



# CUDA COMPILER DRIVER NVCC

TRM-06721-001\_v8.0 | May 2016

## Reference Guide



## CHANGES FROM PREVIOUS VERSION

- Major update to the document to reflect recent **nvcc** changes.

# TABLE OF CONTENTS

<b>Chapter 1. Introduction.....</b>	<b>1</b>
1.1. Overview.....	1
1.1.1. CUDA Programming Model.....	1
1.1.2. CUDA Sources.....	1
1.1.3. Purpose of NVCC.....	2
1.2. Supported Host Compilers.....	2
<b>Chapter 2. Compilation Phases.....</b>	<b>3</b>
2.1. NVCC Identification Macro.....	3
2.2. NVCC Phases.....	3
2.3. Supported Input File Suffixes.....	4
2.4. Supported Phases.....	4
<b>Chapter 3. NVCC Command Options.....</b>	<b>7</b>
3.1. Command Option Types and Notation.....	7
3.2. Command Option Description.....	8
3.2.1. Options for Specifying the Compilation Phase.....	8
3.2.2. File and Path Specifications.....	9
3.2.3. Options for Specifying Behavior of Compiler/Linker.....	10
3.2.4. Options for Passing Specific Phase Options.....	11
3.2.5. Options for Guiding the Compiler Driver.....	12
3.2.6. Options for Steering CUDA Compilation.....	13
3.2.7. Options for Steering GPU Code Generation.....	13
3.2.8. Generic Tool Options.....	17
3.2.9. Phase Options.....	17
3.2.9.1. Ptxas Options.....	17
3.2.9.2. NVLINK Options.....	19
<b>Chapter 4. The CUDA Compilation Trajectory.....</b>	<b>20</b>
<b>Chapter 5. GPU Compilation.....</b>	<b>22</b>
5.1. GPU Generations.....	22
5.2. GPU Feature List.....	22
5.3. Application Compatibility.....	23
5.4. Virtual Architectures.....	23
5.5. Virtual Architecture Feature List.....	24
5.6. Further Mechanisms.....	25
5.6.1. Just-in-Time Compilation.....	25
5.6.2. Fatbinaries.....	26
5.7. NVCC Examples.....	26
5.7.1. Base Notation.....	26
5.7.2. Shorthand.....	26
5.7.2.1. Shorthand 1.....	26
5.7.2.2. Shorthand 2.....	26

5.7.2.3. Shorthand 3.....	27
5.7.3. Extended Notation.....	27
5.7.4. Virtual Architecture Identification Macro.....	28
<b>Chapter 6. Using Separate Compilation in CUDA.....</b>	<b>29</b>
6.1. Code Changes for Separate Compilation.....	29
6.2. NVCC Options for Separate Compilation.....	29
6.3. Libraries.....	30
6.4. Examples.....	31
6.5. Potential Separate Compilation Issues.....	32
6.5.1. Object Compatibility.....	32
6.5.2. JIT Linking Support.....	33
6.5.3. Implicit CUDA Host Code.....	33
6.5.4. Using <code>__CUDA_ARCH__</code> .....	33
6.5.5. Device Code in Libraries.....	34
<b>Chapter 7. Miscellaneous NVCC Usage.....</b>	<b>35</b>
7.1. Cross Compilation.....	35
7.2. Keeping Intermediate Phase Files.....	35
7.3. Cleaning Up Generated Files.....	35
7.4. Printing Code Generation Statistics.....	36

## LIST OF FIGURES

Figure 1	CUDA Whole Program Compilation Trajectory .....	21
Figure 2	Two-Staged Compilation with Virtual and Real Architectures .....	24
Figure 3	Just-in-Time Compilation of Device Code .....	25
Figure 4	CUDA Separate Compilation Trajectory .....	30



# Chapter 1.

## INTRODUCTION

### 1.1. Overview

#### 1.1.1. CUDA Programming Model

The CUDA Toolkit targets a class of applications whose control part runs as a process on a general purpose computing device, and which use one or more NVIDIA GPUs as coprocessors for accelerating *single program, multiple data* (SPMD) parallel jobs. Such jobs are self-contained, in the sense that they can be executed and completed by a batch of GPU threads entirely without intervention by the host process, thereby gaining optimal benefit from the parallel graphics hardware.

The GPU code is implemented as a collection of functions in a language that is essentially C++, but with some annotations for distinguishing them from the host code, plus annotations for distinguishing different types of data memory that exists on the GPU. Such functions may have parameters, and they can be called using a syntax that is very similar to regular C function calling, but slightly extended for being able to specify the matrix of GPU threads that must execute the called function. During its life time, the host process may dispatch many parallel GPU tasks.

For more information on the CUDA programming model, consult the [CUDA C Programming Guide](#).

#### 1.1.2. CUDA Sources

Source files for CUDA applications consist of a mixture of conventional C++ host code, plus GPU device functions. The CUDA compilation trajectory separates the device functions from the host code, compiles the device functions using the proprietary NVIDIA compilers and assembler, compiles the host code using a C++ host compiler that is available, and afterwards embeds the compiled GPU functions as fatbinary images in the host object file. In the linking stage, specific CUDA runtime libraries are added for supporting remote SPMD procedure calling and for providing explicit GPU manipulation such as allocation of GPU memory buffers and host-GPU data transfer.

### 1.1.3. Purpose of NVCC

The compilation trajectory involves several splitting, compilation, preprocessing, and merging steps for each CUDA source file. It is the purpose of **nvcc**, the CUDA compiler driver, to hide the intricate details of CUDA compilation from developers. It accepts a range of conventional compiler options, such as for defining macros and include/library paths, and for steering the compilation process. All non-CUDA compilation steps are forwarded to a C++ host compiler that is supported by **nvcc**, and **nvcc** translates its options to appropriate host compiler command line options.

## 1.2. Supported Host Compilers

A general purpose C++ host compiler is needed by **nvcc** in the following situations:

- ▶ During non-CUDA phases (except the run phase), because these phases will be forwarded by **nvcc** to this compiler.
- ▶ During CUDA phases, for several preprocessing stages and host code compilation (see also [The CUDA Compilation Trajectory](#)).

**nvcc** assumes that the host compiler is installed with the standard method designed by the compiler provider. If the host compiler installation is non-standard, the user must make sure that the environment is set appropriately and use relevant **nvcc** compile options.

The following documents provide detailed information about supported host compilers:

- ▶ [NVIDIA CUDA Installation Guide for Linux](#)
- ▶ [NVIDIA CUDA Installation Guide for Mac OS X](#)
- ▶ [NVIDIA CUDA Installation Guide for Microsoft Windows](#)

On all platforms, the default host compiler executable (**gcc** and **g++** on Linux, **clang** and **clang++** on Mac OS X, and **cl.exe** on Windows) found in the current execution search path will be used, unless specified otherwise with appropriate options (see [File and Path Specifications](#)).



# Chapter 2.

## COMPILATION PHASES

### 2.1. NVCC Identification Macro

**nvcc** predefines the following macros:

- \_\_NVCC\_\_**  
Defined when compiling C/C++/CUDA source files.
- \_\_CUDACC\_\_**  
Defined when compiling CUDA source files.
- \_\_CUDACC\_RDC\_\_**  
Defined when compiling CUDA sources files in relocatable device code mode (see [NVCC Options for Separate Compilation](#)).
- \_\_CUDACC\_RELAXED\_CONSTEXPR\_\_**  
Defined when the `--expt-relaxed-constexpr` flag is specified on the command line. Refer to [CUDA C Programming Guide](#) for more details.
- \_\_CUDACC\_EXTENDED\_LAMBDA\_\_**  
Defined when the `--expt-extended-lambda` flag is specified on the command line. Refer to [CUDA C Programming Guide](#) for more details.
- \_\_CUDACC\_VER\_MAJOR\_\_**  
Defined with the major version number of **nvcc**.
- \_\_CUDACC\_VER\_MINOR\_\_**  
Defined with the minor version number of **nvcc**.
- \_\_CUDACC\_VER\_BUILD\_\_**  
Defined with the build version number of **nvcc**.
- \_\_CUDACC\_VER\_\_**  
Defined with the full version number of **nvcc**, represented as  
$$\text{__CUDACC\_VER\_MAJOR__} * 10000 + \text{__CUDACC\_VER\_MINOR__} * 100 + \text{__CUDACC\_VER\_BUILD__}.$$

### 2.2. NVCC Phases

A compilation phase is the a logical translation step that can be selected by command line options to **nvcc**. A single compilation phase can still be broken up by **nvcc** into

smaller steps, but these smaller steps are just implementations of the phase: they depend on seemingly arbitrary capabilities of the internal tools that **nvcc** uses, and all of these internals may change with a new release of the CUDA Toolkit. Hence, only compilation phases are stable across releases, and although **nvcc** provides options to display the compilation steps that it executes, these are for debugging purposes only and must not be copied and used into build scripts.

**nvcc** phases are selected by a combination of command line options and input file name suffixes, and the execution of these phases may be modified by other command line options. In phase selection, the input file suffix defines the phase input, while the command line option defines the required output of the phase.

The following paragraphs will list the recognized file name suffixes and the supported compilation phases. A full explanation of the **nvcc** command line options can be found in the next chapter.

