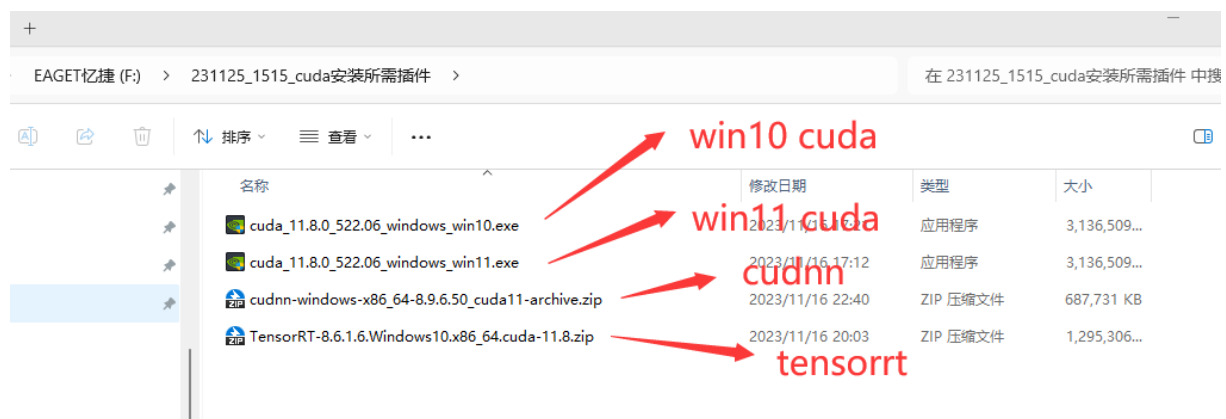


# 231125\_tenserrt 环境安装

范仁义

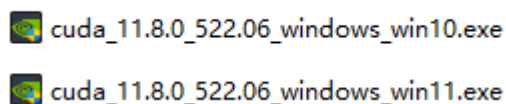
## 一、所需文件

### 1、所需文件



## 二、cuda 安装教程

### 1、cuda 安装



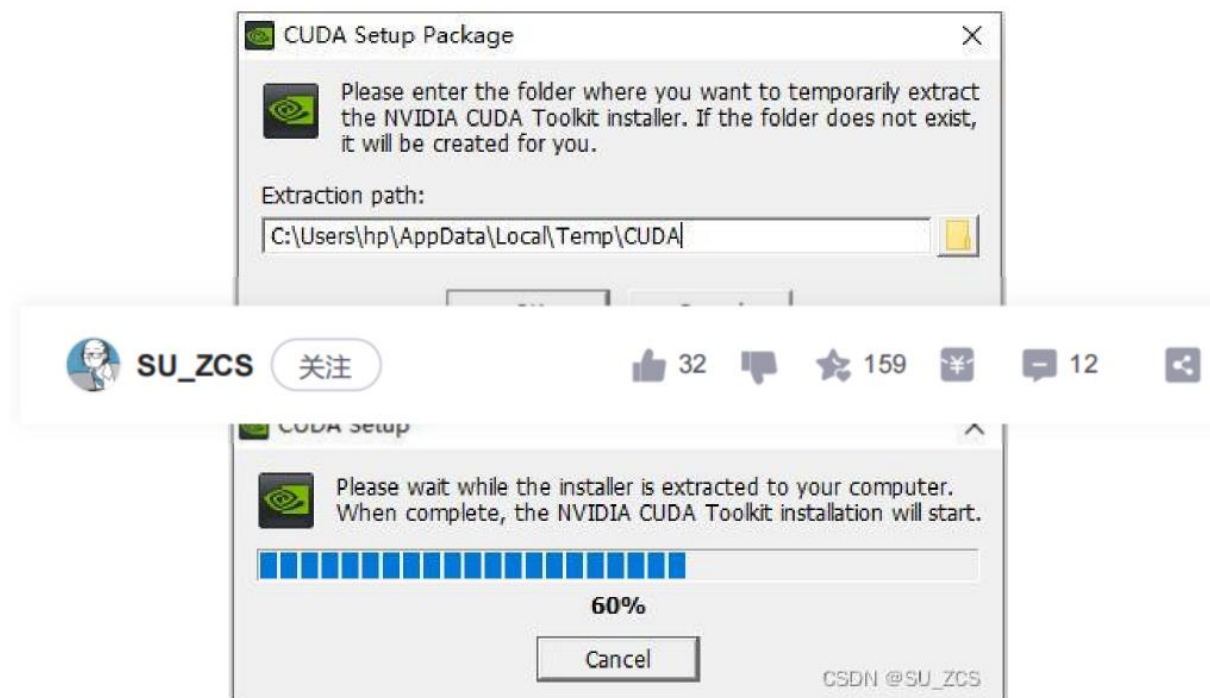
点击上面的 cuda.exe，win10 就装 win10 的，win11 就装 win11 的

