



# SANITIZER API

v2024.2.0 | April 2024

## Reference Manual



# TABLE OF CONTENTS

<b>Chapter 1. Introduction.....</b>	<b>1</b>
1.1. Overview.....	1
<b>Chapter 2. Usage.....</b>	<b>2</b>
2.1. Compatibility and Requirements.....	2
2.2. Callback API.....	2
2.2.1. Driver and Runtime API Callbacks.....	3
2.2.2. Resource Callbacks.....	4
2.2.3. Synchronization Callbacks.....	4
2.2.4. Launch Callbacks.....	5
2.2.5. Memcpy Callbacks.....	5
2.2.6. Memset Callbacks.....	5
2.2.7. Batch Memory Operations Callbacks.....	5
2.3. Patching API.....	5
2.3.1. Writing a Patch.....	6
2.3.2. Insert a Patch.....	6
2.4. Memory API.....	7
2.5. Special cases.....	8
2.5.1. Device graph launch.....	8
<b>Chapter 3. Limitations.....</b>	<b>9</b>

# Chapter 1.

## INTRODUCTION

### 1.1. Overview

The Compute Sanitizer API enables the creation of sanitizing and tracing tools that target CUDA applications. Examples of such tools are memory and race condition checkers. The Compute Sanitizer API is composed of three APIs: the callback API, the patching API and the memory API. It is delivered as a dynamic library on supported platforms.

# Chapter 2.

## USAGE

### 2.1. Compatibility and Requirements

The Compute Sanitizer tools require CUDA 11.0 or newer.

The Compute Sanitizer API requires CUDA 10.1 or newer. Compute Sanitizer API calls will fail with **SANITIZER\_ERROR\_NOT\_INITIALIZED** if the CUDA driver version is not compatible with the Compute Sanitizer version.

### 2.2. Callback API

The Compute Sanitizer Callback API allows you to register a callback into user code. The callback is invoked when the application calls a CUDA runtime or driver function, or when certain events occur in the CUDA driver. The following terminology is used by the Callback API.

- ▶ **Callback domain:** Callbacks are grouped into domains to make it easier to associate callback functions with groups of related CUDA functions or events. The following callback domains are defined by **Sanitizer\_CallbackDomain**.
  1. CUDA driver functions
  2. CUDA runtime functions
  3. CUDA resource tracking
  4. CUDA synchronization notification
  5. CUDA grid launches
  6. CUDA memcpy operations
  7. CUDA memset operations
  8. CUDA batch memory operations
- ▶ **Callback ID:** Each callback is given a unique ID within the corresponding callback domain in order to identify it within the callback function. The CUDA driver API IDs are defined in **sanitizer\_driver\_cbid.h** and the CUDA runtime API IDs are defined in **sanitizer\_runtime\_cbid.h**. Other callback IDs are

defined in `sanitizer_callbacks.h`. All of these headers are included as part of `sanitizer.h`.

- ▶ **Callback Function:** The callback function must be of the type `Sanitizer_CallbackFunc`. This function type has two arguments that specify the callback: the domain and the ID that identifies why the callback is occurring. The type also has a `cbdata` argument that is used to pass data specific to the callback.
- ▶ **Subscriber:** A subscriber is used to associate each of the callback functions with one or more CUDA API functions. There can be at most one subscriber initialized with `sanitizerSubscribe` at any time. Before initializing a new subscriber, the existing one must be finalized with `sanitizerUnsubscribe`.

The subscriber should be initialized prior to making any CUDA API call to ensure correctness of the reported data.

Each callback domain is described in detail below. Unless explicitly stated, it is not supported to call any CUDA runtime or driver API from within a callback function. Doing so may cause the application to hang. However, it is supported to call [Compute Sanitizer Memory APIs](#) from within callback functions.

## 2.2.1. Driver and Runtime API Callbacks

Using the Callback API with the `SANITIZER_CB_DOMAIN_DRIVER_API` or `SANITIZER_CB_DOMAIN_RUNTIME_API` domains, a callback function can be associated with one or more CUDA API functions. When those CUDA functions are called in the application, the callback function is invoked as well. For these domains, the `cbdata` argument to the callback function will be of the type `Sanitizer_CallbackData`.

