# MLOps – NLP – Case Study

**Q1. System design: Based on the above information, describe the KPI that the business should track.**

A) The business should track several key performance indicators (KPIs) related to the automated extraction of diseases and treatments from free-text clinical notes. These KPIs would help assess the effectiveness, efficiency, and impact of the automated solution. In the context of the BeHealthy problem statement, the business goal is to automate the process of extracting diseases and treatments from free text, which would result in reduced man-hours, eliminate the need for a data-entry team, reduce manual errors, and scale up the size of the data without increasing the size of the data-entry team.

This KPI would measure how well the automated process is performing in correctly identifying the diseases and treatments from the clinical notes compared to the manually extracted information. The accuracy of the automated process should be high enough to ensure that there are no errors in the extracted data, and it should be continually monitored and improved. By tracking this KPI, BeHealthy can ensure that their automated process is meeting the desired business goals of reducing man-hours, eliminating the need for a data-entry team, reducing manual errors, and scaling up the size of the data without increasing the size of the data-entry team.We can see this in detail.

Accuracy of Extraction: This KPI measures the accuracy of the extracted diseases and treatments compared to manual extraction by clinical experts. It can be calculated as the percentage of correctly extracted entities over the total number of entities. High accuracy is crucial for ensuring the reliability of the extracted information.

Speed of Processing: Track the time taken to process a batch of clinical notes and extract structured information. This KPI will indicate the efficiency of the solution in handling large volumes of data. Faster processing times would lead to improved productivity and responsiveness.

Reduction in Manual Effort: Measure the reduction in manual effort required for data entry and validation. This can be quantified by comparing the number of man-hours before and after implementing the automated solution. Reduced manual effort signifies increased efficiency and cost savings for the business.

Error Rate: Monitor the error rate in extracted information to ensure the reliability of the system. This KPI measures the proportion of inaccurately extracted entities and helps identify areas for improvement and optimization. Lower error rates indicate better quality of extracted data.

Scalability: Assess the scalability of the solution to handle increasing volumes of clinical notes without significant degradation in performance. This KPI can be measured by analyzing resource utilization and processing times as the dataset size grows. Scalability ensures that the solution can accommodate growing demands without compromising performance.

Cost Savings: Calculate the cost savings achieved by eliminating the need for manual data entry and validation. This includes savings in labor costs and potential revenue gains from improved efficiency. Cost savings demonstrate the business value of the automated solution.

User Satisfaction: Gather feedback from users, including doctors and data validation teams, to assess their satisfaction with the automated system. This qualitative measure provides insights into usability, functionality, and user experience. High user satisfaction indicates successful adoption and acceptance of the automated solution.

Data Drift Detection: Monitor for data drift in the clinical notes dataset to ensure the continued accuracy of the extraction model. Track metrics related to data distribution and concept drift to identify deviations from the training data. Timely detection of data drift allows for proactive maintenance and model retraining maintaining the performance.

These KPIs will help the business evaluate the performance, efficiency, and impact of the automated solution for extracting diseases and treatments from clinical notes.

Q2. System Design: Your Company has decided to build a MLOps system. What advantages would you get by opting to build a MLOps system?

A) Building an Machine Learning Operations(MLOps) system, provides several advantages for companies that are developing and deploying machine learning models. Some of the advantages are:

- Reduction in human errors.
- Reduction in man-hours spent on building and training the system.
- Reduce the cost required to build and maintain such a model.
- Building a model takes a lot of effort. An MLOps system will bridge the gap and also reduce lag in model development and deployment.
- The system will have tools that help teams to collaborate during the model building stage.
- MLOps systems provide robust version control and governance, ensuring that models are developed, deployed, and maintained in a structured and compliant manner.
- By developing and deploying machine learning models more efficiently, businesses can deliver better customer experiences, such as faster response times, more accurate predictions, and personalized recommendations.
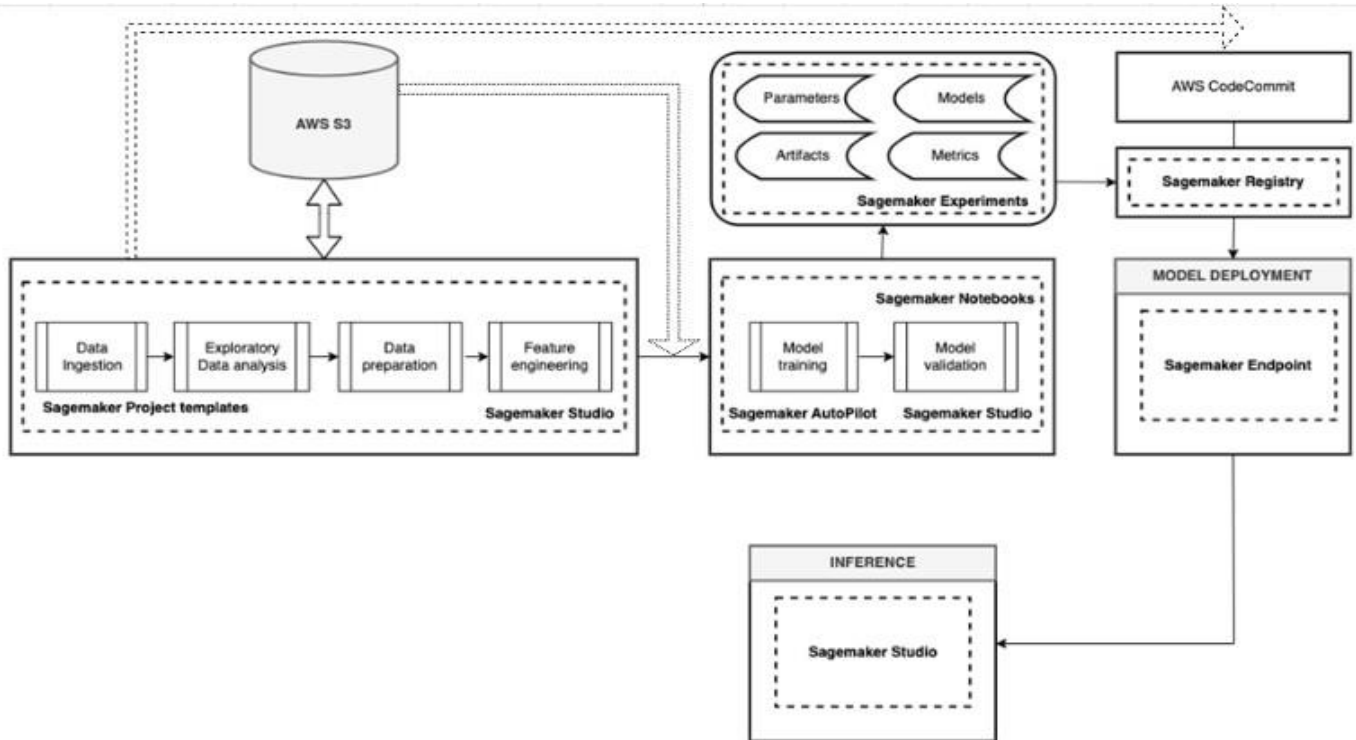
- There will be a good view on what experiments were conducted and their metrics so the best model can be identified.
- Experiment tracking in MLOps system will track all the modelling exercises that were conducted. It will give a good view of what is working and what is not and then it will pick the best model.
- Results of experiments can also be shared between various teams for evaluation and validation.
- MLOps systems enable better collaboration between data scientists, developers, and IT operation teams, resulting in improved communication and teamwork.
- Symmetry can be maintained between development and production environments to ensure that the results got during development are the ones that are got in production.
- With MLOps, companies can automate many aspects of machine learning model development, deployment, and maintenance, resulting in reduced operational costs.
- MLOps systems provide robust security and compliance measures, ensuring that data is handled and stored appropriately. This can help businesses mitigate data privacy and security risks.
- MLOps systems can help businesses gain a competitive advantage by enabling them to quickly and efficiently develop and deploy machine learning models. This can help businesses stay ahead of the competition and improve their market position.
- ML models tend to become stale over time. MLOps systems have monitoring systems to detect data drift or model bias automatically. This will ensure that prediction accuracy is maintained.
- The pipeline can be triggered, models would be continuously trained as and when a drift is detected or new data arrives. It will help avoid model decay.
- API endpoints can be created for real time inference.
- The whole pipeline can be automated with minimal manual intervention.

