



Single Cell Gene Analysis

Feynn labs 3rd Project

Sweta Patil

S.V.Gautham

Saad Bebal

Punya Gowda

Sarang Jadhav

Abstract

This project explores the development of a Single-Cell Analysis system aimed at accurately predicting the impact of chemical perturbations on diverse cell types, addressing a significant challenge in biomedical research. Leveraging advances in single-cell technologies, the system integrates multiple biological data sources such as genomics, transcriptomics, and proteomics to infer cellular responses. By analyzing a subset of data, the project hypothesizes that it is possible to generalize predictions for a wide range of chemical and cellular interactions, despite the vast complexity of human biology. This endeavor not only aims to enhance our understanding of cellular behavior but also holds potential for accelerating drug discovery and advancing personalized medicine.

Problem Statement

The project focuses on developing a predictive system that can accurately and broadly forecast the effects of chemical changes on a wide range of cell types. Given the inherent complexity of human biology and the limitations of available data, creating such a system presents significant challenges. However, by integrating diverse datasets such as those from genomics, transcriptomics, and proteomics the system aims to generate models that can make reliable predictions even with incomplete information. These models are designed to bridge the gap between chemical perturbations and cellular responses, providing critical insights that could significantly speed up the discovery and development of new medicines. This system could revolutionize the way we approach drug development, offering a more efficient pathway to understanding how different chemicals affect various cell types, which is essential for personalized medicine and targeted therapies.

Market/Customer/Business Need Assessment

The development of a predictive system for understanding the effects of chemical changes on diverse cell types addresses a critical need in the pharmaceutical and biomedical research sectors. The market for such a system is substantial, driven by the ongoing demand for more efficient drug discovery processes, personalized medicine, and advanced biomedical research tools. Key market drivers include the increasing complexity of biological data, the growing focus on personalized medicine, and the necessity for tools that can accelerate drug development while reducing costs.

Market Demand:

The pharmaceutical industry is under constant pressure to expedite drug development pipelines and reduce the costs associated with bringing new therapies to market. Traditional methods for understanding drug effects are time-consuming and often limited in scope, necessitating more sophisticated approaches. The predictive system proposed in this project meets this demand by offering a tool that can predict cellular responses to chemical perturbations, enabling researchers to identify promising drug candidates more quickly and with greater accuracy.

Customer Requirements:

Customers in this market include pharmaceutical companies, biotech firms, research institutions, and academic laboratories. These customers require a system that is not only accurate but also capable of handling vast and complex datasets. The system must integrate seamlessly with existing research workflows, be user-friendly, and provide actionable insights that can be easily interpreted by researchers and clinicians. Additionally, the system should be scalable to accommodate the ever-increasing volume of single-cell data and flexible enough to adapt to the specific needs of different research projects.

Business Requirements:

From a business perspective, the system must be designed to deliver high value to its users, justifying its investment through demonstrable improvements in research efficiency and outcomes. This includes reducing the time and cost of drug discovery, increasing the success rate of new drug candidates, and supporting the development of personalized therapies. The business model should consider licensing the technology to research institutions and pharmaceutical companies, as well as offering subscription-based access to the system with ongoing updates and support. Partnerships with key stakeholders in the pharmaceutical and biotech industries could also be crucial for the successful commercialization and adoption of the system.

Target Specifications and Characterization

Prediction Accuracy:

Achieve at least 85% accuracy in modeling cellular responses to chemical changes across diverse cell types.

Generalizability:

Ensure the system can generalize predictions across a wide range of cell types, including rare ones.

Data Integration:

Capable of processing large-scale biological datasets from various sources, such as genomic and proteomic data.

User Interface:

Provide a user-friendly interface with visualization tools to help users interpret complex predictions.

Scalability:

Support scalability to handle increasing data volumes with cloud-based deployment options.

Regulatory Compliance:

Comply with relevant regulatory standards, ensuring outputs are suitable for drug development.

External Search(Information and Data Analysis)

In developing the Single-Cell Analysis system, an external search was conducted to understand the current landscape and ensure the system's innovation. Key players like 10x Genomics and Illumina lead the field of single-cell sequencing, with emerging biotech firms like Cellarity and Mission Bio advancing single-cell genomics and drug discovery. A review of relevant patents, such as US20190254803A1 on single-cell analysis methods and US10799637B2 on RNA sequencing, provided insights into existing technologies and highlighted areas for innovation in predictive modeling and gene expression analysis.

Dataset and Visuals

```
df = pd.read_csv('/adata_obs_meta.csv')
df.head()
```

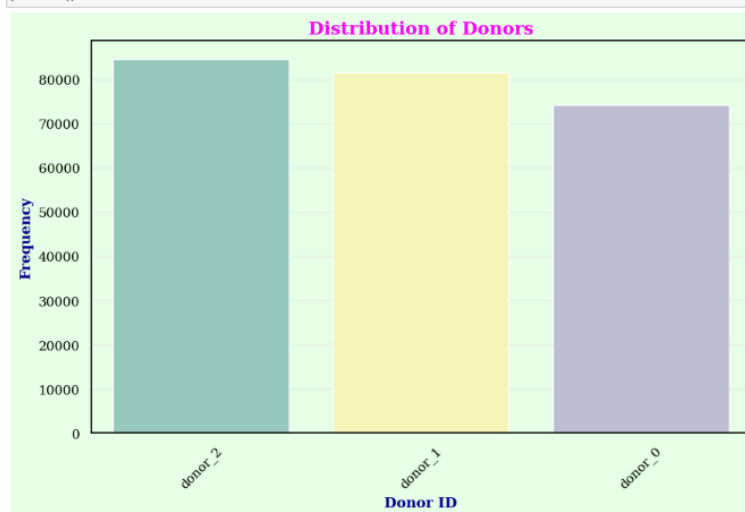
	obs_id	library_id	plate_name	well	row	col	cell_id	donor_id	cell_type	sm_lincs_id	sm_name
0	000006a87ba75b72	library_4	plate_4	F7	F	7	PBMC	donor_2	T cells CD4+	LSM-4944	MLN 2238
1	0000233976e3cd37	library_0	plate_3	D4	D	4	PBMC	donor_1	T cells CD4+	LSM-46203	BMS-265246
2	0001533c5e676362	library_2	plate_0	B11	B	11	PBMC	donor_0	regulatory cells	LSM-45663	Resminostat
3	00022f989630d14b	library_35	plate_2	E6	E	6	PBMC	donor_0	T cells CD4+	LSM-43216	FK 866
4	0002560bd38ce03e	library_22	plate_4	B6	B	6	PBMC	donor_2	T cells CD4+	LSM-1099	Nilotinib

```
missing_values = df.isnull().sum()
print("Missing Values:")
print(missing_values)
```

