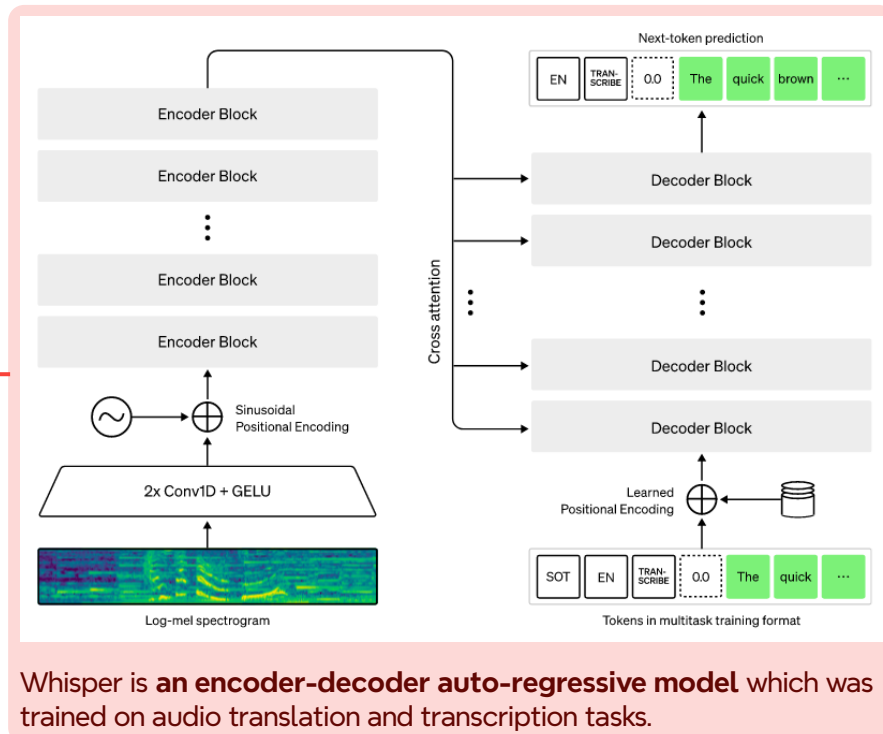


Introduction about Whisper

Introduction:



Given audio data, the model is able to generate the corresponding text.

This notebook will allow you to play around with the model and evaluate it in a few lines of code. If you want more details about Whisper, checkout each 'model card' (see for example [whisper-tiny] (https://huggingface.co/openai/whisper-tiny)) or the [original paper] (https://openai.com/blog/whisper/).

Encoder

A log-mel spectrogram is extracted from a raw audio using a [Processor] (https://huggingface.co/docs/transformers/main_classes/processors), before it is passed to the encoder.

Decoder

The decoder inputs are text tokens, and special tokens such as "<startoftext>", "<|transcribe|>" and "<|en|>" are used to specify the desired task and the language of the audio.

Two different models were trained, an English only and a multilingual model. English only checkpoints have the '.en' suffix on the Hub and use a different tokenizer vocabulary.

package installation

```
!pip install git+https://github.com/huggingface/transformers.git
!pip install datasets
!pip install evaluate
!wget https://cdn-media.huggingface.co/speech_samples/sample1.flac
!wget https://huggingface.co/datasets/openai/whisper-tiny
!pip install gradio
!pip install jwer
```

