

```

# %%capture
# #!unzip Datasets.zip

# from google.colab import drive

# # Mount the Google Drive
# drive.mount('/content/drive')

# %%capture
# !pip install datasets
# !pip install transformers
# !pip install librosa
# !pip install jiwer
# !pip install evaluate

import os
import datasets
import pandas as pd
from sklearn.model_selection import train_test_split
from datasets import Dataset

# Set paths
csv_path = "/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final.csv"
audio_folder = "/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final-w"

# Load the CSV
df = pd.read_csv(csv_path)
df = pd.read_csv(csv_path)
# Ensure the column names match
df.columns = ["Filename", "Transcription"] # Rename columns if needed

# Append '.wav' to the file names
df['Filename'] = df['Filename'].apply(lambda x: f"{x}.wav")

# Add full paths to the audio files
df['file_path'] = df['Filename'].apply(lambda x: os.path.join(audio_folder, x))

# Verify that all audio files exist
missing_files = df[~df['file_path'].apply(os.path.exists)]
if not missing_files.empty:
    print("The following audio files are missing:")
    print(missing_files)
    raise FileNotFoundError("Some audio files listed in the CSV are missing in the folder.")

# Split into train (27) and test (3)
train_df, test_df = train_test_split(df, test_size=0.2, random_state=42)

# Save splits to CSV for reference
train_csv_path = "train_split.csv"
test_csv_path = "test_split.csv"
train_df.to_csv(train_csv_path, index=False)
test_df.to_csv(test_csv_path, index=False)

# Convert to HuggingFace Dataset format
train_dataset = Dataset.from_pandas(train_df)
test_dataset = Dataset.from_pandas(test_df)

# Save HuggingFace datasets
train_dataset_path = "train_dataset"
test_dataset_path = "test_dataset"
train_dataset.save_to_disk(train_dataset_path)
test_dataset.save_to_disk(test_dataset_path)

# Output
print(f"Train set saved to: {train_csv_path} and {train_dataset_path}")
print(f"Test set saved to: {test_csv_path} and {test_dataset_path}")

↗ /home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/tqdm/auto.py:21: TqdmWarning: IProgress not found. Please
from .autonotebook import tqdm as notebook_tqdm
Saving the dataset (1/1 shards): 100%|██████████| 2600/2600 [00:00<00:00, 586017.00 examples/s]
Saving the dataset (1/1 shards): 100%|██████████| 650/650 [00:00<00:00, 376144.81 examples/s]Train set saved to: train_s
Test set saved to: test_split.csv and test_dataset

from datasets import load_from_disk

train_dataset = load_from_disk("train_dataset")
test_dataset = load_from_disk("test_dataset")

```

```
show_random_elements(train_dataset)
```

	Filename	Transcription	file_path	__index_level_0__
0	8140118.wav	چانیہ گج تھہ مکانس نہ اسی وار نہ بر	/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final-waves/8140118.wav	117
1	rafiya-06_13.wav	فارسی شوی نگارو ساس مُنتلق وُنت چھہ تر چھہ مہ در نظر	/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final-waves/rafiya-06_13.wav	3196
2	8140072.wav	پوان لہذا یمن بنزیر ترتیب قائم مچاوی چھہ	/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final-waves/8140072.wav	71
3	8150134.wav	کہ رنگ اکھ سینھار زت اضافہ بلکہ ونہ	/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final-waves/8150134.wav	530
4	farhat-01_16.wav	نِس کئی تام رازن بیس عشوش ونان اسی بسو ومنت	/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final-waves/farhat-01_16.wav	732
5	8140165.wav	تیداد داہے وُنت چھہ تر چھہ واریاہ کم	/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final-waves/8140165.wav	164
6	8140173.wav	تیر ہر گاہ مہ اکاڈمی اندر کانہہ ہڑ روایت	/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final-waves/8140173.wav	172
7	jhon-02_104.wav	ہستی ہڑ ہجو ز تو ہم نہش کتھہ بو کتھہ ہو چھہ از تام فائش حاصل کورم	/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final-waves/jhon-02_104.wav	2553
8	ishrat1-40_6.wav	بر ووتھر ہے اگر شیکھ وونل کُاب ولو	/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final-waves/ishrat1-40_6.wav	1773
9	guddy1-09_3.wav	شمن لہ گُش مہ تر چھہ نہ زمین انہ گُش لہ تر چھہ انہ زمین اسی	/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final-waves/guddy1-09_3.wav	1248