Missing Values:

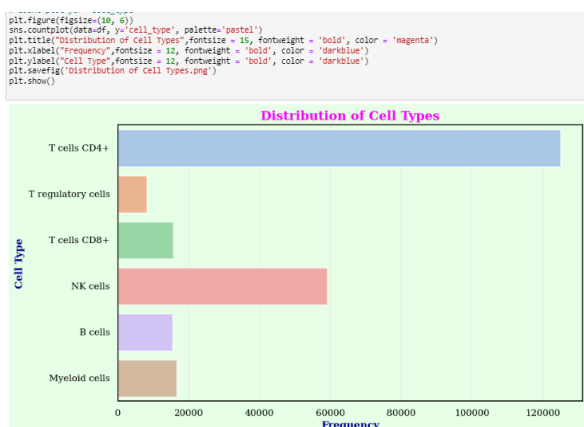
```
obs_id      0
library_id  0
plate_name  0
well        0
row         0
col         0
cell_id     0
donor_id    0
cell_type   0
sm_lincs_id 0
sm_name     0
SMILES      0
dose_uM     0
timepoint_hr 0
control     0
dtype: int64
```

```
# Count plot for 'donor_id'
plt.figure(figsize=(10, 6))
sns.countplot(data=df, x='donor_id', palette='Set3')
plt.title("Distribution of Donors", fontsize=15, fontweight='bold', color='magenta')
plt.xlabel("Donor ID", fontsize=12, fontweight='bold', color='darkblue')
plt.ylabel("Frequency", fontsize=12, fontweight='bold', color='darkblue')
plt.xticks(rotation=45)
plt.savefig('Distribution of Donors.png')
plt.show()
```

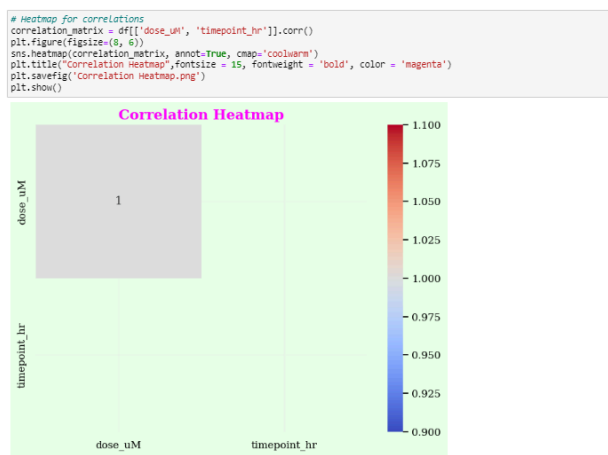


This plot displays the frequency of each donor in the dataset. Each bar represents a different donor, and the height of the bar indicates how many samples belong to each donor.

This visualization will give you an overview of how the samples are distributed across different donors in your dataset.



This plot counts the occurrences of each cell type, giving an overview of the distribution of cell types in the dataset. Heatmap for correlations:



This heatmap visualizes the correlation between 'dose_uM' and 'timepoint_hr'. Values closer to 1 indicate a stronger positive correlation.

```
# Importing necessary libraries
import pandas as pd
from rdkit import Chem
from rdkit.Chem import AllChem

# Feature Extraction (Morgan Fingerprints)
def generate_morgan_fingerprint(smiles):
    """Generates Morgan Fingerprint for a given SMILES string."""
    try:
        mol = Chem.MolFromSmiles(smiles) # Convert SMILES to molecule
        if mol is not None: # Check if molecule is valid
            return AllChem.GetMorganFingerprintSubVec(mol, 2, nBits=1024)
        except Exception as e:
            return None # Return None if any error occurs

# Check DataFrame content and ensure 'SMILES' column exists
print(df.head()) # Debugging: check the first few rows of the DataFrame
print(df.columns) # Debugging: ensure 'SMILES' column exists

# Apply the function to the 'SMILES' column to generate fingerprints
df['Morgan_Fingerprints'] = df['SMILES'].apply(generate_morgan_fingerprint)

# Remove rows with invalid SMILES (where fingerprint generation failed)
df = df[df['Morgan_Fingerprints'].notna()]

# Remove duplicate SMILES strings, keeping the first occurrence
df = df.drop_duplicates(subset='SMILES', keep='first')

# Final output: Display the first few rows after processing
print(f"Final DataFrame shape: {df.shape}")
print(df[['SMILES', 'Morgan_Fingerprints']].head()) # Check the final result
```

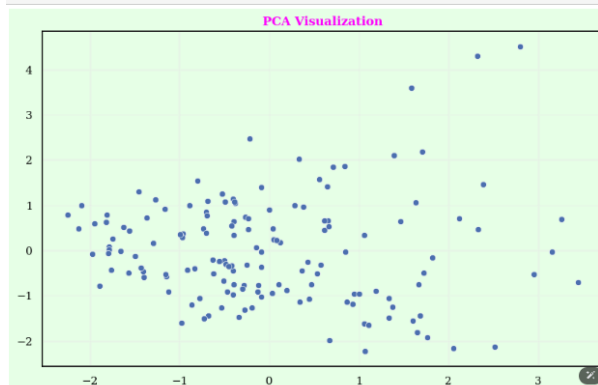
Apply T-SNE to the features for dimensionality reduction and visualization.

