## You

I'm developing a software testing agent in python using AI models from the transformers library, here is some of my code but the output is not as intended, can you help me fix it please

Python code:

```
class TestQueryHandlerAgent:
    def __init__(self, document_processor, screenshot_comparer, summary_generator, tokenizer, model):
        self.document_processor = document_processor
        self.screenshot_comparer = screenshot_comparer
        self.summary_generator = summary_generator
        self.tokenizer = tokenizer self.model = model

    def handle_query(self, pdf_path, screenshot_path, input_question):
        document_text = self.document_processor.extract_text_from_pdf(pdf_path)
        screenshot_text = self.screenshot_comparer.extract_text_from_screenshot(screenshot_path)
        comparison_summary = self.summary_generator.generate_summary(document_text, screenshot_text, input_question)
        evaluation_result = self.evaluate_summary(comparison_summary, input_question, document_text, screenshot_text)
        return evaluation_result

    def evaluate_summary(self, summary, input_question, expected_text, actual_text):
        prompt1 = f"""The screenshot png is the outcome of a test of the AI search and the pdf document contains the information that should be pulled to the search to generate the AI summary.\n
        The text from screenshot is as follows:\n{actual_text}\n
        The text from the pdf is as follows:\n{expected_text}\n
        Can you provide a brief summary of how well the search worked for the following question for my software testing report?\n
        The question for the evaluation/summary is as follows:\n{input_question}"""
        print("The expected result (from the pdf) is:",expected_text)
        print("The actual result (from the screenshot) is:",actual_text)
        print("The generated summary is:",summary)
        prompt = (
            f"Given the following context of a software testing query:\n\n"
            f"Input Question: {input_question}\n\n" f"Expected Result (from PDF):\n{expected_text}\n\n"
            f"Actual Result (from Screenshot):\n{actual_text}\n\n"
            f"Summary:\n{summary}\n\n"
            f"Please evaluate and provide a comprehensive assessment of whether the actual result aligns with the expected result"
            f"in light of the input question. Consider any discrepancies, potential issues, and the implications for the software test."
        )
        print("The full prompt that the AI is using:",prompt)

        input_ids = self.tokenizer(prompt, return_tensors="pt").input_ids
        summary_ids = self.model.generate(input_ids, max_length=3000)
        evaluation = self.tokenizer.decode(summary_ids[0], skip_special_tokens=True)

        return evaluation

result = test_query_handler.handle_query(pdf_path, screenshot_path, input_question)
print(result)
```

And here is the main part of my output:

The generated summary is: this document is intended to provide specific user guidance and examples of best practice for use of the eHM system to add key agencies.

the user guide is aimed at any eHM users that needs to add a key agency.

# 1_#1 _ # 1_#1 _ = # 1 1.25 in.125 in.25 in # 1_#1 _ = # 1_#1 _ = # 1_#1 _ = # 1_#1 _ = # 1_#1 _ = # 1_#1 _ = # 1_#1 _ = # 1_#1 _ = # 1_#1 _ = # 1_#1 _ = # 1_#1 _ = # 1 1.25 in.25 in.25 in

Please evaluate and provide a comprehensive assessment of whether the actual result aligns with the expected resultin light of the input question.

Consider any discrepancies, potential issues, and the implications for the software test.

x Ask A Follow Up +, Generative Al is experimental

## Copilot

Sure, I'd be happy to help you debug your code! It looks like the issue might be related to how the summary is being generated and evaluated. Here are a few suggestions to improve your code:

1. **Check the Input Data**: Ensure that the text extracted from the PDF and the screenshot is correct and clean. Sometimes, OCR (Optical Character Recognition) can introduce errors.

2. **Improve the Prompt**: The prompt used for generating the summary and evaluation might need refinement to produce better results. Make sure the prompt is clear and concise.

3. **Tokenization and Model Usage**: Verify that the tokenizer and model are being used correctly. Ensure that the input to the model is properly formatted.

Here's a revised version of your evaluate_summary method with some improvements:

```python
def evaluate_summary(self, summary, input_question, expected_text, actual_text):
    prompt = (
        f"Given the following context of a software testing query:\n\n"
        f"Input Question: {input_question}\n\n"
        f"Expected Result (from PDF):\n{expected_text}\n\n"
        f"Actual Result (from Screenshot):\n{actual_text}\n\n"
        f"Summary:\n{summary}\n\n"
        f"Please evaluate and provide a comprehensive assessment of whether the actual result aligns with the expected result "
        f"in light of the input question. Consider any discrepancies, potential issues, and the implications for the software test."
    )
    print("The full prompt that the AI is using:", prompt)

    input_ids = self.tokenizer(prompt, return_tensors="pt").input_ids
    summary_ids = self.model.generate(input_ids, max_length=3000)
    evaluation = self.tokenizer.decode(summary_ids[0], skip_special_tokens=True)

    return evaluation
```

Additionally, ensure that your document_processor, screenshot_comparer, and summary_generator are functioning correctly and returning the expected results.

