

```
pip install -i https://mirrors.aliyun.com/pypi/simple/ nvidia-smi
```

```
lscpu
```

- P-core
- E-core

```
lscpu | grep cache
```

```
ls
```

```
conda env list
```

```
conda deactivate
```

```
pwd
```

```
conda create -n videollava python = 3.10 -y
```

```
mv oldname newname
```

```
conda remove --name myenv --all
```

```
rm -rf your_folder_name
```

```
conda search qwen-vl-utils -c conda-forge
```

```
conda install -c conda-forge qwen-vl-utils=0.0.11 -y
```

```
uptime
```

```
uptime -s
```

```
top
```

- 1CPU
- us
- sy
- id

```
history
```

```
df -h
```

- tmfps
- /dev/nvme0n1p2
- /dev/nvme0n1p1
- efivarfs

```
free -h
```

- total
- used
- free
- buff/cache
- available

```
nvcc --verison
```

GPU	Tensor Type	Bit Width	Architecture
RTX 3080(FP32)	torch.float32	32	Tensor Cores
RTX 3080(FP16)	torch.float16	16	Tensor Cores
RTX 3080(BF16)	torch.bfloat16	16	Ampere GPU
TensorRT(32)	NVIDIA Tensor(FP32)	19	NVIDIA Ampere Tensor Processing Unit
INT8	1-bit/FP32	8-bit	Tensor Cores
INT4	0.5-bit/FP32	4-bit	Tensor Cores

dtype	PyTorch Type	Tensor Type
fp32	<code>torch.float32</code>	
fp16	<code>torch.float16</code>	
bf16	<code>torch.bfloat16</code>	
tf32	<code>torch.cuda.FloatTensor</code>	

```
# 8-bit APIs
quant_config = BitsAndBytesConfig(
    load_in_8bit=True,
    bnb_4bit_compute_dtype=torch.bfloat16
)
```

Tensor API bfloat16 と torch.cuda.FloatTensor INT8-bit

CUDA

概要

- **CUDA Core** GPU の計算アーキテクチャ
- **SM** Streaming Multiprocessor
 - GPU の構成要素 SM は CUDA Core の複数の実行エンジン

構成要素

- **Global Memory** GPU メモリ CPU と RAM
- **Shared Memory** SM 内部メモリ CPU と GPU
- **Registers** GPU の内部レジスタ

並列化

- **Kernel** 実行単位
 - GPU の実行単位
- パターン
 - **Thread** 個々の実行単位
 - **Block** 総合的な実行単位 SM
 - **Grid** Block の総合的な実行単位 Kernel

CUDAツール

- **CUDA** GPU の開発環境
- **cuDNN** ネットワーク
- **cuFFT** フーリエ変換
- **cuBLAS** ベクトル演算
- **NVCC** CUDA の GPU の開発環境 GPU