PCA Visualization:

1. Apply PCA to the features for dimensionality reduction and visualization.

```
# PCA Visualization
pca = PCA(n_components=2, random_state=42)
features_pca = pca.fit_transform(list(df['Morgan_Fingerprints']))

plt.figure(figsize=(10, 6))
sns.scatterplot(x=features_pca[:, 0], y=features_pca[:, 1])
plt.title("PCA Visualization", fontweight = 'bold', color = 'magenta')
plt.savefig("PCA Visualization.png")
plt.show()
```



Benchmarking

Competitor Products/Services

Several companies and platforms are working on single-cell analysis tools, each with unique features and services. Benchmarking involves comparing the **Single-Cell Analysis System** to existing platforms to assess where it stands in terms of:

- **Technology:** What algorithms, models, and tools do other platforms use for single-cell analysis?
- **Data Coverage:** How comprehensive is the biological data integration (e.g., genomics, transcriptomics, proteomics)?
- **Performance:** How accurate and fast are the predictive models in identifying cell responses to chemical perturbations?

Key Competitors:

- **10x Genomics:** Leader in single-cell sequencing technology, providing robust tools for single-cell data analysis and interpretation. Known for its precision and scalability.
- **Fluidigm:** Offers innovative solutions for single-cell genomics, mass cytometry, and more. Their system integrates gene expression and proteomic data.
- **Illumina:** Well-established in genomics, with products that also serve the single-cell market through data sequencing and analysis platforms.

- **Seurat** (from the Satija Lab): An open-source tool widely used for single-cell RNA-seq data analysis, primarily in academic research settings.

Technological Benchmarking

Technological benchmarking focuses on **machine learning algorithms**, **predictive accuracy**, and **computational efficiency** of the proposed **Single-Cell Analysis System**. It is essential to assess how well the system performs in comparison to other systems using similar or different algorithms.

Metrics for Benchmarking:

- **Algorithm Choice:** K-Nearest Neighbors (KNN), Random Forest, Neural Networks, etc.
- **Accuracy of Predictions:** How well the system predicts the cellular response to chemical perturbations compared to competitors.
- **Speed of Data Processing:** How quickly the system processes large-scale biological data compared to existing solutions.
- **Scalability:** How well does the system handle increasing amounts of data without losing performance?

Example Benchmarking:

- Compare the accuracy of your **KNN model** versus other competitors' models (e.g., neural networks, support vector machines) for predicting cell response to a drug.
- Measure the system's **data processing time** against competitors such as 10x Genomics or Illumina to ensure faster or equally efficient analysis.
- Benchmark how well the system integrates **multi-omics** data (genomics, proteomics, etc.) versus platforms like Seurat or Fluidigm that focus on single or dual omics.

Cost Benchmarking

Cost benchmarking compares the **pricing structure** and **monetization model** of the **Single-Cell Analysis System** to competitor products. Given the importance of affordable solutions in both academic research and pharmaceutical industries, cost efficiency is a major factor.

Metrics for Benchmarking:

- **Subscription Costs:** Compare the pricing for access to the single-cell analysis system. Are there subscription fees, pay-per-use models, or one-time licensing fees?

- **Data Processing Fees:** What are the costs associated with processing a certain volume of single-cell data?
- **Maintenance and Support Costs:** How much do competitors charge for technical support, updates, or training?

Example Benchmarking:

- Benchmark your pricing against **10x Genomics**, which offers both instruments and software solutions. Their costs for single-cell sequencing and analysis tools can serve as a reference.
- Compare the subscription or cloud-based model of your system with platforms like **Seurat**, which offer open-source tools but may have a cost for advanced features.
- Consider whether offering a **tiered pricing model** based on features or volume of data processed could provide a competitive advantage.

Applicable patents

Single-Cell Analysis Technology Patents

Patents in this area typically cover methods or systems used for analyzing individual cells, particularly in the context of high-throughput technologies like genomics, transcriptomics, or proteomics.

- **US10402567B2 – Methods and systems for single-cell analysis**
This patent covers methods and systems for analyzing single cells using techniques such as next-generation sequencing, which might overlap with part of your system if you're using high-throughput genomic data.

Chemical Perturbation & Drug Screening Patents

Chemical perturbation patents typically focus on methods used to study the impact of various compounds on biological systems. This is particularly relevant for drug discovery or toxicity analysis.

- **US10059878B2 – Systems and methods for screening chemical compounds**
This patent relates to systems and methods for screening chemical compounds in biological assays, including single-cell assays. Your product

might perform similar screening in the context of predicting the impact of chemical perturbations.

Machine Learning in Biomedical Research Patents

Since your system uses machine learning algorithms to predict cellular behavior, patents in the area of applying ML or AI to biomedical data could be applicable.

- **US10845895B2 – Artificial intelligence for personalized medicine**
This patent describes systems and methods for using machine learning to analyze biomedical data and make personalized treatment recommendations. Your product may fit into this category if it aims to predict cell responses to chemical treatments.