## 2.3. Supported Input File Suffixes

The following table defines how **nvcc** interprets its input files:

Input File Prefix	Description
.cu	CUDA source file, containing host code and device functions
.c	C source file
.cc, .cxx, .cpp	C++ source file
.gpu	GPU intermediate file (see <a href="#">Figure 1</a> )
.ptx	PTX intermediate assembly file (see <a href="#">Figure 1</a> )
.o, .obj	Object file
.a, .lib	Library file
.res	Resource file
.so	Shared object file

Note that **nvcc** does not make any distinction between object, library or resource files. It just passes files of these types to the linker when the linking phase is executed.

## 2.4. Supported Phases

The following table specifies the supported compilation phases, plus the option to **nvcc** that enables execution of this phase. It also lists the default name of the output file generated by this phase, which will take effect when no explicit output file name is specified using option **--output-file**:

Phase	nvcc Option		Default Output File Name
	Long Name	Short Name	
CUDA compilation to	<b>--cuda</b>	<b>-cuda</b>	.cpp.i.i appended to source file name, as in x.cu.cpp.i.i. This output file can be compiled

Phase	nvcc Option		Default Output File Name
	Long Name	Short Name	
C/C++ source file			by the host compiler that was used by <code>nvcc</code> to preprocess the <code>.cu</code> file.
C/C++ preprocessing	<code>--preprocess</code>	<code>-E</code>	<code>&lt;result on standard output&gt;</code>
C/C++ compilation to object file	<code>--compile</code>	<code>-c</code>	Source file name with suffix replaced by <code>o</code> on Linux and Mac OS X, or <code>obj</code> on Windows
Cubin generation from CUDA source files	<code>--cubin</code>	<code>-cubin</code>	Source file name with suffix replaced by <code>cubin</code>
Cubin generation from <code>.gpu</code> intermediate files	<code>--cubin</code>	<code>-cubin</code>	Source file name with suffix replaced by <code>cubin</code>
Cubin generation from PTX intermediate files.	<code>--cubin</code>	<code>-cubin</code>	Source file name with suffix replaced by <code>cubin</code>
PTX generation from CUDA source files	<code>--ptx</code>	<code>-ptx</code>	Source file name with suffix replaced by <code>ptx</code>
PTX generation from <code>.gpu</code> intermediate files	<code>--ptx</code>	<code>-ptx</code>	Source file name with suffix replaced by <code>ptx</code>
Fatbinary generation from source, PTX or cubin files	<code>--fatbin</code>	<code>-fatbin</code>	Source file name with suffix replaced by <code>fatbin</code>
GPU C code generation from CUDA source files	<code>--gpu</code>	<code>-gpu</code>	Source file name with suffix replaced by <code>gpu</code>
Linking relocatable device code.	<code>--device-link</code>	<code>-dlink</code>	<code>a_dlink.obj</code> on Windows or <code>a_dlink.o</code> on other platforms
Cubin generation from linked relocatable device code.	<code>--device-link --cubin</code>	<code>-dlink -cubin</code>	<code>a_dlink.cubin</code>
Fatbinary generation from linked	<code>--device-link --fatbin</code>	<code>-dlink -fatbin</code>	<code>a_dlink.fatbin</code>

Phase	nvcc Option		Default Output File Name
	Long Name	Short Name	
relocatable device code			
Linking an executable	<no phase option>		a.exe on Windows or a.out on other platforms
Constructing an object file archive, or library	--lib	-lib	a.lib on Windows or a.a on other platforms
make dependency generation	--generate-dependencies	-M	<result on standard output>
Running an executable	--run	-run	

**Notes:**

- ▶ The last phase in this list is more of a convenience phase. It allows running the compiled and linked executable without having to explicitly set the library path to the CUDA dynamic libraries.
- ▶ Unless a phase option is specified, **nvcc** will compile and link all its input files.

# Chapter 3.

## NVCC COMMAND OPTIONS

### 3.1. Command Option Types and Notation

Each **nvcc** option has a long name and a short name, which are interchangeable with each other. These two variants are distinguished by the number of hyphens that must precede the option name: long names must be preceded by two hyphens, while short names must be preceded by a single hyphen. For example, **-I** is the short name of **--include-path**. Long options are intended for use in build scripts, where size of the option is less important than descriptive value. In contrast, short options are intended for interactive use.

**nvcc** recognizes three types of command options: boolean options, single value options, and list options.

Boolean options do not have an argument; they are either specified on a command line or not. Single value options must be specified at most once, and list options may be repeated. Examples of each of these option types are, respectively: **--verbose** (switch to verbose mode), **--output-file** (specify output file), and **--include-path** (specify include path).

Single value options and list options must have arguments, which must follow the name of the option itself by either one of more spaces or an equals character. When a one-character short name such as **-I**, **-l**, and **-L** is used, the value of the option may also immediately follow the option itself without being separated by spaces or an equal character. The individual values of list options may be separated by commas in a single instance of the option, or the option may be repeated, or any combination of these two cases.

Hence, for the two sample options mentioned above that may take values, the following notations are legal:

```
-o file
-o=file
-I dir1,dir2 -I=dir3 -I dir4,dir5
```

Long option names are used throughout the document, unless specified otherwise, however, short names can be used instead of long names to have the same effect.

## 3.2. Command Option Description

This section presents tables of **nvcc** options. The option type in the tables can be recognized as follows: boolean options do not have arguments specified in the first column, while the other two types do. List options can be recognized by the repeat indicator `, . . .` at the end of the argument.

Long options are described in the first columns of the options tables, and short options occupy the second columns.

### 3.2.1. Options for Specifying the Compilation Phase

Options of this category specify up to which stage the input files must be compiled.

Long Name	Short Name	Description
<code>--cuda</code>	<code>-cuda</code>	Compile all <code>.cu</code> input files to <code>.cu.cpp.i</code> output.
<code>--cubin</code>	<code>-cubin</code>	Compile all <code>.cu/.gpu/.ptx</code> input files to device-only <code>.cubin</code> files. This step discards the host code for each <code>.cu</code> input file.
<code>--fatbin</code>	<code>-fatbin</code>	Compile all <code>.cu/.gpu/.ptx/.cubin</code> input files to device-only <code>.fatbin</code> files. This step discards the host code for each <code>.cu</code> input file.
<code>--ptx</code>	<code>-ptx</code>	Compile all <code>.cu/.gpu</code> input files to device-only <code>.ptx</code> files. This step discards the host code for each <code>.cu</code> input file.
<code>--gpu</code>	<code>-gpu</code>	Compile all <code>.cu</code> input files to device-only <code>.gpu</code> files. This step discards the host code for each <code>.cu</code> input file.
<code>--preprocess</code>	<code>-E</code>	Preprocess all <code>.c/.cc/.cpp/.cxx/.cu</code> input files.
<code>--generate-dependencies</code>	<code>-M</code>	Generate a dependency file that can be included in a <code>make</code> file for the <code>.c/.cc/.cpp/.cxx/.cu</code> input file (more than one are not allowed in this mode).
<code>--compile</code>	<code>-c</code>	Compile each <code>.c/.cc/.cpp/.cxx/.cu</code> input file into an object file.
<code>--device-c</code>	<code>-dc</code>	Compile each <code>.c/.cc/.cpp/.cxx/.cu</code> input file into an object file that contains relocatable device code. It is equivalent to <code>--relocatable-device-code=true --compile</code> .
<code>--device-w</code>	<code>-dw</code>	Compile each <code>.c/.cc/.cpp/.cxx/.cu</code> input file into an object file that contains executable device code. It is

Long Name	Short Name	Description
		equivalent to <code>--relocatable-device-code=false --compile</code> .
<code>--device-link</code>	<code>-dlink</code>	Link object files with relocatable device code and <code>.ptx/.cubin/.fatbin</code> files into an object file with executable device code, which can be passed to the host linker.
<code>--link</code>	<code>-link</code>	This option specifies the default behavior: compile and link all inputs.
<code>--lib</code>	<code>-lib</code>	Compile all input files into object files (if necessary), and add the results to the specified library output file.
<code>--run</code>	<code>-run</code>	This option compiles and links all inputs into an executable, and executes it. Or, when the input is a single executable, it is executed without any compilation or linking. This step is intended for developers who do not want to be bothered with setting the necessary environment variables; these are set temporarily by <code>nvcc</code> .

