

torch.cuda

This package adds support for CUDA tensor types.

It implements the same function as CPU tensors, but they utilize GPUs for computation.

It is lazily initialized, so you can always import it, and use `is_available()` to determine if your system supports CUDA.

CUDA semantics has more details about working with CUDA.

| | |
|-------------------------------------|--|
| <code>StreamContext</code> | Context-manager that selects a given stream. |
| <code>can_device_access_peer</code> | Check if peer access between two devices is possible. |
| <code>current_blas_handle</code> | Return cublasHandle_t pointer to current cuBLAS handle. |
| <code>current_device</code> | Return the index of a currently selected device. |
| <code>current_stream</code> | Return the currently selected <code>Stream</code> for a given device. |
| <code>cuda</code> | Retrieves the CUDA runtime API module. |
| <code>default_stream</code> | Return the default <code>Stream</code> for a given device. |
| <code>device</code> | Context-manager that changes the selected device. |
| <code>device_count</code> | Return the number of GPUs available. |
| <code>device_memory_used</code> | Return usage of (device) memory in bytes as given by <i>nvidia-smi</i> or <i>amd-smi</i> . |
| <code>device_of</code> | Context-manager that changes the current device to that of given object. |
| <code>get_arch_list</code> | Return list of CUDA architectures this library was compiled for. |
| <code>get_device_capability</code> | Get the CUDA capability of a device. |
| <code>get_device_name</code> | Get the name of a device. |
| <code>get_device_properties</code> | Get the properties of a device. |
| <code>get_gencode_flags</code> | Return NVCC generate flags this library was compiled with. |
| <code>get_sync_debug_mode</code> | Return current value of debug mode for CUDA synchronizing operations. |
| <code>init</code> | Initialize PyTorch's CUDA state. |
| <code>ipc_collect</code> | Force collect GPU memory after it has been released by CUDA IPC. |
| <code>is_available</code> | Return a bool indicating if CUDA is currently available. |
| <code>is_initialized</code> | Return whether PyTorch's CUDA state has been initialized. |
| <code>memory_usage</code> | Return the percent of time over the past sample period during which global (device) memory was being read or written as given by <i>nvidia-smi</i> . |
| <code>set_device</code> | Set the current device. |
| <code>set_stream</code> | Set the current stream. This is a wrapper API to set the stream. |
| <code>set_sync_debug_mode</code> | Set the debug mode for CUDA synchronizing operations. |



| | |
|-------------------------------|--|
| <code>synchronize</code> | Wait for all kernels in all streams on a CUDA device to complete. |
| <code>utilization</code> | Return the percent of time over the past sample period during which one or more kernels was executing on the GPU as given by <i>nvidia-smi</i> . |
| <code>temperature</code> | Return the average temperature of the GPU sensor in Degrees C (Centigrades). |
| <code>power_draw</code> | Return the average power draw of the GPU sensor in mW (MilliWatts) |
| <code>clock_rate</code> | Return the clock speed of the GPU SM in Hz Hertz over the past sample period as given by <i>nvidia-smi</i> . |
| <code>OutOfMemoryError</code> | Exception raised when device is out of memory |

Random Number Generator

| | |
|--------------------------------|--|
| <code>get_rng_state</code> | Return the random number generator state of the specified GPU as a ByteTensor. |
| <code>get_rng_state_all</code> | Return a list of ByteTensor representing the random number states of all devices. |
| <code>set_rng_state</code> | Set the random number generator state of the specified GPU. |
| <code>set_rng_state_all</code> | Set the random number generator state of all devices. |
| <code>manual_seed</code> | Set the seed for generating random numbers for the current GPU. |
| <code>manual_seed_all</code> | Set the seed for generating random numbers on all GPUs. |
| <code>seed</code> | Set the seed for generating random numbers to a random number for the current GPU. |
| <code>seed_all</code> | Set the seed for generating random numbers to a random number on all GPUs. |
| <code>initial_seed</code> | Return the current random seed of the current GPU. |

Communication collectives

| | |
|---------------------------------------|--|
| <code>comm.broadcast</code> | Broadcasts a tensor to specified GPU devices. |
| <code>comm.broadcast_coalesced</code> | Broadcast a sequence of tensors to the specified GPUs. |
| <code>comm.reduce_add</code> | Sum tensors from multiple GPUs. |
| <code>comm.scatter</code> | Scatters tensor across multiple GPUs. |
| <code>comm.gather</code> | Gathers tensors from multiple GPU devices. |

Streams and events

| | |
|-----------------------------|---|
| <code>Stream</code> | Wrapper around a CUDA stream. |
| <code>ExternalStream</code> | Wrapper around an externally allocated CUDA stream. |
| <code>Event</code> | Wrapper around a CUDA event. |

Graphs (eta)

| | |
|--|--|
| <code>is_current_stream_capturing</code> | Return True if CUDA graph capture is underway on the current CUDA stream, False otherwise. |
| <code>graph_pool_handle</code> | Return an opaque token representing the id of a graph memory pool. |
| <code>CUDAGraph</code> | Wrapper around a CUDA graph. |

Context manager that captures CUDA work into a graph



`make_graphed_callables`

versions.

Memory management

`empty_cache`

Release all unoccupied cache memory currently held by the caching allocator so that those can be used in other GPU application and visible in *nvidia-smi*.