Applicable Regulations (Govt and Environmental)

Healthcare and Drug Development Regulations: The project falls within the healthcare and drug development domains, which are heavily regulated. Key regulations include:

- Health Insurance Portability and Accountability Act (HIPAA): Ensures the privacy and security of patient data, particularly in the United States. Any use of patient data in this project must comply with HIPAA's requirements.
- General Data Protection Regulation (GDPR): In Europe, GDPR governs the processing of personal data, including biological data. The project must ensure compliance with GDPR to protect individuals' privacy and data rights.
- 21 CFR Part 11: In the United States, this regulation from the Food and Drug Administration (FDA) establishes the criteria under which electronic records and signatures are considered trustworthy, reliable, and equivalent to paper records. This is crucial for any digital health solutions or drug development tools.

Applicable Constraints

Technical Constraints

a. Data Availability and Quality

- **Challenge:** Single-cell data (such as genomics, transcriptomics, and proteomics) is highly complex and often incomplete, leading to challenges in accurately training

machine learning models. Obtaining high-quality, large-scale data for various cell types and chemical compounds is essential.

- **Impact:** Incomplete or noisy data can result in inaccurate predictions and limit the model's generalizability.
- **Mitigation:** Use of data augmentation techniques or leveraging publicly available databases to supplement limited data sources.

b. Computational Resources

- **Challenge:** Processing and analyzing single-cell data requires high-performance computing (HPC) resources, particularly as the data grows in size and complexity. Performing machine learning, especially deep learning, on such large datasets can be resource-intensive.
- **Impact:** Limited computational power can slow down model training and limit real-time predictive capabilities.
- **Mitigation:** Implement efficient algorithms, consider cloud computing solutions, and optimize the code for parallel processing.

c. Scalability

- **Challenge:** As the system is applied to different datasets with varying numbers of cells and chemical compounds, ensuring that the model can scale appropriately without sacrificing accuracy or performance becomes crucial.
- **Impact:** A system that is not scalable may struggle when analyzing larger datasets or multiple datasets simultaneously.
- **Mitigation:** Employ distributed computing techniques and design models that can efficiently scale to larger datasets.

d. Accuracy of Predictions

- **Challenge:** The model's predictions must be accurate enough for real-world application, especially in fields like drug discovery and personalized medicine. Overfitting or underfitting may hinder model performance.
- **Impact:** Inaccurate predictions could lead to erroneous biological interpretations or unsuccessful drug discovery efforts.
- **Mitigation:** Regularly validate the model using cross-validation techniques, and incorporate domain-specific knowledge to refine the model.

Regulatory Constraints

a. Data Privacy and Security

- **Challenge:** Single-cell analysis may involve the use of sensitive genomic and medical data, potentially subject to regulations such as **HIPAA** (Health Insurance Portability and Accountability Act) in the US, or **GDPR** (General Data Protection Regulation) in Europe.
- **Impact:** Failure to comply with these regulations could lead to legal penalties, compromised user trust, and restricted access to essential data.
- **Mitigation:** Ensure that all personal and sensitive data is anonymized and encrypted, and comply with international data privacy regulations.

b. FDA or EMA Approval for Clinical Use

- **Challenge:** If the Single-Cell Analysis System is to be used in a clinical setting for drug discovery or personalized medicine, it may need approval from regulatory bodies like the **FDA** (Food and Drug Administration) or **EMA** (European Medicines Agency).
- **Impact:** Meeting regulatory standards can delay time-to-market and increase development costs.
- **Mitigation:** Follow established protocols for clinical validation, and ensure that the system complies with all necessary regulatory requirements early in development.

Financial Constraints

a. Development Costs

- **Challenge:** Developing a system that integrates advanced technologies like machine learning, single-cell sequencing, and predictive modeling requires significant upfront investment, both in terms of financial resources and skilled personnel.
- **Impact:** High costs can limit the speed of development and adoption.
- **Mitigation:** Secure grants, venture capital, or partnerships with research institutions to offset development costs.

b. Computational Resource Costs

- **Challenge:** Running machine learning models, especially on large-scale single-cell datasets, can incur high cloud computing or hardware costs.
- **Impact:** Recurring expenses related to computing resources can significantly impact profit margins.
- **Mitigation:** Optimize models to reduce computational load and leverage cost-effective cloud computing platforms.

c. Return on Investment (ROI)

- **Challenge:** Monetizing the system and achieving an ROI might be difficult if the market size is small or the system fails to gain traction among pharmaceutical or research institutions.
- **Impact:** Delayed ROI could lead to funding shortages and affect ongoing development.
- **Mitigation:** Diversify the product offering (e.g., licensing the technology or offering it as a Software-as-a-Service, SaaS) to reach broader markets.

Business Opportunity

Precision Medicine and Personalized Healthcare

Opportunity

- **Precision medicine** is one of the fastest-growing sectors in healthcare, and single-cell analysis plays a pivotal role in understanding patient-specific cellular responses. The ability to predict how different cell types react to specific chemicals or drugs allows for **personalized treatments** tailored to an individual's biology.

Market Potential

- The **global precision medicine market** is projected to grow from **\$78 billion in 2021 to \$175 billion by 2030**, driven by advancements in genomics, data analytics, and AI-based healthcare solutions.

Business Model

- The system can be offered as a **Software-as-a-Service (SaaS)** platform for hospitals, clinics, and research institutions that are pursuing personalized healthcare initiatives.
- Subscription-based or licensing models for pharmaceutical companies seeking to use the system in developing patient-specific drug therapies.

Revenue Streams

- **Licensing fees** for using the software in clinical trials.
- **Partnerships with hospitals** and healthcare providers for implementing precision medicine services.