### 3.2.2. File and Path Specifications

Long Name	Short Name	Description
<code>--output-file file</code>	<code>-o</code>	Specify name and location of the output file. Only a single input file is allowed when this option is present in <code>nvcc</code> non-linking/archiving mode.
<code>--pre-include file,...</code>	<code>-include</code>	Specify header files that must be preincluded during preprocessing or compilation.
<code>--library library,...</code>	<code>-l</code>	Specify libraries to be used in the linking stage without the library file extension. The libraries are searched for on the library search paths that have been specified using option <code>--library-path</code> (see <a href="#">Libraries</a> ).
<code>--define-macro def,...</code>	<code>-D</code>	Specify macro definitions for use during preprocessing or compilation.
<code>--undefine-macro def,...</code>	<code>-U</code>	Undefine macro definitions during preprocessing or compilation.
<code>--include-path path,...</code>	<code>-I</code>	Specify include search paths.
<code>--system-include path,...</code>	<code>-isystem</code>	Specify system include search paths.
<code>--library-path path,...</code>	<code>-L</code>	Specify library search paths (see <a href="#">Libraries</a> ).
<code>--output-directory directory</code>	<code>-odir</code>	Specify the directory of the output file. This option is intended for letting the

Long Name	Short Name	Description
		dependency generation step (see <code>--generate-dependencies</code> ) generate a rule that defines the target object file in the proper directory.
<code>--compiler-bindir directory</code>	<code>-ccbin</code>	Specify the directory in which the compiler executable resides. The host compiler executable name can be also specified to ensure that the correct host compiler is selected. In addition, driver prefix options ( <code>--input-drive-prefix</code> , <code>--dependency-drive-prefix</code> , or <code>--drive-prefix</code> ) may need to be specified, if <code>nvcc</code> is executed in a Cygwin shell or a MinGW shell on Windows.
<code>--cudart {none shared static}</code>	<code>-cudart</code>	Specify the type of CUDA runtime library to be used: no CUDA runtime library, shared/dynamic CUDA runtime library, or static CUDA runtime library.  Allowed values for this option: <code>none</code> , <code>shared</code> , <code>static</code> .  Default value: <code>static</code>
<code>--libdevice-directory directory</code>	<code>-ldir</code>	Specify the directory that contains the libdevice library files when option <code>--dont-use-profile</code> is used. Libdevice library files are located in the <code>nvvm/libdevice</code> directory in the CUDA Toolkit.

### 3.2.3. Options for Specifying Behavior of Compiler/Linker

Long Name	Short Name	Description
<code>--profile</code>	<code>-pg</code>	Instrument generated code/executable for use by <code>gprof</code> (Linux only).
<code>--debug</code>	<code>-g</code>	Generate debug information for host code.
<code>--device-debug</code>	<code>-G</code>	Generate debug information for device code.
<code>--generate-line-info</code>	<code>-lineinfo</code>	Generate line-number information for device code.
<code>--optimize level</code>	<code>-O</code>	Specify optimization level for host code.
<code>--ftemplate-backtrace-limit limit</code>	<code>-ftemplate-backtrace-limit</code>	When compiling CUDA source code, set the maximum number of template instantiation notes for a single warning or error to <i>limit</i> . A value of 0 is allowed, and indicates that no limit should be enforced. Note: This flag does not affect the host compiler used for compiling host part of the CUDA source code.



Long Name	Short Name	Description
<code>--shared</code>	<code>-shared</code>	Generate a shared library during linking. Use option <code>--linker-options</code> when other linker options are required for more control.
<code>--x {c c++ cu}</code>	<code>-x</code>	Explicitly specify the language for the input files, rather than letting the compiler choose a default based on the file name suffix.  Allowed values for this option: <code>c</code> , <code>c++</code> , <code>cu</code> .
<code>--std {c++11}</code>	<code>-std</code>	Select a particular C++ dialect. The only value currently supported is <code>c++11</code> . Enabling C++11 mode also turns on C++11 mode for the host compiler.  Allowed value for this option: <code>c++11</code>
<code>--no-host-device-initializer-list</code>	<code>-nohdinitlist</code>	Do not implicitly consider member functions of <code>std::initializer_list</code> as <code>__host__ __device__</code> functions.
<code>--no-host-device-move-forward</code>	<code>-nohdmoveforward</code>	Do not implicitly consider <code>std::move</code> and <code>std::forward</code> as <code>__host__ __device__</code> function templates.
<code>--expt-relaxed-constexpr</code>	<code>-expt-relaxed-constexpr</code>	Experimental flag: Allow host code to invoke <code>__device__ constexpr</code> functions, and device code to invoke <code>__host__ constexpr</code> functions.
<code>--expt-extended-lambda</code>	<code>-expt-extended-lambda</code>	Experimental flag: Allow <code>__host__</code> , <code>__device__</code> annotations in lambda declaration.
<code>--machine {32 64}</code>	<code>-m</code>	Specify 32-bit vs. 64-bit architecture.  Allowed values for this option: <code>32</code> , <code>64</code> .

### 3.2.4. Options for Passing Specific Phase Options

These allow for passing specific options directly to the internal compilation tools that `nvcc` encapsulates, without burdening `nvcc` with too-detailed knowledge on these tools. A table of useful sub-tool options can be found at the end of this chapter.

Long Name	Short Name	Description
<code>--compiler-options options,...</code>	<code>-Xcompiler</code>	Specify options directly to the compiler/preprocessor.
<code>--linker-options options,...</code>	<code>-Xlinker</code>	Specify options directly to the host linker.
<code>--archive-options options,...</code>	<code>-Xarchive</code>	Specify options directly to library manager.
<code>--ptxas-options options,...</code>	<code>-Xptxas</code>	Specify options directly to <code>ptxas</code> , the PTX optimizing assembler.

Long Name	Short Name	Description
<code>--nvlink-options options,...</code>	<code>-Xnvlink</code>	Specify options directly to <code>nvlink</code> .

### 3.2.5. Options for Guiding the Compiler Driver

Long Name	Short Name	Description
<code>--dont-use-profile</code>	<code>-noprof</code>	<code>nvcc</code> uses the <code>nvcc.profiles</code> file for compilation. When specifying this option, the profile file is not used.
<code>--dryrun</code>	<code>-dryrun</code>	Do not execute the compilation commands generated by <code>nvcc</code> . Instead, list them.
<code>--verbose</code>	<code>-v</code>	List the compilation commands generated by this compiler driver, but do not suppress their execution.
<code>--keep</code>	<code>-keep</code>	Keep all intermediate files that are generated during internal compilation steps.
<code>--keep-dir directory</code>	<code>-keep-dir</code>	Keep all intermediate files that are generated during internal compilation steps in this directory.
<code>--save-temps</code>	<code>-save-temps</code>	This option is an alias of <code>--keep</code> .
<code>--clean-targets</code>	<code>-clean</code>	This option reverses the behavior of <code>nvcc</code> . When specified, none of the compilation phases will be executed. Instead, all of the non-temporary files that <code>nvcc</code> would otherwise create will be deleted.
<code>--run-args arguments,...</code>	<code>-run-args</code>	Used in combination with option <code>--run</code> to specify command line arguments for the executable.
<code>--input-drive-prefix prefix</code>	<code>-idp</code>	On Windows, all command line arguments that refer to file names must be converted to the Windows native format before they are passed to pure Windows executables. This option specifies how the current development environment represents absolute paths. Use <code>/cygwin/</code> as <i>prefix</i> for Cygwin build environments and <code>/</code> the <i>prefix</i> for MinGW.
<code>--dependency-drive-prefix prefix</code>	<code>-ddp</code>	On Windows, when generating dependency files (see <code>--generate-dependencies</code> ), all file names must be converted appropriately for the instance of <code>make</code> that is used. Some instances of <code>make</code> have trouble with the colon in absolute paths in the native Windows format, which depends on the environment in which the <code>make</code> instance has been compiled. Use <code>/cygwin/</code> as <i>prefix</i> for a Cygwin <code>make</code> , and <code>/</code> as

Long Name	Short Name	Description
		<i>prefix</i> for MinGW. Or leave these file names in the native Windows format by specifying nothing.
<code>--drive-prefix <i>prefix</i></code>	<code>-dp</code>	Specifies <i>prefix</i> as both <code>--input-drive-prefix</code> and <code>--dependency-drive-prefix</code> .
<code>--dependency-target-name <i>target</i></code>	<code>-MT</code>	Specify the target name of the generated rule when generating a dependency file (see <code>--generate-dependencies</code> ).
<code>--no-align-double</code>		Specifies that <code>-malign-double</code> should not be passed as a compiler argument on 32-bit platforms. <b>WARNING:</b> this makes the ABI incompatible with the CUDA's kernel ABI for certain 64-bit types.
<code>--no-device-link</code>	<code>-nodlink</code>	Skip the device link step when linking object files.

