

Business AI Meeting Companion STT



Introduction

Consider you're attending a business meeting where all conversations are being captured by an advanced AI application. This application not only transcribes the discussions with high accuracy but also provides a concise summary of the meeting, emphasizing the key points and decisions made.

In our project, we'll use OpenAI's Whisper to transform speech into text. Next, we'll use IBM Watson's AI to summarize and find key points. We'll make an app with Hugging Face Gradio as the user interface.

Learning Objectives

After finishing this lab, you will able to:

- Create a Python script to generate text using a model from the Hugging Face Hub, identify some key parameters that influence the model's output, and have a basic understanding of how to switch between different LLM models.
- Use OpenAI's Whisper technology to convert lecture recordings into text, accurately.
- Implement IBM Watson's AI to effectively summarize the transcribed lectures and extract key points.
- Create an intuitive and user-friendly interface using Hugging Face Gradio, ensuring ease of use for students and educators.



Generated by DALL-E-3

Preparing the environment

Let's start with setting up the environment by creating a Python virtual environment and installing the required libraries, using the following commands in the terminal:

```
1. 1
2. 2
3. 3

1. pip3 install virtualenv
2. virtualenv my_env # create a virtual environment my_env
3. source my_env/bin/activate # activate my_env
```

Copied! **Executed!**

Then, install the required libraries in the environment (this will take time 🕒🕒):

```
1. 1
2. 2

1. # installing required libraries in my_env
2. pip install transformers==4.35.2 torch==2.1.1 gradio==4.44.0 langchain==0.0.343 ibm_watson_machine_learning==1.0.335 huggingface-hub==0.19.4
```

Copied! **Executed!**

Have a cup of coffee, it will take a few minutes.

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9

1.      ) (
2.      ( ) )
3.      ) ( (
4.      ) ( (
5.      ) ( (
6.      ) ( (
7.      ) ( (
8.      ) ( (
9.      ) ( (
```

Copied!

We need to install `ffmpeg` to be able to work with audio files in python.

```
1. 1

1. sudo apt update
```

Copied! **Executed!**

Then run:

```
1. 1

1. sudo apt install ffmpeg -y
```

Copied! **Executed!**

Whisper from OpenAI is available in [github](https://github.com/openai/whisper). Whisper's code and model weights are released under the MIT License. See [LICENSE](https://openai.com/whisper/faq) for further details.

Step 1: Speech-to-Text

Initially, we want to create a simple speech-to-text Python file using OpenAI Whisper.

You can test the sample audio file **Sample voice** [link to download](#).

Create and open a Python file and call it `simple_speech2text.py` by clicking the link below:

Open `simple_speech2text.py` in IDE

Let's download the file first (you can do it manually, then drag and drop it into the file environment).

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9
10. 10
11. 11
12. 12
13. 13
14. 14
15. 15
16. 16
17. 17
18. 18
19. 19
20. 20

1. import requests
2.
3. # URL of the audio file to be downloaded
4. url = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMSkillsNetwork-GPXX04C6EN/Testing%20speech%20to%20text.mp3"
5.
6. # Send a GET request to the URL to download the file
7. response = requests.get(url)
8.
9. # Define the local file path where the audio file will be saved
10. audio_file_path = "downloaded_audio.mp3"
11.
12. # Check if the request was successful (status code 200)
13. if response.status_code == 200:
14.     # If successful, write the content to the specified local file path
15.     with open(audio_file_path, "wb") as file:
16.         file.write(response.content)
17.     print("File downloaded successfully")
18. else:
19.     # If the request failed, print an error message
20.     print("Failed to download the file")
```

Copied!

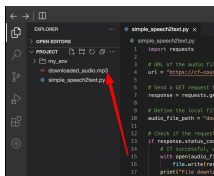
Run the Python file to test it.

```
1. 1
```

```
1. python3 simple_speech2text.py
```

[Copied!](#) [Executed!](#)

You should see the downloaded audio file in the file explorer:



Next, implement OpenAI Whisper for transcribing voice to speech.

You can override the previous code in the Python file.

