# [Read Me] Deriving Relevant Sentences for Patient Summaries in Medical Transcriptions Using BERT

## Introduction

In many different applications, it can be useful to determine whether a particular statement is accurately reflected in a larger corresponding document. Throughout this blog post, we will be considering the benefits of this type of application within the medical field. Often times, doctors are tasked with writing a quick summary statement that will summarize their patient's entire visit. Later on, if that patient revisits the doctor or transfers to a new doctor, then those summary statements will be reviewed by their new doctor in order to get a better understanding of the patient's medical history. That being said, the accuracy of these summary statements can have direct impact on the patient's treatment in the future. This application will provide a way to check that the summary statement is an accurate representation of the information from the corresponding document (aka the full details from the patient visit).

Example
If the summary statement is "The patient is healthy" and the supporting documentation is "The patient's health is good. The patient does not have a fever. The patient is a 28 year old female." the most accurate summary statement would be "The patient's health is good."

Note: From this point forward, we will be referring to the summary statement as a "restatement" to avoid any confusion with the corresponding document that summarizes the patient visit. That being said, the restatement may be a summary statement of the document, a high level description, an exact quote from the document, or bear some other relationship to the document.
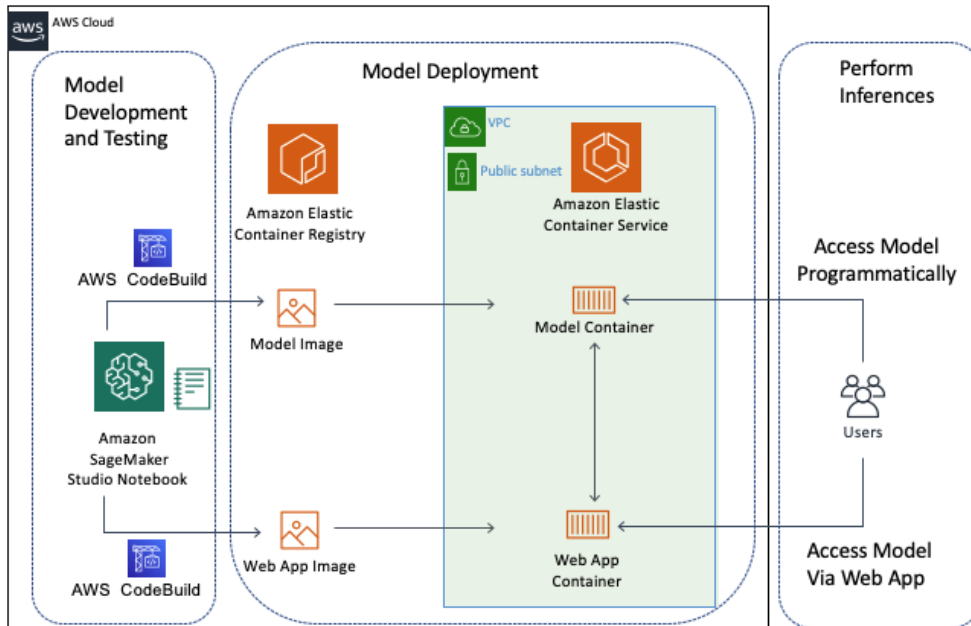
## Approach

As briefly mentioned above, we will deploy a solution that takes in a restatement and a corresponding document. The corresponding document will first be broken up into individual sentences, and then Bio_ClinicalBERT will be used to embed each of those sentences as well as the restatement sentence itself. The script will then compute the distances (using cosine similarity) between the restatement and each of the individual sentences from the corresponding document. The top 5 sentences (i.e. the ones with the lowest distances) will be returned along with their respective scores.

## Machine Learning Method

This core machine learning method used relies on BERT, a language/word embedding model published by Google (the original paper can be found here). BERT is a context aware embedder that can be used for a number of Natural Language Processing tasks (a brief overview of BERT can be found here). Please note that for this specific use case, we are using a pre-trained BERT model, called Bio_ClinicalBERT, that was developed specifically for medical/clinical tests. This pre-trained model was developed by `Alsentzer et al`, and the full paper published by that group can be found here.