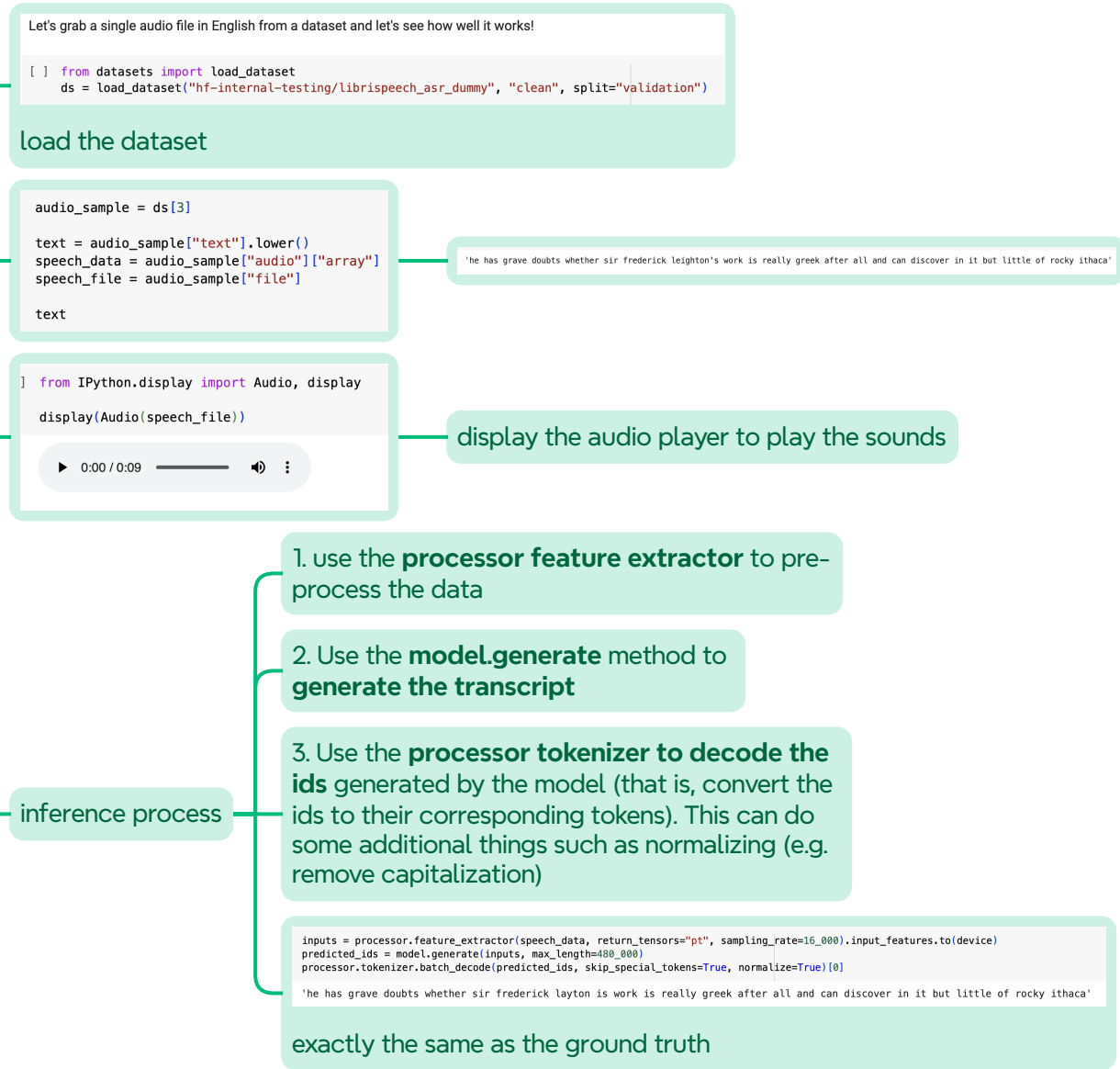
Loading the model and the processor

```
import torch
from transformers import WhisperForConditionalGeneration, WhisperProcessor

device = "cuda" if torch.cuda.is_available() else "cpu"
model = WhisperForConditionalGeneration.from_pretrained("openai/whisper-small").to(device)
processor = WhisperProcessor.from_pretrained("openai/whisper-small")

load the model and processor
```

Simple Inference



Running Inference on a dataset using a base Whisper Model

```
ds = load_dataset("hf-internal-testing/libraspeech_asr_dummy", "clean", split="validation")

def eval_fn(batch):
    arrays = list(batch)
    inputs = processor.feature_extractor(audio_arrays, return_tensors="pt")
    outputs = model.generate(inputs, max_length=400, return_tensors="pt")
    results = processor.tokenizer.batch_decode(outputs, skip_special_tokens=True, normalize=True)
    batch_predictions = [result for result in results]
    batch_references = [reference for reference in batch["text"]]
    return batch_predictions, batch_references

ds = ds.map(eval_fn, batch_size=1, remove_columns=batched=True)
```

obtained batched inference results

```
from evaluate import load
wer = load("wer")
wer_score = wer.compute(predictions=predictions, references=references)
print(f"WER: {wer_score * 100:.2f} %")
WER: 6.16 %
```

Calculating the word error rate

Inference in Other Languages

```
from datasets import load_dataset
sample = load_dataset("osanseviero/dummy_a_audio") ["train"] ["audio"] [0]
speech_data = sample["array"]

inputs = processor.feature_extractor(speech_data, return_tensors="pt")
sampling_rate=8000, input_features.to(device)
forced_decoder_id = processor.get_decoder_prompt_ids(language="ja", task="transcribe")
predictions_ids = model.generate(inputs, max_length=400, forced_decoder_id=forced_decoder_id)
processor.tokenizer.batch_decode(predictions_ids, skip_special_tokens=True, normalize=True)

"カキハチ 8000 Hz 16-bit Little Endian"
```

identify text and translate it to other languages at the same time

processor to extract the feature
set up forced_decoder_id -> this can force the model to generate japanese characters

