. Support for **target data** construct clauses is as follows.

- ▶ The **if**, **device**, **use\_device\_ptr**, and **use\_device\_addr** clauses are supported
- ▶ The **map** clause is supported as described in 2.19.7

#### 2.12.3 Target Enter Data Construct

The **target enter data** construct is supported for GPU targets. Support for **enter data** construct clauses is as follows.

- ▶ The **if**, **device**, and **nowait** clauses are supported
- ▶ The **map** clause is supported as described in 2.19.7.
- ▶ The **depend**(**[dependmodifier]**, **dependtype** : **list**) clause is supported as described in 2.17.11

#### 2.12.4 Target Exit Data Construct

The **target exit data** construct is supported for GPU targets. Support for **exit data** construct clauses is as follows.

- ▶ The **if**, **device**, and **nowait** clauses are supported
- ▶ The **map** clause is supported as described in 2.19.7.
- ▶ The **depend**(**[dependmodifier]**, **dependtype** : **list**) clause is supported as described in 2.17.11

#### 2.12.5 Target Construct

The **target** construct is supported for GPU targets. If there is no GPU or GPU offload is otherwise disabled, execution falls back to CPU mode. Support for **target** construct clauses is as follows:

- ▶ The **if**, **private**, **firstprivate**, **is\_device\_ptr**, and **nowait** clauses are supported
- ▶ The **device** clause is supported without the device-modifier **ancestor** keyword
- ▶ The **map** clause is supported as described in 2.19.7
- ▶ The **defaultmap** clause is supported using OpenMP 5.0 semantics
- ▶ The **depend**(**[dependmodifier,] dependtype : list**) clause is supported as described in 2.17.11
- ▶ The **allocate** and **uses\_allocate** clauses are ignored

### 2.12.6 Target Update Construct

The **target update** construct is supported for GPU targets. Support for **target update** construct clauses is as follows.

- ▶ The **if**, **device**, and **nowait** clauses are supported.
- ▶ The **to** and **from** clauses are supported without **mapper** or **mapid**
- ▶ The **depend**(**[dependmodifier,] dependtype : list**) clause is supported as described in 2.17.11

Array sections are supported in **to** and **from** clauses, including noncontiguous array sections. Array section strides are not supported. If the array section is noncontiguous, the OpenMP runtime may have to use multiple host-to-device or device-to-host data transfer operations, which increases the overhead. If the host data is in host-pinned memory, then **update** data transfers with the **nowait** clause are asynchronous. This means the data transfer for a **target update to nowait** may not occur immediately or synchronously with the program thread, and any changes to the data may affect the transfer, until a synchronizing operation is reached. Similarly, a **target update from nowait** may not occur immediately or synchronously with the program thread, and the downloaded data may not be available until a synchronizing operation is reached. If the host data is not in host-pinned memory, then **update** data transfers with the **nowait** clause may require that the data transfer operation use an intermediate pinned buffer managed by the OpenMP runtime library, and that a memory copy operation on the host between the program memory and the pinned buffer may be needed before starting or before finishing the transfer operation, which affects overhead and performance. To learn more about the pinned buffer, please refer to [Staging Memory Buffer](#).

### 2.12.7 Declare Target Construct

The **declare target** construct is supported for GPU targets.

- ▶ **declare target ... end declare target** is supported
- ▶ **declare target(list)** is supported
- ▶ The **to(list)** clause is supported

- ▶ The **device\_type** clause is supported for C/C++

A function or procedure that is referenced in a function or procedure that appears in a **declare target to** clause (explicitly or implicitly) is treated as if its name had implicitly appeared in a **declare target to** clause.

### 2.13 Combined Constructs

Combined constructs are supported to the extent that the component constructs are themselves supported.

### 2.14 Clauses on Combined and Composite Constructs

Clauses on combined constructs are supported to the extent that the clauses are supported on the component constructs.

### 2.16 Master Construct

The **master** construct is supported for CPU and GPU targets.

### 2.17 Synchronization Constructs and Clauses

#### 2.17.1 Critical Construct

The **critical** construct is supported only for CPU targets; the compiler emits an error for GPU targets.

#### 2.17.2 Barrier Construct

The **barrier** construct is supported.

#### 2.17.3 Implicit Barriers

Implicit barriers are implemented.

#### 2.17.4 Implementation-Specific Barriers

There may be implementation-specific barriers, and they may be different for CPU targets than for GPU targets.

#### 2.17.5 Taskwait Construct

The **taskwait** construct is supported only for CPU targets; it is ignored for GPU targets.

- ▶ The **depend([dependmodifier,] dependtype : list)** clause is supported as described in 2.17.11

#### 2.17.6 Taskgroup Construct

The **taskgroup** construct is supported only for CPU targets. It is ignored for GPU targets.

#### 2.17.7 Atomic Construct

Support for **atomic** construct clauses is as follows.

- ▶ The **read**, **write**, **update**, and **capture** clauses are supported.

- ▶ The memory order clauses **seq\_cst**, **acq\_rel**, **release**, **acquire**, **relaxed** are not supported
- ▶ The **hint** clause is ignored

### 2.17.8 Flush Construct

The **flush** construct is supported only for CPU targets.

### 2.17.9 Ordered Construct and Ordered Directive

The **ordered** block construct is supported only for CPU targets.

### 2.17.11 Depend Clause

The **depend** clause is supported on CPU targets. It is not supported on GPU targets. The dependence types **in**, **out**, and **inout** are supported. The dependence types **mutexinoutset** and **depobj**, dependence modifier **iterator(iters)**, **depend(source)**, and **depend(sink:vector)** are not supported.

## 2.19 Data Environment

### 2.19.2 Threadprivate Directive

The **threadprivate** directive is supported only for CPU targets. It is not supported for GPU targets; references to **threadprivate** variables in device code are not supported.

### 2.19.5 Reduction Clauses and Directives

The **reduction** clause is supported. The optional modifier is not supported.

### 2.19.6 Data Copying Clauses

The data copying **copyin** and **copyprivate** clauses are supported only for CPU targets; the compiler emits a compile-time error for GPU targets.

### 2.19.7 Data Mapping Attribute Rules, Clauses, and Directives

- ▶ The **map([mapmod[,]...] maptype:) datalist)** clause is supported. Of the map-type-modifiers, **always** is supported, **close** is ignored, and **mapper(mapid)** is not supported.
- ▶ The **defaulttmap** clause is supported using OpenMP 5.0 semantics.

## 2.20 Nesting of Regions

For constructs supported in this subset, restrictions on nesting of regions is observed. Additionally, nested parallel regions on CPU are not supported and nested teams or parallel regions in a target region are not supported.

## Runtime Library Routines

### 3.2 Execution Environment Routines

The following execution environment runtime API routines are supported.



- ▶ `omp_set_num_threads, omp_get_num_threads,`  
`omp_get_max_threads, omp_get_thread_num, omp_get_thread_limit,`  
`omp_get_supported_active_levels, omp_set_max_active_levels,`  
`omp_get_max_active_levels, omp_get_level,`  
`omp_get_ancestor_thread_num, omp_get_team_size, omp_get_num_teams,`  
`omp_get_team_num, omp_is_initial_device`

The following execution environment runtime API routines are supported only on the CPU.

- ▶ `omp_get_num_procs, omp_set_dynamic, omp_get_dynamic,`  
`omp_set_schedule, omp_get_schedule, omp_in_final, omp_get_proc_bind,`  
`omp_get_num_places, omp_get_affinity_format, omp_set_default_device,`  
`omp_get_default_device, omp_get_num_devices, omp_get_device_num,`  
`omp_get_initial_device`

The following execution environment runtime API routines have limited support.

- ▶ `omp_get_cancellation, omp_get_nested`; supported only on the CPU; the value returned is always false
- ▶ `omp_display_affinity, omp_capture_affinity`; supported only on the CPU; the format specifier is ignored
- ▶ `omp_set_nested`; supported only on the CPU, the value is ignored

The following execution environment runtime API routines are not supported.

- ▶ `omp_get_place_num_procs, omp_get_place_proc_ids, omp_get_place_num,`  
`omp_get_partition_num_places, omp_get_partition_place_nums,`  
`omp_set_affinity_format, omp_get_max_task_priority,`  
`omp_pause_resource, omp_pause_resource_all`

### 3.3 Lock Routines

Lock runtime API routines are not supported on the GPU. The following lock runtime API routines are supported on the CPU.

- ▶ `omp_init_lock, omp_init_nest_lock, omp_destroy_lock,`  
`omp_destroy_nest_lock, omp_set_lock, omp_set_nest_lock,`  
`omp_unset_lock, omp_unset_nest_lock, omp_test_lock,`  
`omp_test_nest_lock`

The following lock runtime API routines are not supported.

- ▶ `omp_init_lock_with_hint, omp_init_nest_lock_with_hint`

### 3.4 Timing Routines

The following timing runtime API routines are supported.

- ▶ `omp_get_wtime, omp_get_wtick`

### 3.6 Device Memory Routines

The following device memory routines are supported only on the CPU.

- ▶ `omp_target_is_present`, `omp_target_associate_ptr`,  
`omp_target_disassociate_ptr`
- ▶ `omp_target_memcpy` and `omp_target_memcpy_rect` are only supported when copying to and from the same device.

The following device memory routines are supported on the CPU; we extend OpenMP to support these in target regions on a GPU, but only allocation and deallocation on the same device is supported.

- ▶ `omp_target_alloc`, `omp_target_free`

### 3.7 Memory Management Routines

The following memory management routines are supported.

- ▶ `omp_alloc`, `omp_free`

The following memory management routines are not supported.

- ▶ `omp_init_allocator`, `omp_destroy_allocator`,  
`omp_set_default_allocator`, `omp_get_default_allocator`

## 6 Environment Variables

The following environment variables have limited support.

- ▶ `OMP_SCHEDULE`, `OMP_NUM_THREADS`, `OMP_NUM_TEAMS`, `OMP_DYNAMIC`,  
`OMP_PROC_BIND`, `OMP_PLACES`, `OMP_STACKSIZE`, `OMP_WAIT_POLICY`,  
`OMP_MAX_ACTIVE_LEVELS`, `OMP_NESTED`, `OMP_THREAD_LIMIT`,  
`OMP_TEAMS_THREAD_LIMIT`, `OMP_DISPLAY_ENV`, `OMP_DISPLAY_AFFINITY`,  
`OMP_DEFAULT_DEVICE`, and `OMP_TARGET_OFFLOAD` are supported on CPU.
- ▶ `OMP_CANCELLATION` and `OMP_MAX_TASK_PRIORITY` are ignored.
- ▶ `OMP_AFFINITY_FORMAT`, `OMP_TOOL`, `OMP_TOOL_LIBRARIES`, `OMP_DEBUG`, and  
`OMP_ALLOCATOR` are not supported

## 7.5. Using metadirective

This section contains limitations affecting **metadirective** along with a few guidelines for its use.

The Fortran compiler does not support variants leading to an OpenMP directive for which a corresponding **end** directive is required.

Nesting **user** conditions, while legal, may create situations that the HPC Compilers do not handle gracefully. To avoid potential problems, use **device** traits inside **user** conditions instead. The following example illustrates this best practice.

Avoid nesting dynamic **user** conditions like this:

```
#pragma omp metadirective \
  when( user={condition(use_offload)} : target teams distribute) \
  default( parallel for schedule(static) )
  for (i = 0; i < N; i++) {
    ...
#pragma omp metadirective \
  when( user={condition(use_offload)} : parallel for)
    for (j = 0; j < N; j++) {
      ...
    }
    ...
  }
```

Instead, use **target\_device** and **device** traits within dynamic **user** conditions like this:

```
#pragma omp metadirective \
  when( target_device={kind(gpu)}, user={condition(use_offload)} : target teams
  distribute) \
  default( parallel for schedule(static) )
  for (i = 0; i < N; i++) {
    ...
#pragma omp metadirective \
  when( device={kind(gpu)} : parallel for)
    for (j = 0; j < N; j++) {
      ...
    }
    ...
  }
```

The HPC compilers do not support nesting **metadirective** inside a **target** construct applying to a syntactic block leading to a **teams** variant. Some examples:

The compilers will emit an error given the following code:

```
#pragma omp target map(to:v1,v2) map(from:v3)
{
#pragma omp metadirective \
when( device={arch("nvptx")} : teams distribute parallel for) \
default( parallel for)
  for (int i = 0; i < N; i++) {
    v3[i] = v1[i] * v2[i];
  }
}
```

The compilers will always match **device={arch("nvptx")}** given the following code:

```
#pragma omp target map(to:v1,v2) map(from:v3)
#pragma omp metadirective \
when( device={arch("nvptx")} : teams distribute parallel for) \
default( parallel for)
  for (int i = 0; i < N; i++) {
    v3[i] = v1[i] * v2[i];
  }
```

The compilers match **device={"arch"}** for GPU code, and **default** for host fallback, given the following code:

```
#pragma omp target teams distribute map(to:v1,v2) map(from:v3)
for (...)
```

```
{
#pragma omp metadirective \
when( device={arch("nvptx")} ) : parallel for) \
default( simd )
  for (int i = 0; i < N; i++) {
    v3[i] = v1[i] * v2[i];
  }
}
```

## 7.6. Mapping target constructs to CUDA streams

An OpenMP target task generating construct is executed on the GPU in a CUDA stream. The following are target task generating constructs:

- ▶ **target enter data**
- ▶ **target exit data**
- ▶ **target update**
- ▶ **target**

This section explains how these target constructs are mapped to CUDA streams. The relationship with the OpenACC queues is also explained below.

Keep in mind that the **target data** construct does not generate a task and is not necessarily executed in a CUDA stream. It also cannot have the **depend** and **nowait** clauses, thus its behavior cannot be directly controlled by the user application. The rest of this section does not cover the behavior of the **target data** construct.

Any task-generating target construct can have **depend** and **nowait** clauses. The NVIDIA OpenMP Runtime takes these clauses as a guidance for how to map the construct to a specific CUDA stream. Below is a breakdown of how the clauses affect the mapping decisions.

'target' without 'depend', without 'nowait'

For these constructs, the per-thread default CUDA stream is normally used. The stream is unique for each host thread, so target regions created by different host threads will execute independently in different streams according to the CUDA rules described in [CUDA Runtime API](#); see the rules in the "Per-thread default stream" section.

The OpenACC queue **acc\_async\_sync** is initially associated with the same per-thread default CUDA stream. The user is allowed to change the association by calling **acc\_set\_cuda\_stream(acc\_async\_sync, stream)**. This will change accordingly the stream used for **target** without **nowait**.

The CUDA stream handle can be directly obtained via the **omp\_x\_get\_cuda\_stream(int device, int nowait)** function, with the **nowait** parameter set to 0. The per-thread default stream can be obtained with the CUDA handle **CU\_STREAM\_PER\_THREAD** or **cudaStreamPerThread**.

Here is an example of how a custom CUDA stream can be used to substitute the default stream:

```
extern __global__ void kernel(int *data);

    CUSTream stream;
    cuStreamCreate(&stream, CU_STREAM_DEFAULT);
    acc_set_cuda_stream(acc_async_sync, stream);
#pragma omp target enter data map(to:data[:N])
#pragma omp target data use_device_ptr(data)
    kernel<<<N/32, 32, 0, stream>>>(data);
#pragma omp target teams distribute parallel for
    for (int i = 0; i < N; i++) {
        data[i]++;
    }
#pragma omp target exit data map(from:data[:N])
```

Note there is no explicit stream synchronization after the CUDA **kernel** is launched. The stream is synchronized automatically at the **target** constructs that follow.

'target' with 'depend', without 'nowait'

For this construct, the runtime will block the current thread until all dependencies listed in the **depend** clause are resolved. Then, the **target** construct will be executed in the default per-thread CUDA stream as described in the previous section (that is, as if there is no **depend** clause).

'target' with 'nowait', without 'depend'

By default, the runtime will select a CUDA stream for each new **target nowait** construct. The selected stream may be the same that was used for a prior **target nowait** construct. That is, there is no guarantee of uniqueness of the selected stream.

This is different from the OpenACC model that uses the same CUDA stream associated with the **acc\_async\_noval** queue for any asynchronous construct with the **async** clause without an argument. To change this behavior, the user can call the **omp\_set\_cuda\_stream\_auto(int enable)** function with the **enable** parameter set to 0. In this case, the CUDA stream associated with the **acc\_async\_noval** OpenACC queue will be used for all OpenMP **target nowait** constructs. Another way to enable this behavior is to set the environment variable **NVCOMPILER\_OMP\_AUTO\_STREAMS** to **FALSE**.

To access the stream used for the next **target nowait** construct, the user can call the **omp\_get\_cuda\_stream(int device, int nowait)** function, with the **nowait** parameter set to 1.

'target' with both 'depend' and 'nowait'

The decision on which CUDA stream to use in this case relies on previously scheduled target and host tasks sharing a subset of the dependencies listed in the **depend** clause:

- ▶ If the **target** construct has only one dependency, which is of the type **inout** or **out**, and that dependency maps to a previously scheduled **target depend(...)** **nowait** construct, and the same device is used for both target constructs, then the CUDA stream which the previous target task was scheduled to will be used.
- ▶ Otherwise, a CUDA stream will be selected for this target construct according to the stream selection policy.

Note that target constructs with a single **in** dependency can be scheduled on a newly selected CUDA stream. This is to allow parallel execution of multiple **target nowait** constructs that depend on data produced by another previously scheduled **target nowait** construct.

Here is a simplified example of how a **target** construct, a CUDA library function and a CUDA kernel can be executed on the GPU in the same stream asynchronously with respect to the host thread:

```
extern __global__ void kernel(int *data);

cudaStream_t stream =
(cudaStream_t)omp_get_cuda_stream(omp_get_default_device(), 1);
cufftSetStream(cufft_plan, stream);

#pragma omp target enter data map(to:data[:N]) depend(inout:stream) nowait
#pragma omp target data use_device_ptr(data)
{
    kernel<<<N/32, 32, 0, stream>>>(data);
    cufftExecC2C(cufft_plan, data, data, CUFFT_FORWARD);
}
#pragma omp target teams distribute parallel for depend(inout:stream) nowait
for (int i = 0; i < N; i++) {
    data[i]++;
}
#pragma omp target exit data map(from:data[:N]) depend(inout:stream) nowait
```

Note that the **stream** variable holds the CUDA stream handle and also serves as the dependency for the **target** constructs. This dependency enforces the order of execution and also guarantees the target constructs are on the same stream that was returned from the **omp\_get\_cuda\_stream** function call.

NVIDIA OpenMP API to access and control CUDA streams

NVIDIA OpenMP Runtime provides the following API to access CUDA streams and to control their use.

```
void *omp_get_cuda_stream(int device, int nowait);
```

This function returns the handle of the CUDA stream that will be used for the next **target** construct:

- ▶ If the **nowait** parameter is set to 0, it returns the CUDA stream associated with the OpenACC queue **acc\_async\_sync**, which is initially mapped to the default per-thread CUDA stream;

- Otherwise, it returns a CUDA stream which will be used for the next **target nowait** construct that cannot be mapped to an existing stream according to the rules for the **depend** clause.

```
void ompx_set_cuda_stream_auto(int enable);
```

This function sets the policy for how CUDA streams are selected for **target nowait** constructs:

- If the **enable** parameter is set to a non-zero value, an internally selected CUDA stream will be used for each **target nowait** construct that follows. This is the default behavior;
- Otherwise, the CUDA stream associated with the OpenACC queue **acc\_async\_noval** will be used for all **target nowait** constructs that follow. This becomes the default behavior if the environment variable **NVCOMPILER\_OMP\_AUTO\_STREAMS** is set to **FALSE**.

The setting is done only for the host thread which calls this function.

## 7.7. Noncontiguous Array Sections

Array sections can be used in **to** and **from** clauses, including noncontiguous array sections. The noncontiguous array section must be specified in a single **map** clause; it cannot be split between multiple directives. Although this feature may become a part of a future OpenMP specification, at this time it is an NVIDIA HPC compilers extension.

## 7.8. OpenMP with CUDA Unified Memory

This section will focus on OpenMP unified shared memory programming, and assume users are familiar with Separate, Managed, and Unified Memory Modes explained in the [Memory Model](#) and [Managed and Unified Memory Modes](#) sections. OpenMP unified shared memory corresponds to Unified Memory Mode in NVHPC Compilers and it can be enabled with **-gpu=mem:unified** flag. Source code with **requires unified\_shared\_memory** directive is accepted but requires **-gpu=mem:unified** flag to activate Unified Memory Mode.

In Unified Memory Mode, **map** clauses on **target** constructs are optional. Additionally, **declare target** directives are optional for variables with static storage duration accessed inside functions to which such directive is applied. The OpenMP unified shared memory eases accelerator programming on the GPUs removing the need for data management and only requiring to express the parallelism in the compute regions.

In Unified Memory Mode, all data is managed by the CUDA runtime. Explicit data **map** clauses which manage the data movement across the host and devices become optional. All variables are accessible from the OpenMP offload compute regions executing on the GPU. The **map** clause with **alloc**, **to,from**, and **tofrom** type will not result in any device allocation or data transfer. The OpenMP runtime, however, may leverage such

clauses to communicate preferable data placement to the CUDA runtime by means of memory hint APIs as elaborated in the following blog post on the NVIDIA website: [Simplifying GPU Application Development with Heterogeneous Memory Management](#). Device memory can be allocated or deallocated in OpenMP programs in Unified Memory Mode by using the `omp_target_alloc` and `omp_target_free` API calls. Please, note that the memory allocated through `omp_target_alloc` cannot be accessed by the host.

## Understanding Data Movement

When the compiler encounters a compute construct without visible **target data** directives or **map** clauses, it attempts to determine what data is required for correct execution of the region on the GPU. When the compiler is unable to determine the size and shape of data needing to be accessible on the device, it behaves as follows:

- ▶ In Separate Memory Mode, the compiler may not be able to alert you to the need for an explicit data clause specifying size and/or shape of data being copied to/from the GPU. In this case, the default length of one may be used. This may cause illegal memory access errors at runtime on the GPU devices.
- ▶ In Managed Memory Mode (`-gpu=mem:managed`), the compiler assumes the data is allocated in managed memory and thus is accessible from the device; if this assumption is wrong, for example, if the data was defined globally or is located on the CPU stack, the program may fail at runtime.
- ▶ In Unified Memory Mode (`-gpu=mem:unified`), all data is accessible from the device making information about size and shape unnecessary.

Take the following example in C:

```
#pragma omp declare target
void set(int* ptr, int i, int j, int dim){
    int idx = i * dim + j;
    return ptr[idx] = someval(i, j);
}
#pragma omp end declare target

void fill1d(int* ptr, int dim){
#pragma omp target teams distribute parallel for
    for (int i = 0; i < dim; i++)
        for (int j = 0; j < dim; j++)
            set(ptr, i, j, dim);
}
```

In Separate Memory Mode, the only way to guarantee correctness for this example is to specify an array section in the **target** construct as follows:

```
#pragma omp target teams distribute parallel for map(from: ptr[0:dim*dim])
```

This change explicitly instructs the OpenMP implementation about the precise data segment used within the target for loop.

In Unified Memory Mode, the **map** clause is not required.



The next example, in Fortran, illustrates how a global variable can be accessed in an OpenMP routine without requiring any explicit annotation.

```
module m
integer :: globmin = 1234
contains
subroutine findmin(a)
!$omp declare target
integer, intent(in) :: a(:)
integer :: i
do i = 1, size(a)
if (a(i) .lt. globmin) then
globmin = a(i)
endif
end do
end subroutine
end module m
```

Compile the example above for Unified Memory Mode:

```
nvfortran -mp=gpu -gpu=mem:unified example.f90
```

The source does not need any OpenMP directives to access module variable **globmin**, to either read or update its value, in the routine invoked from CPU and GPU. Moreover, any access to **globmin** will be made to the same exact instance of the variable from CPU and GPU; its value is synchronized automatically. In Separate or Managed Memory Modes, such behavior can only be achieved with a combination of OpenMP **declare target** and **target update** directives in the source code.

Migrating existing OpenMP applications written for Separate Memory Mode should, in most cases, be a seamless process requiring no source changes. Some data access patterns, however, may lead to different results produced during application execution in Unified Memory Mode. Applications which rely on having separate data copies in GPU memory to conduct temporary computations on the GPU -- without maintaining data synchronization with the CPU -- pose a challenge for migration to unified memory. For the following Fortran example, the value of variable **c** after the last loop will differ depending on whether the example is compiled with or without **-gpu=mem:unified**.

```
b(:) = ...
c = 0

!$omp target data map(to: b) map(from: a)
!$omp target distribute teams parallel for
do i = 1, N
b(i) = b(i) * i
end do
!$omp target distribute teams parallel for
do i = 1, N
a(i) = b(i) + i
end do
!$omp end target data

do i = 1, N
c = c + a(i) + b(i)
end do
```

Without Unified Memory, array **b** is copied into the GPU memory at the beginning of the OpenMP **target data** region. It is then updated in the GPU memory and used to compute elements of array **a**. As instructed by the data clause **map(to:b)**, **b** is not copied back to the CPU memory at the end of the **target data** region and therefore its initial value is used in the computation of **c**. With **-mp=gpu -gpu=mem:unified**, the updated value of **b** in the first loop is automatically visible in the last loop leading to a different value of **c** at its end.

