

# Blackwell Compatibility Guide Release 12.8

**NVIDIA Corporation** 

## **Contents**

1	About	t this Document	3	
2	2 Application Compatibility on Blackwell Architecture			
3	<b>Verify</b> 3.1 3.2	ring Blackwell Compatibility for Existing Applications Applications Built Using CUDA Toolkit 12.8 or Earlier	<b>7</b> 7 8	
4	<b>Buildi</b> 4.1 4.2 4.3	ng Applications with Blackwell Architecture Support  Building Applications Using CUDA Toolkit 12.7 or Earlier	10	
5	Revisi	ion History	13	
6	<b>Notic</b> 6.1 6.2 6.3	es           Notice           OpenCL           Trademarks	16	

#### **Blackwell Compatibility Guide for CUDA Applications**

The guide to building CUDA applications for Blackwell GPUs

Contents 1

2 Contents

## Chapter 1. About this Document

This application note, Blackwell Architecture Compatibility Guide for CUDA Applications, is intended to help developers ensure that their NVIDIA® CUDA® applications will run on the NVIDIA® Blackwell architecture based GPUs. This document provides guidance to developers who are familiar with programming in CUDA C++ and want to make sure that their software applications are compatible with Blackwell architecture.

# Chapter 2. Application Compatibility on Blackwell Architecture

A CUDA application binary (with one or more GPU kernels) can contain the compiled GPU code in two forms, binary cubin objects and forward-compatible PTX assembly for each kernel. Both cubin and PTX are generated for a certain target compute capability. A cubin generated for a certain compute capability is supported to run on any GPU with the same major revision and same or higher minor revision of compute capability. For example, a cubin generated for compute capability 8.0 is supported to run on a GPU with compute capability 8.6, however a cubin generated for compute capability 8.6 is *not* supported to run on a GPU with compute capability 8.0, and a cubin generated with compute capability 8.x is *not* supported to run on a GPU with compute capability 9.0.

Kernel can also be compiled to a PTX form. PTX is compiled at runtime to cubin and the cubin is used for kernel execution. Unlike cubin, PTX is forward-compatible. Meaning PTX is supported to run on any GPU with compute capability higher than the compute capability assumed for generation of that PTX. For example, PTX code generated for compute capability 9.x is supported to run on compute capability 9.x or any higher revision (major or minor), including compute capability 10.0. Therefore although it is optional, it is recommended that all applications should include PTX of the kernels to ensure forward-compatibility. To read more about cubin and PTX compatibilities see Compilation with NVCC from the CUDA C++ Programming Guide.

When a CUDA application launches a kernel on a GPU, the CUDA Runtime determines the compute capability of the GPU in the system and uses this information to find the best matching cubin or PTX version of the kernel. If a cubin compatible with that GPU is present in the binary, the cubin is used as-is for execution. Otherwise, the CUDA Runtime first generates compatible cubin by JIT-compiling the PTX and then the cubin is used for the execution. If neither compatible cubin nor PTX is available, kernel launch results in a failure.

Application binaries that include PTX version of kernels, should work as-is on the Blackwell GPUs. In such cases, rebuilding the application is not required. However application binaries which do not include PTX (only include cubins), need to be rebuilt to run on the Blackwell GPUs. To know more about building compatible applications read *Building Applications with Blackwell Architecture Support*.

Application binaries that include PTX version of kernels with architecture conditional features using sm\_100a or compute\_100a in order to take full advantage of Blackwell GPU architecture, are not forward or backward compatible. For example, PTX compiled for compute\_90a (Hopper) are not supported on the Blackwell architecture.

<sup>&</sup>lt;sup>1</sup> Just-in-time compilation.

Blackwell Compatibility Guide, Release 12.8						
	Ch	A	0 4! -! !4-	Dl l	. A	

# Chapter 3. Verifying Blackwell Compatibility for Existing Applications

The first step towards making a CUDA application compatible with Blackwell architecture is to check if the application binary already contains compatible GPU code (at least the PTX). The following sections explain how to accomplish this for an already built CUDA application.