基本一步步默认安装即可

以下图是从别的位置截取到的图：

### 1.3 cuda toolkit安装

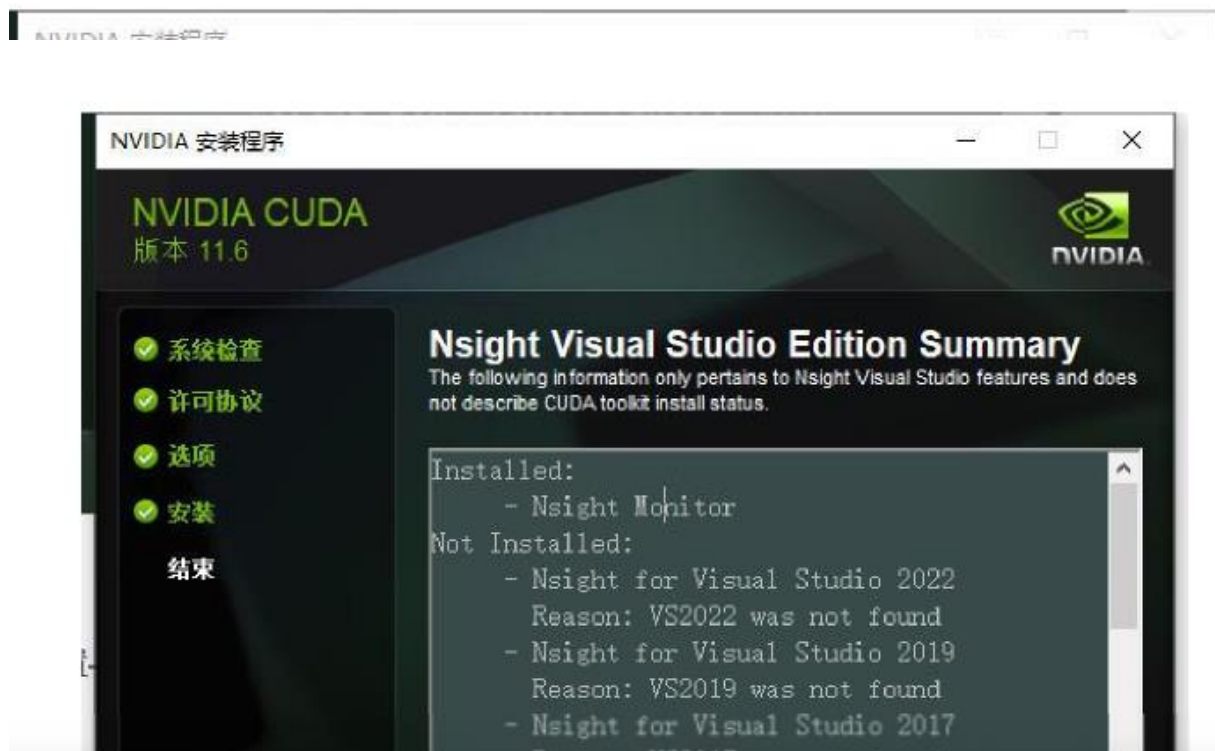
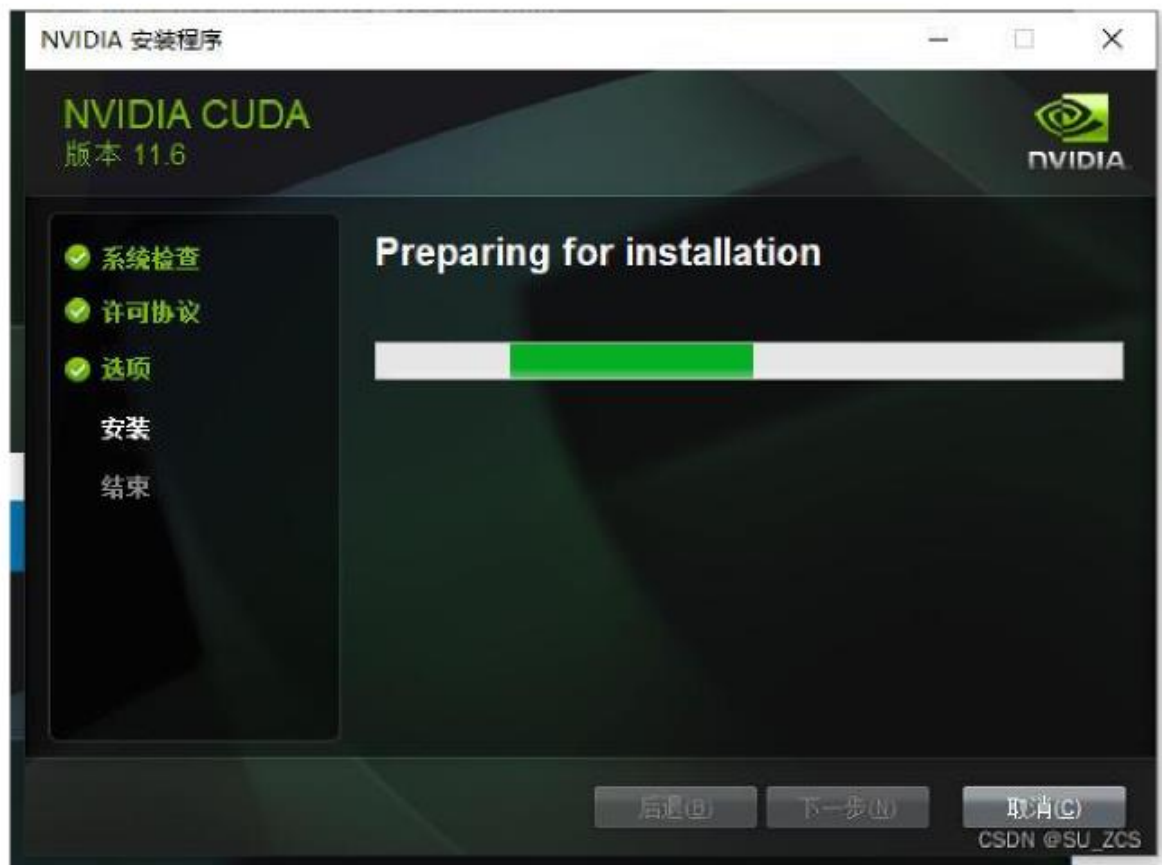
双击exe文件进行安装即可