## Architecture Diagram

As you can see, there are three sections.

1. In the first section `Model Development and Testing` we will use the SageMaker Studio Integrated Development Environment (IDE) for building the application. SageMaker Studio is an IDE developed specifically for machine learning.

2. In the second section `Model Deployment` we will deploy both the model and a web app (powered by Streamlit for querying the model. Both the model and the web app will be packaged as Docker containers using AWS CodeBuild .They will be first registered with Amazon Elastic Container Registry (ECR) and then deployed on Amazon Elastic Container Service(ECS) using AWS Fargate. For a more in-depth tutorial on deploying Streamlit apps to ECS; checkout out this blog post.

3. In the third section `Perform Inferences` we will perform some inferences using our model and web app.

## Contents

The content of the workshop is as follows:

```
docker_containers
├── cloudformation_template.yaml
├── model_container
│   ├── Dockerfile
│   ├── build_and_push_docker.sh
│   ├── data
│   │   └── mtsamples.csv
│   ├── install_packages.sh
│   ├── local_tests
│   ├── model_scripts
│   │   ├── download_embedding_data.py
│   │   ├── hosted_model.py
│   │   ├── nginx.conf
│   │   ├── null_hypothesis.txt
│   │   ├── predictor.py
│   │   ├── sample_local_testing.txt
│   │   ├── serve
│   │   └── wsgi.py
│   └── nlp_bert_medical_workshop.ipynb
└── web_app_container
    ├── Dockerfile
    ├── build_and_push_docker.sh
    └── web_app.py
```

## Workshop Requirements

You will need an AWS account to use this solution. Sign up for an account here.

You will also need to have permission to use AWS CloudFormation and to create all the resources detailed in the architecture section. All AWS permissions can be managed through AWS IAM. Admin users will have the required permissions, but please contact your account's AWS administrator if your user account doesn't have the required permissions.
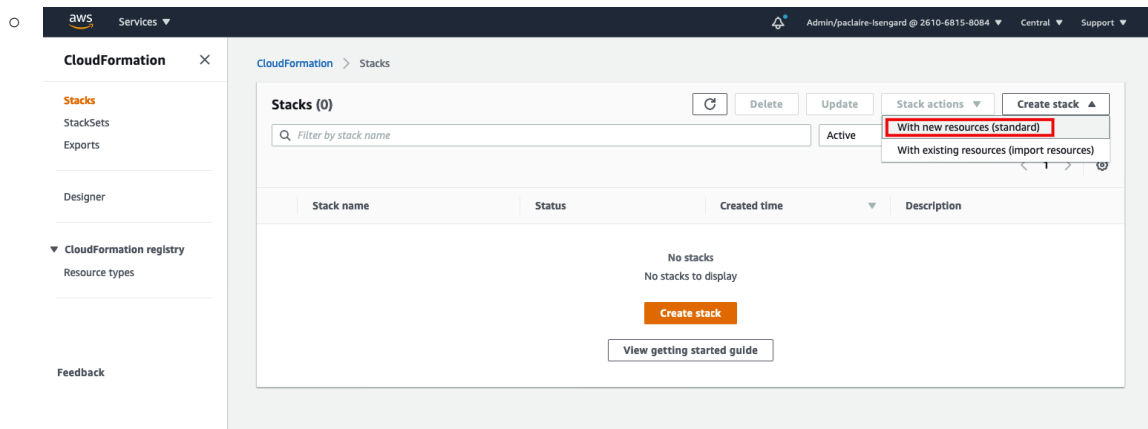
## Steps for deploying CloudFormation Template

To run this workshop, you will need to first deploy the CloudFormation Template

- Open AWS Console
- Select "CloudFormation" from the list of services
- Select "Stacks" from the sidebar

- Select "Create stack" and choose "With new resources (standard)"



