

CuPy and CuPy-Xarray:

Presenter: Giacomo Capodaglio Oct 14th 2025: AMD @ CASTIEL



What is CuPy

- NumPy is a Python interface to optimized routines written in C that provide arrays, multi-dimensional arrays and common numerical operations on them. These are much faster than operating on Python lists
- SciPy provides fundamental algorithms common in scientific and numerical computing. The underlying code is a mixture of Fortran, C and C++
- CuPy is a NumPy/SciPy-compatible array library for GPU-accelerated computing with Python. It is not
 an Nvidia product, despite the Nvidia sounding name.
- CuPy acts as a drop-in replacement to run existing NumPy/SciPy code on Nvidia CUDA or AMD ROCm™ platforms
- CuPy provides the N-dimensional array (ndarray), sparse matrices, and the associated routines for GPU devices, most having the same API as NumPy and SciPy
- CuPy provides interfaces to GPU optimized libraries such as rocBLAS, rocSPARSE, rocFFT, and RCCL

source: <u>cupy documentation</u>



CuPy and HIP

- CuPy uses HIP as backhand to run on AMD GPUs
- HIP: Heterogeneous-compute Interface for Portability
 - C++ runtime API and kernel language
 - Works on AMD and Nvidia GPUs
- The CPU is often referred to as the host, and the GPU as the device
- In HIP, launching a kernel is non-blocking for the host
 - After sending instructions/data, the host continues to do more work while the device executes the kernel. This means GPU execution and CPU activity can overlap
- What it means for CuPy: appropriate synchronization calls have to be made after a kernel call:
 - cupy.cuda.Device(0).synchronize()
 - cupy.cuda.Stream.synchronize()
- In HIP, memory copies such as hipMemcpy is blocking for the host
 - All activity on the host stops until the copy has completed.
- What it means for CuPy: no need to sync if calling a memory copy right after a kernel.

click **here** for differences between CuPy and NumPy

CuPy functions

CuPy vs NumPy API

CuPy-specific functions

Comparison Table # Here is a list of NumPy / SciPy APIs and its corresponding CuPy implementations. in CuPy column denotes that CuPy implementation is not provided yet. We welcome contributions for these functions. NumPy / CuPy APIs Module-Level CuPy NumPy cupy.DataSource (alias of numpy.DataSource) numpy.DataSource numpy.ScalarType cupy.abs numpy.abs numpy.absolute cupy.absolute cupy.add numpy.add numpy.all cupy.all

CuPy-specific functions CuPy-specific functions are placed under cupyx namespace.	
cupyx.rsqrt	Returns the reciprocal square root.
cupyx.scatter_add (a, slices, value)	Adds given values to specified elements of an array.
cupyx.scatter_max (a, slices, value)	Stores a maximum value of elements specified by indices to an array.
<pre>cupyx.scatter_min (a, slices, value)</pre>	Stores a minimum value of elements specified by indices to an array.
<pre>cupyx.empty_pinned (shape[, dtype, order])</pre>	Returns a new, uninitialized NumPy array with the given shape and dtype.
<pre>cupyx.empty_like_pinned (a[, dtype, order,])</pre>	Returns a new, uninitialized NumPy array with the same shape and dtype as those of the given array.