 Show hidden output

```

vocab_dict["|"] = vocab_dict[" "]
del vocab_dict[" "]

vocab_dict["[UNK]"] = len(vocab_dict)
vocab_dict["[PAD]"] = len(vocab_dict)
len(vocab_dict)

↩ 60

import json
with open('vocab.json', 'w') as vocab_file:
    json.dump(vocab_dict, vocab_file)

from transformers import Wav2Vec2CTCTokenizer

tokenizer = Wav2Vec2CTCTokenizer.from_pretrained("./", unk_token="[UNK]", pad_token="[PAD]", word_delimiter_token="|", clear

from transformers import Wav2Vec2FeatureExtractor

feature_extractor = Wav2Vec2FeatureExtractor(feature_size=1, sampling_rate=16000, padding_value=0.0, do_normalize=True, retu

from transformers import Wav2Vec2Processor

processor = Wav2Vec2Processor(feature_extractor=feature_extractor, tokenizer=tokenizer)

train_dataset[0]["file_path"]

↩ ' /home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/final-
waves/farhat-03_21.wav'

```

Replacing the File Path with Actual Audio.

```

from datasets import load_from_disk, Audio

# Load datasets
train_dataset = load_from_disk("train_dataset") # Adjust to your actual path
test_dataset = load_from_disk("test_dataset")

# Rename 'file_path' to 'audio'
train_dataset = train_dataset.rename_column("file_path", "audio")
test_dataset = test_dataset.rename_column("file_path", "audio")

# # Cast the 'audio' column to use the Audio feature
train_dataset = train_dataset.cast_column("audio", Audio(sampling_rate=16_000))
test_dataset = test_dataset.cast_column("audio", Audio(sampling_rate=16_000))

# # Drop unnecessary columns if needed
train_dataset = train_dataset.remove_columns(["__index_level_0__"])
test_dataset = test_dataset.remove_columns(["__index_level_0__"])

# # Verify the dataset structure
print(train_dataset)
print(test_dataset)

# # Inspect the first example
print(train_dataset[0])

↩ Dataset({
  features: ['Filename', 'Transcription', 'audio'],
  num_rows: 2600
})
Dataset({
  features: ['Filename', 'Transcription', 'audio'],
  num_rows: 650
})
{'Filename': 'farhat-03_21.wav', 'Transcription': 'زأنم شاه صأبن بدشابس\مرا قېم پتم تم ييئلم', 'audio': {'path': '/ho
0.0256958 , 0.02392578}], 'sampling_rate': 16000}}

#print(test_dataset[0]['audio'])

rand_int = random.randint(0, len(train_dataset))

print("Target text:", train_dataset[rand_int]["Transcription"])
print("Input array shape:", train_dataset[rand_int]["audio"]["array"].shape)
print("Sampling rate:", train_dataset[rand_int]["audio"]["sampling_rate"])

```

↗ Target text: یم میون کٹھ تم یم دوشوے چھ اکھ اُکس کھوژان
 Input array shape: (64000,) Sampling rate: 16000

```
def prepare_dataset(batch):
    audio = batch["audio"]

    # batched output is "un-batched"
    batch["input_values"] = processor(audio["array"], sampling_rate=audio["sampling_rate"]).input_values[0]
    batch["input_length"] = len(batch["input_values"])

    batch["labels"] = processor(text=batch["Transcription"]).input_ids

    return batch
```

```
train_dataset = train_dataset.map(prepare_dataset, remove_columns=train_dataset.column_names)
test_dataset = test_dataset.map(prepare_dataset, remove_columns=test_dataset.column_names)
```