This system would eliminate the need to have a trained doctor analysing each and every clinical note to identify the mapping as well as reduce the need of a data entry team. This automation would result in reduced human errors and reduced man-hours effort
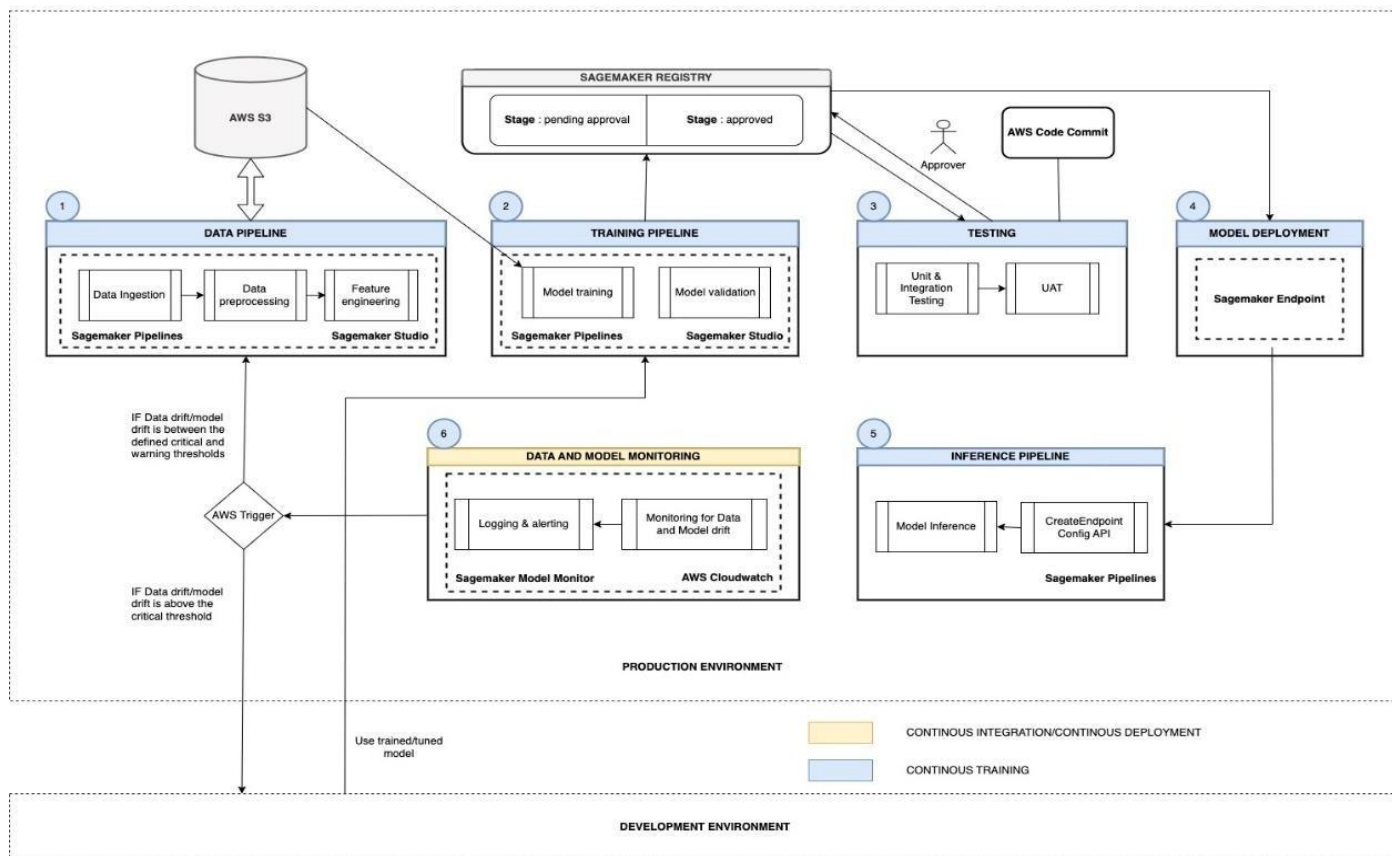- Increase productivity of the team
- Ensure high quality model production

