# Configuration Manual

MSc Research Project
MSc in Cloud Computing

## Sravanthi Challa
Student ID: 21156239

School of Computing
National College of Ireland

Supervisor: Shaguna Gupta

## National College of Ireland

## MSc Project Submission Sheet

## School of Computing

| | |
|---|---|
| **Student Name:** | Sravanthi Challa |
| **Student ID:** | 21156239 |
| **Programme:** | MSc in Cloud Computing          **Year:**  2023 |
| **Module:** | MSc Research Project |
| **Lecturer:** | Shaguna Gupta |
| **Submission Due Date:** | 14 December 2023 |
| **Project Title:** | Machine Learning & PBFT Blockchain Methodology on AWS for Proteomics Analytics |
| **Word Count:** | …………… **Page Count:** …………… |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project.  All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section.  Students are required to use the Referencing Standard specified in the report template.  To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

**Signature:** ………………………………………………………………………………………………………………

**Date:** ………………………………………………………………………………………………………………

### PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies) | ☐ |
| **Attach a Moodle submission receipt of the online project submission,** to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid.  It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

Sravanthi Challa
Student ID: 21156239

# 1  Introduction

Proteomics has transformed the study of proteins and provided an immense amount of new information about biological systems and disease states. However, proteomics data analysis is a difficult and computationally demanding task that frequently calls for advanced machine learning techniques. This configuration manual presents a thorough method of proteomics data analytics on Amazon Web Services (AWS) that makes use of PBFT blockchain technology and machine learning to overcome these obstacles.

The objective of this project is to improve protein identification and quantification accuracy by utilising machine learning, and to guarantee data confidentiality and integrity by integrating blockchain technology. For the purpose of deploying and managing the blockchain application and machine learning model, AWS offers a stable and scalable platform. The integration of these technologies has the potential to revolutionise proteomics research by facilitating more dependable and effective data analysis.

## 1.1  Target Audience

The purpose of this configuration manual is to assist researchers and developers who want to implement an AWS proteomics data analytics solution based on machine learning. From data collection and pre-processing to model training, deployment, and blockchain integration, it offers detailed guidance for every step of the procedure.

## 1.2  Prerequisites

Before proceeding with this configuration manual, ensure that you have the following prerequisites:
- Basic understanding of machine learning, Blockchain and proteomics
- Familiarity with AWS services, such as ECS, Docker, and Kubernetes
- Experience with Python

# 2  Data Collection and Pre-Processing

## 2.1  Data Collection

Proteomics data collection requires collecting information from a variety of experimental methods, including chromatography, NMR (Nuclear Magnetic Resonance), and mass spectrometry. These methods produce a variety of data formats, such as raw output files, chromatograms, and spectra. The data collection process is meticulous and crucial because it often involves large volumes, a variety of biological samples, and multi-dimensional features.

Data will be collected and stored in AWS S3 bucket.

1. First create AWS S3 bucket in AWS.



2. Then check wheather you have access for SageMaker in AWS

3. Go to SageMaker and create the machine learning algorithm file and upload PDB file. Here, the SageMaker will train the data for different models.



4. We will be declaring the best model prediction in SageMaker in our algorithm.

5. Write down the model training code.



6. Data pre-processing declaration is done, Prediction of protiens will be done by X and Y-axis data

7. Based on the macromolecule, Sequence and residuecount the prediction is done. The column can be seen in the data file.