full list here: cupy API vs numpy API

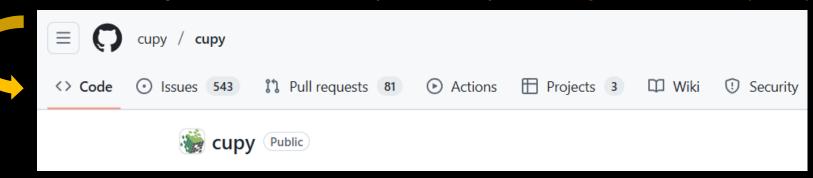
cupy.allclose

full list here: cupy documentation

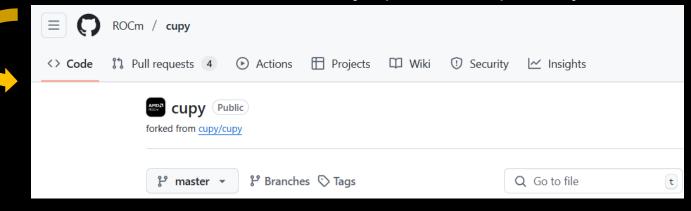
numpy.allclose

CuPy Installation – GitHub Repos

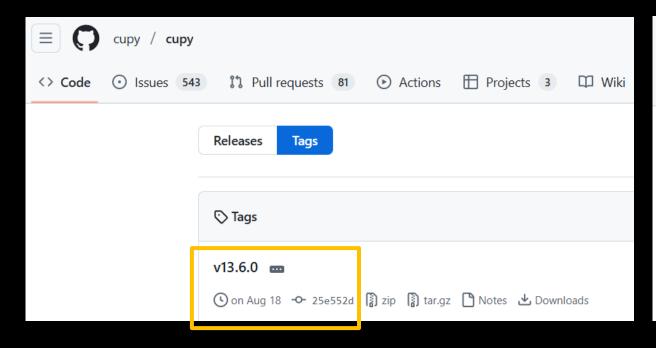
- There are two GitHub repos to take the CuPy source code from to run on AMD GPUs
- We are using the upstream CuPy repository: https://github.com/cupy/cupy

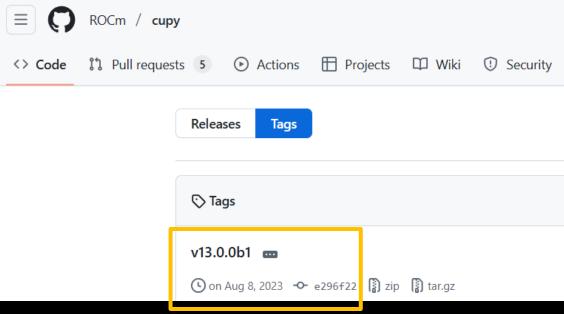


There is also a fork of the CuPy upstream repository in the ROCm github: https://github.com/rocm/cupy



