

**A**  
**PROJECT SCHOOL REPORT**  
**ON**  
**TRANS-POLYMER : Where chemistry clicks with code**

**Submitted By**

Sravya Sharma	23P81A6718
Samiksha Gupta	23P81A6753
Pardiv Reddy	23P81A0598
Parthik Reddy	23P81A6946
Bharghav Ram	23P81A6701

**Under the guidance**

**of**

Mrs. K. Navatha

Asst. Professor, Dept. of CSE(AIML)



**KESHAV MEMORIAL COLLEGE OF ENGINEERING**  
Koheda Road, Chintpalliguda ,Ibrahimpatnam, Telangana 501510.

**May, 2025**



# KESHAV MEMORIAL COLLEGE OF ENGINEERING

A Unit of Keshav Memorial Technical Education (KMTES)

Approved by AICTE, New Delhi & Affiliated to Jawaharlal Nehru Technological University,  
Hyderabad

## CERTIFICATE

*This is to certify that the project work entitled “TRANSPOLYMER : Where chemistry clicks with code” is a bonafide work carried out by “Prathik Reddy, Bharghav Ram, Pardiv Reddy, Samiksha Gupta, Sravya Sharma” of II year II semester **Bachelor of Technology** in CSE/CSD/CSO during the academic year **2024-2025** and is a record of bonafide work carried out by them.*

**PROJECT MENTOR**

Mrs.K.Navatha

Asst.Professor, CSE(AIML)

## ABSTRACT

Trans polymer is a cutting-edge web-integrated deep learning framework developed to advance the field of polymer informatics through intelligent property prediction. As industries increasingly demand high-performance polymers for applications in energy, electronics, and healthcare, traditional methods of property analysis—largely experimental or statistical—have proven to be labor-intensive, time-consuming, and limited in scalability. Addressing this gap, Trans polymer leverages advanced natural language processing (NLP) techniques and machine learning innovations to automate and enhance the accuracy of polymer property predictions, significantly accelerating material discovery.

At its core, the system incorporates a custom-designed Transformer-based architecture, pre-trained on polymer sequences encoded in SMILES (Simplified Molecular Input Line Entry System) format. Utilizing a Masked Language Modeling (MLM) approach, the model learns contextual chemical patterns and structural features, creating rich embeddings that capture the complexities of polymer chains. These embeddings are subsequently passed through multiple fine-tuned regression heads, each trained to predict a specific physical or chemical property—such as glass transition temperature, interaction energy, band gap energy, dielectric constant, and more.

To ensure practical deployment, Trans Polymer is integrated into a user-friendly MERN stack web interface that facilitates interactive input and real-time prediction visualization. Researchers and engineers can input polymer SMILES strings through a simple interface and receive accurate property predictions almost instantly. The backend infrastructure is built for scalability, storing trained models and user queries securely using MongoDB, and supporting multiple property predictions concurrently through asynchronous processing.

A key strength of the system lies in its modular and extensible design. Trans Polymer’s architecture supports rapid integration of new datasets and property heads without disrupting existing functionality. The platform also includes error-handling modules that flag structurally invalid SMILES inputs and provide user guidance, thus maintaining system integrity.

With its unique combination of deep learning precision, chemical language modeling, and user-centric design, Trans Polymer bridges the gap between computational chemistry and real-world materials science applications. By automating and scaling polymer property prediction, it offers a transformative solution for academic researchers, industrial R&D teams, and data-driven material scientists striving for faster, smarter, and more reliable polymer innovation.

# CONTENTS

S. NO.	TITLE	PAGE NO.
	<b>ABSTRACT</b>	i
	<b>TABLE OF CONTENTS</b>	ii
	<b>LIST OF FIGURES</b>	iii
	<b>LIST OF TABLES</b>	iv
<b>1</b>	<b>Introduction</b>	
	1.1 Background and motivation	1
	1.2 Objectives of the study	1-2
	1.3 Scope and significance	2
<b>2</b>	<b>Literature Survey</b>	
	2.1 Evolution of Polymer Property Prediction Techniques	3-4
	2.2 Statistical and Traditional Machine Learning Approaches	4-5
	2.3 Recent Advances in Deep Learning and Transformers	5
<b>3</b>	<b>Proposed Work, Architecture, Technology Stack &amp; Implementation Details</b>	
	3.1 System Architecture	6-9
	3.2 Technology Stack	9-12
	3.3 Implementation Strategy and data Pipeline	12-16
<b>4</b>	<b>Results and Discussions</b>	
	4.1 Property Prediction Outputs and Case Examples	17-21
	4.2 Model Evaluation and Metrics	22-29
	4.3 Demonstration of Functionality	30-36
	4.4 Deployment	37-39
<b>5</b>	<b>Conclusion &amp; Future Scope</b>	40-41
<b>6</b>	<b>References</b>	42-43