Additional complications may arise from the asynchronous execution as the use of unified shared memory may require extra synchronizations to avoid data races.

## 7.9. Multiple Device Support

A program can use multiple devices on a single node.

This functionality is supported using the **omp\_set\_default\_device** API call and the **device()** clause on the **target** constructs. Our experience is that most programs use MPI parallelism with each MPI rank selecting a single GPU to which to offload. Some programs assign multiple MPI ranks to each GPU, in order to keep the GPU fully occupied, though the fixed memory size of the GPU limits how effective this strategy can be. Similarly, other programs use OpenMP thread parallelism on the CPU, with each thread selecting a single GPU to which to offload.

## 7.10. Interoperability with CUDA

The HPC Compilers support interoperability of OpenMP and CUDA to the same extent they support CUDA interoperability with OpenACC.

If OpenMP and CUDA code coexist in the same program, the OpenMP runtime and the CUDA runtime use the same CUDA context on each GPU. To enable this coexistence, use the compilation and linking option **-cuda**. CUDA-allocated data is available for use inside OpenMP target regions with the OpenMP analog **is\_device\_ptr** to OpenACC's **deviceptr()** clause.

OpenMP-allocated data is available for use inside CUDA kernels directly if the data was allocated with the **omp\_target\_alloc()** API call; if the OpenMP data was created with a **target data map** clause, it can be made available for use inside CUDA kernels using the **target data use\_device\_addr()** clause. Calling a CUDA device function inside an OpenMP target region is supported, as long as the CUDA function is a scalar function, that is, does not use CUDA shared memory or any inter-thread synchronization. Calling an OpenMP **declare target** function inside a CUDA kernel is supported as long as the **declare target** function has no OpenMP constructs or API calls.

## 7.11. Interoperability with Other OpenMP Compilers

OpenMP CPU-parallel object files compiled with NVIDIA's HPC compilers are interoperable with OpenMP CPU-parallel object files compiled by other compilers using the KMPC OpenMP runtime interface. Compilers supporting KMPC OpenMP include Intel and CLANG. The HPC compilers support a GNU OpenMP interface layer as well which provides OpenMP CPU-parallel interoperability with the GNU compilers.

For OpenMP GPU computation, there is no similar formal or informal standard library interface for launching GPU compute constructs or managing GPU memory. There is also no standard way to manage the device context in such a way as to interoperate between multiple offload libraries. The HPC compilers therefore do not support interoperability of device compute offload operations and similar operations generated with another compiler.

## 7.12. GNU STL

When using `nvc++` on Linux, the GNU STL is thread-safe to the extent listed in the GNU documentation as required by the C++11 standard. If an STL thread-safe issue is suspected, the suspect code can be run sequentially inside of an OpenMP region using `#pragma omp critical` sections.

# Chapter 8.

## USING STDPAR

This chapter describes the NVIDIA HPC Compiler support for standard language parallelism, also known as Stdpar:

- ▶ ISO C++ standard library parallel algorithms with **nvc++**
- ▶ ISO Fortran **do concurrent** loop construct with **nvfortran**

Use the **-stdpar** compiler option to enable parallel execution with standard parallelism. The sub-options to **-stdpar** are the following:

- ▶ **gpu**: compile for parallel execution on GPU; this sub-option is the default. This feature is supported on the NVIDIA Pascal architecture and newer.
- ▶ **multicore**: compile for multicore CPU execution.

By default, NVC++ auto-detects and generates GPU code for the type of GPU that is installed on the system on which the compiler is running. To generate code for a specific GPU architecture, which may be necessary when the application is compiled and run on different systems, add the **-gpu=ccXX** command-line option. More details can be found in [Compute Capability](#).

### Predefined Macros

The following macros corresponding to the parallel execution target compiled for are added implicitly:

- ▶ **\_\_NVCOMPILER\_STDPAR\_GPU** for parallel execution on GPU.
- ▶ **\_\_NVCOMPILER\_STDPAR\_MULTICORE** for parallel execution on multicore CPU.

## 8.1. GPU Memory Modes

When compiling for GPU execution, Stdpar utilizes [Managed and Unified Memory Modes](#) for managing data accessed from the sequential code running on CPU and from the parallel code on GPU.

The compiler detects the memory capability of the system on which the compiler is running and uses that information to enable the correct memory mode as follows:

- ▶ When compiled on the platform with full CUDA Unified Memory capability, `-stdpar` implies `-gpu=mem:unified`.
- ▶ When compiled on the platform with CUDA Managed Memory capability only, `-stdpar` implies `-gpu=mem:managed`.

To compile code for a specific Memory Mode regardless of the memory capability of the system on which you are compiling, add the desired `-gpu=mem:unified` or `-gpu=mem:managed` option.

Stdpar with Separate Memory Mode can only be supported when the data are fully managed through features of other programming models e.g. OpenACC.

All restrictions on variables used on the GPU in standard language parallel code in Managed Memory Mode have been removed when using Unified Memory Mode.

If the compiler utilises CUDA Managed Memory automatically, the interception of deallocations is enabled implicitly at runtime. This is to prevent deallocating the data with unmatching API which may lead to undefined behavior. The interception incurs some runtime overhead and may be unnecessary if allocations and deallocations for all data in the application are performed using the matching APIs. The interception can be disabled using dedicated command-line options detailed in [Interception of Deallocations](#). More details about the memory modes supported by the NVIDIA HPC Compilers and dedicated command-line options can be found in [Memory Model](#).

## 8.2. Stdpar C++

The NVIDIA HPC C++ compiler, NVC++, supports C++ Standard Language Parallelism (Stdpar) for execution on NVIDIA GPUs and multicore CPUs. As mentioned previously, use the NVC++ command-line option `-stdpar` to enable GPU accelerated C++ Parallel Algorithms. The following sections go into more detail about the NVC++ support for the ISO C++ Standard Library Parallel Algorithms.

### 8.2.1. Introduction to Stdpar C++

The C++17 Standard introduced higher-level parallelism features that allow users to request parallelization of Standard Library algorithms.

This higher-level parallelism is expressed by adding an execution policy as the first parameter to any algorithm that supports execution policies. Most of the existing Standard C++ algorithms were enhanced to support execution policies. C++17 defined several new parallel algorithms, including the useful `std::reduce` and `std::transform_reduce`.

C++17 defines three [execution policies](#):

- ▶ `std::execution::seq`: Sequential execution. No parallelism is allowed.
- ▶ `std::execution::par`: Parallel execution on one or more threads.
- ▶ `std::execution::par_unseq`: Parallel execution on one or more threads, with each thread possibly vectorized.

When you use an execution policy other than `std::execution::seq`, you are communicating two important things to the compiler:

- ▶ You prefer but do not require that the algorithm be run in parallel. A conforming C++17 implementation may ignore the hint and run the algorithm sequentially, but a performance-oriented implementation takes the hint and executes in parallel when possible and prudent.
- ▶ The algorithm is safe to run in parallel. For the `std::execution::par` and `std::execution::par_unseq` policies, any user-provided code—such as iterators, lambdas, or function objects passed into the algorithm—must not introduce data races if run concurrently on separate threads. For the `std::execution::par_unseq` policy, any user-provided code must not introduce data races or deadlocks if multiple calls are interleaved on the same thread, which is what happens when a loop is vectorized. For more information about potential deadlocks, see the [forward progress guarantees](#) provided by the parallel policies or watch [CppCon 2018: Bryce Adelstein Lelbach “The C++ Execution Model”](#).

The C++ Standard grants compilers great freedom to choose if, when, and how to execute algorithms in parallel as long as the forward progress guarantees the user requests are honored. For example, `std::execution::par_unseq` may be implemented with vectorization and `std::execution::par` may be implemented with a CPU thread pool. It is also possible to execute parallel algorithms on a GPU, which is a good choice for invocations with sufficient parallelism to take advantage of the processing power and memory bandwidth of NVIDIA GPU processors.

## 8.2.2. NVC++ Compiler Parallel Algorithms Support

NVC++ supports C++ Standard Language Parallelism with the parallel execution policies `std::execution::par` or `std::execution::par_unseq` for execution on GPUs or multicore CPUs.

Lambdas, including generic lambdas, are fully supported in parallel algorithm invocations. No language extensions or non-standard libraries are required to enable GPU acceleration. All data movement between host memory and GPU device memory is performed implicitly and automatically under the control of [Managed and Unified Memory Modes](#).

It's straightforward to automatically GPU accelerate C++ Parallel Algorithms with NVC++. However, there are some restrictions and limitations you need to be aware of as explained below.

### 8.2.2.1. Enabling Parallel Algorithms with the -stdpar Option

GPU acceleration of C++ Parallel Algorithms is enabled with the `-stdpar=gpu` command-line option to NVC++. If `-stdpar=gpu` is specified (or `-stdpar` without an argument), almost all algorithms that use a parallel execution policy are compiled for offloading to run in parallel on an NVIDIA GPU:

```
nvc++ -stdpar=gpu program.cpp -o program
```

```
nvc++ -stdpar program.cpp -o program
```

Acceleration of C++ Parallel Algorithms with multicore CPUs is enabled with the `-stdpar=multicore` command-line option to NVC++. If `-stdpar=multicore` specified, all algorithms that use a parallel execution policy are compiled to run on a multicore CPU:

```
nvc++ -stdpar=multicore program.cpp -o program
```

### 8.2.3. Stdpar C++ Simple Example

Here are a few simple examples to get a feel for how the C++ Parallel Algorithms work.

From the early days of C++, sorting items stored in an appropriate container has been relatively easy using a single call like the following:

```
std::sort(employees.begin(), employees.end(),
          CompareByLastName());
```

Assuming the comparison class **CompareByLastName** is thread-safe, which is true for most comparison functions, parallelizing this sort is simple with C++ Parallel Algorithms. Include `<execution>` and add an execution policy to the function call:

```
std::sort(std::execution::par,
          employees.begin(), employees.end(),
          CompareByLastName());
```

Calculating the sum of all the elements in a container is also simple with the **std::accumulate** algorithm. Prior to C++17, transforming the data in some way while taking the sum was somewhat awkward. For example, to compute the average age of your employees, you might write the following code:

```
int ave_age =
    std::accumulate(employees.begin(), employees.end(), 0,
                    [](int sum, const Employee& emp){
                        return sum + emp.age();
                    })
    / employees.size();
```

The **std::transform\_reduce** algorithm introduced in C++17 makes it simple to parallelize this code. It also results in cleaner code by separating the reduction operation, in this case **std::plus**, from the transformation operation, in this case **emp.age()**:

```
int ave_age =
    std::transform_reduce(std::execution::par_unseq,
                          employees.begin(), employees.end(),
                          0, std::plus<int>(),
                          [](const Employee& emp){
                              return emp.age();
                          })
    / employees.size();
```

### 8.2.4. Coding Guidelines for GPU-accelerating Parallel Algorithms

GPUs are not simply CPUs with more threads. To effectively take advantage of the massive parallelism and memory bandwidth available on GPUs, it is typical for GPU programming models to put some limitations on code executed on the GPU. The NVC++ implementation of C++ Parallel Algorithms is no exception in this regard. The sections which follow detail the limitations that apply in the current release.

### 8.2.4.1. Parallel Algorithms and Device Function Annotations

Functions to be executed on the GPU within parallel algorithms do not need any `__device__` annotations or other special markings to be compiled for GPU execution. The NVCC++ compiler walks the call graph for each source file and automatically infers which functions must be compiled for GPU execution.

However, this only works when the compiler can see the function definition in the same source file where the function is called. This is true for most inline functions and template functions but may fail when functions are defined in a different source file or linked in from an external library. You need to be aware of this when formulating parallel algorithms invocations that you expect to be offloaded and accelerated on NVIDIA GPUs.

When calling an externally defined function from within a parallel algorithm region, such functions require some form of device annotations from other GPU programming models e.g. OpenACC routine directive (refer to [External Device Function Annotations](#) for more information).

### 8.2.4.2. Data Management in Parallel Algorithms

When offloading parallel algorithms to a GPU, it's essential to consider how data is accessed from the parallel region. Some GPUs may not access certain segments of the CPU's address space. Developers targeting platforms without unified shared memory or those seeking to optimize performance must be aware of these memory distinctions, as they may affect the following types of data accessed in parallel algorithm regions:

- ▶ Pointer data passed into lambda functions within the parallel algorithm.
- ▶ Data captured by reference in lambda functions or pointer data captured by value.
- ▶ Variables with static storage duration referenced inside the parallel algorithm.

To avoid memory access violations, developers must ensure that all of the above data is accessible to the GPU before the parallel algorithm is executed.

Stdpar C++ only supports [Managed and Unified Memory Modes](#) which allow data being accessed from CPU and GPU. Through support in both the CUDA device driver and the NVIDIA GPU hardware, the CUDA Unified Memory manager automatically moves some types of data based on usage.

Stdpar with Separate Memory Mode can only be supported when the data are fully managed through the OpenACC data directives, refer to [Interoperability with OpenACC](#).

Since object-oriented design is fundamental to C++, special consideration must be given to composite data types with pointer or reference members. The data referenced or pointed to may not be stored contiguously within the composite data type. Moreover, such data might not even be allocated in the same memory segment as the composite type itself. As a result, when accessing both the composite data type and its referenced or pointed-to data from parallel algorithms, the developer must ensure that the member data is also made accessible to the GPU. These considerations should also be taken into



account when standard library containers are used in the parallel algorithms as the containers frequently contain member pointers to their elements.

The discussion in this section assumes familiarity with the Managed and Unified Memory Modes covered in [Memory Model](#) and [Managed and Unified Memory Modes](#). The code executing within the parallel algorithm is referred to as the accelerator subprogram. In contrast to the code executing outside of the parallel algorithm which is referred to as the host subprogram.

### Managed Memory Mode

When Stdpar code is compiled with Managed Memory Mode (as default mode or by passing `-gpu=mem:managed`) only data dynamically allocated on the heap in CPU code can be managed automatically. CPU and GPU automatic storage (stack memory) and static storage (global or static data) cannot be automatically managed. Likewise, data that is dynamically allocated in program units not compiled by `nvc++` with the `-stdpar` option is not automatically managed by CUDA Unified Memory even though it is on the CPU heap. The compiler utilizes CUDA Managed Memory for dynamic allocations to make data accessible from CPU and GPU. As managed memory allocation calls can incur higher runtime overhead than standard allocator calls, the implementation uses memory pools for performance reasons by default as detailed in [Memory Pool Allocator](#).

The Managed Memory Mode is intended for binaries run on targets with CUDA Managed Memory capability only. Any pointer that is dereferenced and any C++ object that is referenced within a parallel algorithm invocation must refer to data on the CPU heap that is allocated in a program unit compiled by `nvc++` with `-stdpar`. Dereferencing a pointer to a CPU stack or a global object will result in a memory violation in GPU code.

### Unified Memory Mode

When Unified Memory is the default memory mode or is selected explicitly on the command line by passing `-gpu=mem:unified`, there are no restrictions on variables accessed in the parallel algorithms. Therefore, all CPU data (either residing on stack, heap, or globally) are simply accessible in the parallel algorithm functions. Note that memory dynamically allocated in GPU code is only visible from GPU code and can never be accessed by the CPU regardless of the CUDA Unified Memory capability.

When compiling a binary for platforms with full CUDA Unified Memory capability, only those source files using features from the standard parallel algorithms library must be compiled by `nvc++` with the `-stdpar` option. There is no requirement that the code dynamically allocating memory accessed on GPU is also compiled in such a way.

Unified Memory Mode may utilize CUDA Managed Memory for dynamic allocation, more details can be found in [Managed and Unified Memory Modes](#).

## Summary

The following table provides a key summary of important command-line options selecting memory modes and the impact of memory modes on the Stdpar features.

Table 21 Stdpar C++ Feature Differences for Memory Modes

Command-line options	Dynamically allocated variables outside of parallel algorithm region	Automatic or static storage variables outside of parallel algorithm region	Dynamic allocator
No memory-specific flags passed, compiling on target with CUDA Managed Memory only	Can be accessed within parallel region code	Cannot be accessed within parallel algorithm code	cudaMallocManaged
No memory-specific flags passed, compiling on target with full CUDA Unified Memory	Can be accessed within parallel region code	Can be accessed within parallel algorithm code	cudaMallocManaged or system allocators: new/malloc (compiler picks the most suitable allocator)
<b>-gpu=mem:managed</b>	Can be accessed within parallel region code	Cannot be accessed within parallel algorithm code	cudaMallocManaged
<b>-gpu=mem:unified</b>	Can be accessed within parallel region code	Can be accessed within parallel algorithm code	cudaMallocManaged or system allocators: new/malloc (compiler picks the most suitable allocator)
<b>-gpu=mem:unified:managedalloc</b>	Can be accessed within parallel region code	Can be accessed within parallel algorithm code	cudaMallocManaged
<b>-gpu=mem:unified:nomanagedalloc</b>	Can be accessed within parallel region code	Can be accessed within parallel algorithm code	System allocators: new/malloc

## Examples

For example, `std::vector` uses dynamically allocated memory, which is accessible from the GPU when using Stdpar. Iterating over the contents of a `std::vector` in a parallel algorithm works as expected when compiling with either `-gpu=mem:managed` or `-gpu=mem:unified`:

```
std::vector<int> v = ...;
std::sort(std::execution::par,
          v.begin(), v.end()); // Okay, accesses heap memory.
```

On the other hand, `std::array` performs no dynamic allocations. Its contents are stored within the `std::array` object itself, which is often on a CPU stack. Iterating over the contents of a `std::array` will not work on systems with only CUDA Managed Memory support unless the `std::array` itself is allocated on the heap and the code is compiled with `-gpu=mem:managed`:

```
std::array<int, 1024> a = ...;
std::sort(std::execution::par,
          a.begin(), a.end()); // Fails on targets with CUDA Managed
                               // Memory capability only, array is on
                               // a CPU stack inaccessible from GPU.
                               // Works correctly on targets with full
                               // CUDA Unified Memory support.
```

The above example works as expected when run on a target supporting full CUDA Unified Memory capability.

When executing on targets with CUDA Managed Memory capability only, pay particular attention to lambda captures, especially capturing data objects by reference, which may contain non-obvious pointer dereferences:

```
void saxpy(float* x, float* y, int N, float a) {
    std::transform(std::execution::par_unseq, x, x + N, y, y,
                  [&](float xi, float yi){ return a * xi + yi; });
}
```

In the earlier example, the containing function parameter `a` is captured by reference. The code within the body of the lambda, which is running on the GPU, tries to access `a`, which is in the CPU stack memory. This attempt results in a memory violation and undefined behavior. In this case, the problem can easily be fixed by changing the lambda to capture by value:

```
void saxpy(float* x, float* y, int N, float a) {
    std::transform(std::execution::par_unseq, x, x + N, y, y,
                  [=](float xi, float yi){ return a * xi + yi; });
}
```

With this one-character change, the lambda makes a copy of `a`, which is then copied to the GPU, and there are no attempts to reference CPU stack memory from GPU code.

Such code will run correctly without requiring modifications on targets with full CUDA Unified Memory capability.

If `std::vector` is accessed through a subscript operator from the device this would require such a vector object to be accessible from the parallel code executing on the GPU. This means that the `std::vector` needs to be allocated dynamically in order to make it accessible from the GPU when compiled for the systems with only CUDA Managed Memory support.

```
std::vector<int> v = ...;
std::for_each(std::execution::par,
              idx.begin(), idx.end(), [&](auto i)
              {v[i] = 1;}); // Fails on targets with CUDA Managed
                          // Memory capability only, vector object is on
                          // a CPU stack inaccessible from GPU.
                          // Works correctly on targets with full
                          // CUDA Unified Memory support.
```

An alternative approach to managing the content of the `std::vector` on systems with CUDA Managed Memory support only would be to obtain a pointer to its elements data region using `data()` member.

```
std::vector<int> v = ...;
int* vdataptr = v.data();
std::for_each(std::execution::par,
              idx.begin(), idx.end(), [&](auto i)
              {vdataptr[i] = 1;}); // Works, vector elements are in heap
                                  // memory
```

Whether `-gpu=mem:unified` is enabled by default or passed explicitly on the command line, parallel algorithms can access global variables and accesses to global variables from CPU and GPU are kept in sync. Extra care should be taken when accessing global variables within parallel algorithms, as simultaneous updates in different iterations running on the GPU can lead to data races. The following example illustrates the safe update of a global variable in the parallel algorithm since the update only occurs in one iteration.

```
int globvar = 123;
void foo() {
    auto r = std::views::iota(0, N);
    std::for_each(std::execution::par_unseq, r.begin(), r.end(),
                  [](auto i) {
                      if (i == N - 1)
                          globvar += 345;
                  });
    // globvar is equal to 468.
}
```

### 8.2.4.3. Parallel Algorithms and Function Pointers

Functions compiled to run on either the CPU or the GPU must be compiled into two different versions, one with the CPU machine instructions and one with the GPU machine instructions.

In the current implementation, a function pointer either points to the CPU or the GPU version of the functions. This causes problems if you attempt to pass function pointers between CPU and GPU code. You might inadvertently pass a pointer to the CPU version of the function to GPU code. In the future, it may be possible to automatically and seamlessly support the use of function pointers across CPU and GPU code boundaries, but it is not supported in the current implementation.

Function pointers can't be passed to Parallel Algorithms to be run on the GPU, and functions may not be called through a function pointer within GPU code. For example, the following code example won't work correctly:

```
void square(int& x) { x = x * x; }
void square_all(std::vector<int>& v) {
    std::for_each(std::execution::par_unseq,
                  v.begin(), v.end(), &square);
}
```

It passes a pointer to the CPU version of the function `square` to a parallel `for_each` algorithm invocation. When the algorithm is parallelized and offloaded to the GPU, the program fails to resolve the function pointer to the GPU version of `square`.

You can often solve this issue by using a function object, which is an object with a function call operator. The function object's call operator is resolved at compile time to the GPU version of the function, instead of being resolved at run time to the incorrect CPU version of the function as in the previous example. For example, the following code example works:

```
struct squared {
    void operator()(int& x) const { x = x * x; }
};
void square_all(std::vector<int>& v) {
    std::for_each(std::execution::par_unseq,
                  v.begin(), v.end(), squared{});
}
```

Another possible workaround is to change the function to a lambda, because a lambda is implemented as a nameless function object:

```
void square_all(std::vector<int>& v) {
    std::for_each(std::execution::par_unseq, v.begin(), v.end(),
                  [](int& x) { x = x * x; });
}
```

If the function in question is too big to be converted to a function object or a lambda, then it should be possible to wrap the call to the function in a lambda:

```
void compute(int& x) {
    // Assume lots and lots of code here.
}
void compute_all(std::vector<int>& v) {
    std::for_each(std::execution::par_unseq, v.begin(), v.end(),
                  [](int& x) { compute(x); });
}
```

No function pointers are used in this example.

The restriction on calling a function through a function pointer unfortunately means passing polymorphic objects from CPU code to GPU-accelerated Parallel Algorithms is not currently supported, as virtual tables are implemented using function pointers.

#### 8.2.4.4. Random Access Iterators

The C++ Standard requires that the iterators passed to most C++ Parallel Algorithms be forward iterators. However, C++ Parallel Algorithms on GPUs only works with random access iterators. Passing a forward iterator or a bidirectional iterator to a GPU/CPU-accelerated Parallel Algorithm results in a compilation error. Passing raw pointers or Standard Library random access iterators to the algorithms has the best performance, but most other random-access iterators work correctly.

#### 8.2.4.5. Interoperability with the C++ Standard Library

Large parts of the C++ Standard Library can be used with stdpar on GPUs.

- ▶ `std::atomic<T>` objects within GPU code work provided that `T` is a four-byte or eight-byte integer type.
- ▶ Math functions that operate on floating-point types—such as `sin`, `cos`, `log`, and most of the other functions declared in `<cmath>`—can be used in GPU code and resolve to the same implementations that are used in CUDA C++ programs.
- ▶ `std::complex`, `std::tuple`, `std::pair`, `std::optional`, `std::variant`, and `<type_traits>`, are supported and work as expected in GPU code.

The parts of the C++ Standard Library that aren't supported in GPU code include I/O functions and in general any function that accesses the CPU operating system. As a special case, basic `printf` calls can be used within GPU code and leverage the same implementation that is used in NVIDIA CUDA C++.

#### 8.2.4.6. No Exceptions in GPU Code

As with most other GPU programming models, throwing and catching C++ exceptions is not supported within Parallel Algorithm invocations that are offloaded to the GPU.

Unlike some other GPU programming models where try/catch blocks and throw expressions are compilation errors, exception code does compile but with non-standard behavior. Catch clauses are ignored, and throw expressions abort the GPU kernel if actually executed. Exceptions in CPU code work without restrictions.