- Step 1: Specify template
  - Under the Prepare template section select "Template is ready"
  - Then under the Specify template section select "Upload a template file"
  - Click the "Choose file" button and upload your cloudformation template
  - Then select "Next"

- Step 2: Specify Stack Details
  - Enter Medical-Text-Analysis as the "Stack name" and as the "EnvironmentName"
  - Click "Next"



- Step 3: Configure Stack Options
  - Optional: Add any tags that you want to be applied to the resources in your stack
  - Leave everything else as default and select "Next"
- Step 4: Review
  - Review over all of the information and at the bottom of the page select the checkbox to acknowledge that this cloudformation template will be creating an IAM resource.
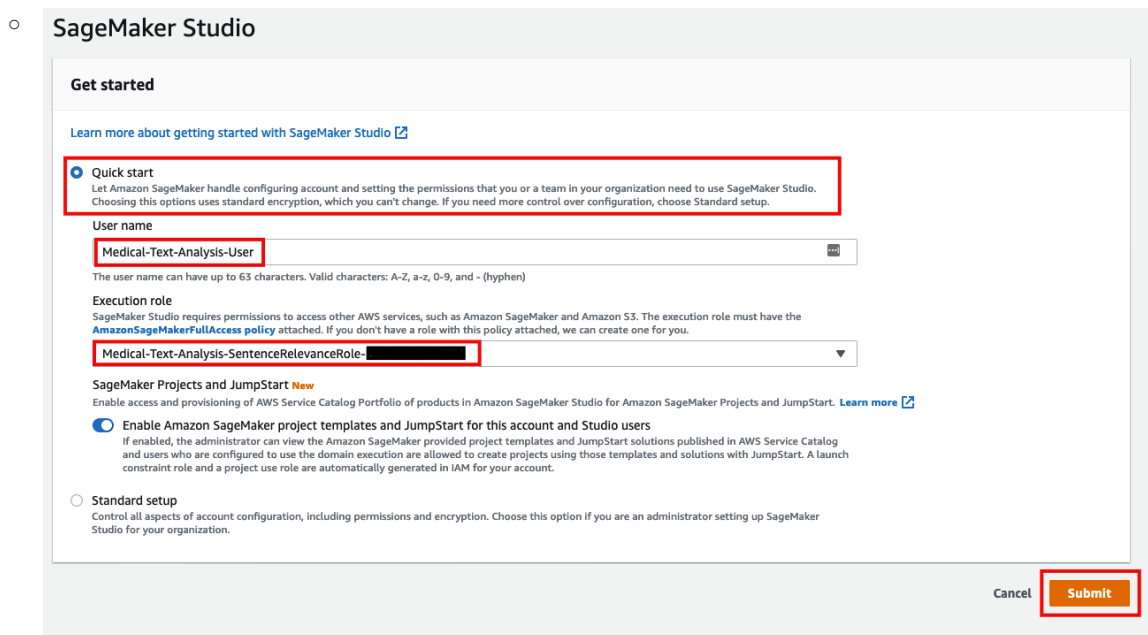
# Steps for opening SageMaker Studio

  - Then select "Create Stack"

