

## 微调

映射:

```

after_run:
(BELOW_NORMAL) LoggerHook

Downloading data files: 100% [REDACTED] 2/
Extracting data files: 100% [REDACTED] 2
Generating train split: 9846 examples [00:00, 147868.35 examples/s]
Generating test split: 518 examples [00:00, 82906.57 examples/s]
Map: 100% [REDACTED] 9846/9846 [00:00<
Map: 100% [REDACTED] 9846/9846 [00:00<
Filter: 100% [REDACTED] 9846/9846 [00:00<
Map: 100% [REDACTED] 9846/9846 [00:1
Flattening the indices: 100% [REDACTED] 9846/9846 [00:00<
Map: 100% [REDACTED] 9846/9846 [00:03
03/28 12:51:57 - mmengine - WARNING - Dataset Dataset has no meta info. ``dataset_meta`` in visualizer will be None.
quantization_config convert to <class 'transformers.utils.quantization_config.BitsAndBytesConfig'>
Loading checkpoint shards: 100% [REDACTED] 8
03/28 12:53:09 - mmengine - INFO - dispatch internlm attn forward
03/28 12:53:09 - mmengine - WARNING - Due to the implementation of the PyTorch version of flash attention, even when the `output_attentions` flag is
ssible to return the `attn_weights`.
Traceback (most recent call last):

```

加速:

```
ft-oasst1# xtuner train ./internlm_chat_7b_qlora_oasst1_e3_copy.py --deepspeed deepspeed_zero2
```

```
System environment:
  sys.platform: linux
  Python: 3.10.13 (main, Sep 11 2023, 13:44:35) [GCC 11.2.0]
  CUDA available: True
  MUSA available: False
  numpy.random_seed: 1156061394
  GPU 0: NVIDIA A100-SXM4-80GB
  CUDA_HOME: /usr/local/cuda
  NVCC: Cuda compilation tools, release 11.7, V11.7.99
  GCC: gcc (Ubuntu 9.4.0-1ubuntu1) 20.04.2 9.4.0
  PyTorch: 2.0.1
  PyTorch compiling details: PyTorch built with:
  - GCC 9.3
  - C++ Version: 201703
  - Intel(R) oneAPI Math Kernel Library Version 2023.1-Product Build 20230303 for Intel(R) 64 architecture applications
  - Intel(R) MKL-DNN v2.7.3 (Git Hash 6dbffba1ef23cbbae17ad77b5b13f1f37c080e)
  - OpenMP 201511 (a.k.a. OpenMP 4.5)
  - LAPACK is enabled (usually provided by MKL)
  - NNPACK is enabled
  - CPU capability usage: AVX2
  - CUDA Runtime 11.7
  - NVCC architecture flags: -gencode=arch=compute_37,code=sm_37;-gencode=arch=compute_50,code=sm_50;-gencode=arch=compute_60,code=sm_60;-gencode=arch=compute_61,code=sm_61;-gencode=arch=compute_70,code=sm_70;-gencode=arch=compute_75,code=sm_75;-gencode=arch=compute_80,code=sm_80;-gencode=arch=compute_86,code=sm_86;-gencode=arch=compute_37,code=compute_37
  - CuDNN 8.5
  - Magma 2.6.1
  - Build settings: BLAS_INFO=mkl, BUILD_TYPE=Release, CUDA_VERSION=11.7, CuDNN_VERSION=8.5.0, CXX_COMPILER=/opt/rh/devtoolset-9/root/usr/bin/c++, CXX_FLAGS=-D_GLIBCXX_USE_CXX11_ABI=0 -fabi-version=11 -fno-deprecated -fvvisability=inline-hidden -DUSE_PTHREADPOOL -DNDEBUG -DUSE_KINETO -DLIBKINETO_NOROCTRACER -DUSE_FBGEMM -DUSE_QNNPACK -DUSE_PYTORCH_QNNPACK -DUSE_NNPACK -DSYMBOLICATE_MOBILE_DEBUG_HANDLE -O2 -fPIC -Wall -Wextra -Werror-return-type -Werror-non-virtual-dtor -Werror=booll-operation -Wnarrowing -Wno-mismatched-gif-initializers -Wno-type-limits -Wno-array-bounds -Wno-unknown-pragmas -Wno-unused-local-type-definitions -Wno-unused-parameter -Wno-unused-function -Wno-unused-result -Wno-strict-overflow -Wno-strict-aliasing -Wno-error=deprecated-declarations -Wno-stringop-overflow -Wno-psabi -Wno-error=pedantic -Wno-error=redundant-decls -Wno-error=old-style-cast -fdiagnostics-color=always -faligned-new -Wno-unused-but-set-variable -Wno-maybe-uninitialized -fno-math-errno -fno-trapping-math -Werror=format -Werror=cast-function-type -Wno-stringop-overflow, LAPACK_INFO=mkl, PERF_WITH_AVX=1, PERF_WITH_AVX2=1, PERF_WITH_AVX512=1, TORCH_DISABLE_GPU_ASSERTS=ON, TORCH_VERSION=2.0.1, USE_CUDA=ON, USE_CUDNN=ON, USE
```