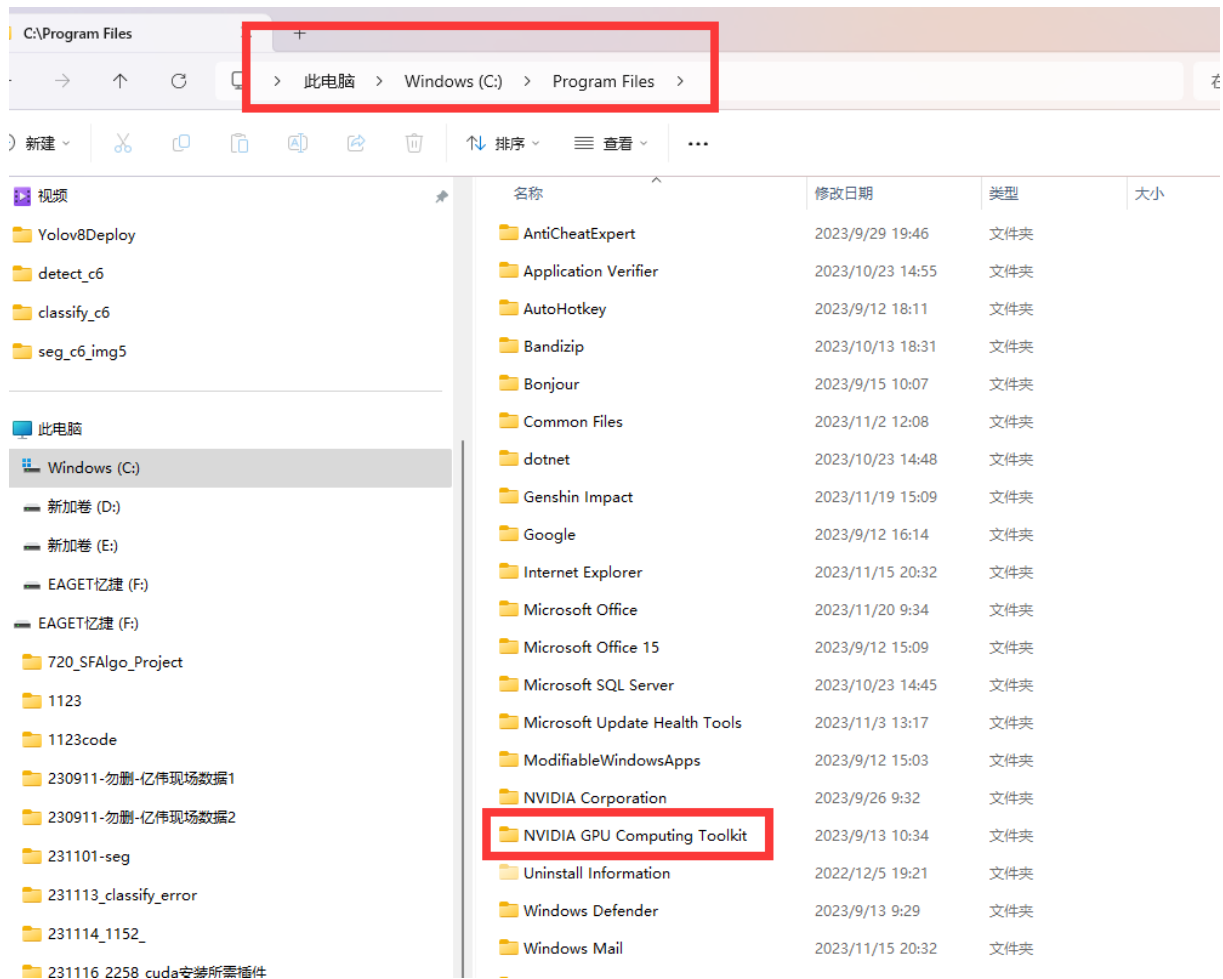


## 2、cudnn 安装

 [cudnn-windows-x86\\_64-8.9.6.50\\_cuda11-archive.zip](#)