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9
10. 10
11. 11
12. 12
13. 13
14. 14
15. 15
16. 16
17. 17
18. 18
19. 19
20. 20
21. 21
22. 22

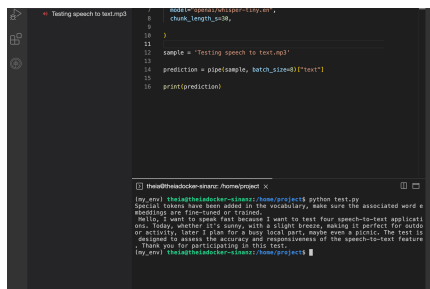
1. import torch
2. from transformers import pipeline
3.
4. # Initialize the speech-to-text pipeline from Hugging Face Transformers
5. # This uses the "openai/whisper-tiny.en" model for automatic speech recognition (ASR)
6. # The 'chunk_length_s' parameter specifies the chunk length in seconds for processing
7. pipe = pipeline(
8.     "automatic-speech-recognition",
9.     model="openai/whisper-tiny.en",
10.     chunk_length_s=30,
11. )
12.
13. # Define the path to the audio file that needs to be transcribed
14. sample = 'downloaded_audio.mp3'
15.
16. # Perform speech recognition on the audio file
17. # The 'batch_size=8' parameter indicates how many chunks are processed at a time
18. # The result is stored in 'prediction' with the key "text" containing the transcribed text
19. prediction = pipe(sample, batch_size=8)["text"]
20.
21. # Print the transcribed text to the console
22. print(prediction)
```

[Copied!](#)

Run the Python file and you will get the output.

```
1. 1
1. python3 simple_speech2text.py
```

[Copied!](#) [Executed!](#)



In the next step, we will utilize Gradio for creating interface for our app.

Gradio interface

Creating a simple demo

Through this project, we will create different LLM applications with Gradio interface. Let's get familiar with Gradio by creating a simple app:

Still in the `project` directory, create a Python file and name it `hello.py`.

Open `hello.py`, paste the following Python code and save the file.

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8

1. import gradio as gr
2.
3. def greet(name):
4.     return "Hello " + name + "!"
5.
6. demo = gr.Interface(fn=greet, inputs="text", outputs="text")
7.
8. demo.launch(server_name="0.0.0.0", server_port= 7868)
```

[Copied!](#)

The above code creates a **gradio.Interface** called `demo`. It wraps the `greet` function with a simple text-to-text user interface that you could interact with.

The **gradio.Interface** class is initialized with 3 required parameters:

- `fn`: the function to wrap a UI around
- `inputs`: which component(s) to use for the input (e.g. "text", "image" or "audio")
- `outputs`: which component(s) to use for the output (e.g. "text", "image" or "label")

The last line `demo.launch()` launches a server to serve our `demo`.

Launching the demo app

Now go back to the terminal and make sure that the `my_env` virtual environment name is displayed at the beginning of the line

Next, run the following command to execute the Python script.

```
1. 1
1. python3 hello.py
```

[Copied!](#) [Executed!](#)

As the Python code is served by a local host, click on the button below and you will be able to see the simple application we just created. Feel free to play around with the input and output of the web app!

Click here to see the application:

[Web application](#)

You should see the following, here we entered the name Bob:

name

bob

Clear

Submit

output

Hello bob!

Flag

If you finish playing with the app and want to exit, press `ctrl+c` in the terminal and close the application tab.

If you wish to learn a little bit more about customization in Gradio, you are invited to take the guided project called **Bring your Machine Learning model to life with Gradio**. You can find it under **Courses & Projects** on [cognitiveclass.ai!](#)

For the rest of this project, we will use Gradio as an interface for LLM apps.

Step 2: Creating audio transcription app

Create a new python file `speech2text_app.py`

Open `speech2text_app.py` in IDE

Exercise: Complete the `transcript_audio` function.

From the step1: fill the missing parts in `transcript_audio` function.