- Once the cloudformation template has finished deploying all of the resources, open SageMaker from the main AWS console
- We will be using SageMaker Studio for this workshop. Here's the list of regions that SageMaker Studio is available in.
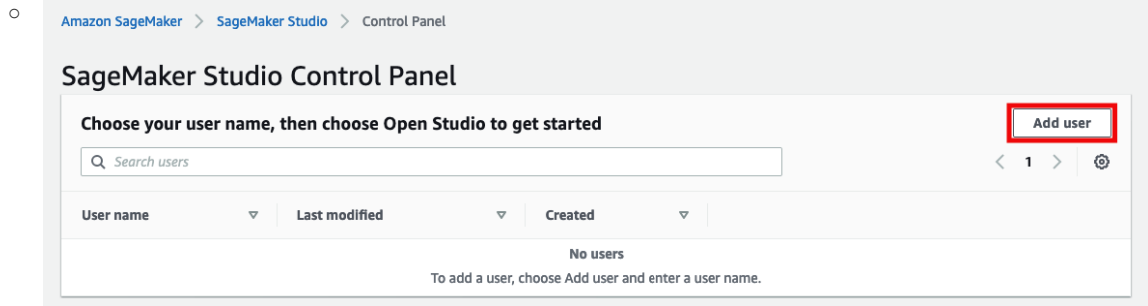
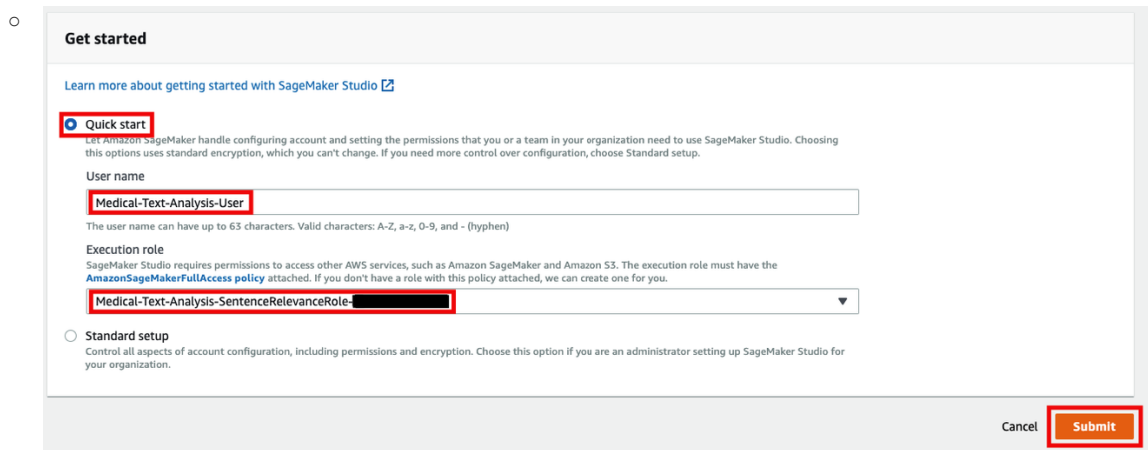- Go to SageMaker Studio by selecting the tab on the side bar

- If this is your first time using the SageMaker Studio console, you will need to first set it up.
  - Choose the "Quick start" option
  - Set the "User name" to "Medical-Text-Analysis-User" (or any other user name will do!)
  - For the Execution role, select the IAM role that was created when the CloudFormation template was deployed. This will look something like "Medical-Text-Analysis-Sentence-RelevanceRole-XXXXXXXXXXXXX"
  - Select Submit



- If this is not your first time using SageMaker Studio, go ahead and just create a new user
  - Select "Add User"