In summary, building an MLOps system can provide businesses with a range of benefits, including increased model accuracy and reliability, enhanced data security and compliance, improved customer experience, competitive advantage, and better resource management.

**Q3. System design: You must create an ML system that has the features of a complete production stack, from experiment tracking to automated model deployment and monitoring. For the given problem, create an ML system design (diagram).**



*Development Pipeline for MLOps systems(Experimentation and Trials*

*Production Pipeline for MLOps system*

**Q4.System design: After creating the architecture, please specify your reasons for choosing the specific tools you chose for the use case.**

**A)Based on the use case and the fact that the BeHealthy is a late-stage startup, we have chosen managed cloud services (AWS PageMaker) as our primary tool for this system. Here are some of the reasons why:**

- ☐ **Easy deployment and management: Managed cloud services like AWS SageMaker provide easy deployment and management of the entire system. This is especially beneficial for a startup like BeHealthy, as it reduces the need for dedicated DevOps personnel and saves time and resources.**
- ☐ **Cost-effective: Managed cloud services offer a pay-as-you-go pricing model, which can be more cost-effective than traditional on-premises infrastructure. This is important for a startup like BeHealthy, which may have limited resources.**
- ☐ **Scalability: Managed cloud services offer scalability, allowing the system to scale up or down based on the volume of data and user traffic. This ensures that BeHealthy can handle an increase in demand for its services without experiencing downtime or performance issues.**

- **Availability and reliability:** Managed cloud services provide high availability and reliability, ensuring that the system is always up and running. This is essential for a healthcare platform like BeHealthy, which needs to provide uninterrupted services to its users.
- **Security:** Managed cloud services like AWS SageMaker offer robust security features, including encryption, access control, and compliance with industry standards. This ensures that user data is secure and that BeHealthy meets its regulatory requirements.

Overall, we believe that managed cloud services like AWS SageMaker provide the right combination of ease of deployment, cost-effectiveness, scalability, availability and reliability, and security to meet the needs of BeHealthy's system.

Here are some of the key tools and services offered by AWS SageMaker which can be leveraged by BeHealthy:

- **Jupyter Notebook:** SageMaker offers Jupyter Notebook, an open-source web application that allows users to create and share documents that contain live code, equations, visualizations, and narrative text. This tool is particularly useful for exploratory data analysis and rapid prototyping of machine learning models.
- **Data Wrangler:** This tool helps data scientists to clean, transform and aggregate data quickly and easily. It has built-in data cleaning and feature engineering capabilities and can handle a variety of data formats, including CSV, JSON, and Parquet.
- **Amazon SageMaker Studio:** This is an integrated development environment (IDE) that provides a web-based interface to build, train, and deploy machine learning models. It offers a range of tools and services, including Jupyter Notebook, AWS Glue DataBrew, and built-in machine learning algorithms.
- **Built-in algorithms:** SageMaker offers a range of built-in machine learning algorithms for common tasks such as classification, regression, and clustering. These algorithms are pre-optimized and can be used with large datasets.
- **Custom algorithms:** SageMaker also allows users to develop and deploy custom machine learning algorithms using frameworks such as TensorFlow and PyTorch.
- **Automatic Model Tuning:** This tool automates the process of hyperparameter tuning by running multiple training jobs in parallel, selecting the best-performing model based on the user-specified objective metric.
- **Model hosting and deployment:** SageMaker makes it easy to deploy trained machine learning models as web services, which can be accessed by other applications through a REST API.