CuPy Installation – Versions





Upstream versions are more recent

The one above is the one we have installed

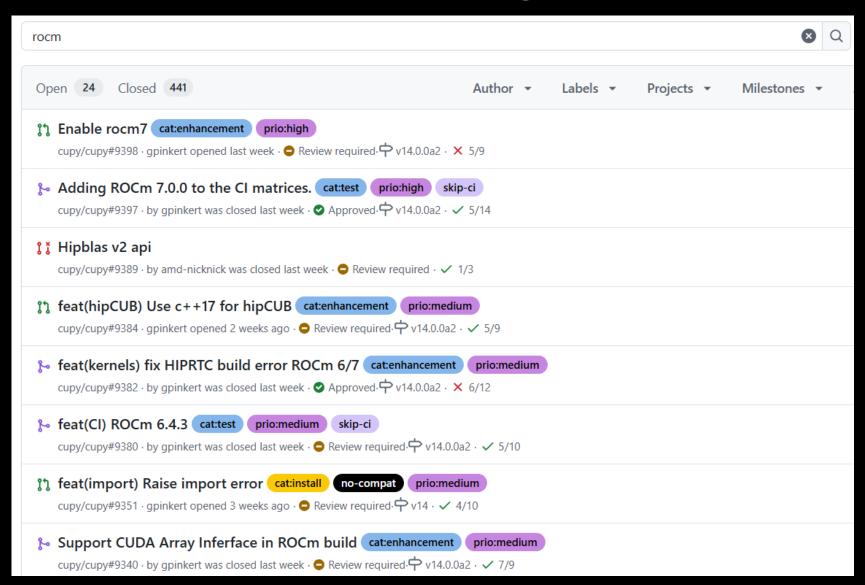
As of September 23rd 2025

ROCm repo versions tend to be behind

There is work from AMD to get changes pushed directly to the upstream repo

The ROCm/cupy will soon be updated too

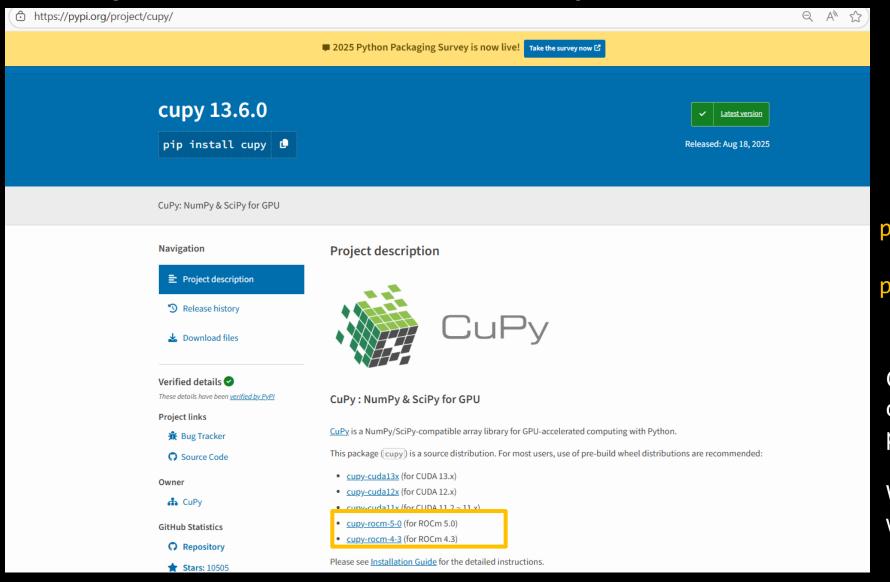
Current AMD work to be integrated into upstream repo



As of September 23rd 2025



CuPy – Installation with pip3 (pre-built wheel for Linux x86_64)



As of September 23rd 2025

pip3 install cupy-rocm-5-0

pip3 install cupy-rocm-4-3

Only old versions of ROCm currently available as pre-built wheels

Wheels for ROCm 6.4 and 7 will soon be available

CuPy – Simple Installation from Source (latest version)

```
export CUPY_INSTALL_USE_HIP=1
export ROCM_HOME=${ROCM_PATH}
export HIPCC=${ROCM_HOME}/bin/hipcc
export HCC_AMDGPU_ARCH=${AMDGPU_GFXMODEL}}
pip3 install cupy --target=$CUPY PATH
```

NOTE: it will **not** work on Ubuntu 24.04 need to use a virtual environment

error: externally-managed-environment

x This environment is externally managed
To install Python packages system-wide, try apt install python3-xyz, where xyz is the package you are trying to install.



If you don't specify the HCC_AMDGPU_TARGET environment variable, CuPy will be built for the GPU architectures available on the build host. This behavior is specific to ROCm builds; when building CuPy for NVIDIA CUDA, the build result is not affected by the host configuration.

source: <u>cupy docs</u>



CuPy – Robust Installation from Source

Installation from source script available in our model installation repository: https://github.com/amd/HPCTrainingDock/blob/main/extras/scripts/cupy_setup.sh

also installs **numpy-allocator** to leverage unified shared memory and **cupy-xarray**

... most relevant part reported below...

```
export CUPY INSTALL USE HIP=1
export ROCM HOME=${ROCM PATH}
export HIPCC=${ROCM HOME}/bin/hipcc
export HCC AMDGPU ARCH=${AMDGPU GFXMODEL}
python3 -m venv cupy build
source cupy build/bin/activate
pip3 install -v --target=$CUPY PATH pytest mock xarray[complete] dask build numpy-allocator --no-cache
export PYTHONPATH=$PYTHONPATH:$CUPY PATH
# Get source from the upstream repository of CuPy.
git clone -q --depth 1 -b v$CUPY VERSION --recursive https://github.com/cupy/cupy.git
cd cupy
python3 -m build --wheel
pip3 install -v --upgrade --target=$CUPY PATH dist/*.whl
pip3 install -v --target=$CUPY PATH cupy-xarray --no-deps
deactivate
```

Basics of CuPy

Must import the CuPy Python™ module in your Python code:

import cupy as cp

- To create an array on the device use cp.array:
- To copy data from GPU to CPU, use cp.asnumpy:
- To copy back from CPU to GPU use cp.asarray:
- Operations between GPU arrays will be done on the GPU:
- CuPy has the concept of a current device usually GPU device 0:

Note that the device will be called <CUDA Device 0> even if you are on AMD GPUs.

gpu array = cp.array(cpu array)

- cpu_array = cp.asnumpy(gpu_array)
- gpu_array2 = cp.asarray(cpu_array)

result_gpu = gpu_array + gpu_array2

gpu_array.device

NumPy – CuPy Interoperability

- CuPy implements a subset of the NumPy interface by implementing cupy.ndarray, a counterpart to NumPy ndarrays
- The cupy.ndarray object implements the __array_ufunc__ interface. This enables NumPy universal functions (ufunc) to be applied to CuPy arrays. Note that the return type of these operations is still consistent with the initial type.

```
>>> import cupy as cp
>>> import numpy as np
>>> gpu_arr = cp.random.randn(1, 2, 3, 4).astype(cp.float32)
>>> result = np.sum(gpu_arr)
>>> print(type(result))
<class 'cupy._core.core.ndarray'>
```

cupy.ndarray also implements the __array_function__ interface, meaning it is possible to do operations such as

```
a = np.random.randn(100, 100)
a_gpu = cp.asarray(a)
qr_gpu = np.linalg.qr(a_gpu)
Oct 13-16, 2025
```

source: numpy-documentation

Simple CuPy code example

First get the example to run from the training examples repository
 git clone https://github.com/amd/HPCTrainingExamples
 cd HPCTrainingExamples/Python/cupy

- Set up the environment: note the "module" below is not the Python™ module module load cupy
- Run the example python3 cupy_array_sum.py

Output should be:

```
CuPy Array: [1 2 3 4 5]
Squared CuPy Array: [ 1 4 9 16 25]
NumPy Array: [5 6 7 8 9]
CuPy Array from NumPy: [5 6 7 8 9]
Addition Result on GPU: [ 6 8 10 12 14]
Result on CPU: [ 6 8 10 12 14]
```

Simple CuPy code example: a closer look

```
import cupy as cp
import numpy as np
# Create a CuPy array
gpu_array = cp.array([1, 2, 3, 4, 5]) 			 Creates an array on the device
print("CuPy Array:", gpu array)
# Perform operations on the GPU
gpu_array_squared = gpu_array ** 2   Operations occur on the GPU
print("Squared CuPy Array:", gpu array squared)
# Create a NumPy array
cpu array = np.array([5, 6, 7, 8, 9])
print("NumPy Array:", cpu array)
# Transfer NumPy array to GPU
print("CuPy Array from NumPy:",
gpu_array_from cpu)
# Perform element-wise addition
result_gpu = gpu_array + gpu_array_from_cpu 	— Operations occur on the GPU
print("Addition Result on GPU:", result gpu)
# Transfer result back to CPU
                                          Returns an array on the host memory from an
result_cpu = cp.asnumpy(result_gpu) 
print("Result on CPU:", result cpu)
                                          arbitrary source array (device in this case)
```

Verifying that CuPy code example runs on the AMD GPU

```
Now run with the AMD LOG LEVEL environment variable set
 export AMD LOG LEVEL=3
 python3 cupy array sum.py
Lots of output now – showing just a little bit:
 hipMemcpyAsync ( 0x559ea98f65f0, 0x7f4556800000, 40, hipMemcpyDeviceToHost, stream:<null> )
  Signal = (0x7f4d5efff280), Translated start/end = 1083534945452078 / 1083534945453358,
   Elapsed = 1280 ns, ticks start/end = 27091222405615 / 27091222405647, Ticks elapsed = 32
 Host active wait for Signal = (0x7f4d5efff200) for -1 ns
 Set Handler: handle(0x7f4d5efff180), timestamp(0x559eaabead90)
 Host active wait for Signal = (0x7f4d5efff180) for -1 ns
 hipMemcpyAsync: Returned hipSuccess : : duration: 5948d us
 hipStreamSynchronize ( stream:<null> )
 Handler: value(0), timestamp(0x559eaa7e7350), handle(0x7f4d5efff180)
 hipStreamSynchronize: Returned hipSuccess :
 hipSetDevice ( 0 )
 hipSetDevice: Returned hipSuccess :
 CuPy Array: [1 2 3 4 5]
```