### 3.2.6. Options for Steering CUDA Compilation

Long Name	Short Name	Description
<code>--default-stream {<i>legacy</i> <i>null</i> <i>per-thread</i>}</code>	<code>-default-stream</code>	<p>Specify the stream that CUDA commands from the compiled program will be sent to by default.</p> <p>Allowed values for this option:</p> <p><b>legacy</b> The CUDA legacy stream (per context, implicitly synchronizes with other streams)</p> <p><b>per-thread</b> A normal CUDA stream (per thread, does not implicitly synchronize with other streams)</p> <p><b>null</b> is a deprecated alias for <b>legacy</b>.</p> <p>Default value: <b>legacy</b></p>

### 3.2.7. Options for Steering GPU Code Generation

Long Name	Short Name	Description
<code>--gpu-architecture <i>arch</i></code>	<code>-arch</code>	<p>Specify the name of the class of NVIDIA <i>virtual</i> GPU architecture for which the CUDA input files must be compiled.</p> <p>With the exception as described for the shorthand below, the architecture specified with this option must be a <i>virtual</i> architecture (such as <code>compute_20</code>). Normally, this option alone does not trigger assembly of the generated PTX for a <i>real</i> architecture (that is the role of <code>nvcc</code> option <code>--gpu-</code></p>

Long Name	Short Name	Description
		<p><code>code</code>, see below); rather, its purpose is to control preprocessing and compilation of the input to PTX.</p> <p>For convenience, in case of simple <code>nvcc</code> compilations, the following shorthand is supported. If no value for option <code>--gpu-code</code> is specified, then the value of this option defaults to the value of <code>--gpu-architecture</code>. In this situation, as only exception to the description above, the value specified for <code>--gpu-architecture</code> may be a <i>real</i> architecture (such as a <code>sm_20</code>), in which case <code>nvcc</code> uses the specified <i>real</i> architecture and its closest <i>virtual</i> architecture as effective architecture values. For example, <code>nvcc --gpu-architecture=sm_20</code> is equivalent to <code>nvcc --gpu-architecture=compute_20 --gpu-code=sm_20,compute_20</code>.</p> <p>See <a href="#">Virtual Architecture Feature List</a> for the list of supported <i>virtual</i> architectures and <a href="#">GPU Feature List</a> for the list of supported <i>real</i> architectures.</p>
<code>--gpu-code code,...</code>	<code>-code</code>	<p>Specify the name of the NVIDIA GPU to assemble and optimize PTX for.</p> <p><code>nvcc</code> embeds a compiled code image in the resulting executable for each specified <i>code</i> architecture, which is a true binary load image for each <i>real</i> architecture (such as <code>sm_20</code>), and PTX code for the <i>virtual</i> architecture (such as <code>compute_20</code>).</p> <p>During runtime, such embedded PTX code is dynamically compiled by the CUDA runtime system if no binary load image is found for the <i>current</i> GPU.</p> <p>Architectures specified for options <code>--gpu-architecture</code> and <code>--gpu-code</code> may be <i>virtual</i> as well as <i>real</i>, but the <i>code</i> architectures must be compatible with the <i>arch</i> architecture. When the <code>--gpu-code</code> option is used, the value for the <code>--gpu-architecture</code> option must be a <i>virtual</i> PTX architecture.</p> <p>For instance, <code>--gpu-architecture=compute_35</code> is not compatible with <code>--gpu-code=sm_30</code>, because the earlier compilation stages will assume the availability of <code>compute_35</code> features that are not present on <code>sm_30</code>.</p>

Long Name	Short Name	Description
		See <a href="#">Virtual Architecture Feature List</a> for the list of supported <i>virtual</i> architectures and <a href="#">GPU Feature List</a> for the list of supported <i>real</i> architectures.
<code>--generate-code</code> <i>specification</i>	<code>-gencode</code>	<p>This option provides a generalization of the <code>--gpu-architecture=arch</code> <code>--gpu-code=code,...</code> option combination for specifying <code>nvcc</code> behavior with respect to code generation. Where use of the previous options generates code for different <i>real</i> architectures with the PTX for the same <i>virtual</i> architecture, option <code>--generate-code</code> allows multiple PTX generations for different <i>virtual</i> architectures. In fact, <code>--gpu-architecture=arch</code> <code>--gpu-code=code,...</code> is equivalent to <code>--generate-code arch=arch,code=code,...</code>.</p> <p><code>--generate-code</code> options may be repeated for different virtual architectures.</p> <p>See <a href="#">Virtual Architecture Feature List</a> for the list of supported <i>virtual</i> architectures and <a href="#">GPU Feature List</a> for the list of supported <i>real</i> architectures.</p>
<code>--relocatable-device-code</code> <code>{true false}</code>	<code>-rdc</code>	<p>Enable (disable) the generation of relocatable device code. If disabled, executable device code is generated. Relocatable device code must be linked before it can be executed.</p> <p>Allowed values for this option: <code>true</code>, <code>false</code>.</p> <p>Default value: <code>false</code></p>
<code>--entries entry,...</code>	<code>-e</code>	Specify the global entry functions for which code must be generated. By default, code will be generated for all entries.
<code>--maxrregcount amount</code>	<code>-maxrregcount</code>	<p>Specify the maximum amount of registers that GPU functions can use.</p> <p>Until a function-specific limit, a higher value will generally increase the performance of individual GPU threads that execute this function. However, because thread registers are allocated from a global register pool on each GPU, a higher value of this option will also reduce the maximum thread block size, thereby reducing the amount of thread parallelism. Hence, a good <code>maxrregcount</code> value is the result of a trade-off.</p>

Long Name	Short Name	Description
		<p>If this option is not specified, then no maximum is assumed.</p> <p>Value less than the minimum registers required by ABI will be bumped up by the compiler to ABI minimum limit.</p>
<code>--use_fast_math</code>	<code>-use_fast_math</code>	<p>Make use of fast math library. <code>--use_fast_math</code> implies <code>--ftz=true</code> <code>--prec-div=false</code> <code>--prec-sqrt=false</code> <code>--fmad=true</code>.</p>
<code>--ftz {true false}</code>	<code>-ftz</code>	<p>This option controls single-precision denormals support. <code>--ftz=true</code> flushes denormal values to zero and <code>--ftz=false</code> preserves denormal values.</p> <p><code>--use_fast_math</code> implies <code>--ftz=true</code>.</p> <p>Allowed values for this option: <code>true</code>, <code>false</code>.</p> <p>Default value: <code>false</code></p>
<code>--prec-div {true false}</code>	<code>-prec-div</code>	<p>This option controls single-precision floating-point division and reciprocals. <code>--prec-div=true</code> enables the IEEE round-to-nearest mode and <code>--prec-div=false</code> enables the fast approximation mode.</p> <p><code>--use_fast_math</code> implies <code>--prec-div=false</code>.</p> <p>Allowed values for this option: <code>true</code>, <code>false</code>.</p> <p>Default value: <code>true</code></p>
<code>--prec-sqrt {true false}</code>	<code>-prec-sqrt</code>	<p>This option controls single-precision floating-point square root. <code>--prec-sqrt=true</code> enables the IEEE round-to-nearest mode and <code>--prec-sqrt=false</code> enables the fast approximation mode.</p> <p><code>--use_fast_math</code> implies <code>--prec-sqrt=false</code>.</p> <p>Allowed values for this option: <code>true</code>, <code>false</code>.</p> <p>Default value: <code>true</code></p>
<code>--fmad {true false}</code>	<code>-fmad</code>	<p>This option enables (disables) the contraction of floating-point multiplies and adds/subtracts into floating-point multiply-add operations (FMAD, FFMA, or DFMA).</p> <p><code>--use_fast_math</code> implies <code>--fmad=true</code>.</p> <p>Allowed values for this option: <code>true</code>, <code>false</code>.</p> <p>Default value: <code>true</code></p>

### 3.2.8. Generic Tool Options

Long Name	Short Name	Description
<code>--disable-warnings</code>	<code>-w</code>	Inhibit all warning messages.
<code>--source-in-ptx</code>	<code>-src-in-ptx</code>	Interleave source in PTX. May only be used in conjunction with <code>--device-debug</code> or <code>--generate-line-info</code> .
<code>--restrict</code>	<code>-restrict</code>	Programmer assertion that all kernel pointer parameters are restrict pointers.
<code>--Wno-deprecated-gpu-targets</code>	<code>-Wno-deprecated-gpu-targets</code>	Suppress warnings about deprecated GPU target architectures.
<code>--Werror kind,...</code>	<code>-Werror</code>	Make warnings of the specified kinds into errors. The following is the list of warning kinds accepted by this option: <code>cross-execution-space-call</code> Be more strict about unsupported cross execution space calls. The compiler will generate an error instead of a warning for a call from a <code>__host__ __device__</code> to a <code>__host__</code> function.
<code>--resource-usage</code>	<code>-res-usage</code>	Show resource usage such as registers and memory of the GPU code.  This option implies <code>--nvlink-options=--verbose</code> when <code>--relocatable-device-code=true</code> is set. Otherwise, it implies <code>--ptxas-options=--verbose</code> .
<code>--help</code>	<code>-h</code>	Print help information on this tool.
<code>--version</code>	<code>-v</code>	Print version information on this tool.
<code>--options-file file,...</code>	<code>-optf</code>	Include command line options from specified file.

### 3.2.9. Phase Options

The following sections lists some useful options to lower level compilation tools.

#### 3.2.9.1. Ptxas Options

The following table lists some useful **ptxas** options which can be specified with **nvcc** option **-Xptxas**.

Long Name	Short Name	Description
<code>--allow-expensive-optimizations</code>	<code>-allow-expensive-optimizations</code>	Enable (disable) to allow compiler to perform expensive optimizations using maximum available resources (memory and compile-time).