↗ Map: 0%| | 0/2600 [00:00<?, ? examples/s]2025-04-05 00:09:12.277485: E external/local_xla/xla/stream_executor WARNING: All log messages before absl::InitializeLog() is called are written to STDERR
 E0000 00:00:1743791952.328213 12371 cuda_dnn.cc:8310] Unable to register cuDNN factory: Attempting to register factory
 E0000 00:00:1743791952.343200 12371 cuda_blas.cc:1418] Unable to register cuBLAS factory: Attempting to register facto
 2025-04-05 00:09:12.460245: I tensorflow/core/platform/cpu_feature_guard.cc:210] This TensorFlow binary is optimized to
 To enable the following instructions: AVX2 FMA, in other operations, rebuild TensorFlow with the appropriate compiler fl
 Map: 100%| | 2600/2600 [00:05<00:00, 478.89 examples/s]
 Map: 100%| | 650/650 [00:00<00:00, 753.23 examples/s]

```
import torch
```

```
from dataclasses import dataclass, field
from typing import Any, Dict, List, Optional, Union
```

```
@dataclass
class DataCollatorCTCWithPadding:
    """
    Data collator that will dynamically pad the inputs received.
    Args:
        processor (:class:`~transformers.Wav2Vec2Processor`)
            The processor used for processing the data.
        padding (:obj:`bool`, :obj:`str` or :class:`~transformers.tokenization_utils_base.PaddingStrategy`, `optional`, defa
            Select a strategy to pad the returned sequences (according to the model's padding side and padding index)
            among:
            * :obj:`True` or :obj:`'longest'`: Pad to the longest sequence in the batch (or no padding if only a single
              sequence if provided).
            * :obj:`'max_length'`: Pad to a maximum length specified with the argument :obj:`max_length` or to the
              maximum acceptable input length for the model if that argument is not provided.
            * :obj:`False` or :obj:`'do_not_pad'` (default): No padding (i.e., can output a batch with sequences of
              different lengths).
    """

    processor: Wav2Vec2Processor
    padding: Union[bool, str] = True

    def __call__(self, features: List[Dict[str, Union[List[int], torch.Tensor]]]) -> Dict[str, torch.Tensor]:
        # split inputs and labels since they have to be of different lengths and need
        # different padding methods
        input_features = [{"input_values": feature["input_values"]} for feature in features]
        label_features = [{"input_ids": feature["labels"]} for feature in features]

        batch = self.processor.pad(
            input_features,
            padding=self.padding,
            return_tensors="pt",
        )

        with self.processor.as_target_processor():
            labels_batch = self.processor.pad(
                label_features,
                padding=self.padding,
                return_tensors="pt",
            )

        # replace padding with -100 to ignore loss correctly
        labels = labels_batch["input_ids"].masked_fill(labels_batch.attention_mask.ne(1), -100)

        batch["labels"] = labels

        return batch

data_collator = DataCollatorCTCWithPadding(processor=processor, padding=True)
```

```
import evaluate
```

```
wer_metric = evaluate.load("wer")
```

➡ Using the latest cached version of the module from /home/muzaffar/.cache/huggingface/modules/evaluate_modules/metrics/ev

```
from evaluate import load
```

```
cer_metric = load("cer")
```

➡ Using the latest cached version of the module from /home/muzaffar/.cache/huggingface/modules/evaluate_modules/metrics/ev

includes both WER and CER

```
def compute_metrics(pred):
    pred_logits = pred.predictions
    pred_ids = np.argmax(pred_logits, axis=-1)

    # Replace padding token (-100) with pad_token_id
    pred.label_ids[pred.label_ids == -100] = processor.tokenizer.pad_token_id

    # Decode predictions and labels to strings
    pred_str = processor.batch_decode(pred_ids)
    label_str = processor.batch_decode(pred.label_ids, group_tokens=False)

    if isinstance(label_str, list):
        if isinstance(pred_str, list) and len(pred_str) == len(label_str):
            for index in random.sample(range(len(label_str)), 3):
                print(f'reference: "{label_str[index]}"')
                print(f'predicted: "{pred_str[index]}"')

        else:
            for index in random.sample(range(len(label_str)), 3):
                print(f'reference: "{label_str[index]}"')
                print(f'predicted: "{pred_str}"')

    # Compute WER
    wer = wer_metric.compute(predictions=pred_str, references=label_str)

    # Compute CER
    cer = cer_metric.compute(predictions=pred_str, references=label_str)

    return {"wer": wer, "cer": cer}
```

```
from transformers import Wav2Vec2ForCTC
```

```
model = Wav2Vec2ForCTC.from_pretrained(
    # "facebook/wav2vec2-xls-r-300m",
    'facebook/wav2vec2-large-xlsr-53',
    attention_dropout=0.05,
    hidden_dropout=0.1,
    feat_proj_dropout=0.1,
    mask_time_prob=0.05,
    layerdrop=0.01377,
    gradient_checkpointing=True,
    ctc_loss_reduction="mean",
    ctc_zero_infinity=True,
    pad_token_id=processor.tokenizer.pad_token_id,
    vocab_size=len(processor.tokenizer),

)
```