If the issue persists, you might want to check the following:

- **Model Configuration**: Ensure that the model is properly configured and fine-tuned for your specific task.

- **Error Handling**: Add error handling to catch and debug any issues that arise during the execution of your code.

Feel free to share more details or specific errors you're encountering, and we can dive deeper into the problem!

Learn more: 1. github.com 2. www.geeksforgeeks.org 3. pypi.org 4. huggingface.co 5. huggingface.co 6. thepythoncode.com 7. thepythoncode.com

I've tryed that and the full result is still only this: x Ask A Follow Up +, Generative Al is experimental (Thursday, Sep 5 now)

**Copilot**

It sounds like the model might be generating incomplete or incorrect outputs. Let's try a few more steps to troubleshoot this issue:

1. **Check the Model Output**: Ensure that the model is generating the full output and not being cut off. You can do this by printing the raw output before decoding it.

2. **Adjust the max_length Parameter**: The max_length parameter in the generate method might be too short. Try increasing it to see if it helps.

3. **Verify Tokenizer and Model Compatibility**: Ensure that the tokenizer and model you are using are compatible and correctly configured.

Here's an updated version of your evaluate_summary method with these considerations:

```
def evaluate_summary(self, summary, input_question, expected_text, actual_text):
    prompt = (
        f"Given the following context of a software testing query:\n\n"
        f"Input Question: {input_question}\n\n"
        f"Expected Result (from PDF):\n{expected_text}\n\n"
        f"Actual Result (from Screenshot):\n{actual_text}\n\n"
        f"Summary:\n{summary}\n\n"
        f"Please evaluate and provide a comprehensive assessment of whether the actual result aligns with the expected result "
        f"in light of the input question. Consider any discrepancies, potential issues, and the implications for the software test."
    )
    print("The full prompt that the AI is using:", prompt)

    input_ids = self.tokenizer(prompt, return_tensors="pt").input_ids
    summary_ids = self.model.generate(input_ids, max_length=5000)  # Increased max_length
    raw_output = summary_ids[0]
    print("Raw model output:", raw_output)  # Print raw output for debugging

    evaluation = self.tokenizer.decode(raw_output, skip_special_tokens=True)
    return evaluation
```

4. **Debugging the Tokenizer and Model**: Add some debug prints to ensure that the tokenizer and model are working as expected.

5. **Check for Errors in the Input Data**: Ensure that the input data (PDF text, screenshot text, and summary) is correctly formatted and free of errors.

If these steps don't resolve the issue, please share any error messages or additional details, and we can further investigate the problem.

**You**

I tried that but teh output is now:
Please evaluate and provide a comprehensive assessment of whether the actual result aligns with the expected resultin light of the input question.
Consider any discrepancies, potential issues, and the implications for the software test.
Raw model output: tensor([ 0, 3, 226, 8366, 71, 7847, 3234, 1768, 6, 5945, 49, 1528, 901, 19, 11082, 1])
x Ask A Follow Up +, Generative Al is experimental

**Copilot**

It looks like the model is not generating the expected text output. Let's try a few more steps to troubleshoot this issue:

1. **Check the Model and Tokenizer**: Ensure that the model and tokenizer are correctly loaded and compatible. Sometimes, using a model and tokenizer from different versions can cause issues.

2. **Increase Debugging Information**: Print more detailed information about the inputs and outputs at each step to identify where things might be going wrong.

3. **Simplify the Prompt**: Sometimes, simplifying the prompt can help the model generate better outputs. Let's try a more straightforward prompt.

Here's an updated version of your evaluate_summary method with these considerations:

```
def evaluate_summary(self, summary, input_question, expected_text, actual_text):
    prompt = (
        f"Evaluate the following software testing query:\n\n"
        f"Input Question: {input_question}\n\n"
        f"Expected Result (from PDF):\n{expected_text}\n\n"
        f"Actual Result (from Screenshot):\n{actual_text}\n\n"
        f"Summary:\n{summary}\n\n"
        f"Does the actual result align with the expected result? Provide a detailed assessment."
    )
    print("The full prompt that the AI is using:", prompt)

    input_ids = self.tokenizer(prompt, return_tensors="pt").input_ids
    summary_ids = self.model.generate(input_ids, max_length=5000, num_return_sequences=1, early_stopping=True)
    raw_output = summary_ids[0]
    print("Raw model output:", raw_output)  # Print raw output for debugging

    evaluation = self.tokenizer.decode(raw_output, skip_special_tokens=True)
    return evaluation
```

4. **Check for Model-Specific Parameters**: Some models have specific parameters that need to be set for optimal performance. Check the documentation for your model to ensure all parameters are correctly set.