Long Name	Short Name	Description
		If unspecified, default behavior is to enable this feature for optimization level $\geq O2$ .
<code>--compile-only</code>	<code>-c</code>	Generate relocatable object.
<code>--def-load-cache</code>	<code>-dlcm</code>	Default cache modifier on global/generic load. Default value: <code>ca</code> .
<code>--def-store-cache</code>	<code>-dscm</code>	Default cache modifier on global/generic store.
<code>--device-debug</code>	<code>-g</code>	Semantics same as <code>nvcc</code> option <code>--device-debug</code> .
<code>--disable-optimizer-constants</code>	<code>-disable-optimizer-consts</code>	Disable use of optimizer constant bank.
<code>--entry entry,...</code>	<code>-e</code>	Semantics same as <code>nvcc</code> option <code>--entries</code> .
<code>--fmad</code>	<code>-fmad</code>	Semantics same as <code>nvcc</code> option <code>--fmad</code> .
<code>--force-load-cache</code>	<code>-flcm</code>	Force specified cache modifier on global/generic load.
<code>--force-store-cache</code>	<code>-fscm</code>	Force specified cache modifier on global/generic store.
<code>--generate-line-info</code>	<code>-lineinfo</code>	Semantics same as <code>nvcc</code> option <code>--generate-line-info</code> .
<code>--gpu-name gpuname</code>	<code>-arch</code>	Specify name of NVIDIA GPU to generate code for. This option also takes virtual compute architectures, in which case code generation is suppressed. This can be used for parsing only.  Allowed values for this option: <code>compute_20</code> , <code>compute_30</code> , <code>compute_35</code> , <code>compute_50</code> , <code>compute_52</code> ; and <code>sm_20</code> , <code>sm_21</code> , <code>sm_30</code> , <code>sm_32</code> , <code>sm_35</code> , <code>sm_50</code> and <code>sm_52</code> .  Default value: <code>sm_20</code> .
<code>--help</code>	<code>-h</code>	Semantics same as <code>nvcc</code> option <code>--help</code> .
<code>--machine</code>	<code>-m</code>	Semantics same as <code>nvcc</code> option <code>--machine</code> .
<code>--maxrregcount amount</code>	<code>-maxrregcount</code>	Semantics same as <code>nvcc</code> option <code>--maxrregcount</code> .
<code>--opt-level N</code>	<code>-O</code>	Specify optimization level. Default value: 3.
<code>--options-file file,...</code>	<code>-optf</code>	Semantics same as <code>nvcc</code> option <code>--options-file</code> .
<code>--output-file file</code>	<code>-o</code>	Specify name of output file. Default value: <code>elf.o</code> .



Long Name	Short Name	Description
<code>--preserve-relocs</code>	<code>-preserve-relocs</code>	This option will make <code>ptxas</code> to generate relocatable references for variables and preserve relocations generated for them in linked executable.
<code>--sp-bound-check</code>	<code>-sp-bound-check</code>	Generate stack-pointer bounds-checking code sequence. This option is turned on automatically when <code>--device-debug</code> or <code>--opt-level=0</code> is specified.
<code>--verbose</code>	<code>-v</code>	Enable verbose mode which prints code generation statistics.
<code>--version</code>	<code>-V</code>	Semantics same as <code>nvcc</code> option <code>--version</code> .
<code>--warning-as-error</code>	<code>-Werror</code>	Make all warnings into errors.
<code>--warn-on-double-precision-use</code>	<code>-warn-double-usage</code>	Warning if double(s) are used in an instruction.
<code>--warn-on-local-memory-usage</code>	<code>-warn-lmem-usage</code>	Warning if local memory is used.
<code>--warn-on-spills</code>	<code>-warn-spills</code>	Warning if registers are spilled to local memory.

### 3.2.9.2. NVLINK Options

The following table lists some useful `nvlink` options which can be specified with `nvcc` option `--nvlink-options`.

Long Name	Short Name	Description
<code>--disable-warnings</code>	<code>-w</code>	Inhibit all warning messages.
<code>--preserve-relocs</code>	<code>-preserve-relocs</code>	Preserve resolved relocations in linked executable.
<code>--verbose</code>	<code>-v</code>	Enable verbose mode which prints code generation statistics.
<code>--warning-as-error</code>	<code>-Werror</code>	Make all warnings into errors.

## Chapter 4.

# THE CUDA COMPILATION TRAJECTORY

The CUDA phase converts a source file coded in the extended CUDA language into a regular ANSI C++ source file that can be handed over to a general purpose C++ host compiler for further compilation and linking. The exact steps that are followed to achieve this are displayed in [Figure 1](#).

CUDA compilation works as follows: the input program is preprocessed for device compilation and is compiled to CUDA binary (**cubin**) and/or PTX intermediate code, which are placed in a fatbinary. The input program is preprocessed once again for host compilation and is synthesized to embed the fatbinary and transform CUDA specific C++ extensions into standard C++ constructs. Then the C++ host compiler compiles the synthesized host code with the embedded fatbinary into a host object.

The embedded fatbinary is inspected by the CUDA runtime system whenever the device code is launched by the host program to obtain an appropriate fatbinary image for the current GPU.

The CUDA compilation trajectory is more complicated in the separate compilation mode. For more information, see [Using Separate Compilation in CUDA](#).

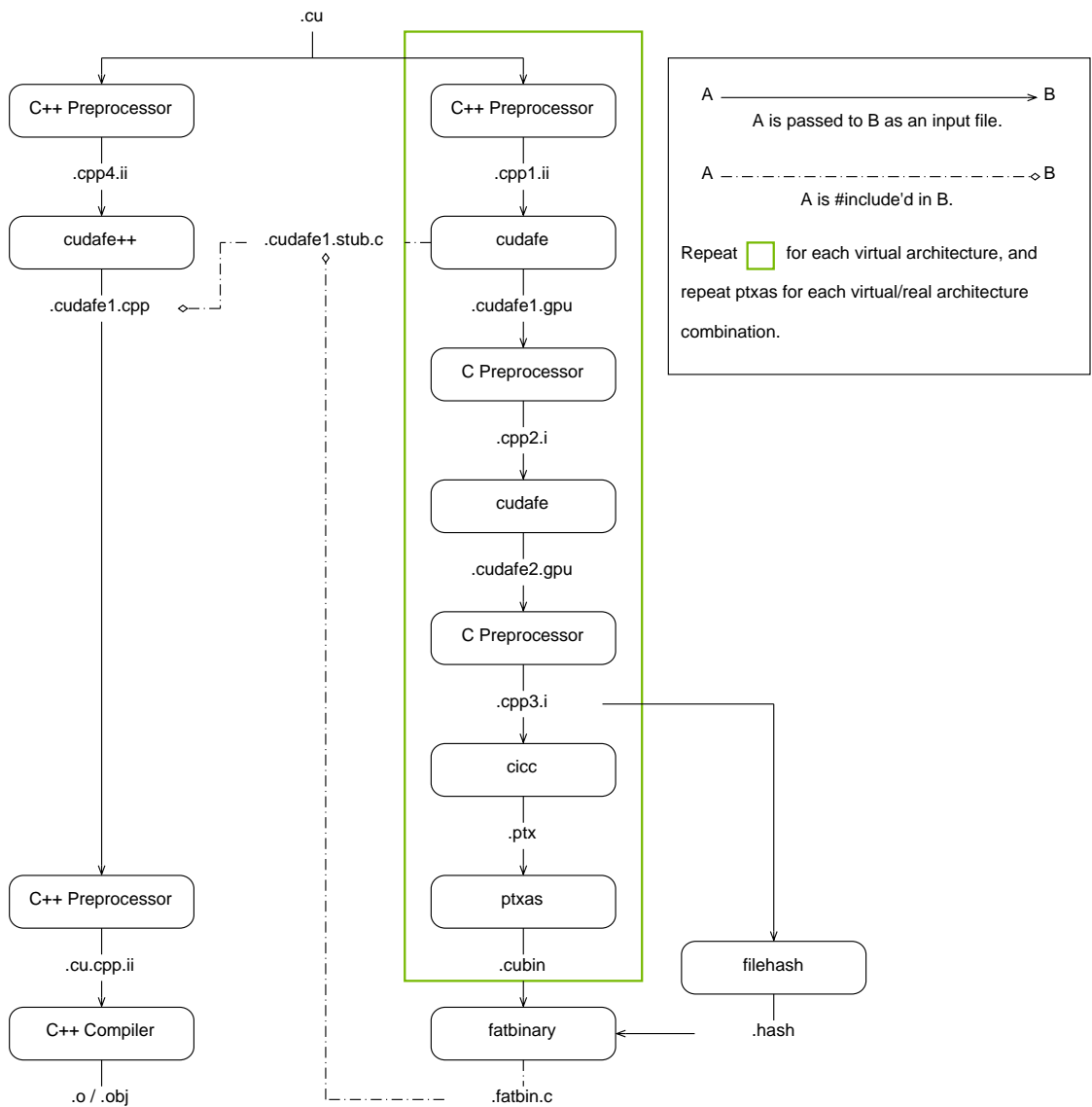


Figure 1 CUDA Whole Program Compilation Trajectory

# Chapter 5.

## GPU COMPILATION

This chapter describes the GPU compilation model that is maintained by **nvcc**, in cooperation with the CUDA driver. It goes through some technical sections, with concrete examples at the end.

### 5.1. GPU Generations

In order to allow for architectural evolution, NVIDIA GPUs are released in different generations. New generations introduce major improvements in functionality and/or chip architecture, while GPU models within the same generation show minor configuration differences that *moderately* affect functionality, performance, or both.

Binary compatibility of GPU applications is not guaranteed across different generations. For example, a CUDA application that has been compiled for a Fermi GPU will very likely not run on a Kepler GPU (and vice versa). This is the instruction set and instruction encodings of a generation is different from those of other generations.