➡ Some weights of Wav2Vec2ForCTC were not initialized from the model checkpoint at facebook/wav2vec2-large-xlsr-53 and are You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```
# from transformers import Wav2Vec2ForCTC
```

```
# model = Wav2Vec2ForCTC.from_pretrained(
#     'facebook/wav2vec2-large-xlsr-53',
#     attention_dropout=0.05,
#     activation_dropout=0.1,
#     hidden_dropout=0.1,
#     feat_proj_dropout=0.01249,
#     final_dropout=0.0,
#     mask_time_prob=0.05,
```

```
# mask_time_length=10,
# mask_feature_prob=0,
# mask_feature_length=10,
# layerdrop=0.01377,
# gradient_checkpointing=True,
# ctc_loss_reduction="mean",
# ctc_zero_infinity=True,
# bos_token_id=processor.tokenizer.bos_token_id,
# eos_token_id=processor.tokenizer.eos_token_id,
# pad_token_id=processor.tokenizer.pad_token_id,
# vocab_size=len(processor.tokenizer.get_vocab())
# )
```

```
model.freeze_feature_encoder()
```

```
# import huggingface_hub
```

```
# huggingface_hub.login()
```

```
#repo_name = "wav2vec2-kashmiri-jhon-data-one"
```

```
save_dir = "/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/training_ε"
```

```
from transformers import TrainingArguments
```

```
training_args = TrainingArguments(
    output_dir=save_dir,
    group_by_length=True,
    per_device_train_batch_size=8,
    per_device_eval_batch_size=8,
    gradient_accumulation_steps=2,
    evaluation_strategy="steps",
    num_train_epochs=30,
    fp16=True,
    save_steps=500,
    eval_steps=500,
    logging_steps=10,
    learning_rate=4e-4,
    warmup_steps=250,
    save_total_limit=2,
    dataloader_num_workers=24
)
```

```
→ /home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/training_args.py:1594: FutureWarning: `eval`
warnings.warn()
```

```
# from transformers import TrainingArguments
```

```
# training_args = TrainingArguments(
#     output_dir=repo_name,
#     group_by_length=True,
#     per_device_train_batch_size=2,
#     gradient_accumulation_steps=2,
#     eval_strategy="steps",
#     num_train_epochs=20,
#     gradient_checkpointing=True,
#     fp16=True,
#     save_steps=20,
#     eval_steps=20,
#     logging_steps=40,
#     learning_rate=3e-4,
#     warmup_steps=50,
#     save_total_limit=2,
#     push_to_hub=True,
# )
```

```
# import numpy as np
# from transformers import Trainer
# trainer = Trainer(
#     model=model,
#     data_collator=data_collator,
#     args=training_args,
#     compute_metrics=compute_metrics,
#     train_dataset=train_dataset,
#     eval_dataset=test_dataset,
#     tokenizer=processor.feature_extractor,
# )
```

```
import numpy as np
from transformers import Trainer

# Assuming processor is an instance of Wav2Vec2Processor (or similar for your model)
trainer = Trainer(
    model=model,
    data_collator=data_collator,
    args=training_args,
    compute_metrics=compute_metrics,
    train_dataset=train_dataset,
    eval_dataset=test_dataset,
    processing_class=processor, # Use the processor directly for feature extraction
)

print("step1")
train_result = trainer.train()
print("step2")

metrics = train_result.metrics
print("step3")
max_train_samples = len(train_dataset)
metrics["train_samples"] = min(max_train_samples, len(train_dataset))
print("step4")
trainer.save_model()
print("model created!")
trainer.log_metrics("train", metrics)
trainer.save_metrics("train", metrics)
trainer.save_state()
```

Step	Training Loss	Validation Loss	Wer	Cer
500	1.442300	1.118181	0.860588	0.301463
1000	0.524800	0.618262	0.568061	0.162705