`get_per_process_memory_fraction`

Get memory fraction for a process.

`list_gpu_processes`

Return a human-readable printout of the running processes and their GPU memory use for a given device.

`mem_get_info`

Return the global free and total GPU memory for a given device using `cudaMemGetInfo`.

`memory_stats`

Return a dictionary of CUDA memory allocator statistics for a given device.

`memory_summary`

Return a human-readable printout of the current memory allocator statistics for a given device.

`memory_snapshot`

Return a snapshot of the CUDA memory allocator state across all devices.

`memory_allocated`

Return the current GPU memory occupied by tensors in bytes for a given device.

`max_memory_allocated`

Return the maximum GPU memory occupied by tensors in bytes for a given device.

`reset_max_memory_allocated`

Reset the starting point in tracking maximum GPU memory occupied by tensors for a given device.

`memory_reserved`

Return the current GPU memory managed by the caching allocator in bytes for a given device.

`max_memory_reserved`

Return the maximum GPU memory managed by the caching allocator in bytes for a given device.

`set_per_process_memory_fraction`

Set memory fraction for a process.

`memory_cached`

Deprecate ; see `memory_reserved()` .

`max_memory_cached`

Deprecate ; see `max_memory_reserved()` .

`reset_max_memory_cached`

Reset the starting point in tracking maximum GPU memory managed by the caching allocator for a given device.

`reset_peak_memory_stats`

Reset the “peak” stats tracked by the CUDA memory allocator.

`caching_allocator_alloc`

Perform a memory allocation using the CUDA memory allocator.

`caching_allocator_delete`

Delete memory allocated using the CUDA memory allocator.

`get_allocator_backend`

Return a string describing the active allocator backend as set by `PYTORCH_CUDA_ALLOC_CONF` .

`CUDAPluggableAllocator`

CUDA memory allocator loaded from a .so file.

`change_current_allocator`

Change the currently used memory allocator to be the one provided .

`MemPool`

MemPool represents a pool of memory in a caching allocator.

`MemPoolContext`

MemPoolContext holds the currently active pool and stashes the previous pool.

`caching_allocator_enable`

Enable or disable the CUDA memory allocator.

CLASS torch.cuda.use_mem_pool(*pool*, *device=None*) [SO RCE]

A context manager that routes allocations to a given pool.

Parameters

- ool** (*torch.cuda.MemPool*) – a MemPool object to be made active so that allocations route to this pool.

NVIDIA Tools Extension (NVTX)

| | |
|------------------------------|--|
| <code>nvtx.mark</code> | Describe an instantaneous event that occurred at some point. |
| <code>nvtx.range_push</code> | Push a range onto a stack of nested range spans. |
| <code>nvtx.range_pop</code> | Pop a range off of a stack of nested range spans. |
| <code>nvtx.range</code> | Context manager / decorator that pushes an NVTX range at the beginning of its scope, and pops it at the end. |

Jiterator (Jeta)

| | |
|--|--|
| <code>jiterator._create_jit_fn</code> | Create a jiterator-generate CUDA kernel for an elementwise op. |
| <code>jiterator._create_multi_output_jit_fn</code> | Create a jiterator-generate CUDA kernel for an elementwise op that supports returning one or more outputs. |

Tuna leO

Some operations could be implemented using more than one library or more than one technique. For example, a GEMM could be implemented for CUDA or ROCm using either the cuDNN/cuBLASLt libraries or hipBLAS/hipBLASLt libraries, respectively. How does one know which implementation is the fastest and should be chosen? That's what Tuna leOp provides. Certain operators have been implemented using multiple strategies as Tuna le Operators. At runtime, all strategies are profiled and the fastest is selected for all subsequent operations.

See the [documentation](#) for information on how to use it.

Stream Sanitizer (Strotty)

CUDA Sanitizer is a prototype tool for detecting synchronization errors between streams in PyTorch. See the [documentation](#) for information on how to use it.

[← Previous](#)

[Next →](#)

Docs

Access comprehensive developer documentation for PyTorch

[View Docs →](#)

Tutorials

Get in-depth tutorials for beginners and advanced developers

[View Tutorials →](#)

Resources

Find development resources and get your questions answered

[View Resources →](#)



PyTorch

[Get Started](#)

[Features](#)

[Ecosystem](#)

[Blog](#)

[Contributing](#)

Resources

[Tutorials](#)

[Docs](#)

[Discuss](#)

[GitHub Issues](#)

[Brand Guidelines](#)

Stay updated

[Facebook](#)

[Twitter](#)

[YouTube](#)

[LinkedIn](#)

PyTorch Podcasts

[Subscribe](#)

[Apple](#)

[Google](#)

[Amazon](#)

[Terms](#) | [Privacy](#)

© Copyright The Linux Foundation. The PyTorch Foundation is a project of The Linux Foundation. For website terms of use, trademark policy and other policies applicable to The PyTorch Foundation, please see www.linuxfoundation.org/policies/. The PyTorch Foundation supports the PyTorch open source project, which has been established as PyTorch Project a Series of LF Projects, LLC. For policies applicable to the PyTorch Project a Series of LF Projects, LLC, please see www.lfprojects.org/policies/.