8. Sagemaker is connected to our streamlit application, and cluster is been created.

9. Once we run the model training it will create a pickle file (PKL) of the that particular model in S3 bucket



10. The trained file, gives us precision, recall, f1score and support

11. Create Github repository and upload all your UI code in it.



12. For creating cluster, Docker and Kubernetes setup has to be done.

```
 Deploying Docker container to AWS and managing it with Amazon Elastic
Kubernetes Service (EKS) involves several steps. Here's a step-by-step guide
to help you through the process:

### Step 1: Prepare Your Docker Image

1. **Build your Docker image** (if not already done):
   ```sh
   docker build -t proteomics-app .
   ```

2. **Tag your Docker image** for Amazon Elastic Container Registry (ECR):
   ```sh
   docker tag proteomics:latest ASIATUYJP7SUEWGXB6WN.dkr.ecr.us-east-
1.amazonaws.com/proteomics:latest
   ```
   Replace `<aws_account_id>` with your AWS account ID and `<region>` with
your AWS region.

### Step 2: Push Your Image to Amazon ECR

1. **Authenticate Docker to your default ECR registry**:
   ```sh
   aws ecr get-login-password --region eu-east-1 | docker login --username AWS
--password-stdin ASIATUYJP7SUEWGXB6WN.dkr.ecr.eu-east-1.amazonaws.com
   ```
```

2. **Create an ECR repository** (if you haven't already):
   ```sh
   aws ecr create-repository --repository-name proteomicsRepo --region us-east-1
   ```

3. **Push your Docker image to ECR**:
   ```sh
   docker push ASIATUYJP7SUEWGXB6WN.dkr.ecr.eu-east-1.amazonaws.com/proteomics-app:latest
   ```

### Step 3: Set Up Amazon EKS

1. **Create an EKS cluster**. This can be done via the AWS Management Console or using AWS CLI. The CLI command is:
   ```sh
   eksctl create cluster --name proteomics-cluster --version 1.28 --region us-east-1 --nodegroup-name standard-workers --node-type t3.medium --nodes 1 --nodes-min 1 --nodes-max 1 --managed
   ```

2. **Configure `kubectl` to communicate with your cluster**:
   ```sh
   aws eks --region us-east-1 update-kubeconfig --name proteomics-cluster
   ```

### Step 4: Deploy Your Application on EKS

1. **Create a Kubernetes deployment**. You need a deployment YAML file (e.g., `proteomics-deployment.yaml`) that references your ECR image. Here's an example of what this file might look like:

   ```yaml
   apiVersion: apps/v1
   kind: Deployment
   metadata:
     name: proteomics
   spec:
     replicas: 2
     selector:
       matchLabels:
         app: proteomics
     template:
       metadata:
         labels:
           app: proteomics
         spec:
   ```

```yaml
      containers:
      - name: proteomics
        image: ASIATUYJP7SUEWGXB6WN.dkr.ecr.us-east-
1.amazonaws.com/proteomics-app:latest
        ports:
        - containerPort: 8501
```

2. **Deploy the application**:
   ```sh
   kubectl apply -f proteomics-deployment.yaml
   ```

3. **Expose the application** (e.g., using a LoadBalancer service):
   ```yaml
   apiVersion: v1
   kind: Service
   metadata:
     name: proteomics-app-service
   spec:
     type: LoadBalancer
     ports:
     - port: 80
       targetPort: 8501
     selector:
       app: proteomics-app
   ```

   Then apply this configuration:
   ```sh
   kubectl apply -f proteomics-service.yaml
   ```

4. **Access your application**. After a few minutes, get the LoadBalancer URL:
   ```sh
   kubectl get services
   ```
   Look for the `EXTERNAL-IP` of your `proteomics-app-service`.

### Step 5: Monitoring and Management

- Use Kubernetes Dashboard or AWS CloudWatch for monitoring.
- Set up autoscaling if needed.
- Regularly update your application with security patches.

working ECR:

1. Retrieve an authentication token and authenticate your Docker client to your registry.

    Use the AWS CLI:

    aws ecr-public get-login-password --region us-east-1 | docker login --username AWS --password-stdin public.ecr.aws/y3d9e0t2

    Note: If you receive an error using the AWS CLI, make sure that you have the latest version of the AWS CLI and Docker installed.

2. Build your Docker image using the following command. For information on building a Docker file from scratch see the instructions [here](http://docs.aws.amazon.com/AmazonECS/latest/developerguide/docker-basics.html) . You can skip this step if your image is already built:

    docker build -t x21156239_proteomics-ecr .