<u>**Benefits/Characteristics of Sagemaker :**</u>

- **Amazon SageMaker is a fully managed machine learning service.**
- **With SageMaker, data scientists and developers can quickly and easily build and train machine learning models, and then directly deploy them into a production-ready hosted environment.**
- **It provides an integrated Jupyter authoring notebook instance for easy access to your data sources for exploration and analysis, so we do not have to manage servers.**
- **It also provides common machine learning algorithms that are optimised to run efficiently against extremely large data in a distributed environment. With native support for bring-your-own-algorithms and frameworks, SageMaker offers flexible distributed training options that adjust to your specific workflows. Deploy a model into a secure and scalable environment by launching it with a few clicks from SageMaker Studio or the SageMaker console.**
- **Training and hosting are billed by minutes of usage, with no minimum fees and no upfront commitments.**

<u>**Tools within AWS SageMaker:**</u>

- **Amazon SageMaker Studio - First fully integrated development environment (IDE) for machine learning.**
- **Amazon SageMaker Notebooks - Enhanced notebook experience with quick-start and easy collaboration**
- **Amazon SageMaker Experiments - Experiment management system to organise, track and compare thousands of experiments.**
- **Amazon SageMaker Debugger - Automatic debugging analysis and alerting**
- **Amazon SageMaker Monitor - Model monitoring to detect deviation in quality and take corrective actions.**
- **Amazon SageMaker Autopilot - Automatic generation of machine learning models with full visibility and control.**

**Q5. Workflow of the solution:**
**We must specify the steps to be taken to build such a system end to end. The steps should mention the tools used in each component and how they are connected with one another to solve the problem. Broadly, the workflow should include the following. Be more comprehensive of each step that is involved here.**
1. **Data and model experimentation**
2. **Automation of data pipeline**
3. **Automation of the training pipeline**
4. **Automation of inference pipeline**
5. **Continuous monitoring pipeline**
**The workflow should also explain the actions to be taken under the following conditions.**
**After you deployed the model, you noticed that there was a sudden increase in the drift due to a shift in data.**

**A)** The system architecture chosen for this case study is divided into multiple layers for the different stages of ML system life cycle. The ML system has been designed to considering all steps involved in a software Lifecyle with emphasis on keeping KPI high for entire Lifecyle with all benefits of MLOps.

The system has been divided into 2 environments for separation between development and client facing system:

- Development environment: The Development environment has been designed to facilitate the evaluation of various models and determine the optimal solution for the given problem statement and dataset. It offers an ideal platform for swift experimentation with both the data and models at hand.
- Production environment: This environment where the best model is deployed after testing and contains pipeline for training, testing, monitoring and inference which fulfills their role to allow system to operate in best manner.

Broadly steps are same as any workflow of MLOPS which have been explained below:

## 1) DATA AND MODEL EXPERIMENTATION

This is done completely in the Development environment which it was designed for. In Development environment Data scientist can collaborate and perform experimentation and trails to find the best model with best features which can be used in future as well in case of any drift.

In the Structure we advised the feature engineering is done in Jupyter notebook given in SageMaker studio which contains template which can be used to speed up the process. Here data scientist does the experimentation and prepare final data which is then used by SageMaker Autopilot to do experimentation across optimized models and algorithms to choose best performing model based on a metric.

All theses are tracked by the SageMaker experiments and trials where it stores all artifacts as well as results, which one is approved can be pushed to endpoint from SageMaker model registry for testing purposes using SageMaker endpoint. Once result is obtained as per threshold model is moved to production.

1. **CAPABILITIES OF SAGEMAKER PIPELINES**

• **Build ML workflows:** Using Python SDK, we can build ML workflows comprising parameters, different steps and data dependencies. We can also orchestrate SageMaker jobs such as the processing job and the training job and can also trigger

the execution of these pipelines.

• **Troubleshoot ML workflows:** We can visualize the execution of the pipeline and the status of each step in the pipeline in real time in SageMaker Studio. We can also view additional information about each of the steps in SageMaker Studio.