Unified Memory Programming on CuPy

Unified memory programming (UMP) support (**experimental!**)

It is possible to make both NumPy and CuPy use/share system allocated memory on Heterogeneous Memory Management (HMM) or Address Translation Services (ATS) enabled systems, such as the NVIDIA Grace Hopper Superchip. To activate this capability, currently you need to:

- 1. Install numpy_allocator
- 2. Set the environment variable CUPY_ENABLE_UMP=1
- 3. Make a memory pool for CuPy to draw system memory (malloc system), for example:

```
import cupy as cp
cp.cuda.set_allocator(cp.cuda.MemoryPool(cp.cuda.memory.malloc_system).malloc)
```

4. Switch to the aligned allocator for NumPy to draw system memory

```
import cupy._core.numpy_allocator as ac
import numpy_allocator
import ctypes
lib = ctypes.CDLL(ac.__file__)
class my_allocator(metaclass=numpy_allocator.type):
    _calloc_ = ctypes.addressof(lib._calloc)
    _malloc_ = ctypes.addressof(lib._malloc)
    _realloc_ = ctypes.addressof(lib._realloc)
    _free_ = ctypes.addressof(lib._free)
my_allocator.__enter__() # change the allocator globally
```

On AMD GPUs, you additionally need: export HSA_XNACK=1

this will enable unified shared memory on MI300A

or managed memory on MI200s and MI300X

```
With this setup change, all the data movement APIs such as get(), asnumpy() and asarray()
become no-op (no copy is done), and the following code is accelerated:

# a, b, c are np.ndarray, d is cp.ndarray
a = np.random.random(100)
b = np.random.random(100)
c = np.add(a, b)
d = cp.matmul(cp.asarray(a), cp.asarray(c))

Essentially, the distinction of CPU/GPU memory spaces is gone, and np/cp are used to represent the execution space (whether the code should run on CPU or GPU).

Apart from the setup configuration for NumPy/CuPy, no user code chage is required.
```

source: <u>cupy docs</u>



Unified Memory Programming on CuPy – benchmark 1/2

Unified memory example in training examples repo:

https://github.com/amd/HPCTrainingExamples/blob/main/Python/cupy/unified_bench.py

How to run the code:

The benchmark will transfer data and run some simple operations comparing a baseline (no UMP) with the UMP case (export CUPY_ENABLE_UMP=1 and export HAS_XNACK=1)

Unified Memory Programming on CuPy – benchmark 2/2

It measures the following:

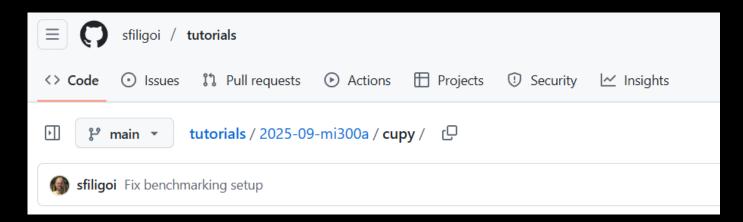
- host to device (HtoD) copies with cp.asarray(np_arr)
- device to host (DtoH) copies with cp_arr.get()
- adding two arrays on device (x+=y)
- pipeline mixing the above three operations (copy to GPU, add on GPU and copy back to CPU)

Results:

```
gcapodag@ppac-pl1-s24-26:~/repos/HPCTrainingExamples/Python/cupy$ python3 cupy_ump_bench.py
=== Running BASELINE (CUPY_ENABLE_UMP=0) ===
=== Running UMP (CUPY_ENABLE_UMP=1) ===
=== RESULTS (avg over repeats) ===
          Mode | HtoD ms | HtoD GiB/s | DtoH ms | DtoH GiB/s | Add ms | Pipeline ms
Size
1.00 GiB |
          BASE | 18.699
                            53.480
                                         116.304 | 8.598
                                                                1.310
                                                                         160.730
                            242413.024 | 125.652 | 7.958
                                                                1.381
1.00 GiB |
[BASELINE] device=AMD Instinct MI300A, SMs=228, mem=94.17 GiB, dtype=float32, repeats=100, warmups=3
[UMP] device=AMD Instinct MI300A, SMs=228, mem=94.17 GiB, dtype=float32, repeats=100, warmups=3
```



