# Mphasis SynthStudio

## Product Description

Mphasis' SynthStudio generates high-quality synthetic data to monetize trustworthy business insights while preserving privacy and protecting data subjects. Its AI-generated synthetic data enables organizations to leverage data's maximum potential to cross-collaborate and build reliable and highly accurate assets with no loss of data privacy and utility.

## Overview

Mphasis' SynthStudio creates synthetic data, based on metadata, constraints, and conditions, for a wide range of requirements - including application testing, creation of POCs, fueling hackathons and innovation initiatives. Synthetic data is created using mathematical models, computer algorithms or other artificial means, and can be used to train Machine Learning models or test software systems in situations where real data is difficult or costly to obtain, or when privacy concerns prevent the use of real data. This synthetic data, artificially constructed, allows businesses to test, learn, and innovate without breaching any real-world data privacy.

The proposed solution incorporates two components: Data Synthetization and Data Enrichment. These components work in tandem to maximize the value and utility of existing data resources. Data Enrichment is a sophisticated process that augments existing datasets with additional, relevant information. This enhancement significantly increases the contextual utility of the data, opening new avenues for innovation and deeper insights. For instance, it can involve generating descriptive tags for images, assigning sentiment labels to text reviews, or adding demographic information to user profiles. By layering these additional dimensions onto the original data, organizations can unlock hidden patterns and correlations, leading to more nuanced understanding and decision-making capabilities. Data Synthetization, on the other hand, is a powerful technique that leverages both representative original data and associated metadata to generate high-quality, synthetic datasets.

Synthetic data that is generated by this solution as a privacy-safe version of the original data unleashing maximum utility for insights. This feature is especially useful when synthetic data is generated from original data while preserving its statistical properties and distributions.

Mphasis' SynthStudio stands out as a synthetic data generation and enrichment solution, utilizing innovative algorithms and proprietary methodologies. This approach establishes a scalable and secure pipeline, ensuring the creation of technically sound and privacy-protected synthetic datasets.

**Highlights**

- Privacy-preserving

- Expert-assisted Synthetic Data Generation
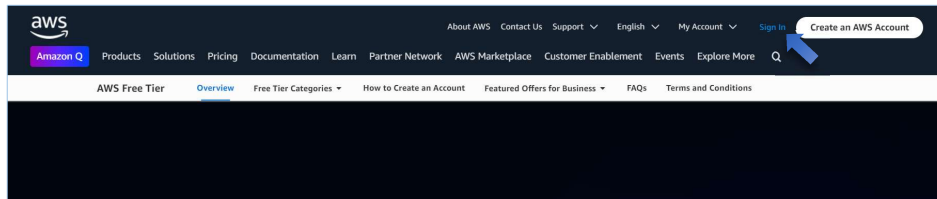
- Deep Learning

## Prerequisites

1. Domain Setup
   If you already have a domain registered or hosted on Route 53, you can use the existing domain and skip the Domain Setup.
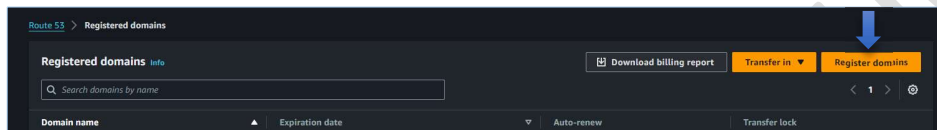   You can get a dedicated URL with AWS Route 53. Follow "Domain Setup" to know how to setup a domain
2. Domain Setup
   a. Create or sign into your AWS account

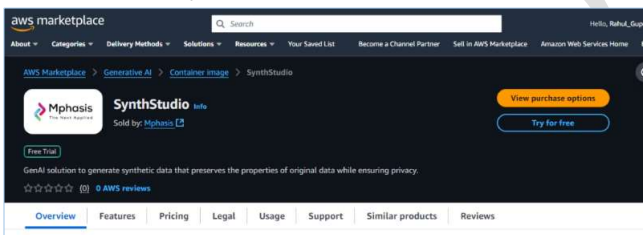   b. Search Route 53 on console and register a domain.

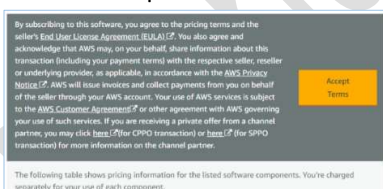   For more details steps on how to register domain click here

## How to Subscribe

1. Navigate to AWS marketplace AWS Marketplace: Homepage
2. Search for SynthStudio Mphasis
3. Subscribe to the product

4. Click on Accept Terms

5. Wait for processing to complete

6. Click on Continue to Configuration to start configuring your product.

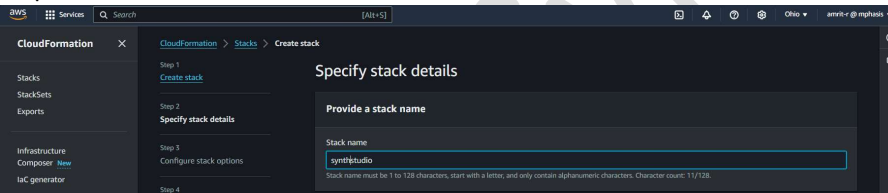7. Click on Continue to Launch to initiate the launch process for your product.