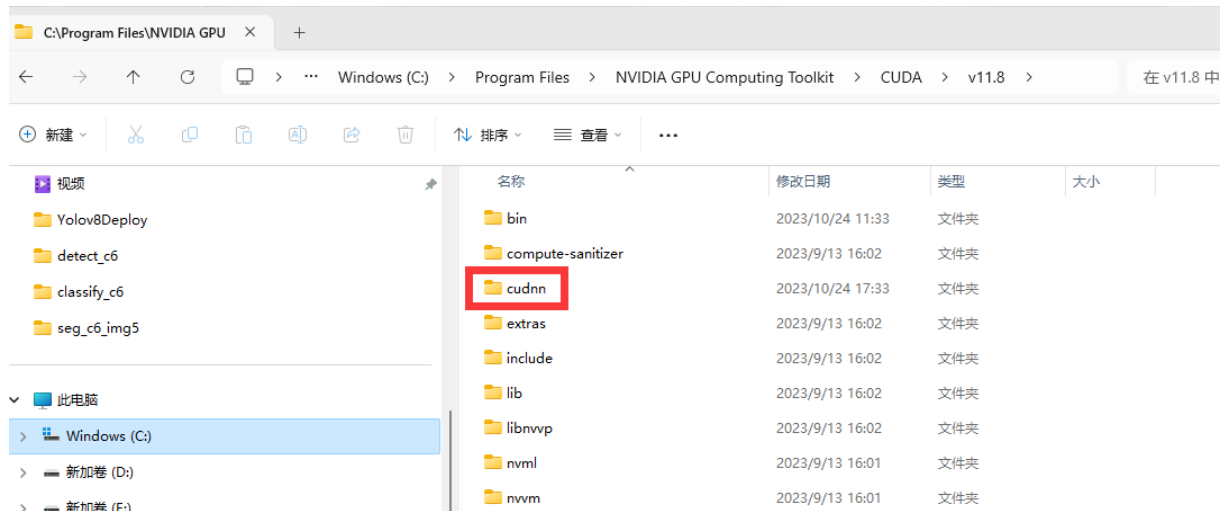
将如下文件复制到 cuda 的安装目录，默认是

C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.8



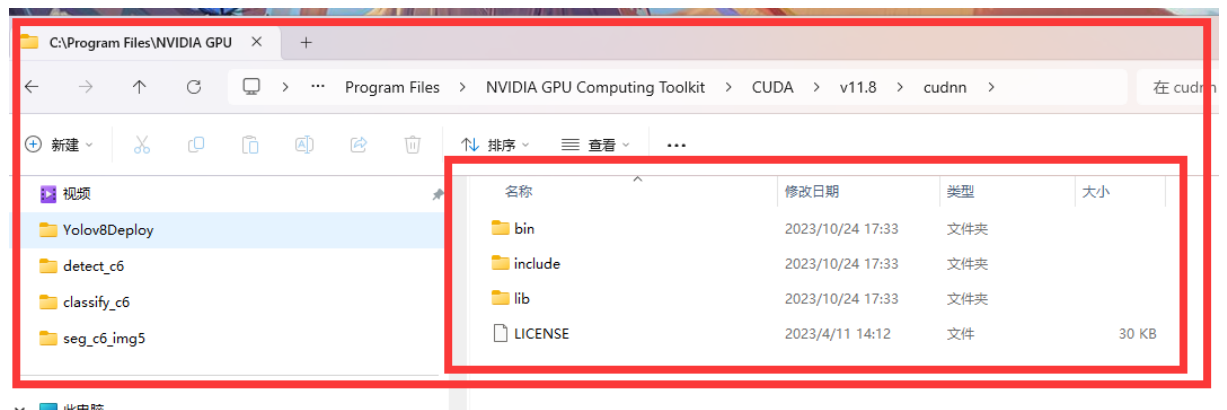


解压这个文件，并且重命名为 cudnn



cudnn 下是这样的

注意：不要在解压的时候多增加一层目录了



特别重要：再将 cudnn 里面的内容直接覆盖到 cuda 目录

就是直接覆盖

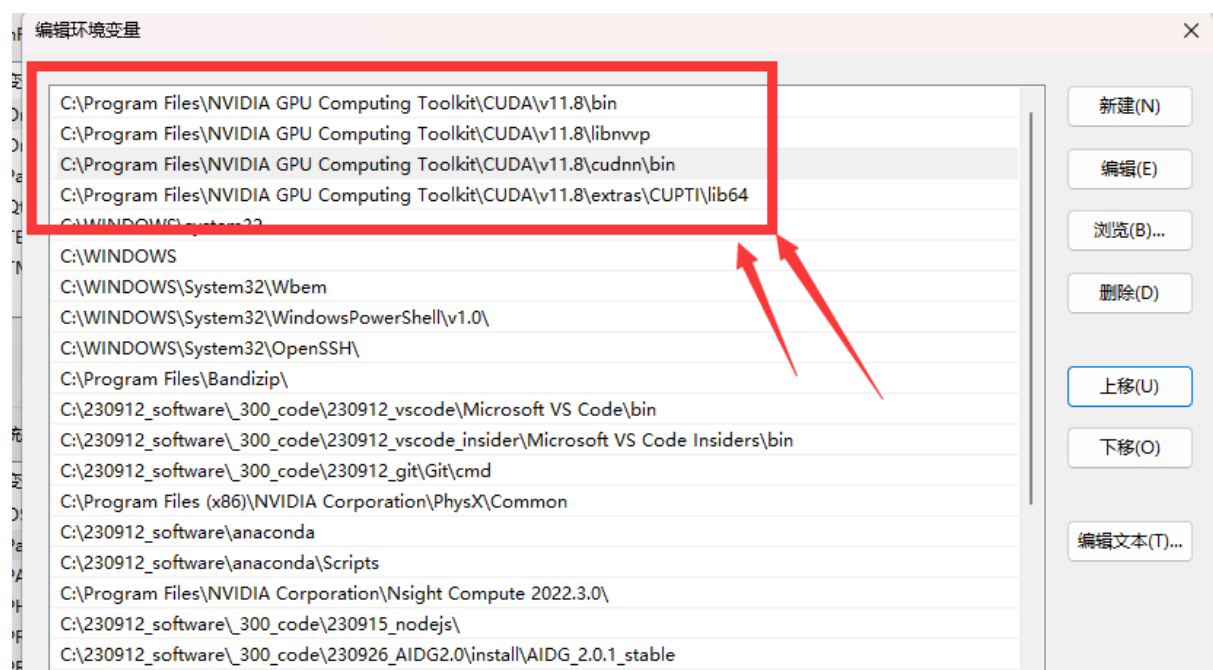


C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.8\bin

C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.8\libnvvp

C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.8\cuda\bin

C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.8\extras\CUPTI\lib64