Additional Examples on CuPy



Code examples by **Igor Sfiligoi**

git clone https://github.com/sfiligoi/tutorials.git
cd tutorials/2025-09-mi300a/cupy

From 2025 ICPP AI tutorial by San Diego Supercomputing Center and AMD.

Compare python3 center_matrix_naive.py with CUPY_ENABLE_UMP=1\

HSA_XNACK=1 python3 center_matrix_naive_apu.py

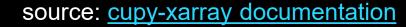
Small on CPU: 0.027914762496948242 Small on GPU: 0.011859893798828125 Medium matrix shape: (8382, 8382) Medium on CPU: 0.21364235877990723 Medium on GPU: 0.11696290969848633 Large matrix shape: (25145, 25145) Large on CPU: 1.9389822483062744 Large on GPU: 1.1117446422576904 ... Retrying Large on CPU: 1.8522789478302002 Large on GPU: 0.890662670135498

Small on CPU: 0.033936262130737305 Small on GPU: 0.06000709533691406 Medium matrix shape: (8382, 8382) 0.22278118133544922 Medium on CPU: Medium on GPU: 0.4624929428100586 Large matrix shape: (25145, 25145) Large on CPU: 1.8617160320281982 Large on GPU: 4.33487343788147 ... Retrying Large on CPU: 1.800774097442627 Large on GPU: 0.9342784881591797

AMD together we advance_

CuPy-Xarray: Xarray on GPUs

- Xarray: Python™ library to work with labelled multi-dimensional arrays
 - Popular for applications where multi-dimensional data needs to be handled (such as climate modeling)
 - Built on top of NumPy
 - Has built-in support for NetCDF
 - Can wrap custom duck array objects (i.e. NumPy-like arrays) that follow specific protocols.
- When used together, Xarray and CuPy can provide an easy way to take advantage of GPU acceleration for scientific computing tasks.
- CuPy-Xarray provides an interface for using CuPy in Xarray, providing accessors on the Xarray objects.
 - CuPy-Xarray relies on an existing CuPy installation, install CuPy first
- Cupy-Xarray github repo: https://github.com/xarray-contrib/cupy-xarray
 - Install with pip install cupy-xarray --no-deps after installing CuPy
- Issue with dask: https://github.com/xarray-contrib/cupy-xarray/pull/62
 - Did not make it into the latest release
 - Make sure to install dask with pip install dask





Simple CuPy-Xarray code example

- First get the example to run from the training examples repository
 git clone https://github.com/amd/HPCTrainingExamples
 cd HPCTrainingExamples/Python/cupy
- Set up the environment: note the "module" below is not the Python™ module module load cupy
- Run the example python cupy xarray test.py

```
Is the array used to create da_np on device? False

Is the array used to create da_cp on device? True

da_cp.data is of type: <class 'cupy.ndarray'>
check that arr_gpu and cupy_array are the same with CuPy: True

check the arr_gpu and cupy_array are the same with NumPy (interoperability): True

arr_gpu is on device: <CUDA Device 0>
arr_cpu is on device: cpu

total number of available devices: 8

arr gpu2 is on device: <CUDA Device 2>
```