8. Click on CloudFormation Template link under Deployment template to access the installation.
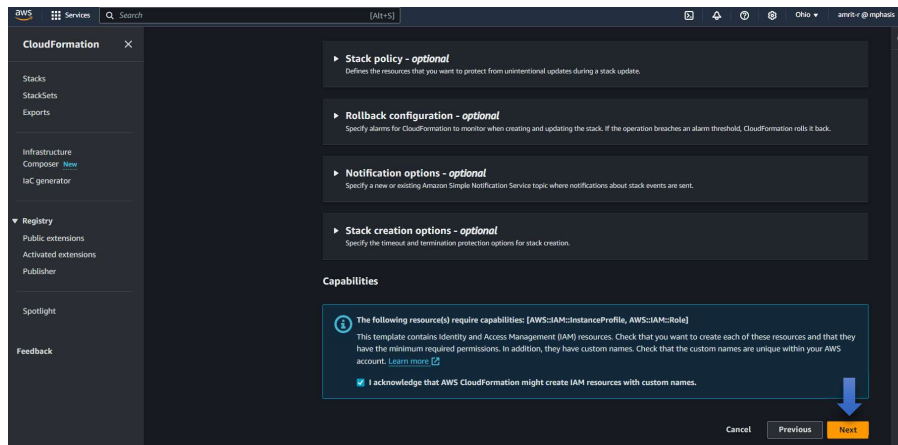


## How to Setup

1. Launch the cloud formation template. You will be redirected to AWS CloudFormation service from the previous step



2. Fill in the necessary **Parameters** fields required in the cloud formation template.
   a. **Clientname**: The name you give to your application. This name is a unique identifier, and all resources created will be using this name.
   b. **DNSRecord**: Your domain name. This domain name is part of AWS Route 52 Hosted zone.
   c. **InstanceTypeInference**: Options to select your EC2 Instance type for inference. Picking a GPU instance is recommended. Default g4dn.xlarge.
   d. **InstanceTypeTrain**: Options to select your EC2 Instance type while training the model. Picking a GPU instance is recommended. Default g4dn.xlarge.
   e. **Route53HostedZoneID**: The hosted zone ID from route53 where your domain is hosted.
   f. **SSLCertificateARN**: A valid certificate ID from AWS certificate manager which is issued against your domain name.
3. For advanced users only (Optional Step): Fill in the optional parameters if need and click next.
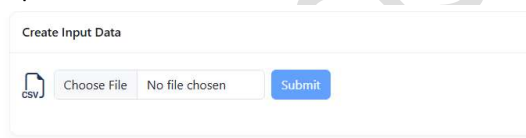
4. Wait till the resources are created and CloudFormation status goes to "CREATE_COMPLETE".
5. (Optional): Modify the load balancer. The Security group attached to the load balancer is open to the public by default and its inbound rules can be modified according to organization policies.
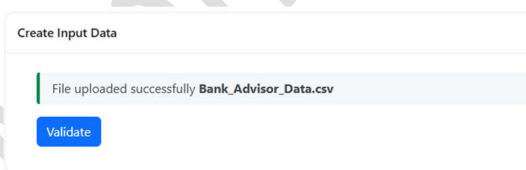6. Go to your browser and open https://clientname.DNSRecordName.
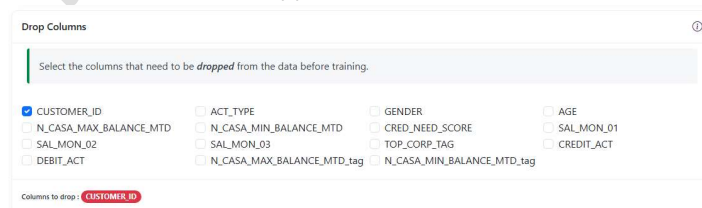
## How to Run

1. Open Homepage



2. Upload a CSV file



3. Click on Validate



4. Select columns to be dropped

5. Select variable categories if not identified correctly. There should be at least one categorical and at least one numerical variable.



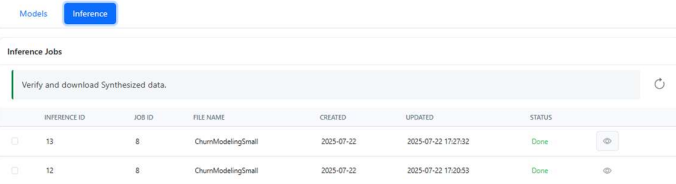6. Add constraints based on requirement



7. Add Constraint details



8. Proceed for training

9. Go to Models page to check the training progress. After training is Done, add number of rows to be generated and proceed



10. Under Inference tab, check the inference progress. After inference is Done, proceed to look at the data



11. On the results page you can have a quick look at the generated data, evaluation metrics and download the results