Too much code!

```
[ ] from transformers import pipeline

pipe = pipeline("automatic-speech-recognition", model="openai/whisper-small")
pipe(speech_file)
```

pipeline is a very convenient tool provided by OpenAI

```
from transformers import pipeline
import gradio as gr

pipe = pipeline("automatic-speech-recognition", model="openai/whisper-small")

def inference(speech_file):
    return pipe(speech_file) ["text"]

gr.Interface(inference, gr.Audio(type="filepath"), "text").launch()

Colab notebook detected. To show errors in colab notebook, set debug=True in launch()
Notes: opening Chrome Inspector may crash dom inside Colab notebooks.
To create a public link, set 'share=True' in 'launch()'.
```

gradio is a visualization tool

Step by Step

1. audio is represented as an array

```
[ ] print(len(speech_data), speech_data)

150000 [-0.00045776 -0.00039673 -0.00040020 ... -0.00021362 -0.00030311
-0.00027466]
```

2. Using the feature extractor, we pre-process the audio to a format usable by the model (i.e. extract the log-mel spectrogram from the audio).

```
inputs = processor.feature_extractor(speech_data, return_tensors="pt")
sampling_rate=8000, input_features.to(device)
inputs
```

print(inputs.shape)
torch.Size([1, 80, 3000])

3. We pass this input to the model generate methods

```
predicted_ids = model.generate(inputs, max_length=400)
predicted_ids
```

```
tensor([ [50258, 50259, 50359, 50363, 634, 575, 12525, 22618, 1968, 6144,
35617, 20804, 1756, 311, 589, 387, 534, 10201, 934, 439,
293, 395, 4611, 206, 309, 457, 787, 250, 6100, 1,
392, 6628, 13, 50257], device='cuda:0'] )
```

model predict the text

4. Tokenizer decodes the ids to actual readable text

```
processor.tokenizer.batch_decode(predicted_ids)

["He has grave doubts whether Sir Frederick Layton's work is really Greek after all and can discover in it but little of rocky Ithaca."]

decode and normalize to a sentence
```

```
[ ] processor.tokenizer.batch_decode(predicted_ids, skip_special_tokens=True) [0]

'He has grave doubts whether Sir Frederick Layton's work is really Greek after all and can discover in it but little of rocky Ithaca.'
```

get rid of special tokens

```
!python3.10 -m pip install transformers==4.18.0
/usr/local/lib/python3.10/dist-packages/transformers/generation/utils.py:1346: UserWarning: Using 'max_length' to control the generation length is deprecated in v4.9+ and will be removed in v5.0. Please use 'max_new_tokens' instead.
warnings.warn(
("text") "He has grave doubts whether Sir Frederick Layton's work is really Greek after all and can discover in it but little of rocky Ithaca.")
```

just two lines of code and done!

Calculation Example
Suppose the audio is sampled at 16,000 Hz and has a duration of 20 seconds. The total number of samples N would be:
 $N = 16,000 \times 20 = 400,000$ samples

Before Log-Mel Spectrogram
Before creating the log-Mel spectrogram, the raw audio typically has the following dimensions:
1. Channels: Audio is typically 1 (mono) or 2 (stereo) channels.
2. Samples: The length of the audio in terms of samples.
For example, a stereo audio signal with 2 channels and N samples would have a shape of $[2, N]$. A mono audio signal with N samples would have a shape of $[1, N]$.

After Log-Mel Spectrogram
The log-Mel spectrogram has the following dimensions:
1. Frequency Bins: This represents the number of Mel filterbank outputs. In Whisper, this is typically set to 80.
2. Time Frames: This depends on the length of the audio signal and the hop length used when creating the Mel spectrogram.

Now, if we use an FFT size $n_{\text{fft}} = 400$ and a hop length $\text{hop_length} = 160$ for the Mel spectrogram:
Time Frames = $\lceil \frac{400,000 - 400}{160} \rceil + 1 = 2999$
Thus, the resulting log-Mel spectrogram would have a shape of:
Shape = $[80, 2999]$

FFT means Fast Fourier Transform (an advanced and quite complex concept in signal processing, optional to understand)