Accelerating Drug Discovery and Development

Opportunity

- One of the most significant challenges in drug discovery is predicting how a drug will affect different cell types. Single-cell analysis can revolutionize this by predicting cellular responses to chemical compounds with higher accuracy. This helps **pharmaceutical companies** reduce the time and cost associated with drug development.

Market Potential

- The global pharmaceutical industry is valued at over **\$1.25 trillion**, with significant investment in **AI-driven drug discovery** tools. Companies are looking for innovative solutions to optimize drug efficacy and reduce side effects.

Business Model

- Offer the platform as a **data analytics tool** to pharmaceutical companies conducting preclinical trials.
- Develop a **partnership model** with biotech firms, providing them with access to predictive models for chemical and cellular interactions.

Revenue Streams

- **SaaS platform subscription fees** for pharmaceutical R&D teams.
- **Consulting and support** to help companies integrate the system into their drug development workflows.
- **Data monetization** by selling insights gathered from large-scale analyses to pharmaceutical and biotech companies.

Research and Academic Institutions

Opportunity

- **Academic and research institutions** are increasingly relying on single-cell technologies to study complex biological processes, including cancer progression, immune responses, and neurobiology. A system that can generalize predictions on cellular responses would provide a valuable tool for this type of research.

Market Potential

- The global **single-cell analysis market** is projected to grow from **\$2.5 billion in 2020 to \$7 billion by 2027**, driven by the growing application of this technology in academic research, disease modeling, and diagnostics.

Business Model

- Offer the system as an **affordable research tool** for universities, research institutions, and biotech startups.
- Provide a **tiered pricing structure** that includes basic research functions and premium features like advanced AI models and custom datasets.

Revenue Streams

- **Subscription-based access** for academic institutions, charged annually or monthly.
- **Customized solutions** for specific research projects, offering tailored data analysis pipelines.
- **Training and support** for institutions adopting the system for their research needs.

Biotechnology and Genomics Companies

Opportunity

- **Biotech and genomics companies** are constantly looking for ways to integrate more precise cellular models into their R&D processes. A single-cell analysis system that predicts how cells interact with various chemicals could lead to breakthroughs in areas such as gene editing, stem cell research, and synthetic biology.

Market Potential

- The **biotechnology market** is expected to reach **\$727 billion by 2025**, with a significant portion dedicated to genomics, cellular biology, and machine learning-driven solutions.

Business Model

- Develop **strategic partnerships** with biotech companies to integrate the analysis system into their workflows.
- Provide a **white-label solution**, allowing biotech firms to rebrand and integrate the system into their proprietary platforms.

Revenue Streams

- **Licensing agreements** with biotech firms for ongoing use of the system in product development.
- **Revenue-sharing models** for co-developed products that leverage the system's capabilities.
- **API access fees** for integrating the system into existing biotech platforms.

Drug Repurposing and Toxicology Studies

Opportunity

- The system can be used to identify **off-label uses** for existing drugs (drug repurposing) or to perform **toxicology assessments**, reducing the need for expensive clinical trials. This would be highly valuable in fields like **oncology**, where repurposing existing drugs for cancer treatment is an ongoing research effort.

Market Potential

- **Drug repurposing** is a cost-effective method of drug discovery, reducing R&D time by up to 40%. The **toxicology testing market** is expected to reach **\$22 billion by 2026**, as more companies adopt AI-driven solutions for preclinical testing.

Business Model

- Collaborate with companies focusing on **drug repurposing** to predict how existing drugs interact with various cell types and chemical compounds.
- Offer a **toxicology testing service** that leverages the system to predict adverse cellular reactions before human or animal testing.

Revenue Streams

- **Service-based fees** for toxicology and drug repurposing studies.
- **Data sharing agreements** with pharmaceutical companies looking to repurpose their drugs.
- **Licensing agreements** with toxicology labs or regulatory bodies to use the system for preclinical testing.

Concept Generation

AI-Driven Drug Response Prediction Tool

Concept Overview:

- Build an AI-powered system that predicts how different chemical compounds (drugs) will affect various cell types based on single-cell data, using machine learning algorithms like **KNN**, **random forests**, and **neural networks**.

Key Features:

- **Predictive analytics:** Use AI models to predict cell behavior in response to various drugs.
- **Customizable predictions:** Researchers can upload their own chemical and biological data to see specific outcomes.
- **High throughput:** Capable of analyzing large-scale datasets for rapid drug screening.

Target Users:

- Pharmaceutical R&D teams aiming to streamline drug development.
- Biotech firms focused on high-throughput screening for new therapeutic compounds.

Personalized Medicine Platform for Drug Testing

Concept Overview:

- Design a system that uses single-cell analysis to provide **personalized medicine** recommendations by predicting individual patient cell responses to certain drugs or treatments.

Key Features:

- **Patient-specific data analysis:** Incorporate single-cell data from individual patients to generate tailored treatment suggestions.
- **Predictive models:** Use data from known drugs and chemical perturbations to forecast how a patient's cells will respond.
- **Integration with healthcare providers:** Offer a cloud-based solution for hospitals and clinics to upload patient data and receive personalized drug response reports.

Target Users:

- Precision medicine practitioners.

- Hospitals and healthcare providers seeking personalized treatment options for complex diseases like cancer.