### 8.2.5. NVC++ Experimental Features

nvc++ experimental features are enabled with the `--experimental-stdpar` compiler flag. Experimental feature headers are exposed via the `<experimental/...>` namespaces and limited support for these features is available in older C++ versions. Table 1 lists all experimental features available and the minimum language version required to use them.

Table 22 Experimental features information

Feature	Recommen	Limited support	Standard proposal	Other notes
Multi-dimensional spans (mdspan)	C++23	C++17	P0009	<a href="https://github.com/NVIDIA/libcudacxx">https://github.com/NVIDIA/libcudacxx</a>

Feature	Recommen	Limited support	Standard proposal	Other notes
Slices of multi-dimensional spans (submdspan)	C++23	C++17	P2630	<a href="https://github.com/NVIDIA/libcudacxx">https://github.com/NVIDIA/libcudacxx</a>
Multi-dimensional arrays (mdarray)	C++23	C++17	P1684	<a href="https://github.com/kokkos/mdspan">https://github.com/kokkos/mdspan</a>
Senders and receivers	C++23	C++20	P2300	<a href="https://github.com/NVIDIA/stdexec">https://github.com/NVIDIA/stdexec</a>
Linear algebra	C++23	C++17	P1673	<a href="https://github.com/kokkos/stdblas">https://github.com/kokkos/stdblas</a>

### 8.2.5.1. Multi-dimensional Spans

Multi-dimensional spans (**std::mdspan**) enable customizable multi-dimensional access to data. This feature was added to C++23 (see [P0009](#) and follow-on papers). [A Gentle Introduction to mdspan](#) gives a tutorial. The reference mdspan implementation <https://github.com/kokkos/mdspan> also has many useful examples.

nvc++ provides an implementation available in the `<experimental/mdspan>` namespace that works with C++17 or newer. It enables applications that are not targeting the C++23 version of the standard to use mdspan.

nvc++ also provides the [P0009R17](#) version of submdspan, which only works for the mdspan layouts in C++23; that is, it does not implement C++26 submdspan ([P2630](#)) yet.

C++23's mdspan uses **operator[]** for array access. For example, if **A** is a rank-2 mdspan, and **i** and **j** are integers, then **A[i, j]** accesses the element of **A** at row **i** and column **j**. Before C++23, **operator[]** was only allowed to take one argument. C++23 changed the language to permit any number of arguments (zero or more). nvc++ does not support this new language feature. As a result, the implementation of mdspan provided by nvc++ permits use of **operator()** as a fall-back (e.g., **A(i, j)** instead of **A[i, j]**). Users may enable this fall-back manually, by defining the macro **MDSPAN\_USE\_PAREN\_OPERATOR** to **1** before including any mdspan headers.

The following example ([godbolt](#)):

```
#include <experimental/mdspan>
#include <iostream>

namespace stdex = std::experimental;

int main() {
    std::array d{
        0, 5, 1,
        3, 8, 4,
        2, 7, 6,
    };

    stdex::mdspan m{d.data(), stdex::extents{3, 3}};
    static_assert(m.rank()==2, "Rank is two");
}
```

```

for (std::size_t i = 0; i < m.extent(0); ++i)
    for (std::size_t j = 0; j < m.extent(1); ++j)
        std::cout << "m(" << i << ", " << j << ") == " << m(i, j) << "\n";

return 0;
}

```

is compiled as follows

```
nvc++ -std=c++17 -o example example.cpp
```

and outputs

```

m(0, 0) == 0
m(0, 1) == 5
m(0, 2) == 1
m(1, 0) == 3
m(1, 1) == 8
m(1, 2) == 4
m(2, 0) == 2
m(2, 1) == 7
m(2, 2) == 6

```

### 8.2.5.2. Senders and Receivers

**P2300 - `std::execution`** proposes a model of asynchronous programming for adoption into the C++26 Standard. For an introduction to this feature, see [Design - user side](#) section of the proposal. The NVIDIA implementation of Senders and receivers is [open source](#) and its repository contains many [useful examples](#). `nvc++` provides access to the NVIDIA implementation which works in C++20 or newer. Since the proposal is still evolving, our implementation is not stable. It is experimental in nature and will change to follow the proposal closely without any warning. The NVIDIA implementation is structured as follows:

Includes	Namespace	Description
<code>&lt;stdexec/...&gt;</code>	<code>::stdexec</code>	Approved for C++ standard
<code>&lt;sexec/...&gt;</code>	<code>::exec</code>	Generic additions and extensions
<code>&lt;nvexec/...&gt;</code>	<code>::nvexec</code>	NVIDIA-specific extensions and customizations

The following example ([godbolt](#)) builds a task graph in which two different vectors, `v0` and `v1`, are concurrently modified in bulk, using a CPU thread pool and a GPU stream context, respectively. This graph then transfers execution to the CPU thread pool, and adds both vectors into `v2` on the CPU, returning the sum of all elements:

```

int main()
{
    // Declare a pool of 8 worker CPU threads:
    exec::static_thread_pool pool(8);

    // Declare a GPU stream context:
    nvexec::stream_context stream_ctx{};

    // Get a handle to the thread pool:
    auto cpu_sched = pool.get_scheduler();
}

```



```

auto gpu_sched = stream_ctx.get_scheduler();

// Declare three dynamic array with N elements
std::size_t N = 5;
std::vector<int> v0 {1, 1, 1, 1, 1};
std::vector<int> v1 {2, 2, 2, 2, 2};
std::vector<int> v2 {0, 0, 0, 0, 0};

// Describe some work:
auto work = stdexec::when_all(
    // Double v0 on the CPU
    stdexec::just()
    | exec::on(cpu_sched,
        stdexec::bulk(N, [v0 = v0.data()] (std::size_t i) {
            v0[i] *= 2;
        })),
    // Triple v1 on the GPU
    stdexec::just()
    | exec::on(gpu_sched,
        stdexec::bulk(N, [v1 = v1.data()] (std::size_t i) {
            v1[i] *= 3;
        })))
)
| stdexec::transfer(cpu_sched)
// Add the two vectors into the output vector v2 = v0 + v1:
| stdexec::bulk(N, [&] (std::size_t i) { v2[i] = v0[i] + v1[i]; })
| stdexec::then([&] {
    int r = 0;
    for (std::size_t i = 0; i < N; ++i) r += v2[i];
    return r;
});
auto [sum] = stdexec::sync_wait(work).value();
// Print the results:
std::printf("sum = %d\n", sum);
for (int i = 0; i < N; ++i) {
    std::printf("v0[%d] = %d, v1[%d] = %d, v2[%d] = %d\n",
        i, v0[i], i, v1[i], i, v2[i]);
}
return 0;
}

```

is compiled as follows:

```
nvc++ --stdpar=gpu --experimental-stdpar -std=c++20 -o example example.cpp
```

and outputs:

```

sum = 40
v0[0] = 2, v1[0] = 6, v2[0] = 8
v0[1] = 2, v1[1] = 6, v2[1] = 8
v0[2] = 2, v1[2] = 6, v2[2] = 8
v0[3] = 2, v1[3] = 6, v2[3] = 8
v0[4] = 2, v1[4] = 6, v2[4] = 8

```

### 8.2.5.3. Linear Algebra

**P1673 - A free function linear algebra interface based on the BLAS** proposes standardizing an idiomatic C++ interface based on `std::mdspan` for a subset of the Basic Linear Algebra Subroutines (BLAS) standard. For an introduction to this feature, see [P1673 \(C++ linear algebra library\) background & motivation](#). There are many useful examples available in `$HPCSDK_HOME/examples/stdpar/stdblas` and in the

repository of the [reference implementation](#). A detailed documentation is available at `$HPCSDK_HOME/compilers/include/experimental/__p1673_bits/README.md`. `nvc++` provides access to the NVIDIA implementation which works in C++17 or newer. Since the proposal is still evolving, our implementation is not stable. It is experimental in nature and will change to follow the proposal closely without any warning. To use the linear algebra library facilities, a suitable linear algebra library must be linked: cuBLAS for GPU execution via the `-cudalib=cublas` flag, and a CPU BLAS library for CPU execution. The HPC SDK bundles OpenBLAS which may be linked using the `-lblas` linker flag.

Execution	BLAS library	Architectures	Compiler flags
Multicore	OpenBLAS	x86_64, aarch64, ppc64l	<code>-stdpar=multicore -lblas</code>
GPU	cuBLAS	All	<code>-stdpar=gpu -cudalib=cublas</code>

The following example ([godbolt](#)):

```
#include <experimental/mdspan>
#include <experimental/linalg>
#include <vector>
#include <array>

namespace stdex = std::experimental;

int main()
{
    constexpr size_t N = 4;
    constexpr size_t M = 2;

    std::vector<double> A_vec(N*M);
    std::vector<double> x_vec(M);
    std::array<double, N> y_vec(N);

    stdex::mdspan A(A_vec.data(), N, M);
    stdex::mdspan x(x_vec.data(), M);
    stdex::mdspan y(y_vec.data(), N);

    for(int i = 0; i < A.extent(0); ++i)
        for(int j = 0; j < A.extent(1); ++j)
            A(i,j) = 100.0 * i + j;

    for(int j = 0; j < x.extent(0); ++j) x(j) = 1.0 * j;
    for(int i = 0; i < y.extent(0); ++i) y(i) = -1.0 * i;

    stdex::linalg::matrix_vector_product(A, x, y); // y = A * x

    // y = 0.5 * y + 2 * A * x
    stdex::linalg::matrix_vector_product(std::execution::par,
    stdex::linalg::scaled(2.0, A), x,
    stdex::linalg::scaled(0.5, y), y);

    // Print the results:
```

```
for (int i = 0; i < N; ++i) std::printf("y[%d] = %f\n", i, y(i));
return 0;
}
```

is compiled as follows for GPU execution:

```
nvc++ -std=c++17 -stdpar=gpu -cudalib=cublas -o example example.cpp
```

And as follows for CPU execution:

```
nvc++ -std=c++17 -stdpar=multicore -o example example.cpp -lblas
```

and produces the same outputs in both cases:

```
y[0] = 2.500000
y[1] = 252.500000
y[2] = 502.500000
y[3] = 752.500000
```

## 8.2.6. Stdpar C++ Larger Example: LULESH

The [LULESH hydrodynamics mini-app](#) was developed at Lawrence Livermore National Laboratory to stress test compilers and model performance of hydrodynamics applications. It is about 9,000 lines of C++ code, of which 2,800 lines are the core computation that should be parallelized.

We ported LULESH to C++ Parallel Algorithms and made the port available on [LULESH's GitHub repository](#). To compile it, install the [NVIDIA HPC SDK](#), check out the 2.0.2-dev branch of the LULESH repository, go to the correct directory, and **run make**.

```
git clone --branch 2.0.2-dev https://github.com/LLNL/LULESH.git
cd LULESH/stdpar/build
make run
```

While LULESH is too large to show the entire source code here, there are some key code sequences that demonstrate the use of stdpar.

The LULESH code has many loops with large bodies and no loop-carried dependencies, making them good candidates for parallelization. Most of these were easily converted into calls to `std::for_each_n` with the `std::execution::par` policy, where the body of the lambda passed to `std::for_each_n` is identical to the original loop body.

The function `CalcMonotonicQRegionForElems` is an example of this. The loop header written for OpenMP looks as follows:

```
#pragma omp parallel for firstprivate(qlc_monoq, qqc_monoq, \
                                   monoq_limiter_mult, monoq_max_slope, ptiny)
for ( Index_t i = 0 ; i < domain.regElemSize(r); ++i ) {
```

This loop header in the C++ Parallel Algorithms version becomes [the following](#):

```
std::for_each_n(
    std::execution::par, counting_iterator(0), domain.regElemSize(r),
    [=, &domain](Index_t i) {
```

The loop body, which in this case is almost 200 lines long, becomes the body of the lambda but is otherwise unchanged from the OpenMP version.

In a number of places, an explicit **for** loop was changed to use C++ Parallel Algorithms that better express the intent of the code, such as the function **CalcPressureForElems**:

```
#pragma omp parallel for firstprivate(length)
for (Index_t i = 0; i < length; ++i) {
    Real_t cls = Real_t(2.0)/Real_t(3.0);
    bvc[i] = cls * (compression[i] + Real_t(1.));
    pbvc[i] = cls;
}
```

This function was rewritten as **as follows**:

```
constexpr Real_t cls = Real_t(2.0) / Real_t(3.0);
std::transform(std::execution::par,
    compression, compression + length, bvc,
    [=](Real_t compression_i) {
        return cls * (compression_i + Real_t(1.0));
    });
std::fill(std::execution::par, pbvc, pbvc + length, cls);
```

## 8.2.7. Interoperability with OpenACC

A subset of OpenACC features can be used when compiling Stdpar code for GPUs. Such a subset is documented in this section. To activate OpenACC directives recognition with Stdpar code add **-acc** command line flag to **nvc++**.

```
nvc++ -stdpar -acc example.cpp
```

OpenACC functionality is detailed in the OpenACC specification and the NVHPC compiler specific differences are detailed in [Using OpenACC](#) of this guide.

Combining OpenACC features with Stdpar offers greater flexibility in how code is written. For instance, it allows external functions to be called from within parallel algorithms. Additionally, it provides opportunities for performance tuning, such as through explicit data management.

### 8.2.7.1. Data Management Directives

C++ parallel algorithms can be offloaded to the GPU when the data accessed in such algorithms is managed through the OpenACC directives. With data fully managed through the OpenACC directives, Stdpar code can run with all GPU Memory Modes including Separate Memory Mode (compiled with **-gpu=mem:separate**).

The following data directives are supported:

- ▶ OpenACC structured data construct directive
- ▶ OpenACC unstructured enter/exit data directives
- ▶ OpenACC host\_data directive
- ▶ OpenACC update directive

Only the data that are captured by reference or pointer-like data captured by values as well as pointer-like data passed as arguments in the parallel algorithm lambdas can be managed through OpenACC. Any non-pointer variables that are captured by value in the parallel algorithm lambda or non-pointer data passed in as lambda arguments are managed by the C++ implementation. A copy of such data is automatically created in the

memory accessible from the GPU. For additional details refer to [Data Management in Parallel Algorithms](#).

OpenACC data management can serve two main purposes:

- ▶ **Explicit Data Management:** This is necessary for data that cannot be managed implicitly, such as on platforms without full CUDA Unified Memory support and when data is not allocated in the CUDA Managed Memory segment.
- ▶ **Performance Tuning:** Even when data is located in the GPU-accessible memory, performance can be optimized via OpenACC features. Many OpenACC data directives and clauses provide hints to the CUDA device driver, which can improve implicit data management.

Data management strategies may differ depending on the specific goals being pursued. These differences are outlined where applicable.

### General Rules

All directives, except **host\_data**, can be used for data management tasks such as allocating memory in the GPU and copying data between the CPU and the GPU. These directives can be used to ensure that the data is present on the device during the execution of parallel algorithms. The **host\_data** construct, on the other hand, is used for address translation between CPU and GPU address spaces when data is accessed in parallel algorithms.

```
int n = get_n();
T* in = new T[nelem];
T* out = new T[nelem];
// Data captured by the lambda are managed explicitly with OpenACC
#pragma acc enter data copyin(n, in[0:nelem]) create(out[0:nelem])
#pragma acc host_data use_device(n, in, out)
{
    std::for_each(std::execution::par_unseq, r.begin(), r.end(),
        [&,in,out](auto i) {
            out[i] = in[i] * n;
        });
}
#pragma acc exit data copyout(out[0:nelem])
```

In the above example all data accessed from **std::for\_each** through the lambda capture are managed explicitly through the OpenACC data directives. Since the data inside the parallel algorithms are either captured by reference or capturing a pointer, the application code must ensure that such data is accessible from the GPU. To make non-GPU resident data accessible in the parallel region, such a region must be enclosed into the **host\_data** construct region with all variables that are managed explicitly via OpenACC runtime listed in the **use\_device** clause. The data need to be present (copied or created) at the time the **host\_data** directive is encountered/executed at runtime and the data must also be present for the duration of parallel algorithm execution. The implications of the above are such that lambdas accessing variables enclosed in **use\_device** regions can not be additionally invoked from the host code (from outside

the parallel region executing on the GPU) because the variable addresses from the GPU obtained through **host\_data** may not be accessible on the CPU.



If the iterator in the above example would be a pointer type it would require explicit data management in addition to the data captured by the lambda.

If the example below is compiled for Separate Memory Mode (-gpu=mem:separate) calling **fn** from within a parallel **std::for\_each** works fine but not from outside of any parallel algorithm function since the data resident on GPU would need to be accessed from the CPU.

```
int n = get_n();
T* in = new T[nelem];
T* out = new T[nelem];
#pragma acc enter data copyin(n, in[0:nelem]) create(out[0:nelem])
#pragma acc host_data use_device(n, in, out)
{
    auto fn = [&in,out](auto i) { out[i] = in[i] * n;};
    std::for_each(std::execution::par_unseq, r.begin(), r.end(), fn);
    // The following line would not be legal, fn accesses variables in GPU memory
    //std::for_each(r.begin(), r.end(), fn);
}
#pragma acc exit data copyout(out[0:nelem])
```



The behavior of using **use\_device** with non-pointer data type is such that all occurrences of non-pointer variables inside the **host\_data** region are converted to using the addresses of the variable in the GPU address space before accessing that variable. This is essentially equivalent to translating original occurrences of such variable **var** into **dvar = \*acc\_device(&var)**.

## Composite Data Types

Composite data types with pointer members can also be managed explicitly but require explicit deep copy to work correctly including pointer attach/detach.

```
struct S {
    float *ptr;
}

int idx[N] = { /*...*/ };
float arr[N];
S s{arr};
// Deep copying ptr member with OpenACC
#pragma acc enter data copyin(s.ptr[0:N])
#pragma acc enter data copyin(s, idx)
#pragma acc data attach(s.ptr)
#pragma acc host_data use_device(s, idx)
{
    std::for_each_n(std::execution::par, idx, N,
        [&](int i) { s.ptr[i] += 5.0; });
}
#pragma acc exit data copyout(s.ptr[0:N])
#pragma acc exit data copyout(s)
```

When variable of struct **s** type in the above example is copied to the device, a deep copy is performed with the content pointed by **s.ptr** copied separately. The pointer attachment is used to ensure the address of the pointer is changed to the device memory equivalent before it is accessed from the GPU. Depending on the order of the copies, the pointer **attach** clause may not be required.



In the above example the pointer-like iterator **idx** is managed through the OpenACC directives in addition to the data captured by the lambda.

## Standard Containers

If the standard containers with non-contiguous storage must be used in host code with explicit data management to GPU memory, the only viable option is to access the raw data directly using the raw pointer to data (e.g. obtained via **data()** member of **std::vector**) unless the iterator over the data can be used.

```
std::vector<T> in(nelem);
std::vector<T> out(nelem);
T *inptr=in.data(),*outptr=out.data();
#pragma acc data copyin(inptr[0:nelem]) copyout(outptr[0:nelem])
#pragma acc host_data use_device(inptr,outptr)
{
    std::for_each(std::execution::par_unseq, r.begin(), r.end(),
        [=](auto i) {
            outptr[i] = inptr[i];
        });
}
```

In the above example vector elements are accessed through raw pointers to their elements obtained through **vector::data()** member, they are explicitly management through the OpenACC data clauses.

## Static Storage Data

Global or static variables can be made accessible in the parallel algorithms using OpenACC data directives similarly to other variables.

```
int glob_arr[N] = { /*...*/ };
void foo(){
#pragma acc data copy(glob_arr)
#pragma acc host_data use_device(glob_arr)
{
    std::for_each_n(std::execution::par, glob_arr, N,
        [](int &e) { e += 1; });
}
}
```

In the above example the global array **glob\_arr** is updated on the GPU with help of OpenACC data directives.

## Member Functions

When the data members are managed inside the member functions the implicit object pointer **this** needs to be explicitly managed for correctness as accessing members is always done through the dereference of the object itself.

```
struct S {
    float *ptr;

    void update_member() {
#pragma acc data copy(ptr[0:N], this)
#pragma acc host_data use_device(ptr, this)
    {
        std::for_each(std::execution::par, ptr, ptr + N,
                      [=](float &e) { ptr[&e - ptr] += 5.0; });
    }
};
```

## GPU Memory Mode Related Differences

In Separate Memory Mode all data must be managed explicitly via extra device allocations and **memcpy** between the host and device and the address translations. This also applies to variables with automatic or static storage duration in Managed Memory Mode.

In Unified Memory Mode all data is automatically managed by the CUDA device driver. Additionally in Managed Memory Mode all dynamic allocations are managed by the CUDA device driver. Use of data clauses and directives can only propagate memory usage hints to the CUDA device driver which are used to improve the data management performance. More details can be found in Memory Model and [OpenACC with CUDA Unified Memory](#).

All the data managed by the CUDA device driver can benefit from the simplified uses of the OpenACC features, particularly:

- Use of **host\_data** directive is not required since the host and device address of data in unified shared memory is identical.
- Use of pointer attach or detach is not required since the host and device pointers in unified shared memory are identical.

The following example illustrates simplified data management with only OpenACC data construct enclosing the **std::for\_each** with Unified Memory Mode.

```
int n = get_n();
T* in = new T[nelem];
T* out = new T[nelem];
#pragma acc data copyin(in[0:nelem]) copyout(out[0:nelem])
{
    std::for_each(std::execution::par_unseq, r.begin(), r.end(),
                  [&](auto i) {
                      out[i] = in[i] * n;
                  });
}
```



In the above example we leverage OpenACC explicit data management construct to indicate how data is used on GPU for the computation executed in **std::for\_each**:

- ▶ **in** is moved into the GPU memory;
- ▶ **out** is moved from the GPU memory.

Both **in** and **out** are captured by reference and therefore their host address is used in the lambda of **std::for\_each**. The scalar variable **n** is not managed. The use of **host\_data** construct is not required.

When standard containers are used in data directives and clauses, the underlying data collection can be managed too. For example, in order to indicate that elements of the **std::vector** are accessed from the GPU the application code must first retrieve the pointer to the array elements using its **data()** member. Then such pointers can be used in the regular data directives.

```
std::vector<T> in(nelem);
std::vector<T> out(nelem);
T *inptr=in.data(), *outptr=out.data();
#pragma acc data copyin(inptr[0:nelem]) copyout(outptr[0:nelem])
{
    std::for_each(std::execution::par_unseq, r.begin(), r.end(),
                  [&](auto i) {
                      out[i] = in[i];
                  });
}
```

The above example demonstrates the use of OpenACC data directives with a raw pointer to elements of **std::vector** which can improve memory performance for data in unified memory and the full deep copy of vector content using attach/detach is not required.

```
int n = get_n();
T* in = new T[nelem];
T* out = new T[nelem];
#pragma acc enter data copyin(n)
#pragma acc host_data use_device(n)
{
    std::for_each(std::execution::par_unseq, r.begin(), r.end(),
                  [&, in, out](auto i) {
                      out[i] = in[i] * n;
                  });
}
#pragma acc enter data delete(n)
```

In the above example, **in** and **out** are dynamically allocated and managed by CUDA device driver with Managed Memory Mode, **n** is on the stack and therefore managed explicitly via OpenACC directives.

### 8.2.7.2. External Device Function Annotations

Using OpenACC routine directive annotations allows calling external device functions.

```
// In file1.cpp
extern int foo();

void bar()
{
    std::for_each(std::execution::par_unseq, r.begin(), r.end(),
                  [=](auto i) {
                      ou[i] = foo();
                  });
}

// In file2.cpp
#pragma acc routine
int foo(){
    return 4;
}
```

The above code can be compiled/linked as follows:

```
nvc++ -stdpar file1.cpp
nvc++ -acc file2.cpp
nvc++ -stdpar -acc file1.o file2.o
```

## 8.2.8. Getting Started with Parallel Algorithms for GPUs

To get started, download and install the [NVIDIA HPC SDK](#) on your x86-64, OpenPOWER, or Arm CPU-based system running a supported version of Linux.

The NVIDIA HPC SDK is freely downloadable and includes a perpetual use license for all NVIDIA Registered Developers, including access to future release updates as they are issued. After you have the NVIDIA HPC SDK installed on your system, the `nvc++` compiler is available under the `/opt/nvidia/hpc_sdk` directory structure.

- ▶ To use the compilers including `nvc++` on a Linux/x86-64 system, add the directory `/opt/nvidia/hpc_sdk/Linux_x86_64/24.9/compilers/bin` to your path.
- ▶ On an OpenPOWER or Arm CPU-based system, replace `Linux_x86_64` with `Linux_ppc64le` or `Linux_aarch64`, respectively.

### 8.2.8.1. Supported NVIDIA GPUs

The NVC++ compiler can automatically offload C++ Parallel Algorithms to NVIDIA GPUs based on the Volta architecture or newer. These architectures include features -- such as independent thread scheduling and hardware optimizations for CUDA Unified Memory -- that were specifically designed to support high-performance, general-purpose parallel programming models like the C++ Parallel Algorithms.