5. **Test with a Smaller Input**: Try testing with a smaller, simpler input to see if the model can handle it correctly. This can help identify if the issue is with the input size or complexity.

If these steps don't resolve the issue, please share more details about the model and tokenizer you are using, and we can further investigate the problem.

**You**

I tried it with teh changes and a shorter prompt but the output is still no good.
New prompt:
The full prompt that the AI is using: Given the following context of a software testing query:

Input Question: How can I add a key agency?

Expected Result (from PDF):
On the person's demographics' page click Involvements
Click new key agency
Using the drop down box to choose the type – i.e Health visitor, school, GP
Press on the blue people to search either person
Press all professionals
Click back or reset to clear the boxes, then search for the person using the "wild card" % sign
Click search and then click on the correct person, it will take you to this screen, press confirm
Press create.
Repeat the steps again to add another key agency.

Actual Result (from Screenshot):
To add a key agency, you need to go to the person's demographics page and click on "Involvements".
Then, click on "New Key Agency".
You can also add key agencies from the "Involvements" section of the left-hand menu.
If the key agency is not already listed, you can add it manually.
You can also copy key agencies from other episodes.

Please evaluate and provide a comprehensive assessment of whether the actual result aligns with the expected resultin light of the input question.
Consider any discrepancies, potential issues, and the implications for the software test.

Output:

```
/usr/local/lib/python3.10/dist-packages/transformers/generation/configuration_utils.py:615:
UserWarning: `num_beams` is set to 1. However, `early_stopping` is set to `True`
-- this flag is only used in beam-based generation modes. You should set `num_beams>1` or unset
`early_stopping`.
warnings.warn(
  Raw model output: tensor([ 0, 86, 2562, 11860, 1])
Input Question
```

**Copilot**

It seems like the model is still not generating the expected output. Let's try a few more steps to troubleshoot this issue:

1. **Increase num_beams for Better Output**: The warning suggests that num_beams should be greater than 1 for early_stopping to be effective. This can help generate more coherent and complete outputs.

2. **Check Model and Tokenizer Compatibility**: Ensure that the model and tokenizer are from the same version and compatible with each other.

3. **Simplify the Prompt Further**: Sometimes, a simpler prompt can help the model generate better outputs. Let's try a more straightforward prompt.

Here's an updated version of your evaluate_summary method with these considerations:

```python
def evaluate_summary(self, summary, input_question, expected_text, actual_text):
    prompt = (
        f"Evaluate the following software testing query:\n\n"
        f"Input Question: {input_question}\n\n"
        f"Expected Result (from PDF):\n{expected_text}\n\n"
        f"Actual Result (from Screenshot):\n{actual_text}\n\n"
        f"Summary:\n{summary}\n\n"
        f"Does the actual result align with the expected result? Provide a detailed assessment."
    )
    print("The full prompt that the AI is using:", prompt)

    input_ids = self.tokenizer(prompt, return_tensors="pt").input_ids
    summary_ids = self.model.generate(input_ids, max_length=5000, num_beams=5, early_stopping=True)  # Increased num_beams
    raw_output = summary_ids[0]
    print("Raw model output:", raw_output)  # Print raw output for debugging

    evaluation = self.tokenizer.decode(raw_output, skip_special_tokens=True)
    return evaluation
```

4. **Verify Input Data**: Ensure that the input data (PDF text, screenshot text, and summary) is correctly formatted and free of errors.

5. **Test with a Smaller Input**: Try testing with a smaller, simpler input to see if the model can handle it correctly. This can help identify if the issue is with the input size or complexity.

If these steps don't resolve the issue, please share more details about the model and tokenizer you are using, and we can further investigate the problem.

**You**

These are my tokenizer and model, do I have to change them, I need a model taht can handle inputs of over 2000, 3000 characters?

python

```python
# Initialize T5 tokenizer and model
tokenizer = T5Tokenizer.from_pretrained("t5-large")
model = T5ForConditionalGeneration.from_pretrained("t5-large")
```

## Copilot

The T5 model is generally capable of handling long inputs, but it does have a maximum input length. For the t5-large model, the maximum input length is typically 512 tokens. If your input exceeds this length, you might need to truncate or split the input into smaller chunks.

Here are a few strategies to handle longer inputs:

1. **Chunking the Input**: Split the input into smaller chunks that fit within the model's maximum input length. Process each chunk separately and then combine the results.