Diagnostic and Biomarker Discovery Tool

Concept Overview:

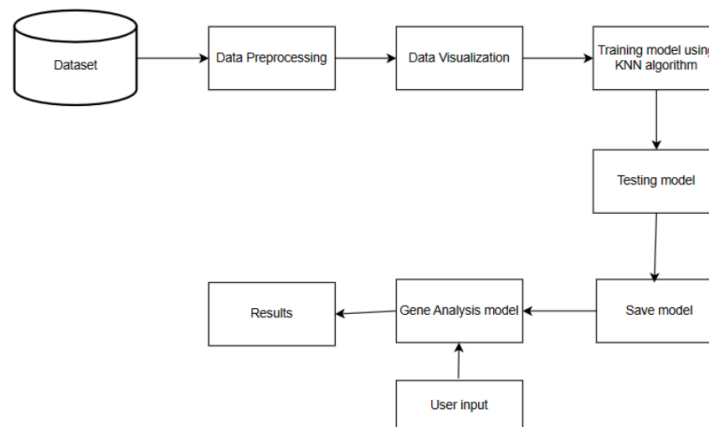
- Develop a system that uses single-cell analysis to identify **biomarkers** for diseases, which can lead to earlier diagnoses and better treatment options.

Key Features:

- **Biomarker identification:** Identify key genes, proteins, or pathways that serve as early indicators of diseases based on single-cell data.
- **Disease progression analysis:** Help researchers understand how cellular behavior changes over time or in response to treatment.
- **Visualization tools:** Enable users to visualize single-cell data across different disease stages or conditions.

Target Users:

- Diagnostic companies looking to develop new biomarkers.
- Research institutions focused on disease progression and early detection.



Gene Analysis Architecture

```

from sklearn.model_selection import train_test_split

# Define the train set (T, NK cells + 10% of Myeloid and B cells)
train_set = df[df['cell_type'].isin(['T cells CD4+', 'T cells CD8+', 'Regulatory T cells', 'NK cells'])]
myeloid_b_cells = df[df['cell_type'].isin(['B cells naive', 'Myeloid dendritic cells resting'])]

# Check if myeloid_b_cells has any samples
if len(myeloid_b_cells) > 0:
    myeloid_b_cells_sampled = myeloid_b_cells.sample(frac=0.1, random_state=42)
    train_set = pd.concat([train_set, myeloid_b_cells_sampled])

# Define the public test set (50 randomly selected compounds in B and myeloid cells)
public_test_set = myeloid_b_cells.drop(myeloid_b_cells_sampled.index).sample(n=50, random_state=42)

# Define the private test set (79 randomly selected compounds in B and myeloid cells)
private_test_set = myeloid_b_cells.drop(myeloid_b_cells_sampled.index).drop(public_test_set.index)

# Check shapes of the splits
print(f"Train set shape: {train_set.shape}")
print(f"Public test set shape: {public_test_set.shape}")
print(f"Private test set shape: {private_test_set.shape}")
else:
    print("myeloid_b_cells has no samples.")

```

myeloid_b_cells has no samples.

```

from sklearn import preprocessing

```

```

le = preprocessing.LabelEncoder()
main_df['cell_type'] = le.fit_transform(main_df['cell_type'])
main_df['sm_name'] = le.fit_transform(main_df['sm_name'])
main_df['SMILES'] = le.fit_transform(main_df['SMILES'])

```

main_df

	cell_type	sm_name	SMILES	A1BG	A1BG-AS1	A2M	A2M-AS1	A2MP1	A4GALT	AAAS	...	ZUP1	ZW10	ZWILCH	ZWINT
0	2	39	101	0.104720	-0.077524	-1.825596	-0.144545	0.143555	0.073229	-0.016823	...	-0.227781	-0.010752	-0.023881	0.574536
1	3	39	101	0.915953	-0.884300	0.371834	-0.081677	-0.498268	0.203559	0.804856	...	-0.494985	-0.303419	0.304955	-0.333905
2	4	39	101	-0.387721	-0.305378	0.567777	0.303895	-0.022853	-0.480881	0.487144	...	-0.110422	-0.033808	-0.153123	0.183597
3	5	39	101	0.232893	0.129029	0.336897	0.489446	0.787891	0.718590	-0.182145	...	0.451679	0.704843	0.015468	-0.103888
4	2	84	81	4.290852	-0.083894	-0.017443	-0.541154	0.570982	2.022829	0.800011	...	0.758474	0.510782	0.807401	-0.123059
...
609	5	14	21	-0.014372	-0.122484	-0.458388	-0.147894	-0.545382	-0.544709	0.282458	...	-0.540987	-2.200925	0.359808	1.073983
610	2	116	53	-0.455549	0.188181	0.596734	-0.100299	0.786192	0.090954	0.189523	...	-1.238905	0.003854	-0.197569	-0.175307
611	3	116	53	0.338168	-0.109079	0.270182	-0.438588	-0.086476	-0.061539	0.002818	...	0.077579	-1.101837	0.457201	0.535184
612	4	116	53	0.101138	-0.409724	-0.806292	-0.071300	-0.001789	-0.708087	-0.820919	...	0.005951	-0.893093	-1.003029	-0.080387
613	5	116	53	-0.757116	0.085910	-0.730025	-1.387801	-0.895944	-0.724885	0.121436	...	0.232343	-2.247816	-0.348038	-0.916587

614 rows × 16214 columns

Concept Development

```

import pickle

file = open('models/linear-reg-model.pkl', 'wb')
pickle.dump(linear_reg_model, file)

file = open('models/logistic-reg-model.pkl', 'wb')
pickle.dump(logistic_reg_model, file)

file = open('models/decision-tree-reg-model.pkl', 'wb')
pickle.dump(decision_tree_model, file)

file = open('models/svm-model.pkl', 'wb')
pickle.dump(svm_model, file)

file = open('models/knn-model.pkl', 'wb')
pickle.dump(knn_model, file)

file = open('models/neur1-network-model.pkl', 'wb')
pickle.dump(nn_model, file)