The NVC++ compiler provides limited support for C++ Parallel Algorithms on the Pascal architecture, which does not have the [independent thread scheduling](#) necessary to properly support the `std::execution::par` policy. When compiling for the Pascal architecture (`-gpu=cc60`), NVC++ compiles algorithms with the

**std::execution::par** policy for serial execution on the CPU. Only algorithms with the **std::execution::par\_unseq** policy will be scheduled to run on Pascal GPUs.

### 8.2.8.2. Supported CUDA Versions

The NVC++ compiler is built on CUDA libraries and technologies and uses CUDA to accelerate C++ Parallel Algorithms on NVIDIA GPUs. A GPU-accelerated system on which NVC++-compiled applications are to be run must have a CUDA 11.2 or newer device driver installed.

The NVIDIA HPC SDK compilers ship with an integrated CUDA toolchain, header files, and libraries to use during compilation, so it is not necessary to have a CUDA Toolkit installed on the system.

When **-stdpar** is specified, NVC++ compiles using the CUDA toolchain version that best matches the CUDA driver installed on the system on which compilation is performed. To compile using a different version of the CUDA toolchain, use the **-gpu=cudaX.Y** option. For example, use the **-gpu=cuda11.8** option to specify that your program should be compiled for a CUDA 11.8 system using the CUDA 11.8 toolchain.

## 8.3. Stdpar Fortran

Fortran 2008 introduced the **do concurrent** (DC) loop construct signaling that loop iterations have no interdependencies. With **-stdpar** such loop iterations will be executed in parallel on the GPU when **-stdpar** (or **-stdpar=gpu**) is passed to **nvfortran** or using CPU threads when **-stdpar=multicore** is passed to **nvfortran**. More details can be found in the following blog post on the NVIDIA website: [Accelerating Fortran DO CONCURRENT with GPUs and the NVIDIA HPC SDK](#).

### 8.3.1. Calling Routines in DO CONCURRENT on the GPU

When compiling for the GPU, calling routines in the body of **do concurrent** loop can be constrained. PURE routines can generally be called inside the **do concurrent** loop body. The compiler detects that such routines are to be compiled for the GPU target. External routines, however, can't be called from within the DC loop unless they are explicitly annotated with the OpenACC routine directive (refer to [Interoperability with OpenACC](#)) or CUDA device attribute (refer to [Interoperability with CUDA Fortran](#)).

The following example will compile successfully.

```
module m
contains
pure subroutine foo()
return
end subroutine
end module m

program dc
use m
implicit none
integer :: i

do concurrent (i=1:10)
```

```

    call foo()
enddo
end program

```

The following example, however, doesn't compile unless **foo** is either

- ▶ annotated with **!\$acc routine**,
- ▶ or attributed with **attributes(device)** and compiled as Stdpar and CUDA Fortran.

```

program dc
implicit none
interface
    pure subroutine foo()
    end subroutine foo
end interface
integer :: i

do concurrent (i=1:10)
    call foo()
enddo
end program

```

### 8.3.2. GPU Data Management

If **-gpu=mem:managed** is enabled by default or is explicitly passed on the command line, some data accesses in **do concurrent** loops are invalid. For example, accessing global variables in the routines called from the **do concurrent** loop does not perform expected value updates in the CPU code.

Additionally, there are rare instances where the compiler cannot accurately determine variable sizes for implicit data movements between CPU and GPU. As demonstrated in the following example, **a** is an assumed-size array, and its access region inside the DC construct cannot be determined at compile time because the element index positions are taken from another array **b** initialized outside of the routine. Such code does not update **a** as expected and may result in a memory violation and undefined behavior.

```

subroutine r(a, b)
    integer :: a(*)
    integer :: b(:)
    do concurrent (i = 1 : size(b))
        a(b(i)) = i
    enddo
end subroutine

```

There are no limitations on the variable accessed in **do concurrent** loops described above when the code is compiled with **-gpu=mem:unified**, whether this option is enabled by default or explicitly via an option on the command line.

### 8.3.3. Interoperability with OpenACC

OpenACC features can be used when compiling Stdpar code for GPUs. To activate OpenACC directives recognition with Stdpar code add **-acc** command line flag to **nvfortran**.

```
nvfortran -stdpar -acc example.f90
```

OpenACC functionality and interoperability with DO-CONCURRENT loop is detailed in the OpenACC specification and the NVIDIA HPC compiler specific differences are detailed in [Using OpenACC](#) of this guide.

Using OpenACC features can enhance functionality of DC-loop for example with the following:

- ▶ Explicit data management to improve performance of CPU-GPU implicit data movements or even leverage separate memory compiling on the GPU when compiling with -gpu=mem:separate passed in.
- ▶ Tuning DC-loop execution on the GPU e.g. GPU kernels launch configuration.
- ▶ Executing DC-loops asynchronously.
- ▶ Calling external routines from within DC-loops.
- ▶ Atomic operations in DC-loops.

#### Examples

Some examples of using OpenACC directives with DC-loops are provided below.

The following example demonstrates how the data accessed inside the DC-loop are fully managed in the OpenACC data construct.

```
!$acc data copyin(b) copyout(a)
do concurrent (j=1:N)
  do i=1,K
    a(j,i) = b(j,i)
  end do
end do
!$acc end data
```

While in the above example the data construct is used for GPU data management, the same effect can be achieved with the use of data clauses on the compute construct enclosing DC-loop.

The following example shows how the scheduling of DC loop on the GPU is controlled through the clauses on the compute construct.

```
!$acc parallel loop num_gangs(50000) vector_length(32)
do concurrent (i=1:K,j=1:N)
  a(j,i) = real(j)
end do
```

Use of OpenACC async clause on the compute constructs can be utilised to perform computations in DC-loop asynchronously.

```
!$acc parallel loop async
do concurrent (j=1:N)
  a(j) = j
end do

b = foo()

#pragma acc wait

c = sum(a) + b
```

In the previous example, array **a** is filled in with values asynchronously in DC-loop.

### 8.3.4. Interoperability with CUDA Fortran

CUDA Fortran features can also be used when compiling Stdpar code for GPUs. To recognize CUDA Fortran features in your source code, compile with the **-cuda** command line flag using **nvfortran**.

```
nvfortran -stdpar -cuda example.f90
```

Using CUDA Fortran extensions can enhance the functionality of a do concurrent (DC) loop and Stdpar program, for several cases:

- ▶ Explicit data locality, accessing CUDA Fortran attributed arrays or other data with the device, managed, unified, or constant attributes from within DC-loops.
- ▶ Tuning DC-loop execution on the GPU e.g. controlling the GPU kernels launch configuration.
- ▶ Executing DC-loops asynchronously using a specific CUDA stream.
- ▶ Calling external, user-defined CUDA device routines from within DC-loops.
- ▶ Using CUDA Atomic operations in DC-loops, or other CUDA-specific device-side runtime library calls.
- ▶ Inserting CUDA Runtime API calls for memory tuning hints outside of DC-loops.

#### Examples

Some examples of using CUDA Fortran features with DC-loops are provided below. The following example demonstrates how a DC-loop can access CUDA Fortran device data, run on a specific CUDA stream, call the CUDA Runtime API for creating a stream, and hide non-standard features behind the CUF sentinel for code portability.

```
!@cuf use cudafor
!@cuf integer(kind=cuda_stream_kind) :: istrm
      real, allocatable :: a(:,,:), b(:,,:)
!@cuf attributes(device) :: a ! A is device array only, not unified/managed
      . . .
!@cuf istat = cudaStreamCreate(istrm)
      . . .
      a(:,,:) = 0.0
      . . .
!$cuf kernel do(1) <<< *, *, stream=istrm>>>
do concurrent (j=1:N)
  do i=1,K
```

```

        a(j,i) = a(j,i) + 2.0 * b(j,i)
    end do
end do

```

This program demonstrates how to call low-level CUDA device functions from within a DC-loop. The function can be written in either CUDA Fortran or CUDA C++, depending on the interface. The CUDA C function must be compiled for relocatable device code. This can be used for accessing features in CUDA and NVIDIA GPUs not readily available in directive-based models or standard languages.

```

module mcuda
contains
    attributes(host,device) pure integer function std_dbg(itype)
    integer, value :: itype
    if (itype.eq.1) then
        std_dbg = threadIdx%x
    else if (itype.eq.2) then
        std_dbg = blockIdx%x
    else
        std_dbg = (blockIdx%x-1)*blockDim%x + threadIdx%x
    end if
end function
end module

program test
use mcuda
integer, parameter :: N = 2000
integer, allocatable :: a(:), b(:), c(:)
allocate(a(N),b(N),c(N))

do concurrent (j=1:N)
    a(j) = std_dbg(1)
    b(j) = std_dbg(2)
    c(j) = std_dbg(3)
end do

print *,a(1),a(N/2),a(N)
print *,b(1),b(N/2),b(N)
print *,c(1),c(N/2),c(N)
end

```

Many functions from the CUDA Fortran **cudadevice** module are available within do concurrent loops, not just atomics. This code snippet shows two uses:

```

real :: tmp(4), x, y
...
block; use cudadevice
do concurrent (i=1:K,j=1:N)
    x = real(j) + a(i,j)
    y = atomicAdd(b(1,j), x)
end do

do concurrent (j=1:N)
    x = real(j)
    tmp(1:4) = __ldca(a(1:4,j))
    tmp(1:4) = tmp(1:4) + x
    call __stwt(b(1:4,j), tmp)
end do
end block

```

# Chapter 9.

## PCAST

Parallel Compiler Assisted Software Testing (PCAST) is a set of API calls and compiler directives useful in testing program correctness. Numerical results produced by a program can diverge when parts of the program are mapped onto a GPU, when new or additional compiler options are used, or when changes are made to the program itself. PCAST can help you determine where these divergences begin, and pinpoint the changes that cause them. It is useful in other situations as well, including when using new libraries, determining whether parallel execution is safe, or porting programs from one ISA or type of processor to another.

### 9.1. Overview

PCAST Comparisons can be performed in two ways. The first saves the initial run's data into a file through the **pcast\_compare** call or directive. Add the calls or directives to your application where you want intermediate results to be compared. Then, execute the program to save the "golden" results where the values are known to be correct. During subsequent runs of the program, the same **pcast\_compare** calls or directives will compare the computed intermediate results to the saved "golden" results and report the differences.

The second approach works in conjunction with the NVIDIA OpenACC implementation to compare GPU computation against the same program running on a CPU. In this case, all compute constructs are performed redundantly, both on the CPU and GPU. GPU results are compared against the CPU results, and differences reported. This is essentially like the first case where the CPU-calculated values are treated as the "golden" results. GPU to CPU comparisons can be done implicitly at the end of data regions with the **autocompare** flag or explicitly after kernels with the **acc\_compare** call or directive.

With the **autocompare** flag, OpenACC regions will run redundantly on the CPU and GPU. On an OpenACC region exit where data is to be downloaded from device to host, PCAST will compare the values calculated on the CPU with those calculated in the GPU. Comparisons done with **autocompare** or **acc\_compare** are handled in memory and do not write results to an intermediate file.



The following table outlines the supported data types that can be used with PCAST. Short, integer, long, and half precision data types are not supported with **ABS**, **REL**, **ULP**, or **IEEE** options; only a bit-for-bit comparison is supported.

For floating-point types, PCAST can calculate absolute, relative, and unit-last-place differences. Absolute differences measures only the absolute value of the difference (subtraction) between two values, i.e.  $abs(A-B)$ . Relative differences are calculated as a ratio between the difference of values,  $A-B$ , and the previous value  $A$ ;  $abs((A-B)/A)$ . Unit-least precision (Unit-last place) is a measure of the smallest distance between two values  $A$  and  $B$ . With the **ULP** option set, PCAST will report if the calculated ULP between two numbers is greater than some threshold.

Table 23 Supported Types for Tolerance Measurements

C/C++ Type	Fortran Type	ABS	REL	ULP	IEEE
float	real, real(4)	Yes	Yes	Yes	Yes
double	double precision, real(8)	Yes	Yes	Yes	Yes
float _Complex	complex, complex(4)	Yes	Yes	Yes	Yes
double _Complex	complex(8)	Yes	Yes	Yes	Yes
-	real(2)	No	No	No	No
(un)signed short	integer(2)	N/A	N/A	N/A	N/A
(un)signed int	integer, integer(4)	N/A	N/A	N/A	N/A
(un)signed long	integer(8)	N/A	N/A	N/A	N/A

## 9.2. PCAST with a "Golden" File

The run-time call **pcast\_compare** highlights differences between successive program runs. It has two modes of operation, depending on the presence of a data file named *pcast\_compare.dat* by default. If the file does not exist, **pcast\_compare** assumes this is the first "golden" run. It will create the file and fill it with the computed data at each call to **pcast\_compare**. If the file exists, **pcast\_compare** assumes it is a test run. It will read the file and compare the computed data with the saved data from the file. The default behavior is to consider the first 50 differences to be a reportable error, no matter how small.

By default, the **pcast\_compare.dat** file is in the same directory as the executable. The behavior of **pcast\_compare**, and other comparison parameters, can be changed at runtime with the PCAST\_COMPARE environment variable discussed in the [Environment Variables](#) section.

The signature of **pcast\_compare** for C++ and C is:

```
void pcast_compare(void*, char*, size_t, char*, char*, char*, int);
```

The signature of **pcast\_compare** for Fortran is:

```

subroutine pcast_compare(a, datatype, len, varname, filename, funcname, lineno)
  type(*), dimension(..) :: a
  character(*) :: datatype, varname, filename, funcname
  integer(8),value :: len
  integer(4),value :: lineno

```

The call takes seven arguments:

1. The address of the data to be saved or compared.
2. A string containing the data type.
3. The number of elements to compare.
4. A string treated as the variable name.
5. A string treated as the source file name.
6. A string treated as the function name.
7. An integer treated as a line number.

For example, the **pcast\_compare** runtime call can be invoked like the following:

```
pcast_compare(a, "float", N, "a", "pcast_compare03.c", "main", 1);
```

```
call pcast_compare(a, 'real', n, 'a', 'pcast_compare1.f90', 'program', 9)
```

The caller should give meaningful names to the last four arguments. They can be anything, since they only serve to annotate the report. It is imperative that the identifiers are not modified between comparisons; comparisons must be called in the same order for each program run. If, for example, you are calling **pcast\_compare** inside a loop, it is reasonable to set the last argument to be the loop index.

There also exists a directive form of the **pcast\_compare**, which is functionally the same as the runtime call. It can be used at any point in the program to compare the current value of data to that recorded in the golden file, same as the runtime call. There are two benefits to using the directive over the API call:

1. The directive syntax is much simpler than the API syntax. Most of what the compare call needs to output data to the user can be gleaned by the compiler at compile-time (The type, variable name, file name, function name, and line number).

```
#pragma nvidia compare(a[0:n])
```

as opposed to:

```
pcast_compare(a, "float", N, "a", "pcast_compare03.c", "main", 1);
```

2. The directive is only enabled when the -Mpcast flag is set, so the source need not be changed when testing is complete. Consider the following usage examples:

```
#pragma nvidia compare(a[0:N]) // C++ and C
!$nvf compare(a(1:N)) ! Fortran
```

The directive interface is given below in C++ or C style, and in Fortran. Note that for Fortran, **var-list** is a variable name, a subarray specification, an array element, or a composite variable member.

```
#pragma nvidia compare (var-list) // C++ and C
!$nvf compare (var-list) ! Fortran
```

Let's look at an example of

```
#include <stdlib.h>
#include <openacc.h>

int main() {
    int size = 1000;
    int i, t;
    float *a1;
    float *a2;

    a1 = (float*)malloc(sizeof(float)*size);
    a2 = (float*)malloc(sizeof(float)*size);

    for (i = 0; i < size; i++) {
        a1[i] = 1.0f;
        a2[i] = 2.0f;
    }

    for (t = 0; t < 5; t++) {
        for (i = 0; i < size; i++) {
            a2[i] += a1[i];
        }
        pcast_compare(a2, "float", size, "a2", "example.c", "main", 23);
    }
    return 0;
}
```

Compile the example using these compiler options:

```
> nvc -fast -o a.out example.c
```

Compiling with redundant or autocompare options are not required to use pcast\_compare. Once again, running the compiled executable using the options below, results in the following output:

```
> PCAST_COMPARE=summary,rel=1 ./out.o
datafile pcast_compare.dat created with 5 blocks, 5000 elements, 20000 bytes
> PCAST_COMPARE=summary,rel=1 ./out.o
datafile pcast_compare.dat compared with 5 blocks, 5000 elements, 20000 bytes
no errors found
relative tolerance = 0.100000, rel=1
```

Running the program for the first time, the data file "pcast\_compare.dat" is created. Subsequent runs compare calculated data against this file. Use the **PCAST\_COMPARE** environment variable to set the name of the file, or force the program to create a new file on the disk with **PCAST\_COMPARE=create**.

The same example above can be written with the compare directive. Notice how much more concise the directive is to the update host and **pcast\_compare** calls.

```
#include <stdlib.h>
#include <openacc.h>

int main() {
    int size = 1000;
    int i, t;
    float *a1;
    float *a2;

    a1 = (float*)malloc(sizeof(float)*size);
    a2 = (float*)malloc(sizeof(float)*size);

    for (i = 0; i < size; i++) {
        a1[i] = 1.0f;
        a2[i] = 2.0f;
    }

    for (t = 0; t < 5; t++) {
        for(i = 0; i < size; i++) {
            a2[i] += a1[i];
        }
        #pragma nvidia compare(a2[0:size])
    }
    return 0;
}
```

With the directive, you will want to add "-Mpcast" to the compilation line to enable the directive. Other than that, the output from this program is identical to the runtime example above.

## 9.3. PCAST with OpenACC

PCAST can also be used with the NVIDIA OpenACC implementation to compare GPU computation against the same program running on a CPU. In this case, all compute constructs are performed redundantly on both the CPU and GPU. The CPU results are considered to be the "golden master" copy which GPU results are compared against.

There are two ways to perform comparisons with GPU-calculated results. The first is with the explicit call or directive **acc\_compare**. To use **acc\_compare**, you must compile with **-acc -gpu=redundant** to force the CPU and GPU to compute results redundantly. Then, insert calls to **acc\_compare** or put an **acc compare** directive at points where you want to compare the GPU-computed values against those computed by the CPU.

The second approach is to turn on autocompare mode by compiling with **-acc -gpu=autocompare**. In autocompare mode, PCAST will automatically perform a comparison at each point where data is moved from the device to the host. It does not require the programmer to add any additional directives or runtime calls; it's a convenient way to do all comparisons at the end of a data region. If there are multiple compute kernels within a data region, and you're only interested in one specific kernel, you should use the previously-mentioned **acc\_compare** to target a specific kernel. Note that autocompare mode implies **-gpu=redundant**.

During redundant execution, the compiler will generate both CPU and GPU code for each compute construct. At runtime, both the CPU and GPU versions will execute redundantly, with the CPU code reading and modifying values in system memory and the GPU reading and modifying values in device memory. Insert calls to **acc\_compare()** calls (or the equivalent **acc compare** directive) at points where you want to compare the GPU-computed values against CPU-computed values. PCAST treats the values generated by the CPU code as the "golden" values. It will compare those results against GPU values. Unlike **pcast\_compare**, **acc\_compare** does not write to an intermediary file; the comparisons are done in-memory.

**acc\_compare** only has two arguments: a pointer to the data to be compared, *hostptr*, and the number of elements to compare, *count*. The type can be inferred in the OpenACC runtime, so it doesn't need to be specified. The C++ and C interface is given below:

```
void acc_compare(void *, size_t);
```

And in Fortran:

```
subroutine acc_compare(a)
subroutine acc_compare(a, len)
  type(*), dimension(*) :: a
  integer(8), value :: len
```

You can call **acc\_compare** on any variable or array that is present in device memory. You can also call **acc\_compare\_all** (no arguments) to compare all values that are present in device memory against the corresponding values in host memory.

```
void acc_compare_all()
```

```
subroutine acc_compare_all()
```

Directive forms of the **acc\_compare** calls exist. They work the same as the API calls and can be used in lieu of them. Similar to PCAST **compare** directives, **acc compare** directives are ignored when redundant or autocompare modes are not enabled on the compilation line.

The **acc compare** directive takes one or more arguments, or the 'all' clause (which corresponds to **acc\_compare\_all()**). The interfaces are given below in C++ or C, and Fortran respectively. Argument "var-list" can be a variable name, a sub-array specification, and array element, or a composite variable member.

```
#pragma acc compare [ (var-list) | all ]
```

```
$!acc compare [ (var-list) | all ]
```

For example:

```
#pragma acc compare(a[0:N])
#pragma acc compare all
```

```
!$acc compare(a, b)
!$acc compare(a(1:N))
!$acc compare all
```

Consider the following OpenACC program that uses the **acc\_compare()** API call and an **acc compare** directive. This Fortran example uses real\*4 and real\*8 arrays.

```
program main
  use openacc
  implicit none
  parameter N = 1000
  integer :: i
  real :: a(N)
  real*4 :: b(N)
  real(4) :: c(N)
  double precision :: d(N)
  real*8 :: e(N)
  real(8) :: f(N)

  d = 1.0d0
  e = 0.1d0

  !$acc data copyout(a, b, c, f) copyin(d, e)

  !$acc parallel loop
  do i = 1,N
    a(i) = 1.0
    b(i) = 2.0
    c(i) = 0.0
  enddo
  !$acc end parallel

  !$acc compare(a(1:N), b(1:N), c(1:N))

  !$acc parallel loop
  do i = 1,N
    f(i) = d(i) * e(i)
  enddo
  !$acc end parallel

  !$acc compare(f)

  !$acc parallel loop
  do i = 1,N
    a(i) = 1.0
    b(i) = 1.0
    c(i) = 1.0
  enddo
  !$acc end parallel

  call acc_compare(a, N)
  call acc_compare(b, N)
  call acc_compare(c, N)

  !$acc parallel loop
  do i = 1,N
    f(i) = 1.0D0
  enddo
  !$acc end parallel

  call acc_compare_all()

  !$acc parallel loop
  do i = 1,N
    a(i) = 3.14;
    b(i) = 3.14;
```

```

c(i) = 3.14;
f(i) = 3.14d0;
enddo
!$acc end parallel

! In redundant mode, no comparison is performed here. In
! autocompare mode, a comparison is made for a, b, c, and f (but
! not e and d), since they are copied out of the data region.

!$acc end data

call verify(N, a, b, c, f)
end program

subroutine verify(N, a, b, c, f)
integer, intent(in) :: N
real, intent(in) :: a(N)
real*4, intent(in) :: b(N)
real(4), intent(in) :: c(N)
real(8), intent(in) :: f(N)
integer :: i, errcnt

errcnt = 0
do i=1,N
if(abs(a(i) - 3.14e0) .gt. 1.0e-06) then
errcnt = errcnt + 1
endif
end do
do i=1,N
if(abs(b(i) - 3.14e0) .gt. 1.0e-06) then
errcnt = errcnt + 1
endif
end do
do i=1,N
if(abs(c(i) - 3.14e0) .gt. 1.0e-06) then
errcnt = errcnt + 1
endif
end do
do i=1,N
if(abs(f(i) - 3.14d0) .gt. 1.0d-06) then
errcnt = errcnt + 1
endif
end do

if(errcnt /= 0) then
write (*, *) "FAILED"
else
write (*, *) "PASSED"
endif
end subroutine verify

```

The program can be compiled with the following command:

```

> nvfortran -fast -acc -gpu=redundant -Minfo=accel example.F90
main:
16, Generating copyout(a(:),b(:))
Generating copyin(e(:))
Generating copyout(f(:),c(:))
Generating copyin(d(:))
18, Generating Tesla code
19, !$acc loop gang, vector(128) ! blockidx%x threadidx%x
26, Generating acc compare(c(:),b(:),a(:))
28, Generating Tesla code
29, !$acc loop gang, vector(128) ! blockidx%x threadidx%x
34, Generating acc compare(f(:))
36, Generating Tesla code

```

```

37, !$acc loop gang, vector(128) ! blockidx%x threadidx%x
48, Generating Tesla code
49, !$acc loop gang, vector(128) ! blockidx%x threadidx%x
56, Generating Tesla code
57, !$acc loop gang, vector(128) ! blockidx%x threadidx%x

```

Here, you can see where the acc compare directives are generated on lines 26 and 34. The program can be run with the following command:

```

> ./a.out
PASSED

```

As you can see, no PCAST output is generated when the comparisons match. We can get more information with the summary option:

```

> PCAST_COMPARE=summary ./a.out
PASSED
compared 13 blocks, 13000 elements, 68000 bytes
no errors found
absolute tolerance = 0.0000000000000000e+00, abs=0

```

There are 13 blocks compared. Let's count the blocks in the compare calls.

```
!$acc compare(a(1:N), b(1:N), c(1:N))
```

Compares three blocks, one each for a, b, and c.

```
!$acc compare(f)
```

Compares one block for f.

```

call acc_compare(a, N)
call acc_compare(b, N)
call acc_compare(c, N)

```

Each call compares one block for their respective array.

```
call acc_compare_all()
```

Compares one block for each array present on the device (a, b, c, d, e, and f) for a total of 6 blocks.