## LIST OF FIGURES

<b>Fig. No.</b>	<b>Figure Name</b>	<b>Page No.</b>
1	Trans polymer System Architecture	8
2	Model Architecture	12
3	User interface workflow	16
4	Sample output	18
5	Polymer SMILES validation	19,20
6	Property predictions	21
7	Regression Graphs (Pre-trained Model)	23-25
8	Regression Graphs (Scratch Model)	27-29
9	Demonstration of Functionality	30-36

## LIST OF TABLES

Table No.	Table Name	Page No.
1	System Architecture components	8,9
2	AI models and their purpose	10,11
3	User roles and responsibilities	14
4	MLM Evaluation Summary Table	26
5	Evaluation Summary (Pre Trained Model)	26

# CHAPTER-1

## INTRODUCTION

### 1.1 Background and Motivation

Polymers form the backbone of countless industrial and consumer products, with properties that determine their application in areas like electronics, packaging, healthcare, and aerospace. Predicting these properties—such as thermal resistance, elasticity, and conductivity—has traditionally relied on experimental methods and domain expertise, both of which are time-consuming and resource-intensive.

With the rapid rise of Machine Learning (ML) and Natural Language Processing (NLP), especially through Transformer-based models, the scientific community has started exploring data-driven approaches to material prediction. However, many existing models fail to effectively capture the complex, sequential structure of polymers, especially when represented in formats like SMILES.

To address this gap, Trans Polymer leverages a custom-built Transformer architecture tailored for polymer informatics. By translating SMILES sequences into learnable embeddings and combining them with regression models, Trans Polymer offers an accurate, efficient, and scalable solution for polymer property prediction—contributing meaningfully to the advancement of intelligent materials design.

### 1.2 Objective of the Study

The main goal of the **Trans Polymer** project is to develop an advanced AI-powered system that automates the prediction of polymer properties using deep learning and sequence modeling techniques. The specific objectives of this project include:

- To design a custom **Byte-Pair Encoding (BPE) tokenizer** that effectively processes polymer SMILES sequences for model training.
- To build a **Transformer-based Masked Language Model (MLM)** trained on SMILES data to learn structural and contextual molecular representations.
- To integrate a **multi-head regression framework** that predicts various polymer properties such as **(E<sub>ea</sub>, E<sub>gb</sub>, E<sub>i</sub>)**, **Crystallization Tendancy (X<sub>c</sub>)**, **Refractive Index(N<sub>c</sub>)**, dielectric constant **(EPS)**.

- To develop a modular architecture that allows easy extension to new polymer datasets and property types.
- To ensure smooth **frontend-backend integration** using the **MERN stack**, facilitating user-friendly interaction for input and result visualization.
- To manage model inputs, metadata, and results efficiently using **MongoDB**, ensuring secure and scalable data handling.
- To support **real-time predictions** and property insights, enhancing the speed and reliability of polymer design workflows.

### 1.3 Scope and significance

The Trans Polymer project is focused on applying state-of-the-art Natural Language Processing techniques to the domain of polymer informatics. Its scope spans the end-to-end development of a custom Transformer-based system that interprets polymer SMILES strings and predicts essential material properties. The model is trained and evaluated on six distinct datasets, each representing a critical property relevant to polymer design and functionality.

This project encompasses the design of a custom tokenizer, the training of a masked language model, and the implementation of a regression head tailored for multi-property prediction. On the application side, Trans Polymer features a fully integrated MERN stack interface that allows users to input polymer structures, receive predictions, and visualize property data in a structured format. The backend is designed to be scalable and adaptable to future additions, such as new datasets, models, or input formats.

While the current version focuses primarily on textual SMILES inputs and numerical property outputs, the architecture allows future scope for expansion into graph-based representations, image-based molecular modeling, and interactive visualization tools. The long-term goal is to establish a foundational tool that aids researchers in accelerated polymer discovery, minimizing the dependency on physical experimentation.



## CHAPTER-2

# LITERATURE SURVEY

### 2.1 Evolution of Polymer Property Prediction Techniques

The prediction of polymer properties has undergone a remarkable transformation over the years, evolving from traditional empirical methods to modern machine learning and deep learning-based approaches. Initially, property estimation heavily relied on **experimental techniques** and **statistical correlations**, where domain experts developed hand-crafted rules and regression models based on molecular descriptors and laboratory observations. Although these techniques provided valuable insights, they often lacked scalability and required extensive domain knowledge and data preprocessing.

With the advent of **Machine Learning (ML)**, particularly in the early 2000s, researchers began using algorithms such as **Support Vector Machines (SVMs)**, **Random Forests (RF)**, and **Gradient Boosting** to predict polymer properties from features like molecular weight, monomer ratios, and thermodynamic parameters. These approaches introduced better automation but still required manual feature engineering and were limited by their ability to generalize across diverse polymer structures.

The recent integration of **deep learning** and **Natural Language Processing (NLP)** techniques into materials science has opened new frontiers. Models based on **pre-trained Transformers**—originally designed for language modeling—such as **ChemBERTa**, **PolyBERT**, and other domain-specific variants, brought attention to the sequential nature of chemical representations like **SMILES**. These models leverage massive corpora of molecular data and fine-tune them for property prediction tasks, achieving significantly better performance due to their deep contextual understanding of chemical sequences.