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9
10. 10
11. 11
12. 12
13. 13
14. 14
15. 15
16. 16
17. 17
18. 18
19. 19
20. 20
21. 21
22. 22
23. 23
24. 24
25. 25

1. import torch
2. from transformers import pipeline
3. import gradio as gr
4.
5. # Function to transcribe audio using the OpenAI Whisper model
6. def transcript_audio(audio_file):
7.     # Initialize the speech recognition pipeline
8.     pipe = #----> Fill here <----
9.
10.    # Transcribe the audio file and return the result
11.    result = #----> Fill here <----
12.    return result
13.
14. # Set up Gradio interface
15. audio_input = gr.Audio(sources="upload", type="filepath") # Audio input
16. output_text = gr.Textbox() # Text output
17.
18. # Create the Gradio interface with the function, inputs, and outputs
19. iface = gr.Interface(fn=transcript_audio,
20.                      inputs=audio_input, outputs=output_text,
21.                      title="Audio Transcription App",
22.                      description="Upload the audio file")
23.
24. # Launch the Gradio app
25. iface.launch(server_name="0.0.0.0", server_port=7860)
```

Copied!

▼ Click here for the answer

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9
10. 10
11. 11
12. 12
13. 13
14. 14
15. 15
16. 16
17. 17
18. 18
19. 19
20. 20
21. 21
22. 22
23. 23
24. 24
25. 25
26. 26
27. 27
28. 28

1. import torch
2. from transformers import pipeline
3. import gradio as gr
4.
5. # Function to transcribe audio using the OpenAI Whisper model
6. def transcript_audio(audio_file):
7.     # Initialize the speech recognition pipeline
8.     pipe = pipeline(
9.         "automatic-speech-recognition",
10.        model="openai/whisper-tiny.en",
11.        chunk_length_s=30,
12.    )
13.    # Transcribe the audio file and return the result
14.    result = pipe(audio_file, batch_size=8)["text"]
15.    return result
16.
17. # Set up Gradio interface
18. audio_input = gr.Audio(sources="upload", type="filepath") # Audio input
19. output_text = gr.Textbox() # Text output
20.
21. # Create the Gradio interface with the function, inputs, and outputs
22. iface = gr.Interface(fn=transcript_audio,
23.                      inputs=audio_input, outputs=output_text,
24.                      title="Audio Transcription App",
25.                      description="Upload the audio file")
26.
27. # Launch the Gradio app
28. iface.launch(server_name="0.0.0.0", server_port=7860)
```

Copied!

Then, run your app:

```
1. 1
1. python3 speech2text_app.py
```

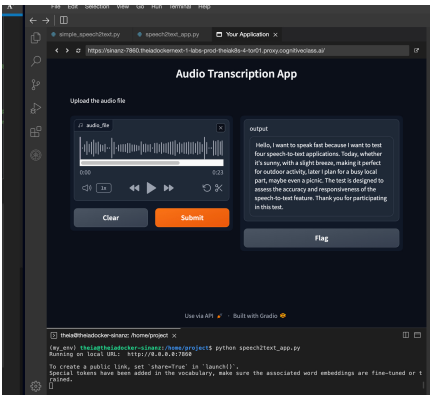
Copied! Executed!

And start the app:

Web application

You can download the sample audio file we've provided by right-clicking on it in the file explorer and selecting "Download." Once downloaded, you can upload this file to the app. Alternatively, feel free to choose and upload any MP3 audio file from your local computer.

The result will be:



Press **Ctrl + C** to stop the application.

Step 3: Integrating LLM: Using Llama 3 in WatsonX as LLM

Running simple LLM

Let's start by generating text with LLMs. Create a Python file and name it `simple_llm.py`. You can proceed by clicking the link below or by referencing the accompanying image.

[Open simple_llm.py in IDE](#)

In case, you want to use Llama 3 as an LLM instance, you can follow the instructions below:

IBM WatsonX utilizes various language models, including Llama 3 by Meta, which is currently the strongest open-source language model.