- Select "Quick Start"
- Set the "User name" to "Medical-Text-Analysis-User" (or any other user name will do!)
- For the Execution role, select the IAM role that was created when the CloudFormation template was deployed. This will look something like "Medical-Text-Analysis-Sentence-RelevanceRole-XXXXXXXXXXXX"
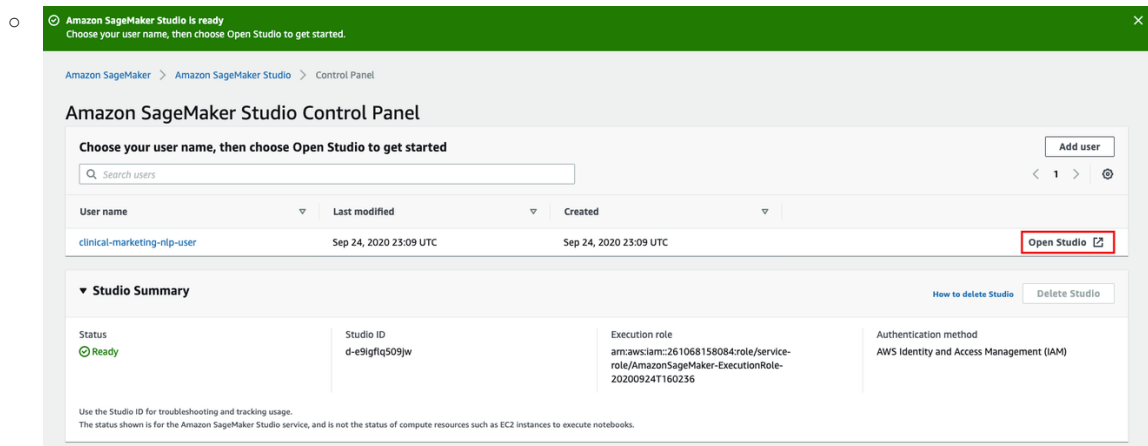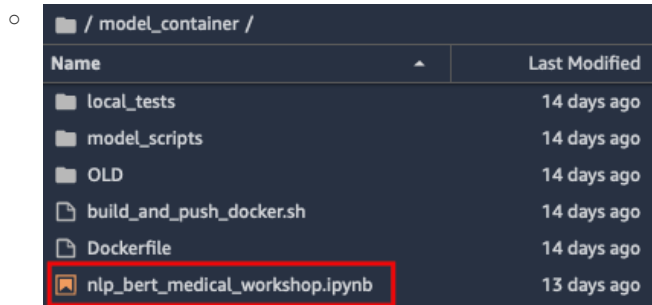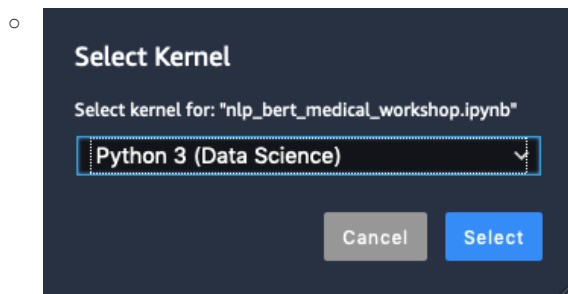- Select Submit



- You will then be prompted to choose a VPC and Subnet you would like to use for SageMaker Studio. Select `sentence_relevance_VPC`. There is only one subnet in that VPC; select it.
- Once you choose your VPC and Subnet, you will then be brought to the Amazon SageMaker Studio Control Panel where you will see a bar at the top explaining that it is "Preparing Amazon SageMaker Studio". This will take a few minutes to create.
- Once SageMaker studio is ready (this will take about 5 minutes), find your Username and select "Open Studio"

- Once the studio opens, clone this repo into the SageMaker Studio console with `git clone <URL> (note: you can also download this repository manually and upload it to SageMaker Studio).
- Click on the directory `Medical_Text_Analysis_Resources`, then `docker_containers` and finally the folder `model_container"`and double click on "nlp_bert_medical_workshop.ipynb" to open the notebook
    - 
- Choose Python 3 (Data Science) as your preferred kernel
    - NOTE: This might take a few minutes to start up
    - 
- In the top right of the notebook click on "Unknown" and set to ml.m5.large

- 



The rest of th... notebook `nlp_bert_medical_workshop.ipynb.`

# Cleaning Up

When you've finished with this solution, make sure that you delete all unwanted AWS resources. AWS CloudFormation can be used to automatically delete all standard resources that have been created by the solution and notebook. Go to the AWS CloudFormation Console, and delete the *Medical_Text_Analysis* stack.

Caution: You need to manually delete the extra resources that you may have created. Examples include extra Amazon S3 buckets (to the solution's default bucket), and extra Amazon SageMaker Console, you can delete the SageMaker Studio domains you created. From the ECS Clusters, you can manually de-deploy the deployed containers.

# Useful Resources

- CloudFormation User Guide
- Amazon SageMaker Developer Guide
- Amazon SageMaker Studio Guide

# Credits

- Bio_ClinicalBERT
- MTSamples