## 4、验证 cuda 是否安装成功

命令行输入：

```
nvcc --version
```

或者

```
nvcc -V
```

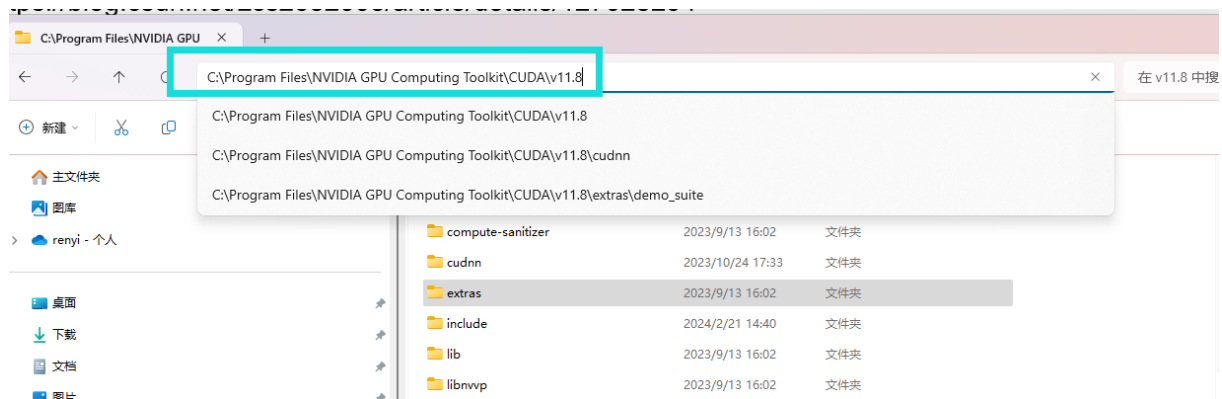
```
C:\WINDOWS\system32\cmd. X + v

Microsoft Windows [版本 10.0.22621.2715]
(c) Microsoft Corporation。保留所有权利。

C:\Users\FanRenyi>nvcc --version
nvcc: NVIDIA (R) Cuda compiler driver
Copyright (c) 2005-2022 NVIDIA Corporation
Built on Wed_Sep_21_10:41:10_Pacific_Daylight_Time_2022
Cuda compilation tools, release 11.8, V11.8.89
Build cuda_11.8.r11.8/compiler.31833905_0

C:\Users\FanRenyi>
```