Here's how the code works:

- Setting up credentials:** The credentials needed to access IBM's services are pre-arranged by the Skills Network team, so you don't have to worry about setting them up yourself.
- Specifying parameters:** The code then defines specific parameters for the language model. 'MAX_NEW_TOKENS' sets the limit on the number of words the model can generate in one go. 'TEMPERATURE' adjusts how creative or predictable the generated text is.
- Setting up Llama 3 model:** Next, the LLAMA3 model is set up using a model ID, the provided credentials, chosen parameters, and a project ID.
- Creating an object for Llama 3:** The code creates an object named `llm`, which is used to interact with the Llama 3 model. A model object, `LLAMA3_model`, is created using the `Model` class, which is initialized with a specific model ID, credentials, parameters, and project ID. Then, an instance of `WatsonxLLM` is created with `LLAMA3_model` as an argument, initializing the language model hub `llm` object.
- Generating and printing response:** Finally, 'llm' is used to generate a response to the question, "How to read a book effectively?" The response is then printed out.

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9
10. 10
11. 11
12. 12
13. 13
14. 14
15. 15
16. 16
17. 17
18. 18
19. 19
20. 20
21. 21
22. 22
23. 23
24. from ibm_watson_machine_learning.foundation_models import Model
25. from ibm_watson_machine_learning.extensions.langchain import WatsonxLLM
26. from ibm_watson_machine_learning.metanames import GenTextParamsMetaNames as GenParams
27.
28. 4.
29. my_credentials = {
30.     "url" : "https://us-south.ml.cloud.ibm.com"
31. }
32.
33. 8.
34. params = {
35.     GenParams.MAX_NEW_TOKENS: 800, # The maximum number of tokens that the model can generate in a single run.
36.     GenParams.TEMPERATURE: 0.1, # A parameter that controls the randomness of the token generation. A lower value makes the generation more deterministic, while a higher value introduces more randomness.
37. }
38.
39. 14. LLAMA2_model = Model(
40.     model_id="meta-llama/llama-3-8b-instruct",
41.     credentials=my_credentials,
42.     params=params,
43.     project_id="skills-network",
44. )
45.
46. 21. llm = WatsonxLLM(LLAMA2_model)
47.
48. 22.
49. 23. print(llm("How to read a book effectively?"))
```

[Copied!](#)

You can then run this script in the terminal using the following command:

```
1. 1
2. python3 simple_llm.py
```

[Copied!](#) [Executed!](#)

Upon running the script, you should see the generated text in your terminal, as shown below:

```
Reading is one of the most efficient ways to gain knowledge and expand your mind. However, not all reading is created equal. Here are some tips to help you read a book effectively:
1. Set goals: Before you start reading, set specific goals for what you want to achieve. Do you want to learn a new skill or gain a deeper understanding of a particular subject? Having clear goals in mind will help you stay focused and motivated.
2. Choose the right book: Not all books are created equal. Choose a book that aligns with your goals and interests. Look for books that have received positive reviews and are written by experts in their fields.
3. Create a reading schedule: Set aside dedicated time each day or week to read consistently. It helps in making progress and retaining information.
4. Take notes: Taking notes while reading can help you retain information and engage with the material on a deeper level. Write down key points, questions, and insights that come to mind as you read.
5. Summarize the material: After finishing a chapter or section, summarize the material in your own words. This will help you understand the material better and
```

You can see how watsonx Llama 2 provides a good answer.

Step 4: Put them all together

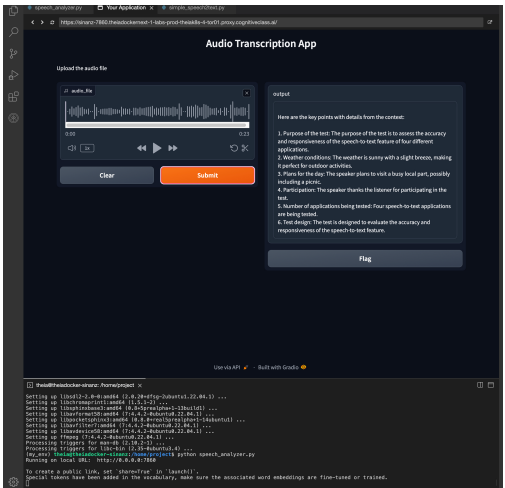
Create a new Python file and call it `speech_analyzer.py`

[Open speech_analyzer.py in IDE](#)

In this exercise, we'll set up a language model (LLM) instance, which could be IBM WatsonxLLM, HuggingFaceHub, or an OpenAI model. Then, we'll establish a prompt template. These templates are structured guides to generate prompts for language models, aiding in output organization (more info in [langchain prompt template](#)).

Next, we'll develop a transcription function that employs the OpenAI Whisper model to convert speech-to-text. This function takes an audio file uploaded through a Gradio app interface (preferably in .mp3 format). The transcribed text is then fed into an LLMChain, which integrates the text with the prompt template and forwards it to the chosen LLM. The final output from the LLM is then displayed in the Gradio app's output textbox.

The output should look:



Notice how the LLM corrected a minor mistake made by the speech-to-text model, resulting in a coherent and accurate output.

Exercise: Fill the missing parts:

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9
10. 10
11. 11
12. 12
13. 13
14. 14
15. 15
16. 16
17. 17
18. 18
19. 19
20. 20
21. 21
22. 22
23. 23
24. 24
25. 25
26. 26
27. 27
28. 28
29. 29
30. 30
31. 31
32. 32
33. 33
34. 34
35. 35
36. 36
37. 37
38. 38
39. 39
40. 40
41. 41
42. 42
43. 43
44. 44
45. 45
46. 46
47. 47
48. 48
49. 49
50. 50
51. 51
52. 52
53. 53
54. 54
55. 55
56. 56
57. 57
58. 58
59. 59
60. 60
61. 61
62. 62
63. 63

1. import torch
2. import os
3. import gradio as gr
4. #from langchain.llms import OpenAI
5. from langchain.llms import HuggingFaceHub
6. from transformers import pipeline
7. from langchain.prompts import PromptTemplate
8. from langchain.chains import LLMChain
9. from ibm_watson_machine_learning.foundation_models.extensions.langchain import WatsonLLM
10. from ibm_watson_machine_learning.foundation_models.utils.enums import DecodingMethods
11. from ibm_watson_machine_learning.metanames import GenTextParamsMetaNames as GenParams
12. from ibm_watson_machine_learning.foundation_models import Model
13.
14. #####----- LLM-----####
15.
16. # initiate LLM instance, this can be IBM WatsonX, huggingface, or OpenAI instance
17.
18. llm = ##### write your code here
19.
20. #####----- Prompt Template-----####
21.
22. # This template is structured based on LLaMA2. If you are using other LLMs, feel free to remove the tags
23. temp = ""
24. <<SYS>>
25. List the key points with details from the context:
26. [INST] The context : {context} [/INST]
27. <<SYS>>
28. ""
29. # here is the simplified version of the prompt template
30. # temp = ""
31. # List the key points with details from the context:
32. # The context : {context}
33. # ""
34.
35. pt = PromptTemplate(
36.     input_variables=["context"],
37.     template=temp)
38.
39. prompt_to_LLaMA2 = LLMChain(llm=llm, prompt=pt)
40.
41. #####----- Speech2Text-----####
42.
43. def transcript_audio(audio_file):
44.     # Initialize the speech recognition pipeline
45.
46.     pipe = ##### write the code here
47.
48.     # Transcribe the audio file and return the result
49.     transcript_txt = pipe(audio_file, batch_size=8)["text"]
50.     # run the chain to merge transcript text with the template and send it to the LLM
51.     result = prompt_to_LLaMA2.run(transcript_txt)
52.
53.     return result
54.
55. #####----- Gradio-----####
56.
57. audio_input = gr.Audio(sources="upload", type="filepath")
58. output_text = gr.Textbox()
59.
60. # Create the Gradio interface with the function, inputs, and outputs
61. iface = ##### write code here
62.
63. iface.launch(server_name="0.0.0.0", server_port=7860)
```

