

```
pip install -i https://mirrors.aliyun.com/pypi/simple/ nvidia-smi 顯示NVIDIAGPU相關資訊  
lscpu 顯示cpu資訊


- P-core 顯示CPU核心數量2個
- E-core 顯示CPU輔助核心


lscpu | grep cache 顯示缓存資訊  
ls 顯示目錄內容  
conda env list 顯示所有虛擬環境  
conda deactivate 退出虛擬環境  
pwd 顯示工作目錄  
conda create -n videollava python = 3.10 -y 建立新虛擬環境videollava  
mv oldname newname 重新命名虛擬環境  
conda remove --name myenv --all 移除所有myenv相關包  
rm -rf your_folder_name 移除文件夾  
conda search qwen-vl-utils -c conda-forge |search| -c conda-forge 搜尋qwen-vl-utils  
conda install -c conda-forge qwen-vl-utils=0.0.11 -y 安裝qwen-vl-utils  
uptime 顯示系統運行時間  
uptime -s 顯示系統啟動時間  
top 顯示CPU資訊


- 1顯示CPU資訊
  - us使用者
  - sy系統
  - id空閒


history 顯示命令歷史  
df -h 顯示磁盤空間資訊  


- tmpfs 顯示tmpfs RAM資訊
- /dev/nvme0n1p2 顯示磁盤
- /dev/nvme0n1p1 顯示磁盤
- efivarfs 顯示efivarfs


free -h 顯示記憶體資訊  


- total 總量
- used 使用量
- free 空閒量
- buff/cache 儲存緩衝區
- available 可用量


nvcc --verison 顯示cuda版本
```

```

tensorboard --logdir=./logs --port=6006 --bind_all

• --logdir : 亂數値logging_dir 亂數値
• --port : 亂數値6006
• --bind_all : 亂數値

ps aux 亂數値

• ps 亂數値Process Status 亂數値
• a 亂數値 亂數値 亂數値 亂數値
• u 亂數値 亂數値 亂數値 亂數値
• x 亂數値 亂數値 亂數値 亂數値

kill -9 1234 亂數値PID1234 亂數値

ssh -L [乱数]:[乱数]:[乱数]@[乱数] 亂數値

亂數値

sudo ufw status 亂數値

sudo ufw allow/deny 8080 亂数/乱数8080 亂数

sudo ufw enable/disable 亂数/乱数

亂數値

sudo systemctl status nginx 亂數値

sudo systemctl start/stop/restart redis 亂数/乱数/乱数

sudo systemctl enable/disable redis 亂数/乱数

tops 亂數値/亂數値/亂數値/亂數値/亂數値/亂數値

亂數値

```

- 1 GOPS = 10GFLOPs
- 1 POPS = 1000GFLOPs

亂數値

亂數値

乱数	乱数	乱数	乱数
乱数(FP32)	torch.float32	32	乱数
乱数(FP16)	torch.float16	16	乱数
乱数(BF16)	torch.bfloat16	16	Ampere GPU
TF(32)	乱数(FP32)	19	NVIDIA Ampere GPU
INT8	1/FP32/FP32/1/4	8-bit	乱数
INT4	0.5/FP32/FP32/1/8)	4-bit	乱数

乱数

乱数乱数乱数乱数乱数乱数乱数

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
- **Threads**
 - **Thread** 実行単位
 - **Block** 総合実行単位 SM
 - **Grid** Block の総合実行単位 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