You can call `cudaDeviceSynchronize`, `cudaStreamSynchronize`, `cuCtxSynchronize` and `cuStreamSynchronize` from within a driver or runtime API callback function.

The following code shows a typical sequence used to associate a callback function with one or more CUDA API functions. For simplicity, error checking code was removed.

```
Sanitizer_SubscriberHandle handle;
MyDataStruct *my_data = ...;
...
sanitizerSubscribe(&handle, my_callback, my_data);
sanitizerEnableDomain(1, handle, SANITIZER_CB_DOMAIN_RUNTIME_API);
```

First, `sanitizerSubscribe` is used to initialize a subscriber with the `my_callback` callback function. Next, `sanitizerEnableDomain` is used to associate that callback with all the CUDA runtime functions. Using this code sequence will cause `my_callback` to be called twice each time any of the CUDA runtime API functions are invoked, once on entry to the CUDA function and once just before the CUDA function exits. Compute Sanitizer callback API functions `sanitizerEnableCallback` and `sanitizerEnableAllDomains` can also be used to associate CUDA API functions with a callback.

The following code shows a typical callback function.

```
void SANITIZERAPI
my_callback(void *userdata,
            Sanitizer_CallbackDomain domain,
            Sanitizer_CallbackId cbid,
            const void *cbdata)
{
    const Sanitizer_CallbackData *cbInfo = (Sanitizer_CallbackData *)cbdata;
    MyDataStruct *my_data = (MyDataStruct *)userdata;

    if ((domain == SANITIZER_CB_DOMAIN_RUNTIME_API) &&
        (cbid == SANITIZER_RUNTIME_TRACE_CBID_cudaMemcpy_v3020) &&
        (cbInfo->callbackSite == SANITIZER_API_ENTER))
    {
        cudaMemcpy_v3020_params *funcParams = (cudaMemcpy_v3020_params *)
(cbInfo->functionParams);
        size_t count = funcParams->count;
        enum cudaMemcpyKind kind = funcParams->kind
        ...
    }
    ...
}
```

In the callback function, **Sanitizer\_CallbackDomain** and **Sanitizer\_CallbackId** parameters can be used to determine which CUDA API function invocation is triggering this callback. In the example above, we are checking for the CUDA runtime **cudaMemcpy** function. The **cbdata** parameter holds a structure of useful information that can be used within the callback. In this case, we use the **callbackSite** member of the structure to detect that the callback is occurring on entry to **cudaMemcpy**, and we use the **functionParams** member to access the parameters to **cudaMemcpy**. To access the parameters, we first cast **functionParams** to a structure type corresponding to the **cudaMemcpy** function. These parameter structures are contained in **generated\_cuda\_runtime\_api\_meta.h**, **generated\_cuda\_meta.h** and a number of other files.

## 2.2.2. Resource Callbacks

Using the Callback API with the **SANITIZER\_CB\_DOMAIN\_RESOURCE** domain, a callback function can be associated with some CUDA resource creation and destruction events. For example, when a CUDA context is created, the callback function is invoked with a callback ID equal to **SANITIZER\_CBID\_RESOURCE\_CONTEXT\_CREATED**. For this domain, the **cbdata** argument is one of the following types:

- ▶ **Sanitizer\_ResourceContextData** for CUDA context creation and destruction
- ▶ **Sanitizer\_ResourceStreamData** for CUDA stream creation and destruction
- ▶ **Sanitizer\_ResourceModuleData** for CUDA module load and unload
- ▶ **Sanitizer\_ResourceMemoryData** for CUDA memory allocation and de-allocation

## 2.2.3. Synchronization Callbacks

Using the Callback API with the **SANITIZER\_CB\_DOMAIN\_SYNCHRONIZE** domain, a callback function can be associated with CUDA context and stream synchronizations. For example, when a CUDA context is synchronized, the callback function is invoked

with a callback ID equal to **SANITIZER\_CBID\_SYNCHRONIZE\_CONTEXT\_SYNCHRONIZED**. For this domain, the **cbdata** argument is of the type **Sanitizer\_SynchronizeData**.

