

```
import torch
!pip install transformers
!pip install accelerate
```

```
Requirement already satisfied: transformers in /usr/local/lib/python3.10/dist-packages (4.46.3)
Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from transformers) (3.16.1)
Requirement already satisfied: huggingface-hub<1.0,>=0.23.2 in /usr/local/lib/python3.10/dist-packages (from transformers) (
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.10/dist-packages (from transformers) (1.26.4)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from transformers) (24.2)
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.10/dist-packages (from transformers) (6.0.2)
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.10/dist-packages (from transformers) (2024.9.11)
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from transformers) (2.32.3)
Requirement already satisfied: tokenizers<0.21,>=0.20 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.20.3)
Requirement already satisfied: safetensors>=0.4.1 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.4.5)
Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.10/dist-packages (from transformers) (4.66.6)
Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub<1.0,>=0.23.
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub<1
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->transform
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (
Requirement already satisfied: accelerate in /usr/local/lib/python3.10/dist-packages (1.1.1)
Requirement already satisfied: huggingface-hub>=0.21.0 in /usr/local/lib/python3.10/dist-packages (from accelerate) (0.26.5)
Requirement already satisfied: numpy<3.0.0,>=1.17 in /usr/local/lib/python3.10/dist-packages (from accelerate) (1.26.4)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from accelerate) (24.2)
Requirement already satisfied: psutil in /usr/local/lib/python3.10/dist-packages (from accelerate) (5.9.5)
Requirement already satisfied: pyyaml in /usr/local/lib/python3.10/dist-packages (from accelerate) (6.0.2)
Requirement already satisfied: safetensors>=0.4.3 in /usr/local/lib/python3.10/dist-packages (from accelerate) (0.4.5)
Requirement already satisfied: torch>=1.10.0 in /usr/local/lib/python3.10/dist-packages (from accelerate) (2.5.1+cu121)
Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=0.21.0->accelerate
Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=0.21.0->ac
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=0.21.0->accelerate
Requirement already satisfied: tqdm>=4.42.1 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=0.21.0->accele
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=
Requirement already satisfied: networkx in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate) (3.4.2)
Requirement already satisfied: Jinja2 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate) (3.1.4)
Requirement already satisfied: sympy==1.13.1 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate) (1.
Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.10/dist-packages (from sympy==1.13.1->torch>=1.1
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from Jinja2->torch>=1.10.0->accel
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->huggingfa
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->huggingface-hub>=0.21
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->huggingface-hub
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->huggingface-hub
```

```
from transformers import AutoModelForSpeechSeq2Seq, AutoProcessor
```

```
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
torch_dtype = torch.float16 if torch.cuda.is_available() else torch.float32
```

```
model_id = "openai/whisper-large-v2"
```

```
model = AutoModelForSpeechSeq2Seq.from_pretrained(
    model_id,
    torch_dtype=torch_dtype,
    low_cpu_mem_usage=True, # fast loading
    use_safetensors=True, # secure (over pickle)
    attn_implementation="sdpa", # Flash Attention speed-up
)
```

```
model.to(device)
```

```
processor = AutoProcessor.from_pretrained(model_id)
```

```
generation_config.json: 100% 4.29k/4.29k [00:00<00:00, 280kB/s]
preprocessor_config.json: 100% 185k/185k [00:00<00:00, 872kB/s]
tokenizer_config.json: 100% 283k/283k [00:00<00:00, 1.35MB/s]
vocab.json: 100% 836k/836k [00:00<00:00, 1.31MB/s]
tokenizer.json: 100% 2.48M/2.48M [00:00<00:00, 3.84MB/s]
merges.txt: 100% 494k/494k [00:00<00:00, 1.14MB/s]
normalizer.json: 100% 52.7k/52.7k [00:00<00:00, 3.22MB/s]
added_tokens.json: 100% 34.6k/34.6k [00:00<00:00, 2.84MB/s]
special_tokens_map.json: 100% 2.19k/2.19k [00:00<00:00, 171kB/s]
```

```
!pip install datasets
from datasets import load_dataset
```

```
dataset = load_dataset("hf-internal-testing/librispeech_asr_dummy", "clean", split="validation")
```

```
Collecting datasets
  Downloading datasets-3.2.0-py3-none-any.whl.metadata (20 kB)
Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from datasets) (3.16.1)
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.10/dist-packages (from datasets) (1.26.4)