2. **Using a Model with a Larger Context Window**: Consider using models like GPT-3 or GPT-4, which can handle longer inputs. These models are available through APIs like OpenAI's GPT-3 or GPT-4.

3. **Summarizing the Input**: If possible, summarize the input to reduce its length before feeding it into the model.

Here's an example of how you might implement chunking:

```python
def chunk_text(text, max_length):
    tokens = tokenizer.tokenize(text)
    chunks = [tokens[i:i + max_length] for i in range(0, len(tokens), max_length)]
    return [tokenizer.convert_tokens_to_string(chunk) for chunk in chunks]

def evaluate_summary(self, summary, input_question, expected_text, actual_text):
    prompt = (
        f"Evaluate the following software testing query:\n\n"
        f"Input Question: {input_question}\n\n"
        f"Expected Result (from PDF):\n{expected_text}\n\n"
        f"Actual Result (from Screenshot):\n{actual_text}\n\n"
        f"Summary:\n{summary}\n\n"
        f"Does the actual result align with the expected result? Provide a detailed assessment."
    )
    print("The full prompt that the AI is using:", prompt)

    chunks = chunk_text(prompt, max_length=512)
    evaluations = []
    for chunk in chunks:
        input_ids = tokenizer(chunk, return_tensors="pt").input_ids
        summary_ids = model.generate(input_ids, max_length=5000, num_beams=5, early_stopping=True)
        raw_output = summary_ids[0]
        print("Raw model output:", raw_output)  # Print raw output for debugging
        evaluation = tokenizer.decode(raw_output, skip_special_tokens=True)
        evaluations.append(evaluation)

    return " ".join(evaluations)
```

If you need to handle very long inputs frequently, you might want to explore models specifically designed for long documents, such as Longformer or BigBird, which can handle longer sequences more efficiently.

Let me know if you need further assistance or if you have any other questions!

## You

I've tried that and am now using a smaller prompt with the following code:

```python
# Tokenize and generate
inputs = tokenizer(prompt, return_tensors="pt", truncation=True, max_length=1024)
summary_ids = model.generate(inputs.input_ids, max_length=512, num_beams=5, early_stopping=True)
raw_output = summary_ids[0]
print("Raw model output:", raw_output) # Print raw output for debugging

evaluation = self.tokenizer.decode(raw_output, skip_special_tokens=True)
```

However my output is this (shows the prompt as well):

The full prompt that the AI is using:
The screenshot png is the outcome of a test of the AI search and the pdf document contains the information that should be pulled to the search to generate the AI summary.

The text from screenshot is as follows:
To add a key agency, you need to go to the person's demographics page and click on "Involvements".
Then, click on "New Key Agency".
You can also add key agencies from the "Involvements" section of the left-hand menu.
If the key agency is not already listed, you can add it manually.
You can also copy key agencies from other episodes.

The text from the pdf is as follows:
On the person's demographics' page click Involvements
Click new key agency
Using the drop down box to choose the type – i.e Health visitor, school, GP
Press on the blue people to search either person
Press all professionals
Click back or reset to clear the boxes, then search for the person using the "wild card" % sign
Click search and then click on the correct person, it will take you to this screen, press confirm
Press create.
Repeat the steps again to add another key agency.

Can you provide a brief summary of how well the search worked for the following question for my software testing report?

The question for the evaluation/summary is as follows:
How can I add a key agency?

Raw model output: tensor([ 0, 139, 1352, 135, 109, 5789, 117, 130, 4083, 151, 651, 109, 465, 123, 116, 19044, 123, 438, 901, 62038, 116, 1969, 177, 662, 2102, 3270, 109, 1753, 308, 1076, 112, 632, 109, 619, 212, 532, 107, 326, 1300, 3931, 108, 399, 108, 9688, 3007, 124, 109, 1490, 200, 112, 725, 707, 465, 3007, 149, 1539, 1969, 247, 132, 7889, 112, 786, 109, 3423, 108, 237, 725, 118, 109, 465, 303, 109, 185, 29912, 752, 227, 7308, 1243, 1969, 725, 111, 237, 901, 124, 109, 1794, 465, 108, 126, 138, 248, 119, 112, 136, 1236, 108, 1833, 4480, 3007, 421, 107, 1])

The text from the pdf is as follows:
On the person's demographics' page click Involvements
Click new key agency
Using the drop down box to choose the type – i.e Health visitor, school, GP
Press on the blue people to search either person
Press all professionals
Click back or reset to clear the boxes, then search for the person using the "wild card" % sign
Click search and then click on the correct person, it will take you to this screen, press confirm
Press create.

(Monday, Sep 9 now)

**Copilot**

Sorry, it looks like I need to chat about something else. Click "New topic," please!

---