# 3.1. Applications Built Using CUDA Toolkit 12.8 or Earlier

CUDA applications built using CUDA Toolkit versions 2.1 through 12.8 are compatible with Blackwell GPUs as long as they are built to include PTX versions of their kernels. This can be tested by forcing the PTX to JIT-compile at application load time with following the steps:

- ▶ Download and install the latest driver from <a href="https://www.nvidia.com/drivers">https://www.nvidia.com/drivers</a>.
- ▶ Set the environment variable CUDA\_FORCE\_PTX\_JIT=1.
- ► Launch the application.

With CUDA\_FORCE\_PTX\_JIT=1, GPU binary code embedded in an application binary is ignored. Instead PTX code for each kernel is JIT-compiled to produce GPU binary code. An application fails to execute if it does not include PTX. This means the application is not Blackwell architecture compatible and needs to be rebuilt for compatibility. On the other hand, if the application works properly with this environment variable set, then the application is Blackwell compatible.

**Note:** Be sure to unset the CUDA\_FORCE\_PTX\_JIT environment variable after testing is done.

### 3.2. Applications Built Using CUDA Toolkit 12.8

CUDA applications built using CUDA Toolkit 12.8 are compatible with Blackwell architecture as long as they are built to include kernels in native cubin (compute capability 10.0) or PTX form or both.

# Chapter 4. Building Applications with Blackwell Architecture Support

Depending on the version of the CUDA Toolkit used for building the application, it can be built to include PTX and/or native cubin for the Blackwell architecture. Although it is enough to just include PTX, including native cubin is can avoid the need to JIT compile the PTX at runtime.<sup>2</sup>

# 4.1. Building Applications Using CUDA Toolkit 12.7 or Earlier

The nvcc compiler included with version 12.7 or earlier (11.8-12.7) of the CUDA Toolkit can generate cubins native to the NVIDIA Hopper GPU architectures (compute capability 9.x). When using CUDA Toolkit 12.7 or earlier, to ensure that nvcc will generate cubin files for all recent GPU architectures as well as a PTX version for forward compatibility with future GPU architectures, specify the appropriate –gencode= parameters on the nvcc command line as shown in the examples below.

#### **Windows**

```
nvcc.exe -ccbin "C:\vs2010\VC\bin"
  -Xcompiler "/EHsc /W3 /nologo /02 /Zi /MT"
  -gencode=arch=compute_52,code=sm_52
  -gencode=arch=compute_60,code=sm_60
  -gencode=arch=compute_61,code=sm_61
  -gencode=arch=compute_70,code=sm_70
  -gencode=arch=compute_75,code=sm_75
  -gencode=arch=compute_80,code=sm_80
  -gencode=arch=compute_90,code=sm_90
  -gencode=arch=compute_90,code=compute_90
  --compile -o "Release\mykernel.cu.obj" "mykernel.cu"
```

#### Linux

```
/usr/local/cuda/bin/nvcc
-gencode=arch=compute_52,code=sm_52
-gencode=arch=compute_60,code=sm_60

(continues on next page)
```

<sup>&</sup>lt;sup>2</sup> The CUDA driver caches the cubins generated as a result of the PTX JIT, so this is often a one-time cost.

(continued from previous page)

```
-gencode=arch=compute_61,code=sm_61
-gencode=arch=compute_70,code=sm_70
-gencode=arch=compute_75,code=sm_75
-gencode=arch=compute_80,code=sm_80
-gencode=arch=compute_90,code=sm_90
-gencode=arch=compute_90,code=compute_90
-02 -o mykernel.o -c mykernel.cu
```

Alternatively, the simplified nvcc command-line option -arch=sm\_XX can be used. It is a shorthand equivalent to the following more explicit -gencode= command-line options used above. -arch=sm\_XX expands to the following:

```
-gencode=arch=compute_XX,code=sm_XX
-gencode=arch=compute_XX,code=compute_XX
```

However, while the <code>-arch=sm\_XX</code> command-line option does result in inclusion of a PTX back-end target binary by default, it can only specify a single target cubin architecture at a time, and it is not possible to use multiple <code>-arch=</code> options on the same nvcc command line, which is why the examples above use <code>-gencode=</code> explicitly.

For CUDA toolkits prior to 11.0, one or more of the -gencode options need to be removed according to the architectures supported by the specific toolkit version (for example, CUDA toolkit 10.x supports architectures up to sm\_72 and sm\_75). The final -gencode to generate PTX also needs to be updated. For further information and examples see the documentation for the specific CUDA toolkit version.