Requirement already satisfied: pyarrow>=15.0.0 in /usr/local/lib/python3.10/dist-packages (from datasets) (17.0.0)
Collecting dill<0.3.9,>=0.3.0 (from datasets)
  Downloading dill-0.3.8-py3-none-any.whl.metadata (10 kB)
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (from datasets) (2.2.2)
Requirement already satisfied: requests>=2.32.2 in /usr/local/lib/python3.10/dist-packages (from datasets) (2.32.3)
Requirement already satisfied: tqdm>=4.66.3 in /usr/local/lib/python3.10/dist-packages (from datasets) (4.66.6)
Collecting xxhash (from datasets)
  Downloading xxhash-3.5.0-cp310-cp310-manylinux_2_17_x86_64_manylinux2014_x86_64.whl.metadata (12 kB)
Collecting multiprocess<0.70.17 (from datasets)
  Downloading multiprocess-0.70.16-py310-none-any.whl.metadata (7.2 kB)
Collecting fsspec<=2024.9.0,>=2023.1.0 (from fsspec[http]<=2024.9.0,>=2023.1.0->datasets)
  Downloading fsspec-2024.9.0-py3-none-any.whl.metadata (11 kB)
Requirement already satisfied: aiohttp in /usr/local/lib/python3.10/dist-packages (from datasets) (3.11.10)
Requirement already satisfied: huggingface-hub>=0.23.0 in /usr/local/lib/python3.10/dist-packages (from datasets) (0.26.5)
Requirement already satisfied: packaging in /usr/local/lib/python3.10/dist-packages (from datasets) (24.2)
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.10/dist-packages (from datasets) (6.0.2)
Requirement already satisfied: aiohappyeyeballs>=2.3.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (1.5.0)
Requirement already satisfied: aiosignal>=1.1.2 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (1.3.2)
Requirement already satisfied: async-timeout<6.0,>=4.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (4.0.3)
Requirement already satisfied: attrs>=17.3.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (24.2.0)
Requirement already satisfied: frozenlist>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (1.5.0)
Requirement already satisfied: multidict<7.0,>=4.5 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (6.1.0)
Requirement already satisfied: propcache>=0.2.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (0.2.1)
Requirement already satisfied: yarl<2.0,>=1.17.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets) (1.18.3)
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub->datasets) (4.12.2)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests>=2.32.2->datasets) (3.4.0)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests>=2.32.2->datasets) (3.10.1)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests>=2.32.2->datasets) (2.2.3)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests>=2.32.2->datasets) (2024.7.4)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.10/dist-packages (from pandas->datasets) (2.9.0)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas->datasets) (2024.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.10/dist-packages (from pandas->datasets) (2024.2)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.2->pandas->datasets) (1.17.0)
  Downloading datasets-3.2.0-py3-none-any.whl (480 kB)
  480.6/480.6 kB 28.3 MB/s eta 0:00:00
  Downloading dill-0.3.8-py3-none-any.whl (116 kB)
  116.3/116.3 kB 11.9 MB/s eta 0:00:00
  Downloading fsspec-2024.9.0-py3-none-any.whl (179 kB)
  179.3/179.3 kB 17.0 MB/s eta 0:00:00
  Downloading multiprocess-0.70.16-py310-none-any.whl (134 kB)
  134.8/134.8 kB 13.4 MB/s eta 0:00:00
  Downloading xxhash-3.5.0-cp310-cp310-manylinux_2_17_x86_64_manylinux2014_x86_64.whl (194 kB)
  194.1/194.1 kB 19.8 MB/s eta 0:00:00
Installing collected packages: xxhash, fsspec, dill, multiprocess, datasets
  Attempting uninstall: fsspec
    Found existing installation: fsspec 2024.10.0
    Uninstalling fsspec-2024.10.0:
      Successfully uninstalled fsspec-2024.10.0
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is
  gcsfs 2024.10.0 requires fsspec==2024.10.0, but you have fsspec 2024.9.0 which is incompatible.
Successfully installed datasets-3.2.0 dill-0.3.8 fsspec-2024.9.0 multiprocess-0.70.16 xxhash-3.5.0

README.md: 100% 520/520 [00:00<00:00, 33.4kB/s]

validation-00000-of-00001.parquet: 100% 9.19M/9.19M [00:00<00:00, 81.9MB/s]

Generating validation split: 100% 73/73 [00:00<00:00, 679.06 examples/s]
```

```
# time taken test
```

```
import time
```

```
def time_gen(model, inputs, **kwargs):
    start = time.time()
    outputs = model.generate(**inputs, **kwargs)
    gen_time = time.time() - start
    return outputs, gen_time
```

```
from tqdm import tqdm
```

```

all_time = 0
pred = []
ref = []

for sample in tqdm(dataset):
    audio = sample["audio"]
    inputs = processor(audio["array"], sampling_rate=audio["sampling_rate"], return_tensors="pt")
    inputs = inputs.to(device=device, dtype=torch.float16)

    output, gen_time = time_gen(model, inputs)
    all_time += gen_time
    pred.append(processor.batch_decode(output, skip_special_tokens=True, normalize=True)[0])
    ref.append(processor.tokenizer._normalize(sample["text"]))

print(f"Total time taken: {all_time:.2f}s")