Binary compatibility within one GPU generation can be guaranteed under certain conditions because they share the basic instruction set. This is the case between two GPU versions that do not show functional differences at all (for instance when one version is a scaled down version of the other), or when one version is functionally included in the other. An example of the latter is the *base* Kepler version **sm\_30** whose functionality is a subset of all other Kepler versions: any code compiled for **sm\_30** will run on all other Kepler GPUs.

### 5.2. GPU Feature List

The following table lists the names of the current GPU architectures, annotated with the functional capabilities that they provide. There are other differences, such as amounts of register and processor clusters, that only affect execution performance.

In the CUDA naming scheme, GPUs are named **sm\_xy**, where **x** denotes the GPU generation number, and **y** the version in that generation. Additionally, to facilitate comparing GPU capabilities, CUDA attempts to choose its GPU names such that if **x<sub>1</sub>y<sub>1</sub> <= x<sub>2</sub>y<sub>2</sub>** then all non-ISA related capabilities of **sm\_x<sub>1</sub>y<sub>1</sub>** are included in those of

**sm\_x2y2**. From this it indeed follows that **sm\_30** is the *base* Kepler model, and it also explains why higher entries in the tables are always functional extensions to the lower entries. This is denoted by the plus sign in the table. Moreover, if we abstract from the instruction encoding, it implies that **sm\_30**'s functionality will continue to be included in all later GPU generations. As we will see next, this property will be the foundation for application compatibility support by **nvcc**.

<b>sm_20</b>	Basic features + Fermi support
<b>sm_30</b> and <b>sm_32</b>	+ Kepler support + Unified memory programming
<b>sm_35</b>	+ Dynamic parallelism support
<b>sm_50</b> , <b>sm_52</b> , and <b>sm_53</b>	+ Maxwell support

## 5.3. Application Compatibility

Binary code compatibility over CPU generations, together with a published instruction set architecture is the usual mechanism for ensuring that distributed applications *out there in the field* will continue to run on newer versions of the CPU when these become mainstream.

This situation is different for GPUs, because NVIDIA cannot guarantee binary compatibility without sacrificing regular opportunities for GPU improvements. Rather, as is already conventional in the graphics programming domain, **nvcc** relies on a two stage compilation model for ensuring application compatibility with future GPU generations.

## 5.4. Virtual Architectures

GPU compilation is performed via an intermediate representation, PTX, which can be considered as assembly for a virtual GPU architecture. Contrary to an actual graphics processor, such a virtual GPU is defined entirely by the set of capabilities, or features, that it provides to the application. In particular, a virtual GPU architecture provides a (largely) generic instruction set, and binary instruction encoding is a non-issue because PTX programs are always represented in text format.

Hence, a **nvcc** compilation command always uses two architectures: a *virtual* intermediate architecture, plus a *real* GPU architecture to specify the intended processor to execute on. For such an **nvcc** command to be valid, the *real* architecture must be an implementation of the *virtual* architecture. This is further explained below.

The chosen virtual architecture is more of a statement on the GPU capabilities that the application requires: using a *smallest* virtual architecture still allows a *widest* range of actual architectures for the second **nvcc** stage. Conversely, specifying a virtual architecture that provides features unused by the application unnecessarily restricts the set of possible GPUs that can be specified in the second **nvcc** stage.

From this it follows that the virtual architecture should always be chosen as *low* as possible, thereby maximizing the actual GPUs to run on. The *real* architecture should be chosen as *high* as possible (assuming that this always generates better code), but this is only possible with knowledge of the actual GPUs on which the application is expected to run. As we will see later, in the situation of just in time compilation, where the driver has this exact knowledge: the runtime GPU is the one on which the program is about to be launched/executed.

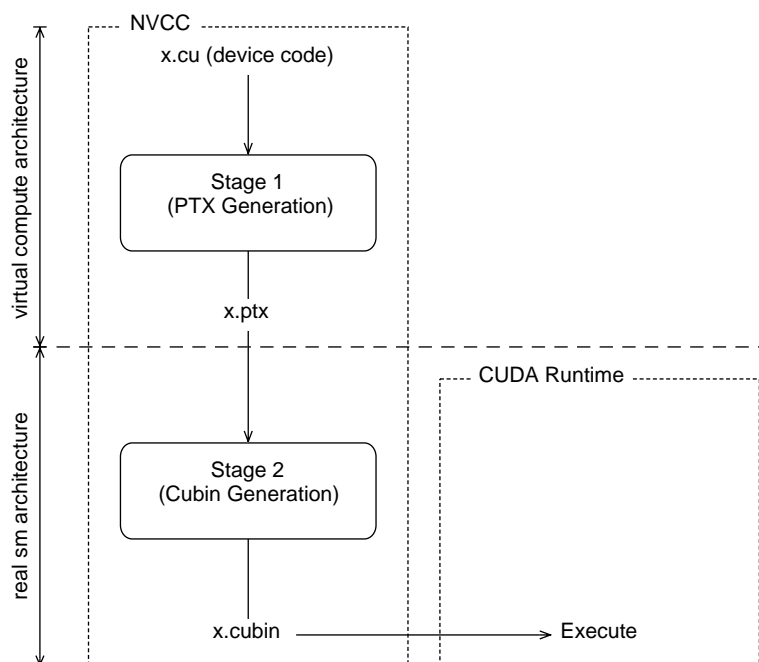


Figure 2 Two-Staged Compilation with Virtual and Real Architectures

## 5.5. Virtual Architecture Feature List

<code>compute_20</code>	Basic features + Fermi support
<code>compute_30</code> and <code>compute_32</code>	+ Kepler support + Unified memory programming
<code>compute_35</code>	+ Dynamic parallelism support
<code>compute_50</code> , <code>compute_52</code> , and <code>compute_53</code>	+ Maxwell support

The above table lists the currently defined virtual architectures. The virtual architecture naming scheme is the same as the real architecture naming scheme shown in Section [GPU Feature List](#).

## 5.6. Further Mechanisms

Clearly, compilation staging in itself does not help towards the goal of application compatibility with future GPUs. For this we need the two other mechanisms by CUDA Samples: just in time compilation (JIT) and fatbinaries.

### 5.6.1. Just-in-Time Compilation

The compilation step to an actual GPU binds the code to one generation of GPUs. Within that generation, it involves a choice between GPU *coverage* and possible performance. For example, compiling to **sm\_30** allows the code to run on all Kepler-generation GPUs, but compiling to **sm\_35** would probably yield better code if Kepler GK110 and later are the only targets.

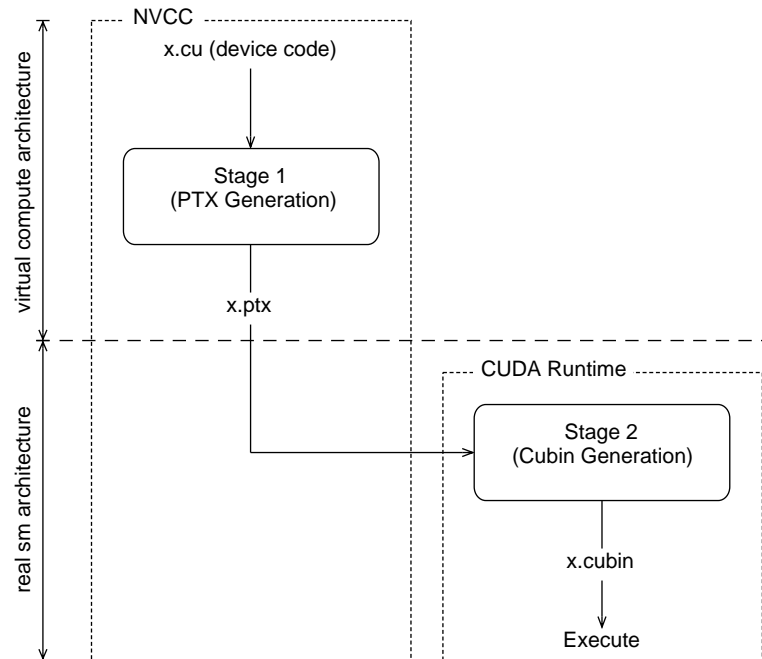


Figure 3 Just-in-Time Compilation of Device Code

By specifying a virtual code architecture instead of a *real* GPU, **nvcc** postpones the assembly of PTX code until application runtime, at which the target GPU is exactly known. For instance, the command below allows generation of exactly matching GPU binary code, when the application is launched on an **sm\_20** or later architecture.

```
nvcc x.cu --gpu-architecture=compute_20 --gpu-code=compute_20
```

The disadvantage of just in time compilation is increased application startup delay, but this can be alleviated by letting the CUDA driver use a compilation cache (refer to "Section 3.1.1.2. Just-in-Time Compilation" of [CUDA C Programming Guide](#)) which is persistent over multiple runs of the applications.