或者进入 cuda 的安装目录：



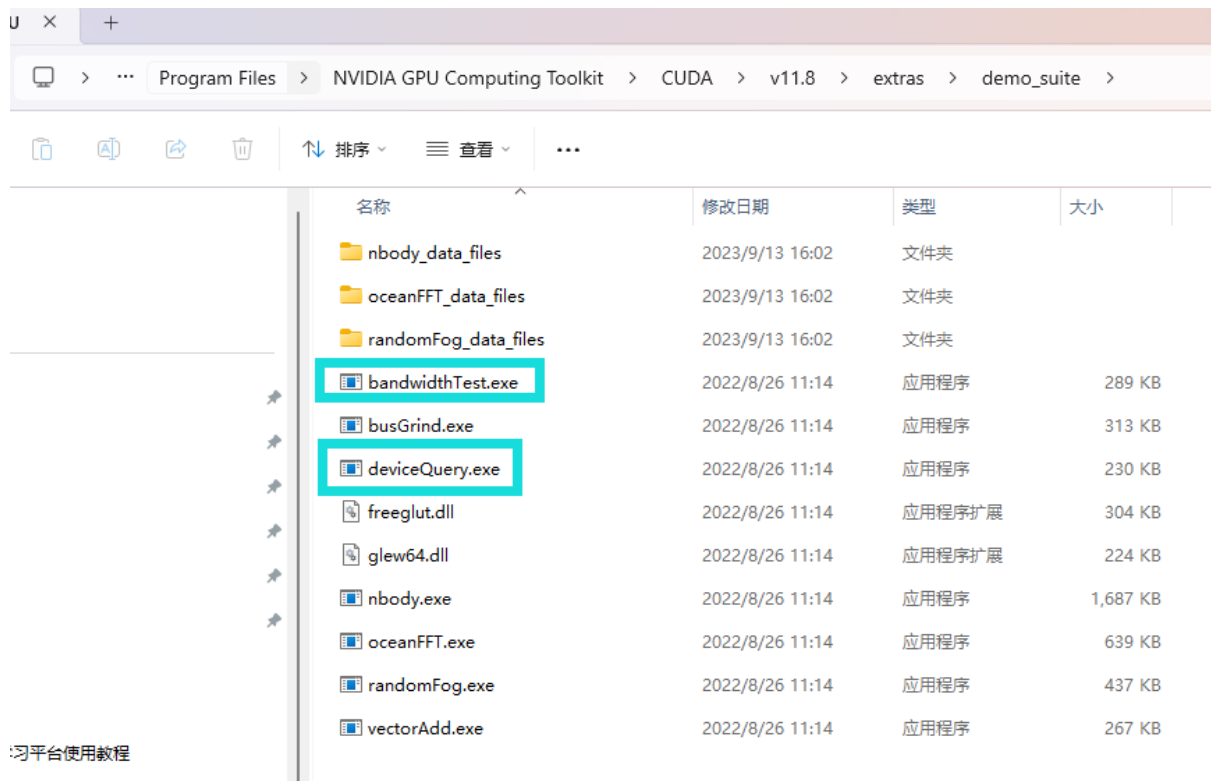
安装的默认地址都在这个目录

C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.8

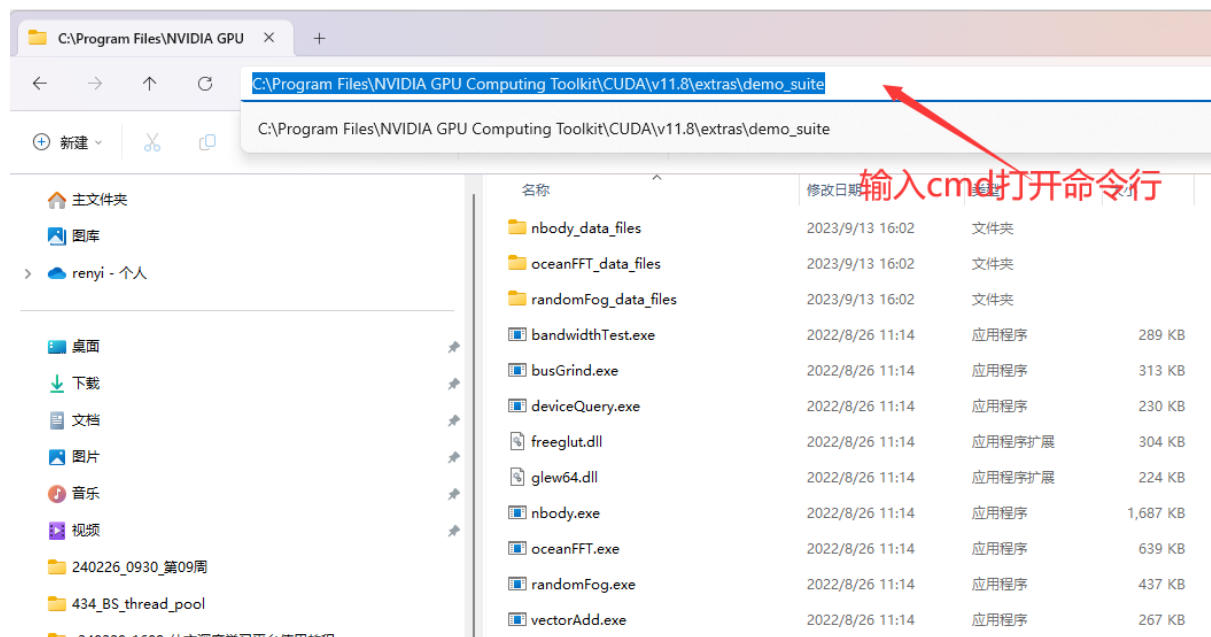
进入

C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.8\extras\demo\_suite

运行 bandwidthTest.exe 和 deviceQuery.exe



习平台使用教程



运行

.\bandwidthTest.exe

```
C:\Windows\System32\cmd.e X + v
Microsoft Windows [版本 10.0.22621.3155]
(c) Microsoft Corporation. 保留所有权利。

C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.8\extras\demo_suite>.\bandwidthTest.exe
[CUDA Bandwidth Test] - Starting...
Running on...

Device 0: NVIDIA GeForce RTX 4060 Laptop GPU
Quick Mode

Host to Device Bandwidth, 1 Device(s)
PINNED Memory Transfers
Transfer Size (Bytes)      Bandwidth(MB/s)
33554432                  12646.9

Device to Host Bandwidth, 1 Device(s)
PINNED Memory Transfers
Transfer Size (Bytes)      Bandwidth(MB/s)
33554432                  12826.6

Device to Device Bandwidth, 1 Device(s)
PINNED Memory Transfers
Transfer Size (Bytes)      Bandwidth(MB/s)
33554432                  225173.7

Result = PASS

NOTE: The CUDA Samples are not meant for performance measurements. Results may vary when GPU Boost is enabled.

C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.8\extras\demo_suite>
```

```
Device to Device Bandwidth, 1 Device(s)
PINNED Memory Transfers
Transfer Size (Bytes)      Bandwidth(MB/s)
33554432                  225173.7

Result = PASS
```

