

# Finetuning Language Models - Can I Patent This?



If you miss a milestone deadline you will be forfeited the corresponding points. Deadlines are quoted here as a week number - check your LMS about actual dates.

## Introduction

The USPTO is the US Patent and Trademark Office. It is the agency that grants patents to inventors and businesses for their inventions.

You are hired by USPTO to work on introducing AI into the organization. You are told that USPTO is facing patent examiner staff and budget shortages and to help the organization you are tasked to create an app that will accept an input patent application and will return its patentability score. The app will be made available to the public and will be used by patent applicants to determine the patentability of their inventions before they file their patent applications, therefore reducing the workload of the patent examiners.

## Milestones

### Milestone 1: (Week 1, 20 points)

In this milestone you will learn the basics of docker and create a development environment. All AI and data science projects are developed in containers.

Learn the basics of docker (the most common container format) by watching the following video:

```
{eval-rst} .. youtube:: pTFZFxd4h0I
```

If you are on Windows you will need to follow [these instructions](#) and install [Docker Desktop](#) and [WSL2](#).

Independent of your OS, you may want to use [VS Code IDE](#) if you have no IDE experience before. Ensure that you are able to debug code in your IDE. It must connect to the [remote container](#).

Submit the github repository URL with a branch titled 'milestone-1' with the README.md file containing the installation instructions you followed and a screenshot of your docker container terminal prompt. Add as collaborator the TA.

### **Milestone 2: Sentiment Analysis App (Week 3, 20 points)**

Merge the earlier branch into the main branch and create a new branch titled 'milestone-2'. Do not delete the milestone-1 branch.

The purpose of this task is to take you through the process of creating a streamlit app. Streamlit is a python library that allows you to create web apps with minimal coding and deploy them in the cloud. This is a necessary step before you can develop your app for the more complex USPTO use case.

After watching this video, you will be able to create a sentiment analysis app using Streamlit and various pretrained models. The pretrained models are available in the [HuggingFace model hub](#).

```
{eval-rst} .. youtube:: 8h0zsFETm4I
```

Develop a streamlit application that allows the user to enter a text, select a pretrained model and get the sentiment analysis of the text. Use the [HuggingFace transformers](#) library just like in the video to do so.

Deploy the streamlite app in [HuggingFace streamlit spaces](#) after you create a free account.

Submit the github repository URL with a branch titled 'milestone-3' with the README.md file containing the link to the deployed HF space where the app must be clearly prepopulated with a sample text and the TA will only have to press the submit button for the app to display the sentiment result.

### **Milestone 3: Finetuning Language Models (Week 5, 40 points)**

Merge the earlier branch into the main branch and create a new branch titled 'milestone-3'. Do not delete the milestone-2 branch.

The following video will show you how to finetune a language model using the [HuggingFace transformers](#) library. Consult the [HuggingFace documentation](#) for more details.

```
{eval-rst} .. youtube:: GSt00_-0ncQ
```

You consult [this publication](#) and [the Harvard USPTO patent dataset](#) to develop the required classifier. You only need to do the finetuning using a small subset of the dataset that corresponds to all patent applications submitted in Jan 2016. You are free to select and test various sections of the patent application but abstract and claims, intuitively, are the most relevant. You can use the GPU in Google Colab to do the training.

Submit the github repository URL with a branch titled ‘milestone-4’ with the README.md file containing the link to the deployed HF space where the app must be clearly prepopulated with a drop down menu to select the application filing number (or any unique identifier that will allow you to retrieve the patent sections you selected), the app must then show in minimally two text boxes the eg abstract and eg claims and the TA will only have to press the submit button for the app to display the patentability score.

#### **Milestone 4: Documentation and Video Production (Week 6, 20 points)**

Merge the earlier branch into the main branch and create a new branch titled ‘milestone-4’. Do not delete the milestone-3 branch.

Document extensively both the code as well as the results (10 points). Use google sites to create a landing page for your USPTO app (5 points), create a video that will demonstrate the app. The video should be no longer than 5 minutes and should be either uploaded to your youtube channel or included in the github repository as an mp4 (5 points).