```
>> Trainer.tokenizer now deprecated. You should use Trainer.processing_class instead.  
Trainer.tokenizer is now deprecated. You should use Trainer.processing_class instead.  
2500 /home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/torch/utils/data/dataloader.py:624: UserWarning: This Da  
warnings.warn(  
3000 /home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
warnings.warn(  
3500 /home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
4000 nings.wa0.713100 0.813449 0.446008 0.125284  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
4500 nings.wa0.72100 0.834218 0.446739 0.125423  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
  
# trainer.push_to_hub()  
  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
~ TESTING WITH CPU KERNL  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:  
/home/muzaffar/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
```



```
import torch
import torchaudio
import librosa
import numpy
from transformers import WavVec2ForCTC, WavVec2Processor
from transformers import WavVec2Processor
```

```

/home/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
Start coding or generate with AI.

```

```

/home/muzaif/anaconda3/envs/trl4/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
model_name_or_path = "/home/muzaif/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/KASHMIRI/experiment5/
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print(model_name_or_path, device)

processor = Wav2Vec2Processor.from_pretrained(model_name_or_path)
model = Wav2Vec2ForCTC.from_pretrained(model_name_or_path).to(device)

```

```
def speech_file_to_array_fn(batch):
    speech_array, sampling_rate = torchaudio.load(batch["file_path"])
    speech_array = speech_array.squeeze().numpy()
    #speech_array = librosa.resample(np.asarray(speech_array), sampling_rate, processor.feature_extractor.sampling_rate)
    speech_array = librosa.resample(y=np.asarray(speech_array), orig_sr=sampling_rate, target_sr=processor.feature_extractor
```

```
batch["speech"] = speech_array
return batch
```

```
def predict(batch):
    features = processor(
        batch["speech"],
        sampling_rate=processor.feature_extractor.sampling_rate,
        return_tensors="pt",
        padding=True
    )

    input_values = features.input_values.to(device)
    #attention_mask = features.attention_mask.to(device)
    attention_mask = features.attention_mask.to(device) if "attention_mask" in features else None

    with torch.no_grad():
        logits = model(input_values, attention_mask=attention_mask).logits

    pred_ids = torch.argmax(logits, dim=-1)

    batch["predicted_N_LM"] = processor.batch_decode(pred_ids)
    return batch
```

```
import torchaudio
import librosa
from datasets import load_dataset
import numpy as np
```

```
dataset = load_dataset("csv", data_files={"/home/muzaffar/Desktop/Research/papers/5-paper Wav2Vec/5. Wave2vec Whisper Paper/"
dataset = dataset.map(speech_file_to_array_fn)
```

[illegible]

10/11

