

نمودج تدقیق

```
import os
import glob
import json
import load_dotenv
import AutoTokenizer, AutoModelForCausalLM, Trainer, TrainingArguments, DataCollatorForLanguageModeling
import Dataset, load_dataset
import torch

load_dotenv()

_ = "deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B" # model name
_= "trained_model"
_= "train.jsonl"
_- _ = 512
_= 8
_= 3

def __(_):
    _ = []
    for _ in os.listdir(_):
        _ = os.path.join(_, _)
        if not os.path.isdir(_):
            for _ in glob.glob(os.path.join(_, "*.md")):
                _ = open(_ , 'r', encoding='utf-8')
                _ = ()read.
                #
                _ = )split.    "___", 2)[-1].strip()
                ()append.

    except Exception as e:
        print(f" {_ : {e}}")

    return

def __(_):
    _, _ = (
        _ , truncation=True, padding=True, max_length=MAX_LENGTH, return_tensors="pt")
        from .__dict__()

def __(_):
    _, _ = (
```

```

        _ = TrainingArguments(
            output_dir=OUTPUT_DIR,
            overwrite_output_dir=True,
            num_train_epochs=EPOCHS,
            per_device_train_batch_size=BATCH_SIZE,
            save_steps=10_000,
            save_total_limit=2,
            prediction_loss_only=True,
            remove_unused_columns=False,
        )

        = AutoModelForCausalLM.from_pretrained(MODEL_NAME, trust_remote_code=True)
        _ = DataCollatorForLanguageModeling(      =      , mlm=False)
        = Trainer(
            model=      ,
            args=      ,
            train_dataset=      ,
            data_collator=      ,
        )

        ()train.
        save._model(OUTPUT_DIR)

def :()
    _ = "_posts"
    =      _ ) _(
        = LlamaTokenizerFast.from_pretrained(MODEL_NAME, trust_remote_code=True, use_fast=True)
    pad.      _token = eos._token
        _ =      _ ) _, (
    ) _ _ _, (
    if __name__ == "__main__":
        ()_

```