## 2.2.4. Launch Callbacks

Using the Callback API with the **SANITIZER\_CB\_DOMAIN\_LAUNCH** domain, a callback function can be associated with CUDA kernel launches. For example, when a CUDA kernel launch has started, the callback function is invoked with a callback ID equal to **SANITIZER\_CBID\_LAUNCH\_BEGIN**. For this domain, the **cbdata** argument is of the type **Sanitizer\_LaunchData**.

## 2.2.5. Memcpy Callbacks

Using the Callback API with the **SANITIZER\_CB\_DOMAIN\_MEMCPY** domain, a callback function can be associated with CUDA memcpy operations. For example, when a **cudaMemcpy** API function is called, the callback function is invoked with a callback ID equal to **SANITIZER\_CBID\_MEMCPY\_STARTING**. For this domain, the **cbdata** argument is of the type **Sanitizer\_MemcpyData**.

## 2.2.6. Memset Callbacks

Using the Callback API with the **SANITIZER\_CB\_DOMAIN\_MEMSET** domain, a callback function can be associated with CUDA memset operations. For example, when a **cudaMemset** API function is called, the callback function is invoked with a callback ID equal to **SANITIZER\_CBID\_MEMSET\_STARTING**. For this domain, the **cbdata** argument is of the type **Sanitizer\_MemsetData**.

## 2.2.7. Batch Memory Operations Callbacks

Using the Callback API with the **SANITIZER\_CB\_DOMAIN\_BATCH\_MEMOP** domain, a callback function can be associated with CUDA batch memory operations. For example, when a **cuStreamWriteValue** API function is called, the callback function is invoked with a callback ID equal to **SANITIZER\_CBID\_BATCH\_MEMOP\_WRITE**. For this domain, the **cbdata** argument is of the type **Sanitizer\_BatchMemopData**.

## 2.3. Patching API

The Compute Sanitizer Patching API allows you to load patch functions and insert them into user code. Patch functions will be invoked when the application's CUDA code executes certain instructions or calls certain CUDA device functions. The following terminology is used by the Patching API:

- ▶ **Instruction ID:** Each patchable event is given a unique ID than can be passed to patching API functions to specify that these events should be patched. Instruction IDs are defined by **Sanitizer\_InstructionId**.
- ▶ **Instrumentation point:** A location in the original CUDA code that is being instrumented by the Compute Sanitizer API. Upon execution, the user code path is

modified so that a patch gets executed either before or after the patched event. All patches are executed prior to the event, with the exception of device-side malloc.

- **Patch:** A CUDA `__device__` function that the Compute Sanitizer will insert into another existing CUDA code. Patch function signatures must match the one expected by the API (see below for the expected signature types).

### 2.3.1. Writing a Patch

The patch must follow the function signature required by the Compute Sanitizer API for a given instruction ID. The mapping of instruction ID to function signature is documented in the comments of `Sanitizer_InstructionId` in `sanitizer_patching.h`. For instance, if we wish to patch a memory access using the instruction ID `SANITIZER_INSTRUCTION_MEMORY_ACCESS`, we need to use the `SanitizerCallbackMemoryAccess` type.

```
extern "C" __device__
SanitizerPatchResult SANITIZERAPI my_memory_access_callback(
    void* userdata,
    uint64_t pc,
    void* ptr,
    uint32_t accessSize,
    uint32_t flags)
{
    MyDeviceDataStruct *my_data = (MyDeviceDataStruct *)userdata

    if ((flags & SANITIZER_MEMORY_DEVICE_FLAG_WRITE) != 0)
        // log write
    else
        // log read

    return SANITIZER_PATCH_SUCCESS;
}
```

In this patch, we log write and read accesses to a structure we allocated previously. **extern "C"** ensures that the patch name will not be mangled, allowing us to use its name as a string directly in calls to `sanitizerPatchInstructions` (see [below](#)).

There can be multiple patches defined in a single CUDA file. This file must then be compiled using the following nvcc options:

```
$ nvcc --cubin --compile-as-tools-patch MySanitizerPatches.cu -o
  MySanitizerPatches.cubin
```

The `--cubin` option can be replaced by `--fatbin` if a fatbin is preferred over a cubin as the output file.