```

0% | 0/73 [00:00<?, ?it/s] Due to a bug fix in <https://github.com/huggingface/transformers/pull/28687> transcripti
 Passing a tuple of 'past_key_values' is deprecated and will be removed in Transformers v4.43.0. You should pass an instance
 The attention mask is not set and cannot be inferred from input because pad token is same as eos token. As a consequence, yo
 /usr/local/lib/python3.10/dist-packages/transformers/models/whisper/tokenization_whisper.py:501: UserWarning: The private me
 warnings.warn(
 100% | 73/73 [01:48<00:00, 1.49s/it] Total time taken: 74.08s

[+ Code](#) [+ Text](#)

```

!pip install evaluate
!pip install jiwer
from evaluate import load

```

```

wer = load("wer")
print(wer.compute(predictions=pred, references=ref))

```

Requirement already satisfied: evaluate in /usr/local/lib/python3.10/dist-packages (0.4.3)
 Requirement already satisfied: datasets>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from evaluate) (3.2.0)
 Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.10/dist-packages (from evaluate) (1.26.4)
 Requirement already satisfied: dill in /usr/local/lib/python3.10/dist-packages (from evaluate) (0.3.8)
 Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (from evaluate) (2.2.2)
 Requirement already satisfied: requests>=2.19.0 in /usr/local/lib/python3.10/dist-packages (from evaluate) (2.32.3)
 Requirement already satisfied: tqdm>=4.62.1 in /usr/local/lib/python3.10/dist-packages (from evaluate) (4.66.6)
 Requirement already satisfied: xxhash in /usr/local/lib/python3.10/dist-packages (from evaluate) (3.5.0)
 Requirement already satisfied: multiprocessing in /usr/local/lib/python3.10/dist-packages (from evaluate) (0.70.16)
 Requirement already satisfied: fsspec>=2021.05.0 in /usr/local/lib/python3.10/dist-packages (from fsspec[http]>=2021.05.0->e
 Requirement already satisfied: huggingface-hub>=0.7.0 in /usr/local/lib/python3.10/dist-packages (from evaluate) (0.26.5)
 Requirement already satisfied: packaging in /usr/local/lib/python3.10/dist-packages (from evaluate) (24.2)
 Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from datasets>=2.0.0->evaluate) (3.16.1)
 Requirement already satisfied: pyarrow>=15.0.0 in /usr/local/lib/python3.10/dist-packages (from datasets>=2.0.0->evaluate) (15.0.2)
 Requirement already satisfied: aiohttp in /usr/local/lib/python3.10/dist-packages (from datasets>=2.0.0->evaluate) (3.11.10)
 Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.10/dist-packages (from datasets>=2.0.0->evaluate) (6.0.2)
 Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=
 Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests>=2.19.0->e
 Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests>=2.19.0->evaluate) (3.10.1)
 Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests>=2.19.0->evaluat
 Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests>=2.19.0->evaluat
 Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.10/dist-packages (from pandas->evaluate) (2.9.0)
 Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas->evaluate) (2024.2)
 Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.10/dist-packages (from pandas->evaluate) (2024.2)
 Requirement already satisfied: aiohappyeyeballs>=2.3.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets>=2
 Requirement already satisfied: aiosignal>=1.1.2 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets>=2.0.0->e
 Requirement already satisfied: async-timeout<6.0,>=4.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets>=2
 Requirement already satisfied: attrs>=17.3.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets>=2.0.0->eval
 Requirement already satisfied: frozenlist>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets>=2.0.0->
 Requirement already satisfied: multidict<7.0,>=4.5 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets>=2.0.0
 Requirement already satisfied: propcache>=0.2.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets>=2.0.0->e
 Requirement already satisfied: yarl<2.0,>=1.17.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets>=2.0.0->
 Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.2->pandas->eva
 Requirement already satisfied: jiwer in /usr/local/lib/python3.10/dist-packages (3.0.5)
 Requirement already satisfied: click<9.0.0,>=8.1.3 in /usr/local/lib/python3.10/dist-packages (from jiwer) (8.1.7)
 Requirement already satisfied: rapidfuzz<4,>=3 in /usr/local/lib/python3.10/dist-packages (from jiwer) (3.11.0)
 0.03507271171941831

```

from transformers import AutoModelForCausalLM