Copied! Click here for the answer

- 1. 1
- 2. 2
- 3. 3
- 4. 4
- 5. 5
- 6. 6
- 7. 7
- 8. 8
- 9. 9
- 10. 10

```
11. 11
12. 12
13. 13
14. 14
15. 15
16. 16
17. 17
18. 18
19. 19
20. 20
21. 21
22. 22
23. 23
24. 24
25. 25
26. 26
27. 27
28. 28
29. 29
30. 30
31. 31
32. 32
33. 33
34. 34
35. 35
36. 36
37. 37
38. 38
39. 39
40. 40
41. 41
42. 42
43. 43
44. 44
45. 45
46. 46
47. 47
48. 48
49. 49
50. 50
51. 51
52. 52
53. 53
54. 54
55. 55
56. 56
57. 57
58. 58
59. 59
60. 60
61. 61
62. 62
63. 63
64. 64
65. 65
66. 66
67. 67
68. 68
69. 69
70. 70
71. 71
72. 72
73. 73

1. import torch
2. import os
3. import gradio as gr
4.
5. #from langchain.llms import OpenAI
6. from langchain.llms import HuggingFaceHub
7.
8. from transformers import pipeline
9. from langchain.prompts import PromptTemplate
10. from langchain.chains import LLMChain
11.
12. from ibm_watson_machine_learning.foundation_models import Model
13. from ibm_watson_machine_learning.extensions.langchain import WatsonxLLM
14. from ibm_watson_machine_learning.metanames import GenTextParamsMetaNames as GenParams
15.
16. my_credentials = {
17.     "url" : "https://us-south.ml.cloud.ibm.com"
18. }
19. params = {
20.     GenParams.MAX_NEW_TOKENS: 800, # The maximum number of tokens that the model can generate in a single run.
21.     GenParams.TEMPERATURE: 0.1, # A parameter that controls the randomness of the token generation. A lower value makes the generation more deterministic, while a higher value introduces more randomness.
22. }
23.
24. LLAMA2_model = Model(
25.     model_id= 'meta-llama/llama-3-8b-instruct',
26.     credentials=my_credentials,
27.     params=params,
28.     project_id="skills-network",
29. )
30.
31. llm = WatsonxLLM(LLAMA2_model)
32.
33. #####----- Prompt Template-----####
34.
35. temp = """
36. <=<<SYS>>
37. List the key points with details from the context:
38. [INST] The context : {context} [/INST]
39. <=</SYS>>
40. """
41.
42. pt = PromptTemplate(
43.     input_variables=["context"],
44.     template= temp)
45.
46. prompt_to_LLAMA2 = LLMChain(llm=llm, prompt=pt)
47.
48. #####----- Speech2text-----####
49.
50. def transcript_audio(audio_file):
51.     # Initialize the speech recognition pipeline
52.     pipe = pipeline(
53.         "automatic-speech-recognition",
54.         model="openai/whisper-tiny.en",
55.         chunk_length_s=30,
56.     )
57.     # Transcribe the audio file and return the result
58.     transcript_txt = pipe(audio_file, batch_size=8)["text"]
59.     result = prompt_to_LLAMA2.run(transcript_txt)
60.
61.     return result
62.
63. #####----- Gradio-----####
64.
65. audio_input = gr.Audio(sources="upload", type="filepath")
66. output_text = gr.Textbox()
67.
68. iface = gr.Interface(fn= transcript_audio,
69.     inputs= audio_input, outputs= output_text,
70.     title= "Audio Transcription App",
71.     description= "Upload the audio file")
72.
73. iface.launch(server_name="0.0.0.0", server_port=7860)
```

Copied!

Run your code:

- 1. 1
- 1. python3 speech_analyzer.py

Copied! Executed!

If there is no error, run the web app:

Web application

Conclusion

Congratulations on completing this project! You have now laid a solid foundation for leveraging powerful Language Models (LLMs) for speech-to-text generation tasks. Here's a quick recap of what you've accomplished:

- Text generation with LLM: You've created a Python script to generate text using a model from the Hugging Face Hub, learned about some key parameters that influence the model's output, and have a basic understanding of how to switch between different LLM models.
- Speech-to-Text conversion: Utilize OpenAI's Whisper technology to convert lecture recordings into text, accurately.
- Content summarization: Implement IBM Watson's AI to effectively summarize the transcribed lectures and extract key points.
- User interface development: Create an intuitive and user-friendly interface using Hugging Face Gradio, ensuring ease of use for students and educators.

Author(s)

Sina Nazeri

© IBM Corporation. All rights reserved.