If the same example is compiled with autocompare, we'll see four additional comparisons, since the four arrays that are copied out (with the copyout clause) are compared at the end of the data region.

```

> nvfortran -fast -acc -gpu=autocompare example.F90
> PCAST_COMPARE=summary ./a.out
PASSED
compared 17 blocks, 17000 elements, 88000 bytes
no errors found
absolute tolerance = 0.0000000000000000e+00, abs=0

```



## 9.4. Limitations

There are currently a few limitations with using PCAST that are worth keeping in mind.

- ▶ Comparisons are not thread-safe. If you are using PCAST with multiple threads, ensure that only one thread is doing the comparisons. This is especially true if you are using PCAST with MPI. If you use **pcast\_compare** with MPI, you must make sure that only one thread is writing to the comparison file. Or, use a script to set PCAST\_COMPARE to encode the file name with the MPI rank.
- ▶ Comparisons must be done with like types; you cannot compare one type with another. It is not possible to, for example, check for differing results after changing from double precision to single. Comparisons are limited to those present in table [Table 23](#). Currently there is no support for structured or derived types.
- ▶ The **-gpu=mem:managed** or **-gpu=mem:unified** options are incompatible with autocompare and **acc\_compare**. Both the CPU and GPU need to calculate result separately and to do so they must have their own working memory spaces.
- ▶ If you do any data movement on the device, you must account for it on the host. For example, if you are using CUDA-aware MPI or GPU-accelerated libraries that modify device data, then you must also make the host aware of the changes. In these cases it is helpful to use the **host\_data** clause, which allows you to use device addresses within host code.

## 9.5. Environment Variables

Behavior of PCAST/Autocompare is controlled through the **PCAST\_COMPARE** variable. Options can be specified in a comma-separated list:

**PCAST\_COMPARE=<opt1>,<opt2>,...**

If no options are specified, the default is to perform comparisons with *abs=0*.

Comparison options are not mutually exclusive. PCAST can compare absolute differences with some *n=3* and relative differences with a different threshold, e.g. *n=5*; *PCAST\_COMPARE=abs=3,rel=5,...*

You can specify either an absolute or relative location to be used with the *datafile* option. The parent directory should be owned by the same user executing the comparisons and the datafile should have the appropriate read/write permissions set.

Table 24 PCAST\_COMPARE Options

Option	Description
abs= <i>n</i>	Compare absolute difference; tolerate differences up to $10^{(-n)}$ , only applicable to floating point types. Default value is 0
create	Specifies that this is the run that will produce the reference file (pcast_compare only)
compare	Specifies that the current run will be compared with a reference file (pcast_compare only)

Option	Description
datafile="name"	Name of the file that data will be saved to, or compared against. If empty will use the default, 'pcast_compare.dat' (pcast_compare only)
disable	Calls to pcast_compare, acc_compare, acc_compare_all, and directives (pcast compare, acc compare, and acc compare) all immediately return from the runtime with no effect. Note that this doesn't disable redundant execution; that will require a recompile.
ieee	Compare IEEE NaN checks (only implemented for floats and doubles)
outputfile="name"	Save comparison output to a specific file. Default behavior is to output to stderr
patch	Patch errors (outside tolerance) with correct values
patchall	Patch all differences (inside and outside tolerance) with correct values
rel=n	Compare relative difference; tolerated differences up to $10^{-n}$ , only applicable to floating point types. Default value is 0.
report=n	Report up to n (default of 50) passes/fails
reportall	Report all passes and fails (overrides limit set in report=n)
reportpass	Report passes; respects limit set with report=n
silent	Suppress output - overrides all other output options, including summary and verbose
stop	Stop at first differences
summary	Print summary of comparisons at end of run
ulp=n	Compare Unit of Least Precision difference (only for floats and doubles)
verbose	Outputs more details of comparison (including patches)
verboseautocompare	Outputs verbose reporting of what and where the host is comparing (autocompare only)

# Chapter 10.

## USING MPI

MPI (the Message Passing Interface) is an industry-standard application programming interface designed for rapid data exchange between processors in a distributed-memory environment. MPI is computer software used in scalable computer systems that allows the processes of a parallel application to communicate with one another.

The NVIDIA HPC SDK includes a pre-compiled version of Open MPI. You can build using alternate versions of MPI with the `-I`, `-L`, and `-l` options.

This section describes how to use Open MPI with the NVIDIA HPC Compilers.

### 10.1. Using Open MPI on Linux

The NVIDIA HPC Compilers for Linux ship with a pre-compiled version of Open MPI that includes everything required to compile, execute and debug MPI programs using Open MPI.

To build an application using Open MPI, use the Open MPI compiler wrappers: `mpicc`, `mpic++` and `mpifort`. These wrappers automatically set up the compiler commands with the correct include file search paths, library directories, and link libraries.

The following MPI example program uses Open MPI.

```
$ cd my_example_dir
$ cp -r /opt/nvidia/hpc_sdk/Linux_x86_64/2024/examples/MPI/samples/mpihello .
$ cd mpihello
$ export PATH=/opt/nvidia/hpc_sdk/Linux_x86_64/24.9/mpi/openmpi/bin:$PATH
$ mpifort mpihello.f -o mpihello
```

```
$ mpiexec mpihello
Hello world! I'm node 0
```

```
$ mpiexec -np 4 mpihello
Hello world! I'm node 0
Hello world! I'm node 2
Hello world! I'm node 1
Hello world! I'm node 3
```

To build an application using Open MPI for debugging, add `-g` to the compiler wrapper command line arguments.

## 10.2. Using MPI Compiler Wrappers

When you use MPI compiler wrappers to build with the `-fpic` or `-mmodel=medium` options, then you must specify `-fortranlibs` to link with the correct libraries. Here are a few examples:

For a static link to the MPI libraries, use this command:

```
% mpifort hello.f
```

For a dynamic link to the MPI libraries, use this command:

```
% mpifort hello.f -fortranlibs
```

To compile with `-fpic`, which, by default, invokes dynamic linking, use this command:

```
% mpifort -fpic -fortranlibs hello.f
```

To compile with `-mmodel=medium`, use this command:

```
% mpifort -mmodel=medium -fortranlibs hello.f
```

## 10.3. Testing and Benchmarking

The `/opt/nvidia/hpc_sdk/Linux_x86_64/2024/examples/MPI` directory contains various benchmarks and tests. Copy this directory into a local working directory by issuing the following command:

```
% cp -r /opt/nvidia/hpc_sdk/Linux_x86_64/2024/examples/MPI .
```

There are several example programs available in this directory.

# Chapter 11.

## CREATING AND USING LIBRARIES

A library is a collection of functions or subprograms that are grouped for reference and ease of linking. This section discusses issues related to NVIDIA-supplied compiler libraries. Specifically, it addresses the use of C++ and C builtin functions in place of the corresponding libc routines, creation of dynamically linked libraries, known as shared objects or shared libraries, and math libraries.



This section does not duplicate material related to using libraries for inlining which are described in [Creating an Inline Library](#).

NVIDIA provides libraries that export C interfaces by using Fortran modules.

### 11.1. Using builtin Math Functions in C++ and C

The name of the math header file is `math.h`. Include the math header file in all of your source files that use a math library routine as in the following example, which calculates the inverse cosine of 3.5.

```
#include <math.h>
#include<stdio.h>
#define PI 3.1415926535
void main()
{
    double x, y;
    x = PI/3.0;
    y = acos(0.5);
    printf('%f %f\n', x, y);
}
```

Including `math.h` causes the NVIDIA C++ and C compilers to use builtin functions, which are much more efficient than library calls. In particular, if you include `math.h`, the following intrinsics calls are processed using builtins:

abs	acosf	asinf	atan	atan2	atan2f
atanf	cos	cosf	exp	expf	fabs
fabsf	fmax	fmaxf	fmin	fminf	log
log10	log10f	logf	pow	powf	sin

[sinf](#)
[sqrt](#)
[sqrtf](#)
[tan](#)
[tanf](#)

## 11.2. Using System Library Routines

Release 24.9 of the NVIDIA HPC Compilers runtime libraries makes use of Linux system libraries to implement, for example, OpenMP and Fortran I/O. The NVIDIA HPC Compilers runtime libraries make use of several additional system library routines.

On 64-bit Linux systems, the system library routines used include these:

<code>aio_error</code>	<code>aio_write</code>	<code>pthread_mutex_init</code>	<code>sleep</code>
<code>aio_read</code>	<code>calloc</code>	<code>pthread_mutex_lock</code>	
<code>aio_return</code>	<code>getrlimit</code>	<code>pthread_mutex_unlock</code>	
<code>aio_suspend</code>	<code>pthread_attr_init</code>	<code>setrlimit</code>	

## 11.3. Creating and Using Shared Object Files on Linux

All of the NVIDIA HPC Fortran, C++ and C compilers support creation of shared object files. Unlike statically-linked object and library files, shared object files link and resolve references with an executable at runtime via a dynamic linker supplied with your operating system. The NVIDIA HPC Compilers must generate position independent code to support creation of shared objects by the linker. However, this is not the default. You must create object files with position independent code and shared object files that will include them.

### 11.3.1. Procedure to create a use a shared object file

The following steps describe how to create and use a shared object file.

1. Create an object file with position independent code.

To do this, compile your code with the appropriate NVIDIA HPC compiler using the `-fpic` option, or one of the equivalent options, such as `-fPIC`, `-Kpic`, and `-KPIC`, which are supported for compatibility with other systems. For example, use the following command to create an object file with position independent code using `nvfortran`:

```
% nvfortran -c -fpic tobeshared.f
```

2. Produce a shared object file.

To do this, use the appropriate NVIDIA HPC compiler to invoke the linker supplied with your system. It is customary to name such files using a `.so` filename extension. On Linux, you do this by passing the `-shared` option to the linker:

```
% nvfortran -shared -o tobeshared.so tobeshared.o
```



Compilation and generation of the shared object can be performed in one step using both the `-fpic` option and the appropriate option for generation of a shared object file.

### 3. Use a shared object file.

To do this, use the appropriate NVIDIA HPC compiler to compile and link the program which will reference functions or subroutines in the shared object file, and list the shared object on the link line, as shown here:

```
% nvfortran -o myprog myprog.f tobeshared.so
```

### 4. Make the executable available.

You now have an executable `myprog` which does not include any code from functions or subroutines in `tobeshared.so`, but which can be executed and dynamically linked to that code. By default, when the program is linked to produce `myprog`, no assumptions are made on the location of `tobeshared.so`. Therefore, for `myprog` to execute correctly, you must initialize the environment variable `LD_LIBRARY_PATH` to include the directory containing `tobeshared.so`. If `LD_LIBRARY_PATH` is already initialized, it is important not to overwrite its contents. If you have placed `tobeshared.so` in directory `/home/myusername/bin`, you can initialize `LD_LIBRARY_PATH` to include that directory and preserve its existing contents, as shown in the following:

```
% setenv LD_LIBRARY_PATH "$LD_LIBRARY_PATH":/home/myusername/bin
```

If you know that `tobeshared.so` always resides in a specific directory, you can create the executable `myprog` in a form that assumes this directory by using the `-R` link-time option. For example, you can link as follows:

```
% nvfortran -o myprog myprog.f tobeshared.so -R/home/myusername/bin
```



As with the `-L` option, there is no space between `-R` and the directory name. If the `-R` option is used, it is not necessary to initialize `LD_LIBRARY_PATH`.

In the previous example, the dynamic linker always looks in `/home/myusername/bin` to resolve references to `tobeshared.so`. By default, if the `LD_LIBRARY_PATH` environment variable is not set, the linker only searches `/usr/lib` and `/lib` for shared objects.

## 11.3.2. ldd Command

The `ldd` command is a useful tool when working with shared object files and executables that reference them. When applied to an executable, as shown in the following example, `ldd` lists all shared object files referenced in the executable along with the pathname of the directory from which they will be extracted.

```
% ldd myprog
```

If the pathname is not hard-coded using the `-R` option, and if `LD_LIBRARY_PATH` is not initialized, the pathname is listed as "not found". For more information on `ldd`, its options and usage, see the online man page for `ldd`.

## 11.4. Using LIB3F

The NVFORTRAN compiler includes support for the de facto standard LIB3F library routines. See the Fortran Language Reference manual for a complete list of available routines in the NVIDIA implementation of LIB3F.

## 11.5. LAPACK, BLAS and FFTs

The NVIDIA HPC SDK includes a BLAS and LAPACK library based on the customized OpenBLAS project source and built with the NVIDIA HPC Compilers. The LAPACK library is called `liblapack.a`. The BLAS library is called `libblas.a`.

To use these libraries, simply link them in using the `-l` option when linking your main program:

```
% nvfortran myprog.f -llapack -lblas
```

## 11.6. Linking with ScaLAPACK

The ScaLAPACK libraries are automatically installed with each MPI library version which accompanies an NVIDIA HPC SDK installation. You can link with the ScaLAPACK libraries by specifying `-Mscalapack` on any of the MPI wrapper command lines. For example:

```
% mpifort myprog.f -Mscalapack
```

A pre-built version of the BLAS library is automatically added when the `-Mscalapack` switch is specified. If you wish to use a different BLAS library, and still use the `-Mscalapack` switch, then you can list the set of libraries explicitly on your link line.

If the `-Mnvpl` switch is also specified in addition to `-Mscalapack`, then the NVPL ScaLAPACK library will be used.

## 11.7. The C++ Standard Template Library

On Linux, the GNU-compatible `nvc++` compiler uses the GNU g++ header files and Standard Template Library (STL) directly. The versions used are dependent on the version of the GNU compilers installed on your system, or specified when `makelocalrc` was run during installation of the NVIDIA HPC Compilers.



## 11.8. NVIDIA Performance Libraries (NVPL)

The NVIDIA Performance Libraries (NVPL) are a suite of high performance mathematical libraries optimized for the NVIDIA Grace Arm architecture. These CPU-only libraries have no dependencies on CUDA or CTK, and are drop in replacements for standard C and Fortran mathematical APIs allowing HPC applications to achieve maximum performance on the Grace platform. They are available for Arm CPUs only. The NVPL includes the following math libraries: BLAS, FFT, LAPACK, RAND, ScaLAPACK, Sparse, and Tensor. Refer to the [NVPL documentation](#) for more information about these math libraries. The following section explains how to use them with the NVHPC compilers.

To use the NVPL libraries, use the `-Mnvpl` option when linking your main program:

```
% nvfortran myprog.f -Mnvpl
```

You can link only the NVPL libraries your application needs using the sub-options to `-Mnvpl`. For example, if you only want the BLAS and FFT libraries from the NVPL, link as follows:

```
% nvfortran myprog.f -Mnvpl=blas,fft
```

Refer to the NVIDIA HPC Compilers Reference Guide for a complete list of supported options for the `-Mnvpl` flag.

### ScaLAPACK

Similar to other ScaLAPACK libraries, the NVPL version is designed to be used with MPI. A straightforward way to access the NVPL ScaLAPACK library is to use an MPI wrapper (i.e., `mpicc`, `mpic++`, `mpifort`) and link with both `-Mnvpl` and `-Mscalapack`. For example:

```
% mpic++ myprog.cpp -Mscalapack -Mnvpl
```

If you choose not to use an MPI wrapper, you can satisfy ScaLAPACK's dependency on `libmpi.so` by explicitly providing this library at link time.

The NVPL ScaLAPACK interfaces are available for the following MPI variants: MPICH, Open MPI 3.x, Open MPI 4.x (including HPC-X), and Open MPI 5.x. The HPC SDK contains builds of Open MPI 3, Open MPI 4, and HPC-X; to take advantage of the NVPL's ScaLAPACK interfaces for MPICH or Open MPI 5.x, you must supply your own build of these MPI libraries.

## 11.9. Linking with the nvmalloc Library

The NVIDIA HPC SDK installation includes a custom host (system) memory allocation library based on the jemalloc memory allocator. This library replaces the system `malloc()`, `free()`, and other related functions used by the `nvc`, `nvc++`, and `nvfortran`

runtime for dynamic heap allocations. You can link with this library by specifying `-nvmmalloc` on any of the compiler command lines used for linking. For example:

```
% nvc main.c -nvmmalloc
```

# Chapter 12.

## ENVIRONMENT VARIABLES

Environment variables allow you to set and pass information that can alter the default behavior of the NVIDIA HPC compilers and the executables which they generate. This section includes explanations of the environment variables specific to the NVIDIA HPC Compilers. .

- ▶ Standard OpenMP environment variables are used to control the behavior of OpenMP programs; these environment variables are described in the OpenMP Specification available online.
- ▶ Several NVIDIA-specific environment variables can be used to control the behavior of OpenACC programs. OpenACC-related environment variables are described in the OpenACC section: [Environment Variables](#) and the [OpenACC Getting Started Guide](#), [docs.nvidia.com/hpc-sdk/compilers/pdf/hpc24%RELEASE\\_VERSION\\_MINORopenacc\\_gs.pdf](https://docs.nvidia.com/hpc-sdk/compilers/pdf/hpc24%RELEASE_VERSION_MINORopenacc_gs.pdf).

### 12.1. Setting Environment Variables

Before we look at the environment variables that you might use with the HPC compilers and tools, let's take a look at how to set environment variables. To illustrate how to set these variables in various environments, let's look at how a user might initialize a Linux shell environment to enable use of the NVIDIA HPC Compilers.

#### 12.1.1. Setting Environment Variables on Linux

Let's assume that you want access to the NVIDIA products when you log in, and that you installed the NVIDIA HPC SDK in `/opt/nvidia/hpc_sdk`. For access at startup, you can add the following lines to your shell startup files on a Linux\_x86\_64 system.

**For csh, use these commands:**

```
% setenv NVHPCSDK /opt/nvidia/hpc_sdk
% setenv MANPATH "$MANPATH":$NVHPCSDK/Linux_x86_64/24.9/compilers/man
% set path = ($NVHPCSDK/Linux_x86_64/24.9/compilers/bin $path)
```

**For bash, sh, zsh, or ksh, use these commands:**

```
$ NVHPCSDK=/opt/nvidia/hpc_sdk; export NVHPCSDK
$ MANPATH=$MANPATH:$NVHPCSDK/Linux_x86_64/24.9/compilers/man; export MANPATH
$ PATH=$NVHPCSDK/Linux_x86_64/24.9/compilers/bin:$PATH; export PATH
```

On a Linux/OpenPOWER system replace **Linux\_x86\_64** with **Linux\_ppc64le**, and on a Linux/Arm Server system replace it with **Linux\_aarch64**.

## 12.2. HPC Compiler Related Environment Variables

The following table provides a listing of environment variables that affect the behavior of the NVIDIA HPC Compilers and the executables they generate.

Table 25 NVIDIA HPC Compilers Environment Variable Summary

Environment Variable	Description
FORTRANOPT	Allows the user to specify that the NVIDIA Fortran compiler should use VAX I/O or other custom I/O conventions.
FORT_FMT_RECL	Allows the user to change the default Fortran stdout (unit 6) line length before a line break occurs. Default: 80 bytes.
GMON_OUT_PREFIX	Specifies the name of the output file for programs that are compiled and linked with the -pg option.
LD_LIBRARY_PATH	Specifies a colon-separated set of directories where libraries should first be searched, prior to searching the standard set of directories.
MANPATH	Sets the directories that are searched for manual pages associated with the command that the user types.
NO_STOP_MESSAGE	If used, the execution of a plain STOP statement does not produce the message <code>FORTRAN STOP</code> .
PATH	Determines which locations are searched for commands the user may type.
NVCOMPILER_FPU_STATE	Manages the initial state of the processor's floating point control and status register at program startup.
NVCOMPILER_TERM	Controls the stack traceback and just-in-time debugging functionality.
NVCOMPILER_TERM_DEBUG	Overrides the default behavior when <code>NVCOMPILER_TERM</code> is set to <code>debug</code> .
PWD	Allows you to display the current directory.
STATIC_RANDOM_SEED	Forces the seed returned by <code>RANDOM_SEED</code> to be constant.
TMP	Sets the directory to use for temporary files created during execution of the HPC compilers and tools; interchangeable with <code>TMPDIR</code> .
TMPDIR	Sets the directory to use for temporary files created during execution of the HPC compilers and tools.

## 12.3. HPC Compilers Environment Variables

Use the environment variables listed in [Table 25](#) to alter the default behavior of the NVIDIA HPC Compilers and the executables which they generate. This section provides more detailed descriptions about the variables in this table.

### 12.3.1. FORTRANOPT

FORTRANOPT allows the user to adjust the behavior of the NVIDIA Fortran compiler.

- ▶ If FORTRANOPT exists and contains the value `vaxio`, the record length in the open statement is in units of 4-byte words, and the \$ edit descriptor only has an effect for lines beginning with a space or a plus sign (+).
- ▶ If FORTRANOPT exists and contains the value `format_relaxed`, an I/O item corresponding to a numerical edit descriptor (such as F, E, I, and so on) is not required to be a type implied by the descriptor.
- ▶ If FORTRANOPT exists and contains the value `no_minus_zero`, an I/O item corresponding to a numerical edit descriptor (such as F, E, I, and so on) equal to negative zero will be output as if it were positive zero.
- ▶ If FORTRANOPT exists and contains the value `crif`, a sequential formatted or list-directed record is allowed to be terminated with the character sequence `\r\n` (carriage return, newline). This approach is useful when reading records from a file produced on a Windows system.

The following example causes the NVIDIA Fortran compiler to use VAX I/O conventions:

```
% setenv FORTRANOPT vaxio
```

### 12.3.2. FORT\_FMT\_RECL

The `FORT_FMT_RECL` environment variable specifies the maximum line in bytes for Fortran formatted output to standard out (unit 6) before a newline will be generated.

If the environment variable `FORT_FMT_RECL` is present, the Fortran runtime library will use the value specified as the number of bytes to output before a newline is generated.

The default value of `FORT_FMT_RECL` is 80.

- ▶ In csh:

```
% setenv FORT_FMT_RECL length-in-bytes
```

- ▶ In bash, sh, zsh, or ksh:

```
$ FORT_FMT_RECL=length-in-bytes
$ export FORT_FMT_RECL
```

### 12.3.3. GMON\_OUT\_PREFIX

`GMON_OUT_PREFIX` specifies the name of the output file for programs that are compiled and linked with the `-pg` option. The default name is `gmon.out`.

If `GMON_OUT_PREFIX` is set, the name of the output file has `GMON_OUT_PREFIX` as a prefix. Further, the suffix is the pid of the running process. The prefix and suffix are separated by a dot. For example, if the output file is `mygmon`, then the full filename may look something similar to this: `mygmon.0012348567`.

The following example causes the NVIDIA Fortran compiler to use `nvout` as the output file for programs compiled and linked with the `-pg` option.

```
% setenv GMON_OUT_PREFIX nvout
```

### 12.3.4. LD\_LIBRARY\_PATH

The `LD_LIBRARY_PATH` variable is a colon-separated set of directories specifying where libraries should first be searched, prior to searching the standard set of directories. This variable is useful when debugging a new library or using a nonstandard library for special purposes.

The following `csch` example adds the current directory to your `LD_LIBRARY_PATH` variable.

```
% setenv LD_LIBRARY_PATH "$LD_LIBRARY_PATH":./"
```

### 12.3.5. MANPATH

The `MANPATH` variable sets the directories that are searched for manual pages associated with the commands that the user types. When using NVIDIA HPC Compilers, it is important that you set your `PATH` to include the location of the compilers and then set the `MANPATH` variable to include the man pages associated with the products.

The following `csch` example targets the `Linux_x86_64` version of the compilers and enables access to the manual pages associated with them. The settings are similar for `Linux_ppc64le` or `Linux_aarch64` targets:

```
% set path = (/opt/nvidia/hpc_sdk/Linux_x86_64/24.9/compilers/bin $path)
% setenv MANPATH "$MANPATH":/opt/nvidia/hpc_sdk/Linux_x86_64/24.9/compilers/man
```

### 12.3.6. NO\_STOP\_MESSAGE

If the `NO_STOP_MESSAGE` variable exists, the execution of a plain `STOP` statement does not produce the message `FORTTRAN STOP`. The default behavior of the NVIDIA Fortran compiler is to issue this message.

### 12.3.7. PATH

The `PATH` variable determines the directories that are searched for commands that the user types. When using the NVIDIA HPC compilers, it is important that you set your `PATH` to include the location of the compilers.

The following `csch` example initializes path settings to use the `Linux_x86_64` versions of the NVIDIA HPC Compilers. Settings for `Linux_ppc64le` and `Linux_aarch64` are done similarly:

```
% set path = (/opt/nvidia/hpc_sdk/Linux_x86_64/24.9/compilers/bin $path)
```

### 12.3.8. NVCOMPILER\_FPU\_STATE

The `NVCOMPILER_FPU_STATE` environment variable manages the initial state of the processor's floating point control and status register. `NVCOMPILER_FPU_STATE` eliminates the need to compile the main entry point (c/c++/Fortran) of programs with -

`M[no]daz`, `-M[no]flushz`, or `-Ktrap=` command line options, as those options can now be specified at runtime.