```

```

assistant_model_id = "distil-whisper/distil-large-v2"

```

```

assistant_model = AutoModelForCausalLM.from_pretrained(
    assistant_model_id,
    torch_dtype=torch_dtype,
    low_cpu_mem_usage=True, # fast loading
    use_safetensors=True, # secure (over pickle)
)

```

```
    attn_implementation="sdpa", # Flash Attention speed-up
)
```

```
assistant_model.to(device)
```



```
config.json: 100% 2.29k/2.29k [00:00<00:00, 43.5kB/s]
model.safetensors: 100% 1.51G/1.51G [00:36<00:00, 41.2MB/s]
generation_config.json: 100% 3.62k/3.62k [00:00<00:00, 229kB/s]
```

```
WhisperForCausalLM(
  (model): WhisperDecoderWrapper(
    (decoder): WhisperDecoder(
      (embed_tokens): Embedding(51865, 1280, padding_idx=50257)
      (embed_positions): WhisperPositionalEmbedding(448, 1280)
      (layers): ModuleList(
        (0-1): 2 x WhisperDecoderLayer(
          (self_attn): WhisperSdpaAttention(
            (k_proj): Linear(in_features=1280, out_features=1280, bias=False)
            (v_proj): Linear(in_features=1280, out_features=1280, bias=True)
            (q_proj): Linear(in_features=1280, out_features=1280, bias=True)
            (out_proj): Linear(in_features=1280, out_features=1280, bias=True)
          )
          (activation_fn): GELUActivation()
          (self_attn_layer_norm): LayerNorm((1280,), eps=1e-05, elementwise_affine=True)
          (encoder_attn): WhisperSdpaAttention(
            (k_proj): Linear(in_features=1280, out_features=1280, bias=False)
            (v_proj): Linear(in_features=1280, out_features=1280, bias=True)
            (q_proj): Linear(in_features=1280, out_features=1280, bias=True)
            (out_proj): Linear(in_features=1280, out_features=1280, bias=True)
          )
          (encoder_attn_layer_norm): LayerNorm((1280,), eps=1e-05, elementwise_affine=True)
          (fc1): Linear(in_features=1280, out_features=5120, bias=True)
          (fc2): Linear(in_features=5120, out_features=1280, bias=True)
          (final_layer_norm): LayerNorm((1280,), eps=1e-05, elementwise_affine=True)
        )
      )
    )
    (layer_norm): LayerNorm((1280,), eps=1e-05, elementwise_affine=True)
  )
  (proj_out): Linear(in_features=1280, out_features=51865, bias=False)
)
```

```
def assisted_generate_with_time(model, inputs, **kwargs):
    start_time = time.time()
    outputs = model.generate(**inputs, assistant_model=assistant_model, **kwargs)
    generation_time = time.time() - start_time
    return outputs, generation_time
```

```
all_time = 0
pred = []
ref = []
```

```
for sample in tqdm(dataset):
    audio = sample["audio"]
    inputs = processor(audio["array"], sampling_rate=audio["sampling_rate"], return_tensors="pt")
    inputs = inputs.to(device=device, dtype=torch.float16)

    out, gen_time = assisted_generate_with_time(model, inputs)
    all_time += gen_time
    pred.append(processor.batch_decode(out, skip_special_tokens=True, normalize=True)[0])
    ref.append(processor.tokenizer._normalize(sample["text"]))
```

```
print(f"Total time taken: {all_time:.2f}s")
```



```
0%|          | 0/73 [00:00<?, ?it/s]From v4.47 onwards, when a model cache is to be returned, `generate` will return a `Ca
/usr/local/lib/python3.10/dist-packages/transformers/models/whisper/tokenization_whisper.py:501: UserWarning: The private me
warnings.warn(
100%|██████████| 73/73 [01:14<00:00, 1.02s/it]Total time taken: 70.96s
```

```
print(wer.compute(predictions=pred, references=ref))
```



```
0.03507271171941831
```

Start coding or [generate](#) with AI.