**Note:** compute\_XX refers to a PTX version and sm\_XX refers to a cubin version. The arch= clause of the -gencode= command-line option to nvcc specifies the front-end compilation target and must always be a PTX version. The code= clause specifies the back-end compilation target and can either be cubin or PTX or both. **Only the back-end target version(s) specified by the code= clause will be retained in the resulting binary; at least one should be PTX to provide compatibility with future architectures.** 

# 4.2. Building Applications Using CUDA Toolkit 12.8

With versions 12.8 of the CUDA Toolkit, nvcc can generate cubin native to the Blackwell architecture (compute capability 10.0). When using CUDA Toolkit 12.8, to ensure that nvcc will generate cubin files for all recent GPU architectures as well as a PTX version for forward compatibility with future GPU architectures, specify the appropriate -gencode= parameters on the nvcc command line as shown in the examples below.

#### Windows

```
nvcc.exe -ccbin "C:\vs2010\VC\bin"
  -Xcompiler "/EHsc /W3 /nologo /02 /Zi /MT"
  -gencode=arch=compute_52,code=sm_52
  -gencode=arch=compute_60,code=sm_60
  -gencode=arch=compute_61,code=sm_61
  -gencode=arch=compute_70,code=sm_70
```

(continues on next page)

(continued from previous page)

```
-gencode=arch=compute_75,code=sm_75
-gencode=arch=compute_75,code=sm_75
-gencode=arch=compute_90,code=sm_90
-gencode=arch=compute_100,code=sm_100
-gencode=arch=compute_100,code=compute_100
--compile -o "Release\mykernel.cu.obj" "mykernel.cu"
```

#### Linux

```
/usr/local/cuda/bin/nvcc
-gencode=arch=compute_52,code=sm_52
-gencode=arch=compute_60,code=sm_60
-gencode=arch=compute_61,code=sm_61
-gencode=arch=compute_70,code=sm_70
-gencode=arch=compute_75,code=sm_75
-gencode=arch=compute_80,code=sm_80
-gencode=arch=compute_90,code=sm_90
-gencode=arch=compute_100,code=sm_100
-gencode=arch=compute_100,code=compute_100
-02 -o mykernel.o -c mykernel.cu
```

**Note:** compute\_XX refers to a PTX version and sm\_XX refers to a cubin version. The arch= clause of the -gencode= command-line option to nvcc specifies the front-end compilation target and must always be a PTX version. The code= clause specifies the back-end compilation target and can either be cubin or PTX or both. **Only the back-end target version(s) specified by the code= clause will be retained in the resulting binary; at least one should be PTX to provide compatibility with future architectures.** 

# 4.3. Independent Thread Scheduling Compatibility

NVIDIA GPUs since Volta architecture have Independent Thread Scheduling among threads in a warp. If the developer made assumptions about warp-synchronicity<sup>3</sup>, this feature can alter the set of threads participating in the executed code compared to previous architectures. Please see Compute Capability 7.x in the CUDA C++ Programming Guide for details and corrective actions. To aid migration to the Blackwell architecture, developers can opt-in to the Pascal scheduling model with the following combination of compiler options.

```
nvcc -gencode=arch=compute_60,code=sm_100 ...
```

<sup>&</sup>lt;sup>3</sup> Warp-synchronous refers to an assumption that threads in the same warp are synchronized at every instruction and can, for example, communicate values without explicit synchronization.

Blackwell Compatibility Guide, Release 12.8						
12	Chantar 4	Duilding A		with Plackwe	all Aughitaat.	una Cumma est

# Chapter 5. Revision History

#### Version 1.0

► Initial public release.

### Chapter 6. Notices

### 6.1. Notice

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

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

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

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

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

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

No license, either expressed or implied, is granted under any NVIDIA patent right, copyright, or other NVIDIA intellectual property right under this document. Information published by NVIDIA regarding third-party products or services does not constitute a license from NVIDIA to use such products or

services or a warranty or endorsement thereof. Use of such information may require a license from a third party under the patents or other intellectual property rights of the third party, or a license from NVIDIA under the patents or other intellectual property rights of NVIDIA.

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

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

### 6.2. OpenCL

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

### 6.3. 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

©2025, NVIDIA Corporation & affiliates. All rights reserved