#### Linux only

If the environment variable `NVCOMPILER_FPU_STATE` is present, all settings from the command line options `-M[no]daz`, `-M[no]flushz`, or `-Ktrap=` are ignored and the FPU is initialized according to the options specified. `NVCOMPILER_FPU_STATE` with no options resets the floating-point control and status register to the system defaults.

The value of `NVCOMPILER_FPU_STATE` is a comma-separated list of options. The commands for setting the environment variable follow.

► In `csh`:

```
% setenv NVCOMPILER_FPU_STATE option[,option...]
```

► In `bash`, `sh`, `zsh`, or `ksh`:

```
$ NVCOMPILER_FPU_STATE=option[,option...]
$ export NVCOMPILER_FPU_STATE
```

Table 26 lists the supported values for `option`.

By default, these options are taken from the compiler command line options `-M[no]daz`, `-M[no]flushz`, and `-Ktrap=`.

Table 26 Supported `NVCOMPILER_FPU_STATE` options

<code>fp</code>	Shorthand for <code>inv,divz,ovf</code>
<code>inv</code>	Raise exception on floating-point invalid operation (infinity - infinity, infinity / infinity, 0 / 0, ...)
<code>invalid</code>	Alias for <code>inv</code>
<code>denorm</code>	Raise exception with floating-point denormalized operands (x86_64 only)
<code>divz</code>	Raise exception on floating-point divide-by-zero
<code>zero</code>	Alias for <code>divz</code>
<code>ovf</code>	Raise exception on floating-point overflow in result
<code>overflow</code>	Alias for <code>ovf</code>
<code>unf</code>	Raise exception on floating-point underflow in result
<code>underflow</code>	Alias for <code>unf</code>
<code>inexact</code>	Raise exception on floating-point inexact result
<code>daz</code>	Convert denormal source operands to zero
<code>nodaz</code>	Do not convert denormal source operands to zero
<code>ftz</code>	Flush underflow results to zero
<code>flushz</code>	Alias for <code>ftz</code>
<code>noftz</code>	Do not flush underflow results to zero
<code>noflushz</code>	Alias for <code>noftz</code>

print	Print to stderr the state of floating point control and status register before and after processing of environment variable NVCOMPILER_FPU_STATE
debug	Alias for print

### 12.3.9. NVCOMPILER\_TERM

The NVCOMPILER\_TERM environment variable controls the stack traceback and just-in-time debugging functionality. The runtime libraries use the value of NVCOMPILER\_TERM to determine what action to take when a program abnormally terminates.

The value of NVCOMPILER\_TERM is a comma-separated list of options. The commands for setting the environment variable follow.

- In csh:

```
% setenv NVCOMPILER_TERM option[,option...]
```

- In bash, sh, zsh, or ksh:

```
$ NVCOMPILER_TERM=option[,option...]
$ export NVCOMPILER_TERM
```

Table 27 lists the supported values for option. Following the table is a complete description of each option that indicates specifically how you might apply the option.

By default, all of these options are disabled.

Table 27 Supported NVCOMPILER\_TERM Values

[no]debug	Enables/disables just-in-time debugging (debugging invoked on error)
[no]trace	Enables/disables stack traceback on error
[no]trace-fp	Enables/disables stack traceback and printing of SIMD registers (ymm/zmm) on error (Linux x86_64 only)
[no]signal	Enables/disables establishment of signal handlers for common signals that cause program termination
[no]abort	Enables/disables calling the system termination routine abort()

[no]debug

This enables/disables just-in-time debugging. The default is nodebug.

When NVCOMPILER\_TERM is set to debug, the command to which NVCOMPILER\_TERM\_DEBUG is set is invoked on error.

[no]trace

This enables/disables stack traceback on error.

[no]trace-fp

This enables/disables stack traceback and printing of SIMD registers (ymm/zmm) on error. (Linux x86\_64 only)



[no]signal

This enables/disables establishing signal handlers for the most common signals that cause program termination. The default is `nosignal`. Setting `trace` and `debug` automatically enables `signal`. Specifically setting `nosignal` allows you to override this behavior.

[no]abort

This enables/disables calling the system termination routine `abort()`. The default is `noabort`. When `noabort` is in effect the process terminates by calling `_exit(127)`.

On Linux, when `abort` is in effect, the `abort` routine creates a core file and exits with code 127.

A few runtime errors just print an error message and call `exit(127)`, regardless of the status of `NVCOMPILER_TERM`. These are mainly errors such as specifying an invalid environment variable value where a traceback would not be useful.

If it appears that `abort()` does not generate core files on a Linux system, be sure to `unlimit` the `coredumpsize`. You can do this in these ways:

- ▶ Using `csh`:

```
% limit coredumpsize unlimited
% setenv NVCOMPILER_TERM abort
```

- ▶ Using `bash`, `sh`, `zsh`, or `ksh`:

```
$ ulimit -c unlimited
$ export NVCOMPILER_TERM=abort
```

To debug a core file with `gdb`, invoke `gdb` with the `--core` option. For example, to view a core file named "core" for a program named "a.out":

```
$ gdb --core=core a.out
```

For more information on why to use this variable, refer to [Stack Traceback and JIT Debugging](#).

## 12.3.10. NVCOMPILER\_TERM\_DEBUG

The `NVCOMPILER_TERM_DEBUG` variable may be set to override the default behavior when `NVCOMPILER_TERM` is set to `debug`.

The value of `NVCOMPILER_TERM_DEBUG` should be set to the command line used to invoke the program. For example:

```
gdb --quiet --pid %d
```

The first occurrence of `%d` in the `NVCOMPILER_TERM_DEBUG` string is replaced by the process id. The program named in the `NVCOMPILER_TERM_DEBUG` string must be found on the current `PATH` or specified with a full path name.

### 12.3.11. PWD

The PWD variable allows you to display the current directory.

### 12.3.12. STATIC\_RANDOM\_SEED

You can use `STATIC_RANDOM_SEED` to force the seed returned by the Fortran 90/95 `RANDOM_SEED` intrinsic to be constant. The first call to `RANDOM_SEED` without arguments resets the random seed to a default value, then advances the seed by a variable amount based on time. Subsequent calls to `RANDOM_SEED` without arguments reset the random seed to the same initial value as the first call. Unless the time is exactly the same, each time a program is run a different random number sequence is generated. Setting the environment variable `STATIC_RANDOM_SEED` to `YES` forces the seed returned by `RANDOM_SEED` to be constant, thereby generating the same sequence of random numbers at each execution of the program.

### 12.3.13. TMP

You can use `TMP` to specify the directory to use for placement of any temporary files created during execution of the NVIDIA HPC Compilers. This variable is interchangeable with `TMPPDIR`.

### 12.3.14. TMPPDIR

You can use `TMPPDIR` to specify the directory to use for placement of any temporary files created during execution of the NVIDIA HPC Compilers.

## 12.4. Using Environment Modules on Linux

On Linux, if you use the Environment Modules package, that is, the `module load` command, the NVIDIA HPC Compilers include a script to set up the appropriate module files. The install script will generate environment module files for you as part of the set up process.

Assuming your installation base directory is `/opt/nvidia/hpc_sdk`, the environment modules will be installed under `/opt/nvidia/hpc_sdk/modulefiles`. There will be three sets of module files:

1. `nvhpc`  
Adds environment variable settings for the NVIDIA HPC Compilers, CUDA libraries, and additional libraries such as MPI, NCCL, and NVSHMEM.
2. `nvhpc-nompi`  
Adds environment variable settings for the NVIDIA HPC Compilers, CUDA libraries, and additional libraries such as NCCL and NVSHMEM. This will not include MPI, if you wish to use an alternate MPI implementation.
3. `nvhpc-byo-compilers`

Adds environment variable settings for the CUDA libraries and additional libraries such as NCCL and NVSHMEM. This will not include the NVIDIA HPC Compilers nor MPI, if you wish to use alternate compilers and MPI.

You can load the `nvhpc` environment module for the 20.11 release as follows:

```
% module load nvhpc/24.9
```

To see what versions of `nvhpc` are available on this system, use this command:

```
% module avail nvhpc
```

The `module load` command sets or modifies the environment variables as indicated in the following table.

This Environment Variable...	Is set or modified by the module load command
<b>CC</b>	Full path to <code>nvcc</code> (nvhpc and nvhpc-nompi only)
<b>CPATH</b>	Prepends the math libraries include directory, the MPI include directory (nvhpc only), and NCCL and NVSHMEM include directories
<b>CPP</b>	C preprocessor, normally <code>cpp</code> (nvhpc and nvhpc-nompi only)
<b>CXX</b>	Path to <code>nvcc++</code> (nvhpc and nvhpc-nompi only)
<b>FC</b>	Full path to <code>nvfortran</code> (nvhpc and nvhpc-nompi only)
<b>F90</b>	Full path to <code>nvfortran</code> (nvhpc and nvhpc-nompi only)
<b>F77</b>	Full path to <code>nvfortran</code> (nvhpc and nvhpc-nompi only)
<b>LD_LIBRARY_PATH</b>	Prepends the CUDA library directory, the NVIDIA HPC Compilers library directory (nvhpc and nvhpc-nompi only), math libraries library directory, MPI library directory (nvhpc only), and NCCL and NVSHMEM library directories
<b>MANPATH</b>	Prepends the NVIDIA HPC Compilers man page directory (nvhpc and nvhpc-nompi only)
<b>OPAL_PREFIX</b>	Full path to the MPI directory (nvhpc only), e.g. <code>/opt/nvidia/hpc_sdk/Linux_x86_64/24.9/comm_libs/mpi</code>
<b>PATH</b>	Prepends the CUDA bin directory, the MPI bin directory (nvhpc only), and the NVIDIA HPC Compilers bin directory (nvhpc and nvhpc-nompi only)



NVIDIA does not provide support for the Environment Modules package. For more information about the package, go to: <http://modules.sourceforge.net>.

## 12.5. Stack Traceback and JIT Debugging

When a programming error results in a runtime error message or an application exception, a program will usually exit, perhaps with an error message. The NVIDIA

HPC Compilers runtime library includes a mechanism to override this default action and instead print a stack traceback, start a debugger, or, on Linux, create a core file for post-mortem debugging.

The stack traceback and just-in-time debugging functionality is controlled by an environment variable, `NVCOMPILER_TERM`, described in [NVCOMPILER\\_TERM](#). The runtime libraries use the value of `NVCOMPILER_TERM` to determine what action to take when a program abnormally terminates.

When the NVIDIA HPC Compilers runtime library detects an error or catches a signal, it calls the routine `nvcompiler_stop_here()` prior to generating a stack traceback or starting the debugger. The `nvcompiler_stop_here()` routine is a convenient spot to set a breakpoint when debugging a program.

# Chapter 13.

## DISTRIBUTING FILES – DEPLOYMENT

Once you have successfully built, debugged and tuned your application, you may want to distribute it to users who need to run it on a variety of systems. This section addresses how to effectively distribute applications built using NVIDIA HPC Compilers. The application must be installed in such a way that it executes accurately on a system other than the one on which it was built, and which may be configured differently.

### 13.1. Deploying Applications on Linux

To successfully deploy your application on Linux, some of the issues to consider include:

- ▶ Runtime Libraries
- ▶ 64-bit Linux Systems
- ▶ Redistribution of Files

#### 13.1.1. Runtime Library Considerations

On Linux systems, the system runtime libraries can be linked to an application either statically or dynamically. For example, for the C runtime library, `libc`, you can use either the static version `libc.a` or the shared object version `libc.so`. If the application is intended to run on Linux systems other than the one on which it was built, it is generally safer to use the shared object version of the library. This approach ensures that the application uses a version of the library that is compatible with the system on which the application is running. Further, it works best when the application is linked on a system that has an equivalent or earlier version of the system software than the system on which the application will be run.



Building on a newer system and running the application on an older system may not produce the desired output.

To use the shared object version of a library, the application must also link to shared object versions of the NVIDIA HPC Compilers runtime libraries. To execute an application built in such a way on a system on which NVIDIA HPC Compilers are

*not* installed, those shared objects must be available. To build using the shared object versions of the runtime libraries, use the `-Bdynamic` option, as shown here:

```
$ nvfortran -Bdynamic myprog.f90
```

### 13.1.2. 64-bit Linux Considerations

On 64-bit Linux systems, 64-bit applications that use the `-mcmodel=medium` option sometimes cannot be successfully linked statically. Therefore, users with executables built with the `-mcmodel=medium` option may need to use shared libraries, linking dynamically. Also, runtime libraries built using the `-fpic` option use 32-bit offsets, so they sometimes need to reside near other runtime `libs` in a shared area of Linux program memory.



If your application is linked dynamically using shared objects, then the shared object versions of the NVIDIA HPC Compilers runtime are required.

### 13.1.3. Linux Redistributable Files

The method for installing the shared object versions of the runtime libraries required for applications built with NVIDIA HPC Compilers is manual distribution.

When the NVIDIA HPC Compilers are installed, there are directories that have a name that begins with `REDIST`; these directories contain the redistributed shared object libraries. These may be redistributed by licensed NVIDIA HPC Compilers users under the terms of the End-User License Agreement.

### 13.1.4. Restrictions on Linux Portability

You cannot expect to be able to run an executable on any given Linux machine. Portability depends on the system you build on as well as how much your program uses system routines that may have changed from Linux release to Linux release. For example, an area of significant change between some versions of Linux is in `libpthread.so` and `libnuma.so`. NVIDIA HPC Compilers use these dynamically linked libraries for the options `-acc` (OpenACC), `-mp` (OpenMP) and `-Mconcur` (multicore auto-parallel). Statically linking these libraries may not be possible, or may result in failure at execution.

Typically, portability is supported for forward execution, meaning running a program on the same or a later version of Linux. But not for backward compatibility, that is, running on a prior release. For example, a user who compiles and links a program under RHEL 7.2 should not expect the program to run without incident on a RHEL 5.2 system, an earlier Linux version. It *may* run, but it is less likely. Developers might consider building applications on earlier Linux versions for wider usage. Dynamic linking of Linux and gcc system routines on the platform executing the program can also reduce problems.

### 13.1.5. Licensing for Redistributable (REDIST) Files

The files in the REDIST directories may be redistributed under the terms of the End-User License Agreement for the product in which they were included.

# Chapter 14.

## INTER-LANGUAGE CALLING

This section describes inter-language calling conventions for C, C++, and Fortran programs using the HPC compilers. Fortran 2003 ISO\_C\_Binding provides a mechanism to support the interoperability with C. This includes the `iso_c_binding` intrinsic module, binding labels, and the `BIND` attribute. Additional interoperability with C is available with Fortran 2018 and the `ISO_Fortran_binding.h` C header file. `nvfortran` supports both the `iso_c_binding` and the `ISO_Fortran_Binding.h` header file. In the absence of these mechanisms, the following sections describe how to call a Fortran function or subroutine from a C or C++ program and how to call a C or C++ function from a Fortran program.

This section provides examples that use the following options related to inter-language calling.

`-c`                      `-Mnomain`                      `-Miface`                      `-Mupcase`

### 14.1. Overview of Calling Conventions

This section includes information on the following topics:

- ▶ Functions and subroutines in Fortran, C, and C++
- ▶ Naming and case conversion conventions
- ▶ Compatible data types
- ▶ Argument passing and special return values
- ▶ Arrays and indexes

The sections [Inter-language Calling Considerations](#) through [Example – C++ Calling Fortran](#) describe how to perform inter-language calling using the Linux or Win64 convention.

### 14.2. Inter-language Calling Considerations

In general, when argument data types and function return values agree, you can call a C or C++ function from Fortran as well as call a Fortran function from C or C++. When



data types for arguments do not agree, you may need to develop custom mechanisms to handle them. For example, the Fortran `COMPLEX` type has a matching type in C99 but does not have a matching type in C89; however, it is still possible to provide inter-language calls but there are no general calling conventions for such cases.



- ▶ If a C++ function contains objects with constructors and destructors, calling such a function from either C or Fortran is not possible unless the initialization in the main program is performed from a C++ program in which constructors and destructors are properly initialized.
- ▶ In general, you can call a C or Fortran function from C++ without problems as long as you use the `extern "C"` keyword to declare the function in the C++ program. This declaration prevents name mangling for the C function name. If you want to call a C++ function from C or Fortran, you also have to use the `extern "C"` keyword to declare the C++ function. This keeps the C++ compiler from mangling the name of the function.
- ▶ You can use the `__cplusplus` macro to allow a program or header file to work for both C and C++. For example, the following defines in the header file `stdio.h` allow this file to work for both C and C++.

```
#ifndef _STDIO_H
#define _STDIO_H
#ifdef __cplusplus
extern "C" {
#endif /* __cplusplus */
.
. /* Functions and data types defined... */
.
#ifdef __cplusplus
}
#endif /* __cplusplus */
#endif
```

- ▶ C++ member functions cannot be declared `extern`, since their names will always be mangled. Therefore, C++ member functions cannot be called from C or Fortran.

## 14.3. Functions and Subroutines

Fortran, C, and C++ define functions and subroutines differently.

For a Fortran program calling a C or C++ function, observe the following return value convention:

- ▶ When a C or C++ function returns a value, call it from Fortran as a function.
- ▶ When a C or C++ function does not return a value, call it as a subroutine.

For a C/C++ program calling a Fortran function, the call should return a similar type. [Table 28, Fortran and C/C++ Data Type Compatibility](#), lists compatible types. If the call is to a Fortran subroutine, or a Fortran `CHARACTER` function, or a Fortran `COMPLEX` function, call it from C/C++ as a function that returns void. The exception to this convention is when a Fortran subroutine has alternate returns; call such a subroutine from C/C++ as a function returning `int` whose value is the value of the integer expression specified in the alternate `RETURN` statement.

## 14.4. Upper and Lower Case Conventions, Underscores

By default on Linux and Win64 systems, all Fortran symbol names are converted to lower case. C and C++ are case sensitive, so upper-case function names stay upper-case. When you use inter-language calling, you can either name your C/C++ functions with lower-case names, or invoke the Fortran compiler command with the option `-Mupcase`, in which case it will not convert symbol names to lower-case.

When programs are compiled using one of the HPC Fortran compilers on Linux and Win64 systems, an underscore is appended to Fortran global names (names of functions, subroutines and common blocks). This mechanism distinguishes Fortran name space from C/C++ name space. Use these naming conventions:

- ▶ If you call a C/C++ function from Fortran, you should rename the C/C++ function by appending an underscore or use `bind(c)` in the Fortran program.
- ▶ If you call a Fortran function from C/C++, you should append an underscore to the Fortran function name in the calling program.

## 14.5. Compatible Data Types

Table 28 shows compatible data types between Fortran and C/C++. Table 29, [Fortran and C/C++ Representation of the COMPLEX Type](#) shows how the Fortran COMPLEX type may be represented in C/C++.



Tip If you can make your function/subroutine parameters as well as your return values match types, you should be able to use inter-language calling.

Table 28 Fortran and C/C++ Data Type Compatibility

Fortran Type (lower case)	C/C++ Type	Size (bytes)
character x	char x	1
character*n x	char x[n]	n
real x	float x	4
real*4 x	float x	4
real*8 x	double x	8
double precision	double x	8
integer x	int x	4
integer*1 x	signed char x	1
integer*2 x	short x	2
integer*4 x	int x	4

Fortran Type (lower case)	C/C++ Type	Size (bytes)
integer*8 x	long long x	8
logical x	int x	4
logical*1 x	char x	1
logical*2 x	short x	2
logical*4	int x	4
logical*8	long x	8

Table 29 Fortran and C/C++ Representation of the `COMPLEX` Type

Fortran Type (lower case)	C/C++ Type	Size (bytes)
complex x	struct {float r,i;} x; float complex x;	8 8
complex*8 x	struct {float r,i;} x; float complex x;	8 8
double complex x	struct {double dr,di;} x; double complex x;	16 16
complex *16 x	struct {double dr,di;} x; double complex x;	16 16



For C/C++, the `complex` type implies C99 or later.

### 14.5.1. Fortran Named Common Blocks

A named Fortran common block can be represented in C/C++ by a structure whose members correspond to the members of the common block. The name of the structure in C/C++ must have the added underscore. For example, here is a Fortran common block:

```
INTEGER I
COMPLEX C
DOUBLE COMPLEX CD
DOUBLE PRECISION D
COMMON /COM/ i, c, cd, d
```

This Fortran Common Block is represented in C with the following equivalent:

```
extern struct {
    int i;
    struct {float real, imag;} c;
    struct {double real, imag;} cd;
    double d;
} com_;
```

This same Fortran Common Block is represented in C++ with the following equivalent:

```
extern "C" struct {
    int i;
    struct {float real, imag;} c;
```

```
struct {double real, imag;} cd;
double d;
} com_;
```



Tip For global or external data sharing, `extern "C"` is not required.

## 14.6. Argument Passing and Return Values

In Fortran, arguments are passed by reference, that is, the address of the argument is passed, rather than the argument itself. In C/C++, arguments are passed by value, except for strings and arrays, which are passed by reference. Due to the flexibility provided in C/C++, you can work around these differences. Solving the parameter passing differences generally involves intelligent use of the `&` and `*` operators in argument passing when C/C++ calls Fortran and in argument declarations when Fortran calls C/C++.

For strings declared in Fortran as type `CHARACTER`, an argument representing the length of the string is also passed to a calling function.

On the following systems, the compiler places the length argument(s) at the end of the parameter list, following the other formal arguments:.

- ▶ On Linux systems
- ▶ On Win64 systems, except when using the option `-Miface=ceref`

The length argument is passed by value, not by reference.

### 14.6.1. Passing by Value (%VAL)

When passing parameters from a Fortran subprogram to a C/C++ function, it is possible to pass by value using the `%VAL` function. If you enclose a Fortran parameter with `%VAL()`, the parameter is passed by value. For example, the following call passes the integer `i` and the logical `bvar` by value.

```
integer*1 i
logical*1 bvar
call cvalue (%VAL(i), %VAL(bvar))
```

### 14.6.2. Character Return Values

[Functions and Subroutines](#) describes the general rules for return values for C/C++ and Fortran inter-language calling. There is a special return value to consider. When a Fortran function returns a character, two arguments need to be added at the beginning of the C/C++ calling function's argument list:

- ▶ The address of the return character or characters
- ▶ The length of the return character

The following example illustrates the extra parameters, `tmp` and `10`, supplied by the caller:

## Character Return Parameters

```
! Fortran function returns a character
CHARACTER*(*) FUNCTION CHF(C1,I)
  CHARACTER*(*) C1
  INTEGER I
END
```

```
/* C declaration of Fortran function */
extern void chf_();
char tmp[10];
char c1[9];
int i;
chf_(tmp, 10, c1, &i, 9);
```

If the Fortran function is declared to return a character value of constant length, for example `CHARACTER*4 FUNCTION CHF()`, the second extra parameter representing the length must still be supplied, but is not used.



The value of the character function is not automatically NULL-terminated.

## 14.6.3. Complex Return Values

When a Fortran function returns a complex value, an argument needs to be added at the beginning of the C/C++ calling function's argument list; this argument is the address of the complex return value. [COMPLEX Return Values](#) illustrates the extra parameter, `cplx`, supplied by the caller.

### COMPLEX Return Values

```
COMPLEX FUNCTION CF(C, I)
  INTEGER I
  . . .
END
```

```
extern void cf_();
typedef struct {float real, imag;} cplx;
cplx c1;
int i;
cf_(&c1, &i);
```

## 14.7. Array Indices

C/C++ arrays and Fortran arrays use different default initial array index values. By default, arrays in C/C++ start at 0 and arrays in Fortran start at 1. If you adjust your array comparisons so that a Fortran second element is compared to a C/C++ first element, and adjust similarly for other elements, you should not have problems working with this difference. If this is not satisfactory, you can declare your Fortran arrays to start at zero.

Another difference between Fortran and C/C++ arrays is the storage method used. Fortran uses column-major order and C/C++ uses row-major order. For one-dimensional arrays, this poses no problems. For two-dimensional arrays, where there are an equal number of rows and columns, row and column indexes can simply be reversed. For

arrays other than single dimensional arrays, and square two-dimensional arrays, inter-language function mixing is not recommended.

## 14.8. Examples

This section contains examples that illustrate inter-language calling.

### 14.8.1. Example – Fortran Calling C



There are other solutions to calling C from Fortran than the one presented in this section. For example, you can use the `iso_c_binding` intrinsic module which NVIDIA does support. For more information on this module and for examples of how to use it, search the web using the keyword `iso_c_binding`.

C function `f2c_func_` shows a C function that is called by the Fortran main program shown in **Fortran Main Program `f2c_main.f`**. Notice that each argument is defined as a pointer, since Fortran passes by reference. Also notice that the C function name uses all lower-case and a trailing `_`.