### 2.3.2. Insert a Patch

Once the patch has been generated, it can be inserted into user code by using the following procedure:

1. **Load the patch.** There are two APIs used to load the patch: `sanitizerAddPatchesFromFile` and `sanitizerAddPatches`. They use the same input format as `cuModuleLoad` and `cuModuleLoadData`, respectively.



2. Select which instructions to patch by using the **sanitizerPatchInstructions** API.
3. Patch user code by using the **sanitizerPatchModule** API.
4. Optionally, set the callback data for patches by using the **sanitizerSetCallbackData** API.

The following code shows a typical sequence using these APIs. For simplicity, error checking was removed.

```
CUcontext ctx = ... // current CUDA context
sanitizerAddPatchesFromFile("MySanitizerPatches.cubin", ctx);

CUmodule module = ... // module containing the user code
sanitizerPatchInstructions(SANITIZER_INSTRUCTION_MEMORY_ACCESS, module,
    "my_memory_access_callback");

sanitizerPatchModule(module);

MyDeviceDataTracker *deviceDataTracker;
cudaMalloc(&deviceDataTracker, sizeof(*deviceDataTracker));

CUfunction function = ... // kernel to be launched for which we want to set the
    callbackdata for the patches
sanitizerSetCallbackData(function, deviceDataTracker);
```

All subsequent launches using code from this CUDA module will be instrumented and **my\_memory\_access\_callback** will be invoked before every memory access. However, the callback data is only set for all subsequent launches of the given kernel. An easy way to have a kernel **CUfunction**, is through the Sanitizer launch callbacks. Instrumentation can be removed by using the **sanitizerUnpatchModule** API.

## 2.4. Memory API

The Compute Sanitizer Memory API provides replacement functions for the CUDA Memory API that can be safely called from within Compute Sanitizer **callbacks**.

- ▶ **sanitizerAlloc** is a replacement for **cudaMalloc**.
- ▶ **sanitizerFree** is a replacement for **cudaFree**.
- ▶ **sanitizerMemcpyHostToDeviceAsync** is a replacement for **cudaMemcpyAsync** for host-to-device copies.
- ▶ **sanitizerMemcpyDeviceToHost** is a replacement for **cudaMemcpy** for device-to-host copies.
- ▶ **sanitizerMemset** is a replacement for **cudaMemset**.

These functions can also be called in normal user code, where they can be mixed with the CUDA API. For instance, memory allocated with **sanitizerAlloc** can be freed with **cudaFree**. However, since only CUDA API calls will cause **callbacks** to be invoked, this can lead to an incoherent tracking state and should be avoided.

## 2.5. Special cases

In some specific scenarios, the Compute Sanitizer API may behave differently from the general cases. These scenarios are listed in this section.

### 2.5.1. Device graph launch

When the application instrumented by the Compute Sanitizer API makes use of device launched CUDA graphs, the following behaviors are to be taken into account:

- ▶ No callbacks will be called from the host when a CUDA graph is launched from the device. However, graph launch-related callbacks are still called when a device launchable graph is launched from the host, or uploaded.
- ▶ Patching API issued callbacks on the device when inside a device launched CUDA graph may have a userdata from a different node if different userdata pointers were set per node.

To compensate these limitations, it is possible to set a device launched graph specific data using **sanitizerSetDeviceGraphData**. It can be called during the launch on the host of a device launchable graph or of a graph containing device graph launches. The data set with **sanitizerSetDeviceGraphData** will then be retrievable from any device graph launched from this host launched graph. To retrieve this data from a patching API device callback, the following code can be used:

```
void* userdata = *((void***)cudaGetCurrentGraphExec());
```

It is however recommended to perform NULL checks at each dereferencing.

## Chapter 3. LIMITATIONS

No known issues at this time.

## **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**

© 2019-2024 NVIDIA Corporation and affiliates. All rights reserved.

This product includes software developed by the Syncro Soft SRL (<http://www.sync.ro/>).