## 5.6.2. Fatbinaries

A different solution to overcome startup delay by JIT while still allowing execution on newer GPUs is to specify multiple code instances, as in

```
nvcc x.cu --gpu-architecture=compute_30 --gpu-code=compute_30,sm_30,sm_35
```

This command generates exact code for two Kepler variants, plus PTX code for use by JIT in case a next-generation GPU is encountered. **nvcc** organizes its device code in fatbinaries, which are able to hold multiple translations of the same GPU source code. At runtime, the CUDA driver will select the most appropriate translation when the device function is launched.

## 5.7. NVCC Examples

### 5.7.1. Base Notation

**nvcc** provides the options **--gpu-architecture** and **--gpu-code** for specifying the target architectures for both translation stages. Except for allowed short hands described below, the **--gpu-architecture** option takes a single value, which must be the name of a virtual compute architecture, while option **--gpu-code** takes a list of values which must all be the names of actual GPUs. **nvcc** performs a stage 2 translation for each of these GPUs, and will embed the result in the result of compilation (which usually is a host object file or executable).

#### Example

```
nvcc x.cu --gpu-architecture=compute_30 --gpu-code=sm_30,sm_35
```

### 5.7.2. Shorthand

**nvcc** allows a number of shorthands for simple cases.

#### 5.7.2.1. Shorthand 1

**--gpu-code** arguments can be virtual architectures. In this case the stage 2 translation will be omitted for such virtual architecture, and the stage 1 PTX result will be embedded instead. At application launch, and in case the driver does not find a better alternative, the stage 2 compilation will be invoked by the driver with the PTX as input.

#### Example

```
nvcc x.cu --gpu-architecture=compute_30 --gpu-code=compute_30,sm_30,sm_35
```

#### 5.7.2.2. Shorthand 2

The **--gpu-code** option can be omitted. Only in this case, the **--gpu-architecture** value can be a non-virtual architecture. The **--gpu-code** values default to the *closest* virtual architecture that is implemented by the GPU specified with **--gpu-**



**architecture**, plus the **--gpu-architecture**, value itself. The *closest* virtual architecture is used as the effective **--gpu-architecture**, value. If the **--gpu-architecture** value is a virtual architecture, it is also used as the effective **--gpu-code** value.

### Example

```
nvcc x.cu --gpu-architecture=sm_35
nvcc x.cu --gpu-architecture=compute_30
```

are equivalent to

```
nvcc x.cu --gpu-architecture=compute_35 --gpu-code=sm_35,compute_35
nvcc x.cu --gpu-architecture=compute_30 --gpu-code=compute_30
```

### 5.7.2.3. Shorthand 3

Both **--gpu-architecture** and **--gpu-code** options can be omitted.

### Example

```
nvcc x.cu
```

is equivalent to

```
nvcc x.cu --gpu-architecture=compute_20 --gpu-code=sm_20,compute_20
```

## 5.7.3. Extended Notation

The options **--gpu-architecture** and **--gpu-code** can be used in all cases where code is to be generated for one or more GPUs using a common virtual architecture. This will cause a single invocation of **nvcc** stage 1 (that is, preprocessing and generation of virtual PTX assembly code), followed by a compilation stage 2 (binary code generation) repeated for each specified GPU.

Using a common virtual architecture means that all assumed GPU features are fixed for the entire **nvcc** compilation. For instance, the following **nvcc** command assumes no warp shuffle support for both the **sm\_20** code and the **sm\_30** code:

```
nvcc x.cu --gpu-architecture=compute_20 --gpu-code=compute_20,sm_20,sm_30
```

Sometimes it is necessary to perform different GPU code generation steps, partitioned over different architectures. This is possible using **nvcc** option **--generate-code**, which then must be used instead of a **--gpu-architecture** and **--gpu-code** combination.

Unlike option **--gpu-architecture** option **--generate-code**, may be repeated on the **nvcc** command line. It takes sub-options **arch** and **code**, which must not be confused with their main option equivalents, but behave similarly. If repeated architecture compilation is used, then the device code must use conditional compilation based on the value of the architecture identification macro **\_\_CUDA\_ARCH\_\_**, which is described in the next section.

For example, the following assumes absence of warp shuffle support for the **sm\_20** and **sm\_21** code, but full support on **sm\_3x**:

```
nvcc x.cu \
  --generate-code arch=compute_20,code=sm_20 \
  --generate-code arch=compute_20,code=sm_21 \
  --generate-code arch=compute_30,code=sm_30
```

Or, leaving actual GPU code generation to the JIT compiler in the CUDA driver:

```
nvcc x.cu \
  --generate-code arch=compute_20,code=compute_20 \
  --generate-code arch=compute_30,code=compute_30
```

The code sub-options can be combined, but for technical reasons must then be quoted, which causes a slightly more complex syntax:

```
nvcc x.cu \
  --generate-code arch=compute_20,code=\"sm_20,sm_21\" \
  --generate-code arch=compute_30,code=\"sm_30,sm_35\"
```

## 5.7.4. Virtual Architecture Identification Macro

The architecture identification macro `__CUDA_ARCH__` is assigned a three-digit value string `xy0` (ending in a literal 0) during each `nvcc` compilation stage 1 that compiles for `compute_xy`.

This macro can be used in the implementation of GPU functions for determining the virtual architecture for which it is currently being compiled. The host code (the non-GPU code) must *not* depend on it.

# Chapter 6.

## USING SEPARATE COMPILE IN CUDA

Prior to the 5.0 release, CUDA did not support separate compilation, so CUDA code could not call device functions or access variables across files. Such compilation is referred to as *whole program compilation*. We have always supported the separate compilation of host code, it was just the device CUDA code that needed to all be within one file. Starting with CUDA 5.0, separate compilation of device code is supported, but the old whole program mode is still the default, so there are new options to invoke separate compilation.

### 6.1. Code Changes for Separate Compilation

The code changes required for separate compilation of device code are the same as what you already do for host code, namely using **extern** and **static** to control the visibility of symbols. Note that previously **extern** was ignored in CUDA code; now it will be honored. With the use of **static** it is possible to have multiple device symbols with the same name in different files. For this reason, the CUDA API calls that referred to symbols by their string name are deprecated; instead the symbol should be referenced by its address.

### 6.2. NVCC Options for Separate Compilation

CUDA works by embedding device code into host objects. In whole program compilation, it embeds executable device code into the host object. In separate compilation, we embed relocatable device code into the host object, and run **nvlink**, the device linker, to link all the device code together. The output of **nvlink** is then linked together with all the host objects by the host linker to form the final executable.

The generation of relocatable vs executable device code is controlled by the **--relocatable-device-code** option.

The **--compile** option is already used to control stopping a compile at a host object, so a new option **--device-c** is added that simply does **--relocatable-device-code=true --compile** .

To invoke just the device linker, the `--device-link` option can be used, which emits a host object containing the embedded executable device code. The output of that must then be passed to the host linker. Or:

```
nvcc <objects>
```

can be used to implicitly call both the device and host linkers. This works because if the device linker does not see any relocatable code it does not do anything.

Figure 4 shows the flow (`nvcc --device-c` has the same flow as Figure 1).

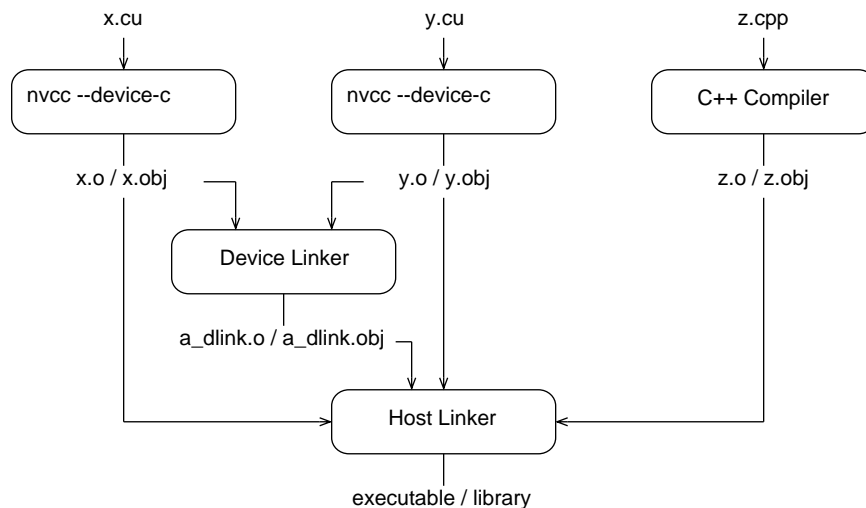


Figure 4 CUDA Separate Compilation Trajectory

## 6.3. Libraries

The device linker has the ability to read the static host library formats (`.a` on Linux and Mac OS X, `.lib` on Windows). It ignores any dynamic (`.so` or `.dll`) libraries. The `--library` and `--library-path` options can be used to pass libraries to both the device and host linker. The library name is specified without the library file extension when the `--library` option is used.

```
nvcc --gpu-architecture=sm_20 a.o b.o --library-path=<path> --library=foo
```

Alternatively, the library name, including the library file extension, can be used without the `--library` option on Windows.

```
nvcc --gpu-architecture=sm_20 a.obj b.obj foo.lib --library-path=<path>
```

Note that the device linker ignores any objects that do not have relocatable device code.