Fortran Main Program `f2c_main.f`

```
logical*1 bool1
character letter1
integer*4 numint1, numint2
real numfloat1
double precision numdoub1
integer*2 numshor1
external f2c_func

call f2c_func(bool1, letter1, numint1, numint2, numfloat1, numdoub1, numshor1)

write( *, "(L2, A2, I5, I5, F6.1, F6.1, I5)")
+ bool1, letter1, numint1, numint2, numfloat1, numdoub1, numshor1

end
```

C function `f2c_func_`

```
#define TRUE 0xff
#define FALSE 0
void f2c_func_( bool1, letter1, numint1, numint2, numfloat1,\
 numdoub1, numshor1, len_letter1)
char *bool1, *letter1;
int *numint1, *numint2;
float *numfloat1;
double *numdoub1;
short *numshor1;
int len_letter1;
{
    *bool1 = TRUE;    *letter1 = 'v';
    *numint1 = 11;    *numint2 = -44;
    *numfloat1 = 39.6 ;
    *numdoub1 = 39.2;
    *numshor1 = 981;
}
```

Compile and execute the program `f2c_main.f` with the call to `f2c_func_` using the following command lines:

```
$ nvc -c f2c_func.c
$ nvfortran f2c_func.o f2c_main.f
```

Executing the `a.out` file should produce the following output:

```
T v 11 -44 39.6 39.2 981
```

## 14.8.2. Example – C Calling Fortran



There are other solutions to calling Fortran from C than the one presented in this section. For example, you can use the `ISO_Fortran_binding.h` C header file which NVIDIA does support. For more information on this header file and for examples of how to use it, search the web using the keyword `ISO_Fortran_binding`.

The example **C Main Program `c2f_main.c`** shows a C main program that calls the Fortran subroutine shown in **Fortran Subroutine `c2f_sub.f`**.

- ▶ Each call uses the `&` operator to pass by reference.
- ▶ The call to the Fortran subroutine uses all lower-case and a trailing `"_"`.

### C Main Program `c2f_main.c`

```
void main () {
    char bool1, letter1;
    int numint1, numint2;
    float numfloat1;
    double numdoub1;
    short numshor1;
    extern void c2f_func();
    c2f_sub (&bool1, &letter1, &numint1, &numint2, &numfloat1, &numdoub1, &numshor1, 1);
    printf(" %s %c %d %d %3.1f %.0f %d\n",
        bool1?"TRUE":"FALSE", letter1, numint1, numint2,
        numfloat1, numdoub1, numshor1);
}
```

### Fortran Subroutine `c2f_sub.f`

```
subroutine c2f_func ( bool1, letter1, numint1, numint2,
+ numfloat1, numdoub1, numshor1)
    logical*1 bool1
    character letter1
    integer numint1, numint2
    double precision numdoub1
    real numfloat1
    integer*2 numshor1

    bool1 = .true.
    letter1 = "v"
    numint1 = 11
    numint2 = -44
    numdoub1 = 902
    numfloat1 = 39.6
    numshor1 = 299
    return
end
```

To compile this Fortran subroutine and C program, use the following commands:

```
$ nvc -c c2f_main.c
$ nvfortran -Mnomain c2f_main.o c2_sub.f
```

Executing the resulting `a.out` file should produce the following output:

```
TRUE v 11 -44 39.6 902 299
```

### 14.8.3. Example – C++ Calling C

[C++ Main Program cp2c\\_main.C Calling a C Function](#) shows a C++ main program that calls the C function shown in [Simple C Function c2cp\\_func.c](#).

C++ Main Program cp2c\_main.C Calling a C Function

```
extern "C" void cp2c_func(int n, int m, int *p);
#include <iostream>
main()
{
    int a,b,c;
    a=8;
    b=2;
    c=0;
    cout << "main: a = "<<a<<" b = "<<b<<" ptr c = "<<hex<<&c<< endl;
    cp2c_func(a,b,&c);
    cout << "main: res = "<<c<<endl;
}
```

Simple C Function c2cp\_func.c

```
void cp2c_func(num1, num2, res)
int num1, num2, *res;
{
    printf("func: a = %d b = %d ptr c = %x\n",num1,num2,res);
    *res=num1/num2;
    printf("func: res = %d\n",*res);
}
```

To compile this C function and C++ main program, use the following commands:

```
$ nvc -c cp2c_func.c
$ nvc++ cp2c_main.C cp2c_func.o
```

Executing the resulting `a.out` file should produce the following output:

```
main: a = 8 b = 2 ptr c = 0xbffffb94
func: a = 8 b = 2 ptr c = bffffb94
func: res = 4
main: res = 4
```

### 14.8.4. Example – C Calling C ++

The example in [C Main Program c2cp\\_main.c Calling a C++ Function](#) shows a C main program that calls the C++ function shown in [Simple C++ Function c2cp\\_func.C with Extern C](#).

C Main Program c2cp\_main.c Calling a C++ Function

```
extern void c2cp_func(int a, int b, int *c);
#include <stdio.h>
main() {
    int a,b,c;
    a=8; b=2;
    printf("main: a = %d b = %d ptr c = %x\n",a,b,&c);
    c2cp_func(a,b,&c);
    printf("main: res = %d\n",c);
}
```



## Simple C++ Function c2cp\_func.C with Extern C

```
#include <iostream>
extern "C" void c2cp_func(int num1,int num2,int *res)
{
    cout << "func: a = "<<num1<<" b = "<<num2<<"ptr c = "<<res<<endl;
    *res=num1/num2;
    cout << "func: res = "<<res<<endl;
}
```

To compile this C function and C++ main program, use the following commands:

```
$ nvc -c c2cp_main.c
$ nvc++ c2cp_main.o c2cp_func.C
```

Executing the resulting a.out file should produce the following output:

```
main: a = 8 b = 2 ptr c = 0xbffffb94
func: a = 8 b = 2 ptr c = bffffb94
func: res = 4
main: res = 4
```



You cannot use the extern "C" form of declaration for an object's member functions.

## 14.8.5. Example – Fortran Calling C++

The Fortran main program shown in [Fortran Main Program f2cp\\_main.f](#) calling a C++ function calls the C++ function shown in [C++ function f2cp\\_func.C](#).

Notice:

- ▶ Each argument is defined as a pointer in the C++ function, since Fortran passes by reference.
- ▶ The C++ function name uses all lower-case and a trailing "\_":

Fortran Main Program f2cp\_main.f calling a C++ function

```
logical*1 bool1
character letter1
integer*4 numint1, numint2
real numfloat1
double precision numdoub1
integer*2 numshor1
external f2cpfunc
call f2cp_func (bool1, letter1, numint1,
+ numint2, numfloat1, numdoub1, numshor1)
write( *, "(L2, A2, I5, I5, F6.1, F6.1, I5)")
+ bool1, letter1, numint1, numint2, numfloat1,
+ numdoub1, numshor1
end
```

## C++ function f2cp\_func.C

```
#define TRUE 0xff
#define FALSE 0
extern "C"
{
extern void f2cp_func_ (
    char *bool1, *letter1,
    int *numint1, *numint2,
    float *numfloat1,
    double *numdoub1,
    short *numshort1,
```

```

int len_letter1)
{
*bool1 = TRUE;      *letter1 = 'v';
*numint1 = 11;      *numint2 = -44;
*numfloat1 = 39.6; *numdoub1 = 39.2;  *numshort1 = 981;
}
}

```

Assuming the Fortran program is in a file `fmain.f`, and the C++ function is in a file `cpfunc.C`, create an executable, using the following command lines:

```

$ nvc++ -c f2cp_func.C
$ nvfortran f2cp_func.o f2cp_main.f -c++libs

```

Executing the `a.out` file should produce the following output:

```
T v 11 -44 39.6 39.2 981
```

## 14.8.6. Example – C++ Calling Fortran

Fortran Subroutine `cp2f_func.f` shows a Fortran subroutine called by the C++ main program shown in C++ main program `cp2f_main.C`. Notice that each call uses the `&` operator to pass by reference. Also notice that the call to the Fortran subroutine uses all lower-case and a trailing `"_"`:

C++ main program `cp2f_main.C`

```

#include <iostream>
extern "C" { extern void cp2f_func_(char *,char *,int *,int *,
    float *,double *,short *); }
main ()
{
    char bool1, letter1;
    int numint1, numint2;
    float numfloat1;
    double numdoub1;
    short numshor1;

    cp2f_func(&bool1,&letter1,&numint1,&numint2,&numfloat1, &numdoub1,&numshor1);
    cout << " bool1 = ";
    bool1?cout << "TRUE ":cout << "FALSE "; cout <<endl;
    cout << " letter1 = " << letter1 <<endl;
    cout << " numint1 = " << numint1 <<endl;
    cout << " numint2 = " << numint2 <<endl;
    cout << " numfloat1 = " << numfloat1 <<endl;
    cout << " numdoub1 = " << numdoub1 <<endl;
    cout << " numshor1 = " << numshor1 <<endl;
}

```

Fortran Subroutine `cp2f_func.f`

```

subroutine cp2f_func ( bool1, letter1, numint1,
+ numint2, numfloat1, numdoub1, numshor1)
logical*1 bool1
character letter1
integer numint1, numint2
double precision numdoub1
real numfloat1
integer*2 numshor1
bool1 = .true. ; letter1 = "v"
numint1 = 11 ; numint2 = -44
numdoub1 = 902 ; numfloat1 = 39.6 ; numshor1 = 299
return
end

```

To compile this Fortran subroutine and C++ program, use the following command lines:

```
$ nvfortran -c cp2f_func.f  
$ nvc++ cp2f_func.o cp2f_main.C -fortranlibs
```

Executing this C++ main should produce the following output:

```
bool1 = TRUE letter1 = v numint1 = 11 numint2 = -44 numfloat1 = 39.6 numdoub1 = 902  
numshor1 = 299
```



You must explicitly link in the NVFORTRAN runtime support libraries when linking nvfortran-compiled program units into C++ or C main programs.

# Chapter 15.

## PROGRAMMING CONSIDERATIONS FOR 64-BIT ENVIRONMENTS

NVIDIA provides 64-bit compilers for 64-bit Linux operating systems running on x86-64 (Linux\_x86\_64), OpenPOWER (Linux\_ppc64) and Arm Server (Linux\_aarch64) architectures. You can use these compilers to create programs that use 64-bit memory addresses. The GNU toolchain on 64-bit Linux systems implements an option to control 32-bit vs 64-bit code generation, as described in [Large Static Data in Linux](#). This section describes the specifics of how to use the NVIDIA compilers to make use of 64-bit memory addressing.



The NVIDIA HPC compilers themselves are 64-bit applications which can only run on 64-bit CPUs running 64-bit Operating Systems.

This section describes how to use the following options related to 64-bit programming.

<code>-fPIC</code>	<code>-mmodel=medium</code>	<code>-Mlarge_arrays</code>
<code>-i8</code>	<code>-Mlargeaddressaware</code>	

### 15.1. Data Types in the 64-Bit Environment

The size of some data types can differ across 64-bit environments. This section describes the major differences.

#### 15.1.1. C++ and C Data Types

On 64-bit Linux operating systems, the size of an `int` is 4 bytes, a `long` is 8 bytes, a `long long` is 8 bytes, and a pointer is 8 bytes.

#### 15.1.2. Fortran Data Types

In Fortran, the default size of the `INTEGER` type is 4 bytes. The `-i8` compiler option may be used to make the default size of all `INTEGER` data in the program 8 bytes.

When using the `-Mlarge_arrays` option, described in [64-Bit Array Indexing](#), any 4-byte INTEGER variables that are used to index arrays are silently promoted by the compiler to 8 bytes. This promotion can lead to unexpected consequences, so 8-byte INTEGER variables are recommended for array indexing when using the option `-Mlarge_arrays`.

## 15.2. Large Static Data in Linux

64-bit Linux operating systems support two different memory models. The default model used by the NVIDIA HPC compilers on `Linux_x86_64` and `Linux_aarch64` targets is the small memory model, which can be specified using `-mcmodel=small`. This is the 32-bit model, which limits the size of code plus statically allocated data, including system and user libraries, to 2GB. The medium memory model, specified by `-mcmodel=medium`, allows combined code and static data areas (`.text` and `.bss` sections) larger than 2GB and is the default on `Linux_ppc64le` targets. The `-mcmodel=medium` option must be used on both the compile command and the link command in order to take effect.

There are implications to using `-mcmodel=medium`. The generated code requires increased addressing overhead to support the large data range. This can affect performance, though the compilers seek to minimize the added overhead through careful instruction selection and optimization.

`Linux_aarch64` does not support `-mcmodel=medium`. If the medium model is specified on the command-line, the compiler driver will automatically select the large model.

## 15.3. Large Dynamically Allocated Data

Dynamically allocated data objects in programs compiled by the NVIDIA HPC compilers can be larger than 2GB. No special compiler options are required to enable this functionality. The size of the allocation is only limited by the system. However, to correctly access dynamically allocated arrays with more than 2G elements you should use the `-Mlarge_arrays` option, described in the following section.

## 15.4. 64-Bit Array Indexing

The NVIDIA Fortran compilers provide an option, `-Mlarge_arrays`, that enables 64-bit indexing of arrays. This means that, as necessary, 64-bit INTEGER constants and variables are used to index arrays.



In the presence of `-Mlarge_arrays`, the compiler may silently promote 32-bit integers to 64 bits, which can have unexpected side effects.

On 64-bit Linux, the `-Mlarge_arrays` option also enables single static data objects larger than 2 GB. This option is the default in the presence of `-mcmodel=medium`.

## 15.5. Compiler Options for 64-bit Programming

The usual switches that apply to 64-bit programmers seeking to increase the data range of their applications are in the following table.

Table 30 64-bit Compiler Options

Option	Purpose	Considerations
-mmodel=medium	Allow for data declarations larger than 2GB. Default on Linux_ppc64le.	Linux_aarch64 does not support -mmodel=medium. If the medium model is specified on the command-line, the compiler driver will automatically select the large model.
-Mlarge_arrays	Perform all array-location-to-address calculations using 64-bit integer arithmetic.	Slightly slower execution. Is implicit with -mmodel=medium. Can be used with option -mmodel=small.
-fpic	Position independent code. Necessary for shared libraries.	Dynamic linking restricted to a 32-bit offset. External symbol references should refer to other shared lib routines, rather than the program calling them.
-i8	All INTEGER functions, data, and constants not explicitly declared INTEGER*4 are assumed to be INTEGER*8.	Users should take care to explicitly declare INTEGER functions as INTEGER*4.

The following table summarizes the limits of these programming models under the specified conditions. The compiler options you use vary by processor.

Table 31 Effects of Options on Memory and Array Sizes


Condition	Addr. Math		Max Size Gbytes		
	A	I	AS	DS	TS
64-bit addr limited by option -mmodel=small	64	32	2	2	2
-fpic incompatible with -mmodel=medium	64	32	2	2	2
Enable full support for 64-bit data addressing	64	64	>2	>2	>2

A	Address Type – size in bits of data used for address calculations, 64-bits.
I	Index Arithmetic -bit-size of data used to index into arrays and other aggregate data structures. If 32-bit, total range of any single data object is limited to 2GB.
AS	Maximum Array Size - the maximum size in gigabytes of any single data object.
DS	Maximum Data Size - max size in gigabytes combined of all data objects in .bss
TS	Maximum Total Size - max size in gigabytes, in aggregate, of all executable code and data objects in a running program.

## 15.6. Practical Limitations of Large Array Programming

The 64-bit addressing capability of 64-bit Linux environments can cause unexpected issues when data sizes are enlarged significantly. The following table describes the most common occurrences of practical limitations of large array programming.

Table 32 64-Bit Limitations

array initialization	Initializing a large array with a data statement may result in very large assembly and object files, where a line of assembler source is required for each element in the initialized array. Compilation and linking can be very time consuming as well. To avoid this issue, consider initializing large arrays in a loop at runtime rather than in a data statement.
stack space	Stack space can be a problem for data that is stack-based. On Linux, stack size is increased in your shell environment. If setting stacksize to unlimited is not large enough, try setting the size explicitly: <pre>limit stacksize new_size ! in csh</pre> <pre>ulimit -s new_size ! in bash</pre>
page swapping	If your executable is much larger than the physical size of memory, page swapping can cause it to run dramatically slower; it may even fail. This is not a compiler problem. Try smaller data sets to determine whether or not a problem is due to page thrashing.
configured space	Be sure your Linux system is configured with swap space sufficiently large to support the data sets used in your application(s). If your memory+swap space is not sufficiently large, your application will likely encounter a segmentation fault at runtime.
support for large address offsets in object file format	Arrays that are not dynamically allocated are limited by how the compiler can express the 'distance' between them when generating code. A field in the object file stores this 'distance' value, which is limited to 32-bits on Linux with -mcmodel=small. It is 64-bits on Linux with -mcmodel=medium. <div>  Without the 64-bit offset support in the object file format, large arrays cannot be declared statically, or locally on the stack. </div>

## 15.7. Medium Memory Model and Large Array in C

Consider the following example, where the aggregate size of the arrays exceeds 2GB.

### Medium Memory Model and Large Array in C

```
% cat bigadd.c
#include <stdio.h>
#define SIZE 600000000 /* > 2GB/4 */
static float a[SIZE], b[SIZE];
int
main()
{
```

```

long long i, n, m;
float c[SIZE]; /* goes on stack */
n = SIZE;
m = 0;
for (i = 0; i < n; i += 10000) {
    a[i] = i + 1;
    b[i] = 2.0 * (i + 1);
    c[i] = a[i] + b[i];
    m = i;
}
printf("a[0]=%g b[0]=%g c[0]=%g\n", a[0], b[0], c[0]);
printf("m=%lld a[%lld]=%g b[%lld]=%gc[%lld]=%g\n",m,m,a[m],m,b[m],m,c[m]);
return 0;
}

```

```
% nvc -mcmmodel=medium -o bigadd bigadd.c
```

When SIZE is greater than 2G/4, and the arrays are of type float with 4 bytes per element, the size of each array is greater than 2GB. With nvc, using the `-mcmmodel=medium` switch, a static data object can now be > 2GB in size. If you execute with these settings in your environment, you may see the following:

```
% bigadd
Segmentation fault
```

Execution fails because the stack size is not large enough. You can most likely correct this error by using the **limit stacksize** command to reset the stack size in your environment:

```
% limit stacksize 3000M
```



The command **limit stacksize unlimited** probably does not provide as large a stack as we are using in the [this example](#).

```
% bigadd
a[0]=1 b[0]=2 c[0]=3
n=599990000 a[599990000]=5.9999e+08 b[599990000]=1.19998e+09
c[599990000]=1.79997e+09
```

## 15.8. Medium Memory Model and Large Array in Fortran

The following example works with the NVFORTRAN compiler. It uses 64-bit addresses and index arithmetic when the `-mcmmodel=medium` option is used.

Consider the following example:

### Medium Memory Model and Large Array in Fortran

```

% cat mat.f
program mat
  integer i, j, k, size, l, m, n
  parameter (size=16000) ! >2GB
  parameter (m=size,n=size)
  real*8 a(m,n),b(m,n),c(m,n),d
  do i = 1, m
    do j = 1, n
      a(i,j)=10000.0D0*dble(i)+dble(j)
    
```



```

        b(i,j)=20000.0D0*dble(i)+dble(j)
    enddo
enddo
!$omp parallel
!$omp do
do i = 1, m
    do j = 1, n
        c(i,j) = a(i,j) + b(i,j)
    enddo
enddo
!$omp do
do i=1,m
    do j = 1, n
        d = 30000.0D0*dble(i)+dble(j)+dble(j)
        if (d .ne. c(i,j)) then
            print *, "err i=", i, "j=", j
            print *, "c(i,j)=", c(i,j)
            print *, "d=", d
            stop
        endif
    enddo
enddo
!$omp end parallel
print *, "M =", M, ", N =", N
print *, "c(M,N) = ", c(m,n)
end

```

When compiled with the NVFORTRAN compiler using -mcmodel=medium:

```

% nvfortran -Mfree -mp -o mat mat.f -i8 -mcmodel=medium
% setenv OMP_NUM_THREADS 2
% mat
M = 16000 , N = 16000
c(M,N) = 480032000.0000000

```

## 15.9. Large Array and Small Memory Model in Fortran

The following example uses large, dynamically-allocated arrays. The code is divided into a main and subroutine so you could put the subroutine into a shared library. Dynamic allocation of large arrays saves space in the size of executable and saves time initializing data.

### Large Array and Small Memory Model in Fortran

```
% cat mat_allo.f90
```

```

program mat_allo
    integer i, j
    integer size, m, n
    parameter (size=16000)
    parameter (m=size,n=size)
    double precision, allocatable::a(:,,:),b(:,,:),c(:,,:)
    allocate(a(m,n), b(m,n), c(m,n))
    do i = 100, m, 1
        do j = 100, n, 1
            a(i,j) = 10000.0D0 * dble(i) + dble(j)
            b(i,j) = 20000.0D0 * dble(i) + dble(j)
        enddo
    enddo
    call mat_add(a,b,c,m,n)
    print *, "M =", m, ", N =", n
end

```

```
    print *, "c(M,N) = ", c(m,n)
end

subroutine mat_add(a,b,c,m,n)
    integer m, n, i, j
    double precision a(m,n),b(m,n),c(m,n)
    do i = 1, m
        do j = 1, n
            c(i,j) = a(i,j) + b(i,j)
        enddo
    enddo
    return
end

% nvfortran -o mat_allo mat_allo.f90 -i8 -Mlarge_arrays -mp -fast
```

# Chapter 16.

## C++ AND C INLINE ASSEMBLY AND INTRINSICS

The examples in this section are shown using x86-64 assembly instructions. Inline assembly is supported on OpenPOWER and Arm Server platforms as well, but is not documented in detail in this section.

### 16.1. Inline Assembly

Inline Assembly lets you specify machine instructions inside a "C" function. The format for an inline assembly instruction is this:

```
{ asm | __asm__ } ("string");
```

The `asm` statement begins with the `asm` or `__asm__` keyword. The `__asm__` keyword is typically used in header files that may be included in ISO "C" programs.

*string* is one or more machine specific instructions separated with a semi-colon (;) or newline (\n) character. These instructions are inserted directly into the compiler's assembly-language output for the enclosing function.

Some simple `asm` statements are:

```
asm ("cli");  
asm ("sti");
```

These `asm` statements disable and enable system interrupts respectively.

In the following example, the `eax` register is set to zero.

```
asm( "pushl %eax\n\t" "movl $0, %eax\n\t" "popl %eax");
```

Notice that `eax` is pushed on the stack so that it is not clobbered. When the statement is done with `eax`, it is restored with the `popl` instruction.

Typically a program uses macros that enclose `asm` statements. The following two examples use the interrupt constructs created previously in this section:

```
#define disableInt __asm__ ("cli");  
#define enableInt __asm__ ("sti");
```

## 16.2. Extended Inline Assembly

**Inline Assembly** explains how to use inline assembly to specify machine specific instructions inside a "C" function. This approach works well for simple machine operations such as disabling and enabling system interrupts. However, inline assembly has three distinct limitations:

1. The programmer must choose the registers required by the inline assembly.
2. To prevent register clobbering, the inline assembly must include push and pop code for registers that get modified by the inline assembly.
3. There is no easy way to access stack variables in an inline assembly statement.

*Extended Inline Assembly* was created to address these limitations. The format for extended inline assembly, also known as *extended asm*, is as follows:

```
{ asm | __asm__ } [ volatile | __volatile__ ]
("string" [: [output operands]] [: [input operands]] [: [clobberlist]]);
```

- ▶ Extended asm statements begin with the *asm* or *\_\_asm\_\_* keyword. Typically the *\_\_asm\_\_* keyword is used in header files that may be included by ISO "C" programs.
- ▶ An optional *volatile* or *\_\_volatile\_\_* keyword may appear after the *asm* keyword. This keyword instructs the compiler not to delete, move significantly, or combine with any other asm statement. Like *\_\_asm\_\_*, the *\_\_volatile\_\_* keyword is typically used with header files that may be included by ISO "C" programs.
- ▶ "string" is one or more machine specific instructions separated with a semi-colon (;) or newline (\n) character. The string can also contain operands specified in the [output operands], [input operands], and [clobber list]. The instructions are inserted directly into the compiler's assembly-language output for the enclosing function.
- ▶ The [output operands], [input operands], and [clobber list] items each describe the effect of the instruction for the compiler. For example:

```
asm( "movl %1, %%eax\n" "movl %%eax, %0" : "=r" (x) : "r" (y) : "%eax" );
```

where

"=r" (x) is an output operand.

"r" (y) is an input operand.

"%eax" is the clobber list consisting of one register, "%eax".

The notation for the output and input operands is a constraint string surrounded by quotes, followed by an expression, and surrounded by parentheses. The constraint string describes how the input and output operands are used in the asm "string". For example, "r" tells the compiler that the operand is a register. The "=" tells the compiler that the operand is write only, which means that a value is stored in an output operand's expression at the end of the asm statement.