运行

.\deviceQuery.exe

```
C:\Windows\System32\cmd.e  X  +  v

NOTE: The CUDA Samples are not meant for performance measurements. Results may vary when GPU Boost is enabled.

C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.8\extras\demo_suite>.\deviceQuery.exe
.\deviceQuery.exe Starting...

CUDA Device Query (Runtime API) version (CUDART static linking)

Detected 1 CUDA Capable device(s)

Device 0: "NVIDIA GeForce RTX 4060 Laptop GPU"
  CUDA Driver Version / Runtime Version      12.2 / 11.8
  CUDA Capability Major/Minor version number:  8.9
  Total amount of global memory:              8188 MBytes (8585216000 bytes)
  MapSMtoCores for SM 8.9 is undefined.  Default to use 128 Cores/SM
  MapSMtoCores for SM 8.9 is undefined.  Default to use 128 Cores/SM
  (24) Multiprocessors, (128) CUDA Cores/MP:  3072 CUDA Cores
  GPU Max Clock rate:                        2100 MHz (2.10 GHz)
  Memory Clock rate:                         8001 Mhz
  Memory Bus Width:                          128-bit
  L2 Cache Size:                             33554432 bytes
  Maximum Texture Dimension Size (x,y,z)      1D=(131072), 2D=(131072, 65536), 3D=(16384, 16384, 16384)
  Maximum Layered 1D Texture Size, (num) layers 1D=(32768), 2048 layers
  Maximum Layered 2D Texture Size, (num) layers 2D=(32768, 32768), 2048 layers
  Total amount of constant memory:             zu bytes
  Total amount of shared memory per block:     zu bytes
  Total number of registers available per block: 65536
  Warp size:                                   32
  Maximum number of threads per multiprocessor: 1536
  Maximum number of threads per block:         1024
```

```
deviceQuery, CUDA Driver = CUDART, CUDA D
A GeForce RTX 4060 Laptop GPU
Result = PASS
C:\Program Files\NVIDIA GPU Computing Tool
```