```

909 drop_last_batch=drop_last_batch,
910 remove_columns=remove_columns,
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
911 keep_in_memory=keep_in_memory,
912 load_from_cache_file=load_from_cache_file
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
913 cache_file_name=cache_file_names[k],
914 writer_batch_size=writer_batch_size
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
915 features=features,
916 disable_nullable=disable_nullable,
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
917 fn_kwargs=fn_kwargs,
918 num_proc=num_proc,
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
919 desc=desc,
920 )
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
921 for k, dataset in self.items()
922 }
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
923 warnings.warn(
File ~/anaconda3/envs/tf14/lib/python3.11/site-packages/datasets/dataset_dict.py:902, in DatasetDict.compact()
903     warnings.warn(
904         "Cache file names is None"
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
905     cache_file_names = {k: None for k in self}
906     return DatasetDict(
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
907     warnings.warn(
908     > 909     function=function,
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
909     with_indices=with_indices,
910     with_rank=with_rank,
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
911     input_columns=input_columns,
912     batched=batched,
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
913     batch_size=batch_size,
914     drop_last_batch=drop_last_batch,
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
915     remove_columns=remove_columns,
916     keep_in_memory=keep_in_memory,
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
917     load_from_cache_file=load_from_cache_file,
918     cache_file_name=cache_file_names[k],
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
919     writer_batch_size=writer_batch_size,
920     features=features,
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
921     disable_nullable=disable_nullable,
922     fn_kwargs=fn_kwargs,
923     num_proc=num_proc,
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
924     desc=desc,
925 )
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
926 for k, dataset in self.items()
927 warnings.warn(
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
928 warnings.warn(
File ~/anaconda3/envs/tf14/lib/python3.11/site-packages/datasets/arrow_dataset.py:562, in DatasetDict._format()
929     warnings.warn(
930     > 931     self._format = {
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
931     warnings.warn(
932     > 933     "format": self._format_type,
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
934     warnings.warn(
935     > 936     "format_kwargs": self._format_kwargs,
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
937     warnings.warn(
938     > 939     "columns": self._format_columns,
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
940     warnings.warn(
941     > 942     "output_all_columns": self._output_all_columns,
943 )
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
944 warnings.warn(
945     > 946     "apply actual function"
947     > 947     out = Union["Dataset", "DatasetDict"] = func(self, *args, **kwargs)
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
948     > 948     dataset = List["Dataset"] = list(out.values()) if isinstance(out, dict) else [out]
949     > 949     # manually format to the output
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
950     warnings.warn(
951     > 951     File ~/anaconda3/envs/tf14/lib/python3.11/site-packages/datasets/arrow_dataset.py:3079, in Dataset.map(self, function,
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
952     > 952     with indices, with_rank, input_columns, batched, batch_size, drop_last_batch, remove_columns, keep_in_memory,
953     > 953     load_from_cache_file, cache_file_name, writer_batch_size, features, disable_nullable, fn_kwargs, num_proc,
954     > 954     suffix_template, new_fingerprint, desc)
955     > 955     warnings.warn(
956     > 956     transformed dataset is None:
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
957     > 957     warnings.warn(
958     > 958     unit="examples"
959     > 959     total_pbar=total,
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
960     > 960     warnings.warn(
961     > 961     desc=desc or "Map",
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
962     > 962     warnings.warn(
963     > 963     as_pbar:
964     > 964     for rank, done, content in Dataset.map_single(**dataset_kwargs):
/home/muza... site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
965     > 965     if done:
966     > 966     shards_done += 1
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
967     > 967     warnings.warn(
968     > 968     File ~/anaconda3/envs/tf14/lib/python3.11/site-packages/datasets/arrow_dataset.py:3546, in Dataset.map_single(shard
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
969     > 969     function, with_indices, with_rank, input_columns, batched, batch_size, drop_last_batch, remove_columns, keep_in_memory,
970     > 970     cache_file_name, writer_batch_size, features, disable_nullable, fn_kwargs, new_fingerprint, rank, offset)
971     > 971     warnings.warn(
972     > 972     if data:
973     > 973     warnings.warn(
974     > 974     if writer is not None:
975     > 975     warnings.warn(
976     > 976     if tmp_file is not None:
/home/muza... anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
977     > 977     warnings.warn(
978     > 978     File ~/anaconda3/envs/tf14/lib/python3.11/site-packages/datasets/arrow_writer.py:636, in ArrowWriter._finalize(self,
979     > 979     close_stream)
980     > 980     warnings.warn(
981     > 981     # Re-initializing to empty list for next batch
982     > 982     warnings.warn(
983     > 983     if hkey_record = 1
984     > 984     self.write_examples_on_file()
985     > 985     warnings.warn(
986     > 986     if schema is known, infer features even if no examples were written
987     > 987     if self.pa_writer is None and self.schema:
988     > 988     warnings.warn(
989     > 989     File ~/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
990     > 990     File ~/anaconda3/envs/tf14/lib/python3.11/site-packages/datasets/arrow_writer.py:495, in
991     > 991     ArrowWriter._write_examples_on_file(self)
992     > 992     File ~/anaconda3/envs/tf14/lib/python3.11/site-packages/transformers/models/wav2vec2/processing_wav2vec2.py:174:
993     > 993     else:

```