Each operand is referenced in the asm "string" by a percent "%" and its number. The first operand is number 0, the second is number 1, the third is number 2, and so on. In the preceding example, "%0" references the output operand, and "%1" references the input operand. The asm "string" also contains "%%eax", which references machine register "%eax". Hard coded registers like "%eax" should be specified in the

clobber list to prevent conflicts with other instructions in the compiler's assembly-language output. *[output operands]*, *[input operands]*, and *[clobber list]* items are described in more detail in the following sections.

### 16.2.1. Output Operands

The *[output operands]* are an optional list of output constraint and expression pairs that specify the result(s) of the asm statement. An output constraint is a string that specifies how a result is delivered to the expression. For example, "*=r*" (*x*) says the output operand is a write-only register that stores its value in the "C" variable *x* at the end of the asm statement. An example follows:

```
int x;
void example()
{
    asm( "movl $0, %0" : "=r" (x) );
}
```

The previous example assigns 0 to the "C" variable *x*. For the function in this example, the compiler produces the following assembly. If you want to produce an assembly listing, compile the example with the `nvc -S` compiler option:

```
example:
..Dcfb0:
    pushq %rbp
..Dcfi0:
    movq %rsp, %rbp
..Dcfi1:
..ENL:
## lineno: 8
    movl $0, %eax
    movl %eax, x(%rip)
## lineno: 0
    popq %rbp
    ret
```

In the generated assembly shown, notice that the compiler generated two statements for the asm statement at line number 5. The compiler generated "*movl \$0, %eax*" from the asm "*string*". Also notice that *%eax* appears in place of "*%0*" because the compiler assigned the *%eax* register to variable *x*. Since item 0 is an output operand, the result must be stored in its expression (*x*).

In addition to write-only output operands, there are read/write output operands designated with a "+" instead of a "=". For example, "*+r*" (*x*) tells the compiler to initialize the output operand with variable *x* at the beginning of the asm statement.

To illustrate this point, the following example increments variable *x* by 1:

```
int x=1;
void example2()
{
    asm( "addl $1, %0" : "+r" (x) );
}
```

To perform the increment, the output operand must be initialized with variable *x*. The *read/write* constraint modifier ("*+*") instructs the compiler to initialize the output operand with its expression. The compiler generates the following assembly code for the `example2()` function:

```
example2:
..Dcfb0:
```

```

pushq %rbp
..Dcfi0:
movq %rsp, %rbp
..Dcfi1:
..EN1:
## lineno: 5
movl x(%rip), %eax
addl $1, %eax
movl %eax, x(%rip)
## lineno: 0
popq %rbp
ret

```

From the `example2()` code, two extraneous moves are generated in the assembly: one `movl` for initializing the output register and a second `movl` to write it to variable `x`. To eliminate these moves, use a memory constraint type instead of a register constraint type, as shown in the following example:

```

int x=1;
void example2()
{
    asm( "addl $1, %0" : "+m" (x) );
}

```

The compiler generates a memory reference in place of a memory constraint. This eliminates the two extraneous moves. Because the assembly uses a memory reference to variable `x`, it does not have to move `x` into a register prior to the `asm` statement; nor does it need to store the result after the `asm` statement. Additional constraint types are found in [Additional Constraints](#).

```

example2:
..Dcfb0:
pushq %rbp
..Dcfi0:
movq %rsp, %rbp
..Dcfi1:
..EN1:
## lineno: 5
addl $1, x(%rip)
## lineno: 0
popq %rbp
ret

```

The examples thus far have used only one output operand. Because extended `asm` accepts a list of output operands, `asm` statements can have more than one result, as shown in the following example:

```

void example4()
{
    int x=1; int y=2;
    asm( "addl $1, %1\n" "addl %1, %0": "+r" (x), "+m" (y) );
}

```

This example increments variable `y` by 1 then adds it to variable `x`. Multiple output operands are separated with a comma. The first output operand is item 0 ("`%0`") and the second is item 1 ("`%1`") in the `asm` "string". The resulting values for `x` and `y` are 4 and 3 respectively.

## 16.2.2. Input Operands

The *[input operands]* are an optional list of input constraint and expression pairs that specify what "C" values are needed by the `asm` statement. The input constraints specify

how the data is delivered to the asm statement. For example, "*r*" (*x*) says that the input operand is a register that has a copy of the value stored in "C" variable *x*. Another example is "*m*" (*x*) which says that the input item is the *memory* location associated with variable *x*. Other constraint types are discussed in [Additional Constraints](#). An example follows:

```
void example5()
{
    int x=1;
    int y=2;
    int z=3;
    asm( "addl %2, %1\n" "addl %2, %0" : "+r" (x), "+m" (y) : "r" (z) );
}
```

The previous example adds variable *z*, item 2, to variable *x* and variable *y*. The resulting values for *x* and *y* are 4 and 5 respectively.

Another type of input constraint worth mentioning here is the *matching constraint*. A matching constraint is used to specify an operand that fills both an input as well as an output role. An example follows:

```
int x=1;
void example6()
{
    asm( "addl $1, %1"
        : "=r" (x)
        : "0" (x) );
}
```

The previous example is equivalent to the *example2()* function shown in [Output Operands](#). The constraint/expression pair, "*0*" (*x*), tells the compiler to initialize output item 0 with variable *x* at the beginning of the *asm* statement. The resulting value for *x* is 2. Also note that "*%1*" in the asm "*string*" means the same thing as "*%0*" in this case. That is because there is only one operand with both an input and an output role.

Matching constraints are very similar to the *read/write* output operands mentioned in [Output Operands](#). However, there is one key difference between *read/write* output operands and *matching constraints*. The *matching constraint* can have an *input expression* that differs from its *output expression*.

The following example uses different values for the input and output roles:

```
int x;
int y=2;
void example7()
{
    asm( "addl $1, %1"
        : "=r" (x)
        : "0" (y) );
}
```

The compiler generates the following assembly for *example7()*:

```
example7:
..Dcfb0:
    pushq %rbp
..Dcfi0:
    movq %rsp, %rbp
..Dcfi1:
..EN1:
## lineno: 8
    movl y(%rip), %eax
    addl $1, %eax
    movl %eax, x(%rip)
```

```
## lineno: 0
popq %rbp
ret
```

Variable *x* gets initialized with the value stored in *y*, which is 2. After adding 1, the resulting value for variable *x* is 3.

Because *matching constraints* perform an input role for an output operand, it does not make sense for the output operand to have the read/write ("+") modifier. In fact, the compiler disallows *matching constraints* with read/write output operands. The output operand must have a write only ("=") modifier.

### 16.2.3. Clobber List

The *[clobber list]* is an optional list of strings that hold machine registers used in the asm "string". Essentially, these strings tell the compiler which registers may be clobbered by the asm statement. By placing registers in this list, the programmer does not have to explicitly save and restore them as required in traditional inline assembly (described in [Inline Assembly](#)). The compiler takes care of any required saving and restoring of the registers in this list.

Each machine register in the [clobber list] is a string separated by a comma. The leading '%' is optional in the register name. For example, "%eax" is equivalent to "eax". When specifying the register inside the asm "string", you must include two leading '%' characters in front of the name (for example, "%eax"). Otherwise, the compiler will behave as if a bad input/output operand was specified and generate an error message.

An example follows:

```
void example8()
{
    int x;
    int y=2;
    asm( "movl %1, %%eax\n"
        "movl %1, %%edx\n"
        "addl %%edx, %%eax\n"
        "addl %%eax, %0"
        : "=r" (x)
        : "0" (y)
        : "eax", "edx" );
}
```

This code uses two hard-coded registers, *eax* and *edx*. It performs the equivalent of  $3*y$  and assigns it to *x*, producing a result of 6.

In addition to machine registers, the clobber list may contain the following special flags:

**"cc"**

The asm statement may alter the control code register.

**"memory"**

The asm statement may modify memory in an unpredictable fashion.

When the "memory" flag is present, the compiler does not keep memory values cached in registers across the asm statement and does not optimize stores or loads to that memory. For example:

```
asm("call MyFunc":::"memory");
```

This asm statement contains a "memory" flag because it contains a call. The callee may otherwise clobber registers in use by the caller without the "memory" flag.



The following function uses extended asm and the "cc" flag to compute a power of 2 that is less than or equal to the input parameter n.

```
#pragma noline
int asmDivideConquer(int n)
{
    int ax = 0;
    int bx = 1;
    asm (
        "LogLoop:n"
        "cmp %2, %1n"
        "jnl Done:n"
        "inc %0n"
        "add %1,%1n"
        "jmp LogLoop:n"
        "Done:n"
        "dec %0n"
        : "+r" (ax), "+r" (bx) : "r" (n) : "cc");
    return ax;
}
```

The 'cc' flag is used because the asm statement contains some control flow that may alter the control code register. The #pragma noline statement prevents the compiler from inlining the asmDivideConquer() function. If the compiler inlines asmDivideConquer(), then it may illegally duplicate the labels LogLoop and Done in the generated assembly.

## 16.2.4. Additional Constraints

Operand constraints can be divided into four main categories:

- ▶ Simple Constraints
- ▶ Machine Constraints
- ▶ Multiple Alternative Constraints
- ▶ Constraint Modifiers

## 16.2.5. Simple Constraints

The simplest kind of constraint is a string of letters or characters, known as *Simple Constraints*, such as the "r" and "m" constraints introduced in [Output Operands](#). [Table 33](#) describes these constraints.

Table 33 Simple Constraints

Constraint	Description
whitespace	Whitespace characters are ignored.
E	An immediate floating point operand.
F	Same as "E".
g	Any general purpose register, memory, or immediate integer operand is allowed.
i	An immediate integer operand.
m	A memory operand. Any address supported by the machine is allowed.
n	Same as "i".
o	Same as "m".

Constraint	Description
p	An operand that is a valid memory address. The expression associated with the constraint is expected to evaluate to an address (for example, "p" (&x) ).
r	A general purpose register operand.
X	Same as "g".
0,1,2,..9	Matching Constraint. See <a href="#">Output Operands</a> for a description.

The following example uses the general or "g" constraint, which allows the compiler to pick an appropriate constraint type for the operand; the compiler chooses from a general purpose register, memory, or immediate operand. This code lets the compiler choose the constraint type for "y".

```
void example9()
{
    int x, y=2;
    asm( "movl %1, %0\n" : "=r"
        (x) : "g" (y) );
}
```

This technique can result in more efficient code. For example, when compiling `example9()` the compiler replaces the load and store of `y` with a constant 2. The compiler can then generate an immediate 2 for the `y` operand in the example. The assembly generated by `nvc` for our example is as follows:

```
example9:
..Dcfb0:
    pushq %rbp
..Dcfi0:
    movq %rsp, %rbp
..Dcfi1:
..EN1:
## lineno: 3
    movl $2, %eax
## lineno: 6
    popq %rbp
    ret
```

In this example, notice the use of \$2 for the "y" operand.

Of course, if `y` is always 2, then the immediate value may be used instead of the variable with the "i" constraint, as shown here:

```
void example10()
{
    int x;
    asm( "movl %1, %0\n"
        : "=r" (x)
        : "i" (2) );
}
```

Compiling `example10()` with `nvc` produces assembly similar to that produced for `example9()`.

## 16.2.6. Machine Constraints

Another category of constraints is *Machine Constraints*. The `x86_64` architectures has several classes of registers. To choose a particular class of register, you can use the `x86_64` machine constraints described in [Table 34](#).

Table 34 x86\_64 Machine Constraints

Constraint	Description
a	a register (e.g., %al, %ax, %eax, %rax)
A	Specifies a or d registers. The d register holds the most significant bits and the a register holds the least significant bits.
b	b register (e.g., %bl, %bx, %ebx, %rbx)
c	c register (e.g., %cl, %cx, %ecx, %rcx)
C	Not supported.
d	d register (e.g., %dl, %dx, %edx, %rdx)
D	di register (e.g., %dil, %di, %edi, %rdi)
e	Constant in range of 0xffffffff to 0x7fffffff
f	Not supported.
G	Floating point constant in range of 0.0 to 1.0.
I	Constant in range of 0 to 31 (e.g., for 32-bit shifts).
J	Constant in range of 0 to 63 (e.g., for 64-bit shifts)
K	Constant in range of 0 to 127.
L	Constant in range of 0 to 65535.
M	Constant in range of 0 to 3 constant (e.g., shifts for lea instruction).
N	Constant in range of 0 to 255 (e.g., for out instruction).
q	Same as "r" simple constraint.
Q	Same as "r" simple constraint.
R	Same as "r" simple constraint.
S	si register (e.g., %sil, %si, %edi, %rsi)
t	Not supported.
u	Not supported.
x	XMM SSE register
y	Not supported.
Z	Constant in range of 0 to 0x7fffffff.

The following example uses the "x" or XMM register constraint to subtract c from b and store the result in a.

```
double example11()
{
    double a;
    double b = 400.99;
    double c = 300.98;
    asm ( "subpd %2, %0;"
        : "=x" (a)
        : "0" (b), "x" (c)
        );
    return a;
}
```

The generated assembly for this example is this:

```
example11:
..Dcfb0:
    pushq %rbp
..Dcfi0:
    movq %rsp, %rbp
..Dcfi1:
..EN1:
## lineno: 4
    movsd .C00128(%rip), %xmm1
    movsd .C00130(%rip), %xmm2
    movapd %xmm1, %xmm0
    subpd %xmm2, %xmm0;
## lineno: 10
## lineno: 11
    popq %rbp
    ret
```

If a specified register is not available, the nvc and nvc++ compilers issue an error message.

## 16.2.7. Multiple Alternative Constraints

Sometimes a single instruction can take a variety of operand types. For example, the x86-64 permits register-to-memory and memory-to-register operations. To allow this flexibility in inline assembly, use *multiple alternative constraints*. An *alternative* is a series of constraints for each operand.

To specify multiple alternatives, separate each alternative with a comma.

Table 35 Multiple Alternative Constraints

Constraint	Description
,	Separates each alternative for a particular operand.
?	Ignored
!	Ignored

The following example uses multiple alternatives for an add operation.

```
void example13()
{
    int x=1;
    int y=1;
    asm( "addl %1, %0\n"
        : "+ab,cd" (x)
        : "db,cam" (y) );
}
```

The preceding *example13()* has two alternatives for each operand: "ab,cd" for the output operand and "db,cam" for the input operand. Each operand must have the same number of alternatives; however, each alternative can have any number of constraints (for example, the output operand in *example13()* has two constraints for its second alternative and the input operand has three for its second alternative).

The compiler first tries to satisfy the left-most alternative of the first operand (for example, the output operand in *example13()*). When satisfying the operand, the compiler starts with the left-most constraint. If the compiler cannot satisfy an alternative with

this constraint (for example, if the desired register is not available), it tries to use any subsequent constraints. If the compiler runs out of constraints, it moves on to the next alternative. If the compiler runs out of alternatives, it issues an error similar to the one mentioned in *example12()*. If an alternative is found, the compiler uses the same alternative for subsequent operands. For example, if the compiler chooses the "c" register for the output operand in *example13()*, then it will use either the "a" or "m" constraint for the input operand.

## 16.2.8. Constraint Modifiers

Characters that affect the compiler's interpretation of a constraint are known as *Constraint Modifiers*. Two constraint modifiers, the "=" and the "+", were introduced in *Output Operands*. The following table summarizes each constraint modifier.

Table 36 Constraint Modifier Characters

Constraint Modifier	Description
=	This operand is write-only. It is valid for output operands only. If specified, the "=" must appear as the first character of the constraint string.
+	This operand is both read and written by the instruction. It is valid for output operands only. The output operand is initialized with its expression before the first instruction in the asm statement. If specified, the "+" must appear as the first character of the constraint string.
&	A constraint or an alternative constraint, as defined in <i>Multiple Alternative Constraints</i> , containing an "&" indicates that the output operand is an early clobber operand. This type operand is an output operand that may be modified before the asm statement finishes using all of the input operands. The compiler will not place this operand in a register that may be used as an input operand or part of any memory address.
%	Ignored.
#	Characters following a "#" up to the first comma (if present) are to be ignored in the constraint.
*	The character that follows the "*" is to be ignored in the constraint.

The "=" and "+" modifiers apply to the operand, regardless of the number of alternatives in the constraint string. For example, the "+" in the output operand of *example13()* appears once and applies to both alternatives in the constraint string. The "&", "#", and "\*" modifiers apply only to the alternative in which they appear.

Normally, the compiler assumes that input operands are used before assigning results to the output operands. This assumption lets the compiler reuse registers as needed inside the asm statement. However, if the asm statement does not follow this convention, the compiler may indiscriminately clobber a result register with an input operand. To prevent this behavior, apply the early clobber "&" modifier. An example follows:

```
void example15()
{
    int w=1;
    int z;
    asm( "movl $1, %0\n"
        "addl %2, %0\n"
```

```
"movl %2, %1"
: "=a" (w), "=r" (z) : "r" (w) );
}
```

The previous code example presents an interesting ambiguity because "w" appears both as an output and as an input operand. So, the value of "z" can be either 1 or 2, depending on whether the compiler uses the same register for operand 0 and operand 2. The use of constraint "r" for operand 2 allows the compiler to pick any general purpose register, so it may (or may not) pick register "a" for operand 2. This ambiguity can be eliminated by changing the constraint for operand 2 from "r" to "a" so the value of "z" will be 2, or by adding an early clobber "&" modifier so that "z" will be 1. The following example shows the same function with an early clobber "&" modifier:

```
void example16()
{
    int w=1;
    int z;
    asm( "movl $1, %0\n"
        "addl %2, %0\n"
        "movl %2, %1"
        : "&a" (w), "=r" (z) : "r" (w) );
}
```

Adding the early clobber "&" forces the compiler not to use the "a" register for anything other than operand 0. Operand 2 will therefore get its own register with its own copy of "w". The result for "z" in *example16()* is 1.

## 16.3. Operand Aliases

Extended asm specifies operands in assembly strings with a percent '%' followed by the operand number. For example, "%0" references operand 0 or the output item "&a" (w) in function *example16()* in the previous example. Extended asm also supports operand aliasing, which allows use of a symbolic name instead of a number for specifying operands, as illustrated in this example:

```
void example17()
{
    int w=1, z=0;
    asm( "movl $1, %[output1]\n"
        "addl %[input], %[output1]\n"
        "movl %[input], %[output2]"
        : [output1] "&a" (w), [output2] "=r"
        (z)
        : [input] "r" (w));
}
```

In *example18()*, "%0" and "%[output1]" both represent the output operand.

## 16.4. Assembly String Modifiers

Special character sequences in the assembly string affect the way the assembly is generated by the compiler. For example, the "%" is an escape sequence for specifying an operand, "%%" produces a percent for hard coded registers, and "\n" specifies a new line. [Table 37](#) summarizes these modifiers, known as *Assembly String Modifiers*.

Table 37 Assembly String Modifier Characters

Modifier	Description
\	Same as \ in printf format strings.
%%	Adds a '%' in the assembly string.
%%	Adds a '%' in the assembly string.
%A	Adds a '*' in front of an operand in the assembly string. (For example, %A0 adds a '*' in front of operand 0 in the assembly output.)
%B	Produces the byte op code suffix for this operand. (For example, %b0 produces 'b' on x86-64.)
%L	Produces the word op code suffix for this operand. (For example, %L0 produces 'l' on x86-64.)
%P	If producing Position Independent Code (PIC), the compiler adds the PIC suffix for this operand. (For example, %P0 produces @PLT on x86-64.)
%Q	Produces a quad word op code suffix for this operand if it is supported by the target. Otherwise, it produces a word op code suffix. (For example, %Q0 produces 'q' on x86-64.)
%S	Produces 's' suffix for this operand. (For example, %S0 produces 's' on x86-64.)
%T	Produces 't' suffix for this operand. (For example, %T0 produces 't' on x86-64.)
%W	Produces the half word op code suffix for this operand. (For example, %W0 produces 'w' on x86-64.)
%a	Adds open and close parentheses ( ) around the operand.
%b	Produces the byte register name for an operand. (For example, if operand 0 is in register 'a', then %b0 will produce '%al'.)
%c	Cuts the '\$' character from an immediate operand.
%k	Produces the word register name for an operand. (For example, if operand 0 is in register 'a', then %k0 will produce '%eax'.)
%q	Produces the quad word register name for an operand if the target supports quad word. Otherwise, it produces a word register name. (For example, if operand 0 is in register 'a', then %q0 produces %rax on x86-64.)
%w	Produces the half word register name for an operand. (For example, if operand 0 is in register 'a', then %w0 will produce '%ax'.)
%z	Produces an op code suffix based on the size of an operand. (For example, 'b' for byte, 'w' for half word, 'l' for word, and 'q' for quad word.)
%+ %C %D %F %O %X %f %h %l %n %s %y are not supported.	

These modifiers begin with either a backslash "\" or a percent "%".

The modifiers that begin with a backslash "\" (e.g., "\n") have the same effect as they do in a printf format string. The modifiers that are preceded with a "%" are used to modify a particular operand.

These modifiers begin with either a backslash "\" or a percent "%". For example, "%b0" means, "produce the byte or 8 bit version of operand 0". If operand 0 is a register, it will produce a byte register such as %al, %bl, %cl, and so on.

## 16.5. Extended Asm Macros

As with traditional inline assembly, described in [Inline Assembly](#), extended asm can be used in a macro. For example, you can use the following macro to access the runtime stack pointer.

```
#define GET_SP(x) \
asm("mov %%sp, %0": "=m" (##x):: "%sp" );
void example20()
{
    void * stack_pointer;
    GET_SP(stack_pointer);
}
```

The GET\_SP macro assigns the value of the stack pointer to whatever is inserted in its argument (for example, stack\_pointer). Another "C" extension known as *statement expressions* is used to write the GET\_SP macro another way:

```
#define GET_SP2 ({ \
void *my_stack_ptr; \
asm("mov %%sp, %0": "=m" (my_stack_ptr) :: "%sp" ); \
my_stack_ptr; \
})
void example21()
{
    void * stack_pointer = GET_SP2;
}
```

The statement expression allows a body of code to evaluate to a single value. This value is specified as the last instruction in the statement expression. In this case, the value is the result of the asm statement, my\_stack\_ptr. By writing an asm macro with a statement expression, the asm result may be assigned directly to another variable (for example, void \* stack\_pointer = GET\_SP2) or included in a larger expression, such as: void \* stack\_pointer = GET\_SP2 - sizeof(long).

Which style of macro to use depends on the application. If the asm statement needs to be a part of an expression, then a macro with a statement expression is a good approach. Otherwise, a traditional macro, like GET\_SP(x), will probably suffice.

## 16.6. Intrinsics

Inline intrinsic functions map to actual x86-64 machine instructions. Intrinsics are inserted inline to avoid the overhead of a function call. The compiler has special knowledge of intrinsics, so with use of intrinsics, better code may be generated as compared to extended inline assembly code.

The NVIDIA HPC Compilers intrinsics library implements MMX, SSE, SS2, SSE3, SSSE3, SSE4a, ABM, and AVX instructions. The intrinsic functions are available to C and C++ programs. Unlike most functions which are in libraries, intrinsics are implemented internally by the compiler. A program can call the intrinsic functions from C/C++ source code after including the corresponding header file.

The intrinsics are divided into header files as follows:



Table 38 Intrinsic Header File Organization

Instructions	Header File		Instructions	Header File
ABM	intrin.h		SSE2	emmintrin.h
AVX	immintrin.h		SSE3	pmmmintrin.h
MMX	mmintrin.h		SSSE3	tmmintrin.h
SSE	xmmmintrin.h		SSE4a	ammintrin.h

The following is a simple example program that calls XMM intrinsics.

```
#include <xmmmintrin.h>
int main(){
    __m128 __A, __B, result;
    __A = _mm_set_ps(23.3, 43.7, 234.234, 98.746);
    __B = _mm_set_ps(15.4, 34.3, 4.1, 8.6);
    result = _mm_add_ps(__A, __B);
    return 0;
}
```

## Notice

ALL NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, "MATERIALS") ARE BEING PROVIDED "AS IS." NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE.

Information furnished is believed to be accurate and reliable. However, NVIDIA Corporation assumes no responsibility for the consequences of use of such information or for any infringement of patents or other rights of third parties that may result from its use. No license is granted by implication or otherwise under any patent rights of NVIDIA Corporation. Specifications mentioned in this publication are subject to change without notice. This publication supersedes and replaces all other information previously supplied. NVIDIA Corporation products are not authorized as critical components in life support devices or systems without express written approval of NVIDIA Corporation.

## Trademarks

NVIDIA, the NVIDIA logo, CUDA, CUDA-X, GPUDirect, HPC SDK, NGC, NVIDIA Volta, NVIDIA DGX, NVIDIA Nsight, NVLink, NVSwitch, and Tesla are trademarks and/or registered trademarks of NVIDIA Corporation in the U.S. and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

## Copyright

© 2013–2024 NVIDIA Corporation. All rights reserved.