3. After the build completes, tag your image so you can push the image to this repository:

    docker tag x21156239_proteomics-ecr:latest public.ecr.aws/y3d9e0t2/x21156239_proteomics-ecr:latest

4. Run the following command to push this image to your newly created AWS repository:

    docker push public.ecr.aws/y3d9e0t2/x21156239_proteomics-ecr:latest

## 13. Kubernetes setup



## 14. ECR repository creation



## 15. After this all is done we have to create EC2 instance in AWS.

16. Open the deployed link. UI will be like below.



Now upload CSV file, After Uploading the bulk CSV data, Single input prediction is done.



The predicted sequence will be stored in prediction history tab

Now we will upload a small CSV file for testing, after uploading below is the view



After saving upload save prediction

After saving the predictions, it will get stored in database



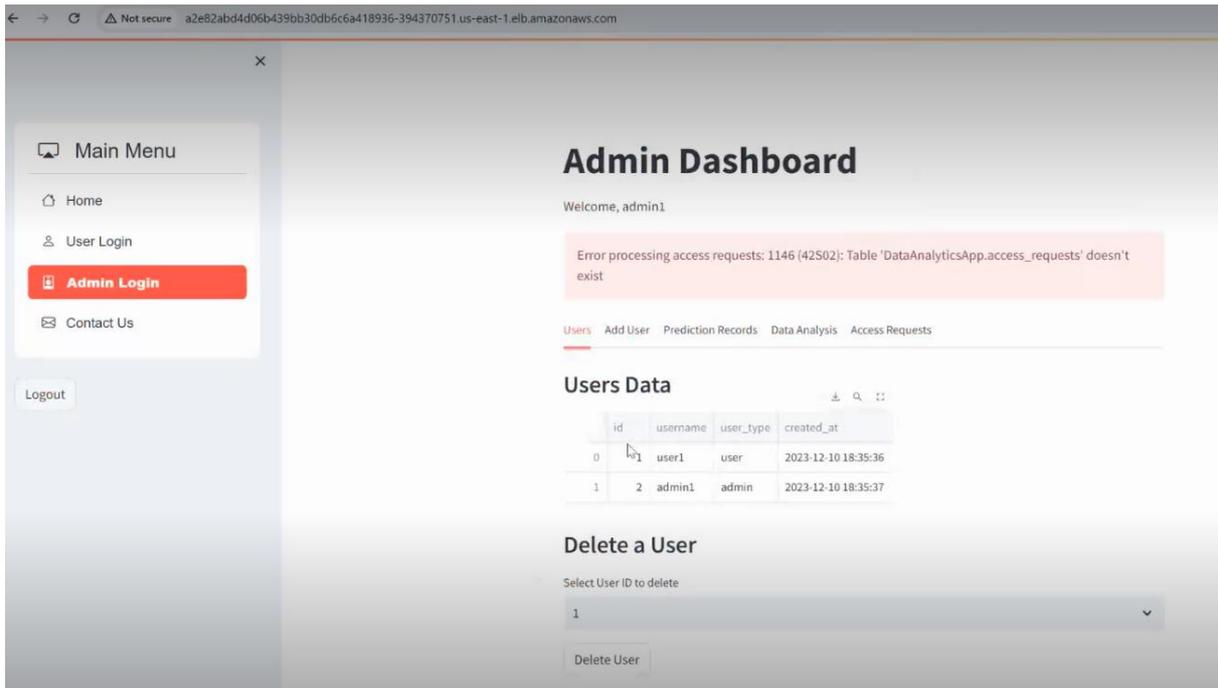While the prediction is getting saved in database, also the data is getting encrypted due to PBFT blockchain algorithm.
Only Admin will be having access to the encrypted key, Only when User request admin for the key, then only it can accessed by user.

Now, Login to Admin panel.



Here, we can create and delete user

On Prediction tab, we can see all the uploaded prediction data.



Even in UI Admin can see the Encryption key. Move the table to right to view encrypted key

In Data Analysis tab, first upload the CSV file.



After uploading file, Visualizations will appear

Bar Chart