• **Manage models:** We can manage different versions of models using the Model Registry. We also have the capability to approve/reject models in the model registry. The model registry consists of different model packages, and each model package consists of multiple versions of the model.

• **Scaling MLOps:** We can create a project in SageMaker Studio and get a code repository, seed code and the MLOps infrastructure set up for me. We are provided with MLOps templates published by SageMaker for building, deploying and establishing end-to-end workflows.

• **Track lineage:** With in-built lineage tracking for SageMaker pipelines, we can track data, models and artifacts. Also, support is provided for tracking custom entities.

## 2. AUTOMATED DATA AND TRAINING PIPELINE

This component focuses on automating model training by converting the code developed in notebooks to Python scripts. With automation coupled with using a feature store, the data and training pipelines can run whenever there is any change in the live data. This helps in the continuous delivery of existing deployed models after it is re-trained on the newly transformed data stored in the feature store. The tool used for automation in this stage is SageMaker Studio and SageMaker Pipelines.

## 3. TESTING

In this step, we will test the different methods used in data preparation, feature extraction and model validation to effectively track whether all the components are working in the desired manner. The tests applied here are unit tests, integration tests, and user acceptance testing (UAT). If the model passes all these tests, it can be moved to production, it can be used for making inferences/predictions. Therefore, testing helps in the continuous integration of models trained on new data.

## 4. INFERENCE PIPELINE

In this stage, once the model/code passes all the tests, we will go ahead and deploy the model for serving predictions. The tool used in this stage is SageMaker Endpoint for deployment and providing end-service.

5. **DATA AND MODEL MONITORING**

Keeping a continuous check on the deployed model is essential for tracking the model performance and ensuring that the model doesn't go stale. It signals what action needs to be performed based on any changes in the live data. The 'trigger' connected to this component decides what action to take: model experimentation or model retraining. The tool used in this stage for monitoring any data drifts is called SageMaker Model Monitor and AWS CloudWatch. The Amazon SageMaker Model Monitor continuously monitors the quality of Amazon SageMaker ML models in production. The Model Monitor provides the following types of monitoring:

• Monitor data quality: Monitor drift in data quality

• Monitor model quality: Monitor drift in model quality metric, such as accuracy

• Monitor bias drift for models in production: Monitor bias in your

model's predictions We can set alerts using AWS CloudWatch when in the

case of deviations in the model quality.

1) What component/pipeline will be triggered if there is any drift detected? What if the drift detected is beyond an acceptable threshold?

A)If the drift is between the defined warning threshold and the critical threshold, then it implies the inference data is not similar to the data model was trained on. This would not give accurate results. Hence, in this case, model should be retrained with the new data, so that we can get optimized model and its parameters.
If the data drift is above the defined critical threshold, then model retraining will not help and in this case, we need to go back to the development environment, perform the process of experimentation again to find the best ML model, which needs to again undergo UAT and then should be good for inference purpose.

2) What component/pipeline will be triggered if you have additional annotated data?

A)To ensure the new patterns in the additional data are captured, the entire pipeline starting from the data ingestion stage must be triggered. This involves ingesting the new data into the pipeline, followed by training, validating, testing, and deploying the model again. By doing so, the model will be updated to incorporate the new patterns in the additional data, resulting in more accurate outcomes.

The complete pipeline starting with the data pipeline should be triggered so that the additional data is ingested and the model is trained, validated, tested and then deployed again. This will make sure that the new patterns in the additional data is captured.

3) How will you ensure the new data you are getting is in the correct format that the inference pipeline takes?

A)In order to prepare the new or test data for accurate predictions by the model, it needs to undergo the data pipeline, which involves preprocessing and formatting the data to match the requirements of the model. As a result, the data pipeline is typically triggered on the test.

The new data/test data will have to undergo the data pipeline, which ensures the data is preprocessed and is available in the format that is expected by the model. Hence,most of the times, data pipeline is triggered on the test dataset and inference pipeline will help to get predictions on the processed test / dataframe datasets. Once the data has been processed, the inference pipeline can be used to generate predictions on the processed test dataset or dataframe.