## 6.4. Examples

Suppose we have the following files:

```
//----- b.h -----
#define N 8

extern __device__ int g[N];

extern __device__ void bar(void);

//----- b.cu -----
#include "b.h"

__device__ int g[N];

__device__ void bar (void)
{
    g[threadIdx.x]++;
}

//----- a.cu -----
#include <stdio.h>
#include "b.h"

__global__ void foo (void) {

    __shared__ int a[N];
    a[threadIdx.x] = threadIdx.x;

    __syncthreads();

    g[threadIdx.x] = a[blockDim.x - threadIdx.x - 1];

    bar();
}

int main (void) {
    unsigned int i;
    int *dg, hg[N];
    int sum = 0;

    foo<<<1, N>>>();

    if(cudaGetSymbolAddress((void**)&dg, g)){
        printf("couldn't get the symbol addr\n");
        return 1;
    }
    if(cudaMemcpy(hg, dg, N * sizeof(int), cudaMemcpyDeviceToHost)){
        printf("couldn't memcpy\n");
        return 1;
    }

    for (i = 0; i < N; i++) {
        sum += hg[i];
    }
    if (sum == 36) {
        printf("PASSED\n");
    } else {
        printf("FAILED (%d)\n", sum);
    }

    return 0;
}
```

These can be compiled with the following commands (these examples are for Linux):

```
nvcc --gpu-architecture=sm_20 --device-c a.cu b.cu
nvcc --gpu-architecture=sm_20 a.o b.o
```

If you want to invoke the device and host linker separately, you can do:

```
nvcc --gpu-architecture=sm_20 --device-c a.cu b.cu
nvcc --gpu-architecture=sm_20 --device-link a.o b.o --output-file link.o
g++ a.o b.o link.o --library-path=<path> --library=cudart
```

Note that a target architecture must be passed to the device linker.

If you want to use the driver API to load a linked cubin, you can request just the cubin:

```
nvcc --gpu-architecture=sm_20 --device-link a.o b.o \
--cubin --output-file link.cubin
```

The objects could be put into a library and used with:

```
nvcc --gpu-architecture=sm_20 --device-c a.cu b.cu
nvcc --lib a.o b.o --output-file test.a
nvcc --gpu-architecture=sm_20 test.a
```

Note that only static libraries are supported by the device linker.

A PTX file can be compiled to a host object file and then linked by using:

```
nvcc --gpu-architecture=sm_20 --device-c a.ptx
```

An example that uses libraries, host linker, and dynamic parallelism would be:

```
nvcc --gpu-architecture=sm_35 --device-c a.cu b.cu
nvcc --gpu-architecture=sm_35 --device-link a.o b.o --output-file link.o
nvcc --lib --output-file libgpu.a a.o b.o link.o
g++ host.o --library=gpu --library-path=<path> \
--library=cudadevrt --library=cudart
```

It is possible to do multiple device links within a single host executable, as long as each device link is independent of the other. This requirement of independence means that they cannot share code across device executables, nor can they share addresses (e.g., a device function address can be passed from host to device for a callback only if the device link sees both the caller and potential callback callee; you cannot pass an address from one device executable to another, as those are separate address spaces).

## 6.5. Potential Separate Compilation Issues

### 6.5.1. Object Compatibility

Only relocatable device code with the same ABI version, same SM target architecture, and same pointer size (32 or 64) can be linked together. Incompatible objects will produce a link error. An object could have been compiled for a different architecture but also have PTX available, in which case the device linker will JIT the PTX to cubin for the desired architecture and then link. Relocatable device code requires CUDA 5.0 or later Toolkit.

If a kernel is limited to a certain number of registers with the `launch_bounds` attribute or the `--maxrregcount` option, then all functions that the kernel calls must not use more than that number of registers; if they exceed the limit, then a link error will be given.

## 6.5.2. JIT Linking Support

CUDA 5.0 does not support JIT linking, while CUDA 5.5 does. This means that to use JIT linking you must recompile your code with CUDA 5.5 or later. JIT linking means doing a relink of the code at startup time. The device linker (`nvlink`) links at the cubin level. If the cubin does not match the target architecture at load time, the driver re-invokes the device linker to generate cubin for the target architecture, by first JIT'ing the PTX for each object to the appropriate cubin, and then linking together the new cubin.

## 6.5.3. Implicit CUDA Host Code

A file like `b.cu` above only contains CUDA device code, so one might think that the `b.o` object doesn't need to be passed to the host linker. But actually there is implicit host code generated whenever a device symbol can be accessed from the host side, either via a launch or an API call like `cudaGetSymbolAddress()`. This implicit host code is put into `b.o`, and needs to be passed to the host linker. Plus, for JIT linking to work all device code must be passed to the host linker, else the host executable will not contain device code needed for the JIT link. So a general rule is that the device linker and host linker must see the same host object files (if the object files have any device references in them—if a file is pure host then the device linker doesn't need to see it). If an object file containing device code is not passed to the host linker, then you will see an error message about the function `__cudaRegisterLinkedBinary_name` calling an undefined or unresolved symbol `__fatbinwrap_name`.

## 6.5.4. Using `__CUDA_ARCH__`

In separate compilation, `__CUDA_ARCH__` must not be used in headers such that different objects could contain different behavior. Or, it must be guaranteed that all objects will compile for the same compute\_arch. If a weak function or template function is defined in a header and its behavior depends on `__CUDA_ARCH__`, then the instances of that function in the objects could conflict if the objects are compiled for different compute arch. For example, if an `a.h` contains:

```
template<typename T>
__device__ T* getptr(void)
{
    #if __CUDA_ARCH__ == 200
        return NULL; /* no address */
    #else
        __shared__ T arr[256];
        return arr;
    #endif
}
```

Then if `a.cu` and `b.cu` both include `a.h` and instantiate `getptr` for the same type, and `b.cu` expects a non-NULL address, and compile with:

```
nvcc --gpu-architecture=compute_20 --device-c a.cu  
nvcc --gpu-architecture=compute_30 --device-c b.cu  
nvcc --gpu-architecture=sm_30 a.o b.o
```

At link time only one version of the `getptr` is used, so the behavior would depend on which version is picked. To avoid this, either `a.cu` and `b.cu` must be compiled for the same compute arch, or `__CUDA_ARCH__` should not be used in the shared header function.

### 6.5.5. Device Code in Libraries

If a device function with non-weak external linkage is defined in a library as well as a non-library object, the device linker will complain about the multiple definitions (this differs from traditional host linkers that may ignore the function definition from the library object, if it was already found in an earlier object).



# Chapter 7.

## MISCELLANEOUS NVCC USAGE

### 7.1. Cross Compilation

Cross compilation is controlled by using the following **nvcc** command line options:

- ▶ **--compiler-bindir** is used for cross compilation, where the underlying host compiler is capable of generating objects for the target platform.
- ▶ **--machine=32**. This option signals that the target platform is a 32-bit platform. Use this when the host platform is a 64-bit platform.

### 7.2. Keeping Intermediate Phase Files

**nvcc** stores intermediate results by default into temporary files that are deleted immediately before it completes. The location of the temporary file directories used are, depending on the current platform, as follows:

#### Windows

Value of environment variable **TEMP** is used. If it is not set, **C:\Windows\temp** is used instead.

#### Other Platforms

Value of environment variable **TMPDIR** is used. If it is not set, **/tmp** is used instead.

Option **--keep** makes **nvcc** store these intermediate files in the current directory or in the directory specified by **--keep-dir** instead, with names as described in [Supported Phases](#).

### 7.3. Cleaning Up Generated Files

All files generated by a particular **nvcc** command can be cleaned up by repeating the command, but with additional option **--clean-targets**. This option is particularly useful after using **--keep**, because the **--keep** option usually leaves quite an amount of intermediate files around.

Because using `--clean-targets` will remove exactly what the original `nvcc` command created, it is important to exactly repeat all of the options in the original command. For instance, in the following example, omitting `--keep`, or adding `--compile` will have different cleanup effects.

```
nvcc acos.cu --keep
nvcc acos.cu --keep --clean-targets
```

## 7.4. Printing Code Generation Statistics

A summary on the amount of used registers and the amount of memory needed per compiled device function can be printed by passing option `--resource-usage` to `nvcc`:

```
$ nvcc --resource-usage acos.cu
ptxas info      : 1536 bytes gmem, 8 bytes cmem[14]
ptxas info      : Compiling entry function 'acos_main' for 'sm_20'
ptxas info      : Function properties for acos_main
    0 bytes stack frame, 0 bytes spill stores, 0 bytes spill loads
ptxas info      : Used 6 registers, 1536 bytes smem, 32 bytes cmem[0]
```

As shown in the above example, the amounts of statically allocated global memory (gmem) and constant memory in bank 14 (cmem) are listed.

Global memory and some of the constant banks are module scoped resources and not per kernel resources. Allocation of constant variables to constant banks is profile specific.

Followed by this, per kernel resource information is printed.

Stack frame is per thread stack usage used by this function. Spill stores and loads represent stores and loads done on stack memory which are being used for storing variables that couldn't be allocated to physical registers.

Similarly number of registers, amount of shared memory and total space in constant bank allocated is shown.

## **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 of 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 and the NVIDIA logo are trademarks 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**

© 2007-2016 NVIDIA Corporation. All rights reserved.