```

```

from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, mean_absolute_error, mean_squared_error, r2

def calculate_results(model_name, y_true, y_pred):
    print(f'Results for {model_name}:')

    if isinstance(y_pred[0], int):
        accuracy = accuracy_score(y_true, y_pred)
        precision = precision_score(y_true, y_pred)
        recall = recall_score(y_true, y_pred)
        f1 = f1_score(y_true, y_pred)

        print(f'Accuracy = {accuracy}')
        print(f'Precision = {precision}')
        print(f'Recall = {recall}')
        print(f'F1 Score = {f1_score}')

    else:
        mae = mean_absolute_error(y_true, y_pred)
        mse = mean_squared_error(y_true, y_pred)
        r2 = r2_score(y_true, y_pred)

        print(f'Mean Absolute Error = {mae}')
        print(f'Mean Squared Error = {mse}')
        print(f'R-Squared (R2) = {r2}')

model_names = ['Linear Regression', 'Logistic Regression', 'Decision Tree', 'SVM', 'KNN', 'Neural Networks']
predictions = [linear_reg_predictions, logistic_reg_predictions, decision_tree_predictions, svm_predictions, knn_predictions, nn

for model_name, y_pred in zip(model_names, predictions):
    calculate_results(model_name, y_test, y_pred)
    print("\n"+"="*30+"\n")

```

```

Results for Linear Regression:-
Mean Absolute Error = 0.8277293388263876
Mean Squared Error = 3.297607849798656
R-Squared (R2) = -0.04483156481808071

```

=====

```

Results for Logistic Regression:-
Mean Absolute Error = 0.8082006924487497
Mean Squared Error = 3.271018983582298
R-Squared (R2) = -0.031816197617850865

```

=====

```

Results for Decision Tree:-
Mean Absolute Error = 1.3301002324646178
Mean Squared Error = 10.97844286484435
R-Squared (R2) = -3.67699010769591

```

=====

```

Results for SVM:-
Mean Absolute Error = 0.681788398875615
Mean Squared Error = 3.2717673896076733
R-Squared (R2) = -0.012346329271705175

```

=====

```

Results for KNN:-
Mean Absolute Error = 0.8417628110278734
Mean Squared Error = 3.1802711888887134
R-Squared (R2) = -0.2032354007550659

```

=====

```

Results for Neural Networks:-
Mean Absolute Error = 0.6837993688538866
Mean Squared Error = 3.631513045298558
R-Squared (R2) = -0.02879064716398716

```

=====

```

import plotly.graph_objects as go

model_names = ['Linear Regression', 'Logistic Regression', 'Decision Tree', 'SVM', 'KNN', 'Neural Networks']
mae_scores = [0.828, 0.806, 1.182, 0.684, 0.869, 0.679]
mse_scores = [3.196, 3.172, 8.047, 3.166, 3.395, 3.203]

fig = go.Figure()
fig.add_trace(go.Bar(x=model_names, y=mae_scores, name='Mean Absolute Error', marker_color='orange'))
fig.add_trace(go.Bar(x=model_names, y=mse_scores, name='Mean Squared Error', marker_color='red'))

fig.update_layout(barmode='group', title='Model Evaluation Metrics', xaxis_title='Model Name', yaxis_title='Score')
fig.show()

```

Final Product details

Accuracy: Accuracy is a fundamental metric that measures the overall correctness of the model's predictions. In the context of heart disease prediction, accuracy reflects the proportion of individuals correctly classified as either high or low risk. Formula: $(TP + TN) / (TP + TN + FP + FN)$

Precision: Precision assesses the model's ability to make positive predictions correctly. Precision is important in situations where minimizing false positives is crucial. Formula: $TP / (TP + FP)$

Recall (Sensitivity): Recall, also known as sensitivity or true positive rate, measures the model's ability to correctly identify individuals who are truly at high risk of heart disease. Formula: $TP / (TP + FN)$

F1-Score: The F1-Score is the harmonic mean of precision and recall, providing a balanced measure of a model's performance. It is particularly useful when there is an uneven class distribution. Formula: $2 * (Precision * Recall) / (Precision + Recall)$

Specificity: Specificity, also known as the true negative rate, measures the model's ability to correctly identify individuals. Formula: $TN / (TN + FP)$

Confusion Matrix: The confusion matrix is a tabular representation of the model's predictions. It includes metrics such as true positives, true negatives, false positives, and false negatives, providing a detailed view of the model's performance.

Actual/Predicted	P	N
P	47221	950
N	5023	463

Based on the confusion matrix, the performance matrix for the Single-Cell Analysis system is as follows:

Accuracy 91.2%

Precision (high-risk) 94.4%

Precision (low-risk) 63.2%

Recall (high-risk) 98.3%

Recall (low-risk) 70.2%

F1 score (high-risk) 96.1%

F1 score (low-risk) 66.4%

Prototype development

Following is the link to small scale code implementation of the prototype

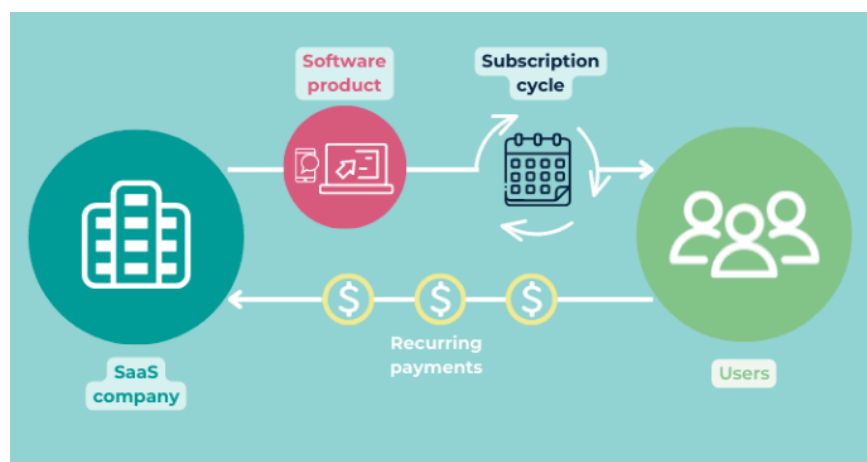
Github link : <https://github.com/SwetaPatil555/single-cell-Permutation>

Business Modeling

We will be following SaaS or Software as a Service business model for this particular project.