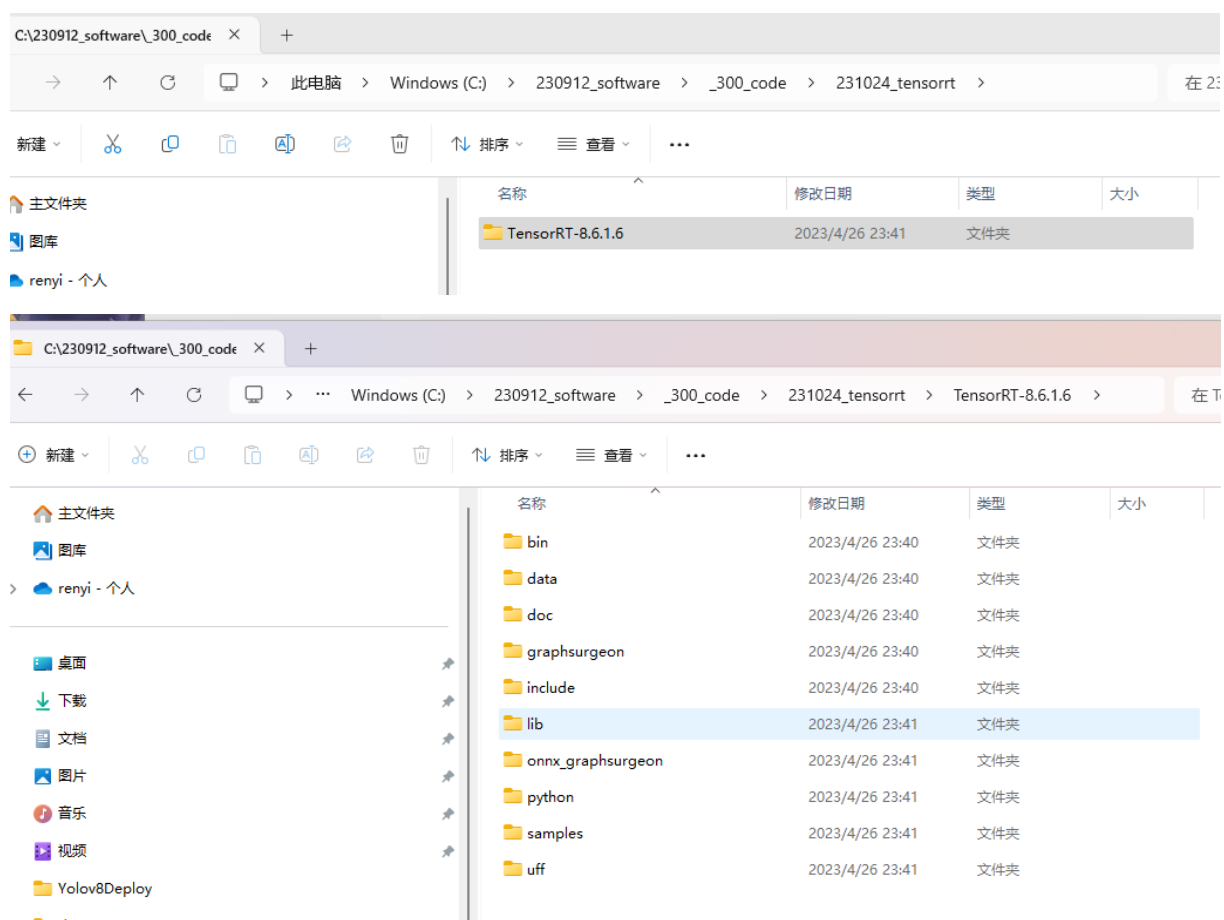
可以看到自己设备的各种信息，以及 是否成功

## 三、安装 tensorrt

### 1、解压文件



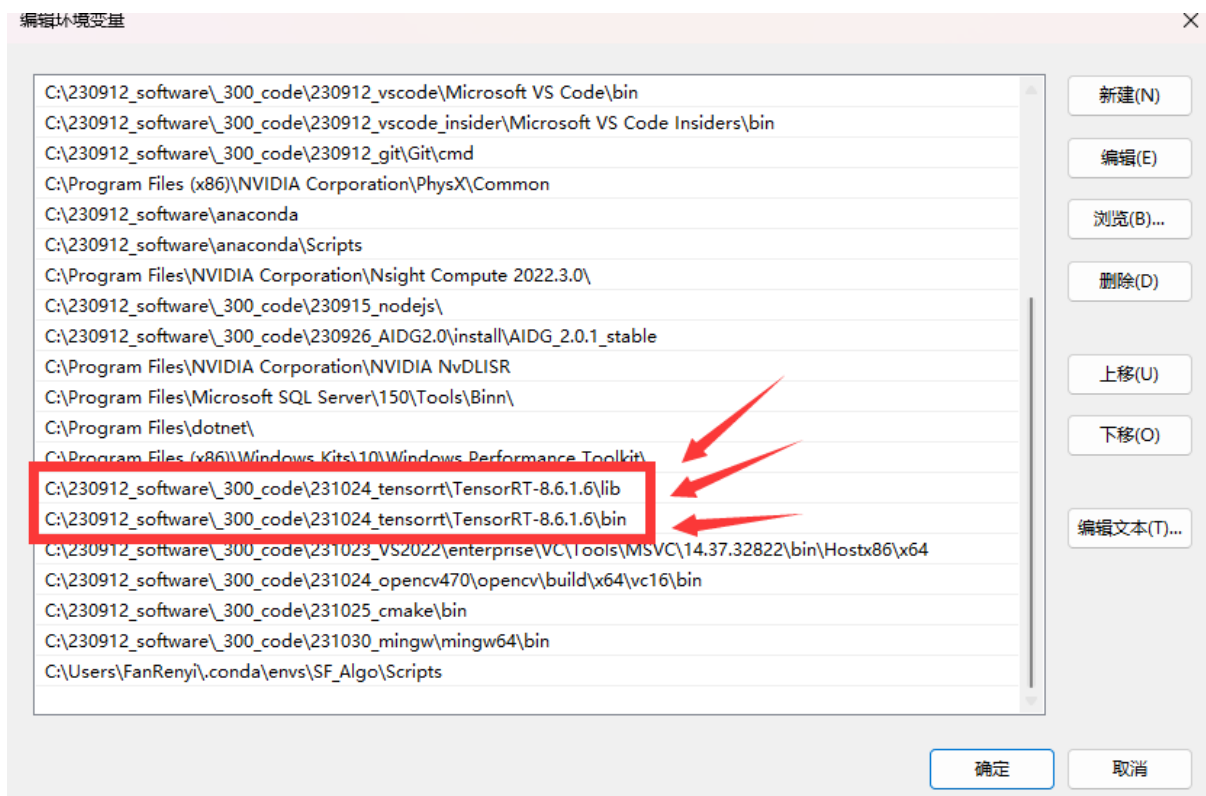
例如我放到了：



## 2、添加环境变量

将如下内容添加到环境变量

C:\230912\_software\\_300\_code\231024\_tensorrt\TensorRT-8.6.1.6\bin  
C:\230912\_software\\_300\_code\231024\_tensorrt\TensorRT-8.6.1.6\lib



### 3、验证是否安装成功

命令行输入:

## trtexec

```
C:\WINDOWS\system32\cmd. X + v
Microsoft Windows [版本 10.0.22621.2715]
(c) Microsoft Corporation。保留所有权利。

C:\Users\FanRenyi>trtexec
&&&& RUNNING TensorRT.trtexec [TensorRT v8601] # trtexec
=== Model Options ===
--uff=<file>          UFF model
--onnx=<file>         ONNX model
--model=<file>        Caffe model (default = no model, random weights used)
--deploy=<file>       Caffe prototxt file
--output=<name>[,<name>]* Output names (it can be specified multiple times); at least one output is required for UFF
and Caffe
--uffInput=<name>,X,Y,Z Input blob name and its dimensions (X,Y,Z=C,H,W), it can be specified multiple times; at l
east one is required for UFF models
--uffNHWC            Set if inputs are in the NHWC layout instead of NCHW (use X,Y,Z=H,W,C order in --uffInput)

=== Build Options ===
--maxBatch           Set max batch size and build an implicit batch engine (default = same size as --bat
ch)
This option should not be used when the input model is ONNX or when dynamic shapes
are provided.
--minShapes=spec     Build with dynamic shapes using a profile with the min shapes provided
--optShapes=spec     Build with dynamic shapes using a profile with the opt shapes provided
--maxShapes=spec     Build with dynamic shapes using a profile with the max shapes provided
--minShapesCalib=spec Calibrate with dynamic shapes using a profile with the min shapes provided
--optShapesCalib=spec Calibrate with dynamic shapes using a profile with the opt shapes provided
--maxShapesCalib=spec Calibrate with dynamic shapes using a profile with the max shapes provided
Note: All three of min, opt and max shapes must be supplied.
However, if only opt shapes is supplied then it will be expanded so
that min shapes and max shapes are set to the same values as opt shapes.
```