source: cupy-xarray documentation



Simple CuPy-Xarray code example: a closer look

```
import cupy as cp
import numpy as np
import xarray as xr
                                      Adds .cupy to Xarray objects
import cupy xarray
arr cpu = np.random.rand(10, 10, 10)
                                                     Creates an array on the CPU with NumPy
arr gpu = cp.random.rand(10, 10, 10)
                                                     Creates an array on the GPU with CuPy
                                                                   Creates a DataArray using NumPy array
da np = xr.DataArray(arr cpu, dims=["x", "y", "time"])
da cp = xr.DataArray(arr gpu, dims=["x", "y", "time"])
                                                                   Creates a DataArray using CuPy array
. . . (some code omitted) . . .
                                             Access the underlying CuPy array used to create the xarray. DataArray
cupy array = da cp.data
print("check that arr gpu and cupy array are the same with CuPy:",
                                                                        Use CuPy to check that the array used to create
cp.allclose(cupy array,arr gpu))
                                                                         the xarray and the one given by xarray are the same
print("check the arr gpu and cupy array are the same with NumPy
                                                                      Use NumPy to check that the array used to create
(interoperability):", np.allclose(cupy_array,arr_gpu))
                                                                      the xarray and the one given by xarray are the same
. . . (some code omitted) . . .
with cp.cuda.Device(2):
                                                                   Use the device context manager to create data on
   arr gpu2 = cp.array([1, 2, 3, 4, 5])
                                                                   a different device
print("arr gpu2 is on device:", arr gpu2.device)
```

Additional Resources

- CuPy vs NumPy speed comparison: https://cupy-xarray.readthedocs.io/latest/examples/01_cupy-basics.html#cupy-vs-numpy-speed-comparison
- Real world example of Cupy-Xarray: https://cupy-xarray.readthedocs.io/latest/examples/06_real-example.html
 - Note: you might need to modify the data read line to this if it is taking too long to get the data: da = xr.open_mfdataset(file_objs, engine="h5netcdf", compat="override", coords='minimal')[var].load()
- Cupy with Xarray vs NumPy with Xarray performance comparison: https://cupy-xarray-vs-numpy-with-xarray

Disclaimer

The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions, and typographical errors. The information contained herein is subject to change and may be rendered inaccurate for many reasons, including but not limited to product and roadmap changes, component and motherboard version changes, new model and/or product releases, product differences between differing manufacturers, software changes, BIOS flashes, firmware upgrades, or the like. Any computer system has risks of security vulnerabilities that cannot be completely prevented or mitigated. AMD assumes no obligation to update or otherwise correct or revise this information. However, AMD reserves the right to revise this information and to make changes from time to time to the content hereof without obligation of AMD to notify any person of such revisions or changes.

THIS INFORMATION IS PROVIDED 'AS IS." AMD MAKES NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE CONTENTS HEREOF AND ASSUMES NO RESPONSIBILITY FOR ANY INACCURACIES, ERRORS, OR OMISSIONS THAT MAY APPEAR IN THIS INFORMATION. AMD SPECIFICALLY DISCLAIMS ANY IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY, OR FITNESS FOR ANY PARTICULAR PURPOSE. IN NO EVENT WILL AMD BE LIABLE TO ANY PERSON FOR ANY RELIANCE, DIRECT, INDIRECT, SPECIAL, OR OTHER CONSEQUENTIAL DAMAGES ARISING FROM THE USE OF ANY INFORMATION CONTAINED HEREIN, EVEN IF AMD IS EXPRESSLY ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

Third-party content is licensed to you directly by the third party that owns the content and is not licensed to you by AMD. ALL LINKED THIRD-PARTY CONTENT IS PROVIDED "AS IS" WITHOUT A WARRANTY OF ANY KIND. USE OF SUCH THIRD-PARTY CONTENT IS DONE AT YOUR SOLE DISCRETION AND UNDER NO CIRCUMSTANCES WILL AMD BE LIABLE TO YOU FOR ANY THIRD-PARTY CONTENT. YOU ASSUME ALL RISK AND ARE SOLELY RESPONSIBLE FOR ANY DAMAGES THAT MAY ARISE FROM YOUR USE OF THIRD-PARTY CONTENT.

© 2025 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, AMD CDNA, AMD ROCm, AMD Instinct, and combinations thereof are trademarks of Advanced Micro Devices, Inc. in the United States and/or other jurisdictions. Other names are for informational purposes only and may be trademarks of their respective owners.

Git and the Git logo are either registered trademarks or trademarks of Software Freedom Conservancy, Inc., corporate home of the Git Project, in the United States and/or other countries

#