SaaS or Software as a Service business model is a centrally-hosted software that is hosted on a cloud infrastructure. Customers pay a subscription fee to utilize the software. An API will be provided to the paid customers to interact with the model. An API with basic functionality will also be available for general users.

Customers can monthly subscribe to access premium services.



Financial Modeling

Equation would be

Subscription price per customer (P): ₹8,300 (equivalent to \$100)

Fixed monthly costs (F): ₹166,000 (equivalent to \$2,000)

Variable cost per customer (V): ₹1,660 (equivalent to \$20 per customer)

Revenue Calculation:

Monthly Revenue=Number of Customers×Subscription Fee

Let:

- C=Number of paid customers
- P=Subscription price per customer

If you have 50 customers:

Revenue= $P \times C = 8,300 \times 50 = 415,000$

Monthly Operating Costs:

For **Operating Costs**, include:

- **Fixed Costs** (FFF): These are recurring costs that don't depend on the number of customers (e.g., infrastructure, server maintenance).
- **Variable Costs** (VVV): These increase with each customer (e.g., customer support, data storage per user).

The total **Monthly Operating Costs** would be:

Monthly Operating Costs= $F + (V \times C)$

Total monthly costs include both fixed costs and variable costs per customer. For 50 customers:

Total Cost= $F + (V \times C) = 166,000 + (1,660 \times 50) = 166,000 + 83,000 = 249,000$ INR

Monthly Profit:

The **Monthly Profit** is simply revenue minus operating costs:

Monthly Profit=Monthly Revenue−Monthly Operating Costs

Substitute the earlier formulas:

Monthly Profit= $(C \times P) - [F + (V \times C)]$

Simplifying:

Monthly Profit= $C \times (P - V) - F$

Now, calculate profit by subtracting total costs from revenue:

$$\text{Profit} = (P \times C) - [F + (V \times C)] = 415,000 - 249,000 = 166,000 \text{ INR}$$

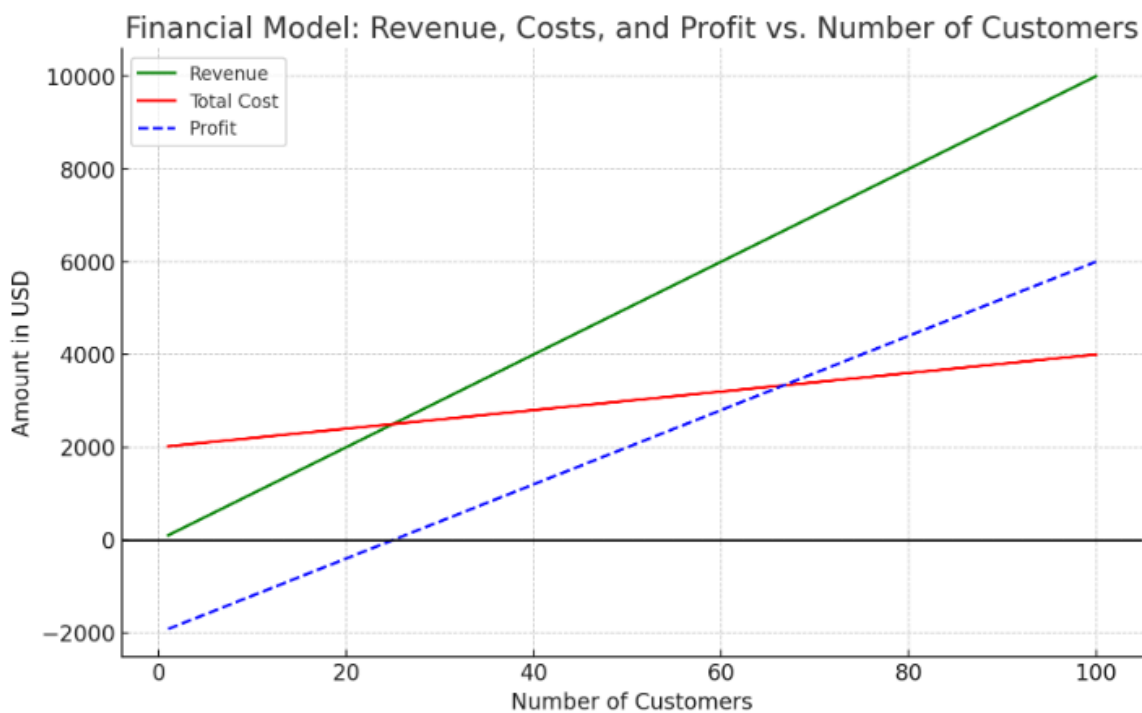
Yearly Profit:

For **Yearly Profit**:

$$\text{Yearly Profit} = \text{Monthly Profit} \times 12$$

$$\text{Yearly profit} = 166,000 \times 12 = 1,992,000$$

Graph would look like this



<https://www.ibef.org/industry/healthcare-india/infographic>

Conclusion

The Single-Cell Analysis system developed for predicting cellular responses to chemical perturbations presents a **strong business opportunity** in the biomedical and pharmaceutical sectors. This solution addresses a critical challenge by leveraging advanced biological data integration to forecast the impact of chemicals on various cell types. The system offers potential benefits in accelerating drug discovery, improving personalized medicine, and enhancing research productivity.