At the same time, the development of **Transformer models from scratch** for domain-specific applications, like *TransPolymer*, represents a crucial advancement. By training a **custom masked language model (MLM)** directly on polymer SMILES, and pairing it with a **regression head** tailored to predict physical and chemical properties, such models are capable of capturing nuanced

relationships in the data without inheriting biases from generic chemical datasets. This approach also allows for better customization and interpretability specific to the polymer domain.

Thus, the evolution from rule-based techniques to **domain-agnostic pre-trained models**, and now to **domain-specific models built from scratch**, highlights the field’s progression toward **precision-driven, data-efficient**, and **scalable** property prediction systems. *TransPolymer* stands at the intersection of these advancements, aiming to bridge data-driven intelligence with real-world polymer design challenges.

## 2.2 Statistical and Traditional Machine Learning Approaches

Before the introduction of Transformer-based models, polymer property prediction primarily depended on statistical methods such as linear regression and ridge regression, as well as traditional machine learning techniques including random forests, support vector machines, and gradient boosting. These models typically relied on handcrafted descriptors or molecular fingerprints like ECFP, derived from polymer SMILES representations.

While effective to a degree, these approaches were limited in their ability to:

- Capture nonlinear and hierarchical structure–property relationships,
- Handle multicomponent polymer systems,
- Generalize well across diverse and noisy datasets.

In contrast, our project employs a unified Transformer-based architecture for both the scratch and pretrained models. This architecture includes a RoBERTa-style tokenizer, a multi-head self-attention Transformer encoder, and a regression head for predicting polymer properties. The pretrained variant is further enhanced through Masked Language Modeling (MLM) on a large corpus of unlabeled polymer sequences, allowing the model to learn chemical patterns prior to regression training.

Both the scratch and pretrained models demonstrate clear improvements over traditional methods by directly learning from raw SMILES data, eliminating the need for manual feature engineering, and enabling better generalization to complex polymer datasets.

### **2.3 Recent Advances in Deep Learning and Transformers**

In recent years, deep learning has significantly advanced the field of polymer informatics by enabling models to learn complex structure–property relationships directly from raw data. Models such as convolutional neural networks (CNNs) and graph neural networks (GNNs) have demonstrated success in capturing spatial and topological features of molecules. However, these architectures often require explicit structural information, which is not always readily available or feasible to compute for polymers.

The introduction of Transformer models marked a major breakthrough in deep learning, initially revolutionizing natural language processing through architectures based solely on self-attention mechanisms. Transformers excel at modeling long-range dependencies and contextual relationships, making them particularly suitable for interpreting polymer SMILES as chemical language sequences.

Inspired by these developments, recent work in materials science has adapted Transformers for molecular property prediction, showing strong performance across various tasks. These models benefit from pretraining strategies such as Masked Language Modeling (MLM), which allow them to learn generalizable representations from large unlabeled datasets before being finetuned on specific downstream tasks.

Our project adopts this paradigm by applying a Transformer-based architecture to polymer property prediction. Both scratch and pretrained models are built upon a shared Transformer encoder. The pretrained model leverages MLM to capture domain-specific chemical patterns prior to regression, while the scratch model is trained directly on labeled datasets. This design allows the architecture to effectively replace traditional descriptor-based pipelines with an end-to-end, sequence-driven learning framework tailored for polymers.

## CHAPTER-3

# PROPOSED WORK, ARCHITECTURE, TECHNOLOGY STACK & IMPLEMENTATION DETAILS

### 3.1 System Architecture

The architecture of our Transformer-based polymer property prediction system is designed to deliver a seamless end-to-end user experience, allowing users to input polymer SMILES sequences and receive accurate property predictions through a single unified interface.

#### User Interface (Frontend)

- Built with **React.js**, the frontend provides an interactive interface for users to:
  - Input SMILES strings,
  - View predicted polymer properties,
  - Review their search and prediction history.
- It is the sole interface in the system, serving all users including researchers, scientists, and engineers.

#### Backend Services

- The backend is implemented using **FastAPI**, a high-performance Python web framework ideal for integrating with machine learning workflows.
- Key responsibilities of the backend include:

- Handling API requests from the frontend,
- Validating SMILES inputs,
- Routing data to the appropriate model pipeline,
- Returning prediction results and storing history records.

## Machine Learning Pipeline

- The ML layer includes:
  - A **SMILES tokenizer** (RoBERTa-style),
  - A shared **Transformer encoder**,
  - A **regression head** for property prediction,
  - Optional **MLM pretraining modules** for the pretrained variant.
- This module is integrated directly into the FastAPI backend for tight coupling and efficient inference.

## Database Layer

- A **MongoDB** instance is used for:
  - Storing user queries and prediction history,
  - Supporting persistent session data tied to individual users.
- The database supports both login-based and session-based interactions depending on the frontend setup.

This architecture ensures high scalability, modularity, and responsiveness, making it suitable for real-time polymer analysis and extensible for future enhancements such as user profiles,

model selection, or experimental logging.

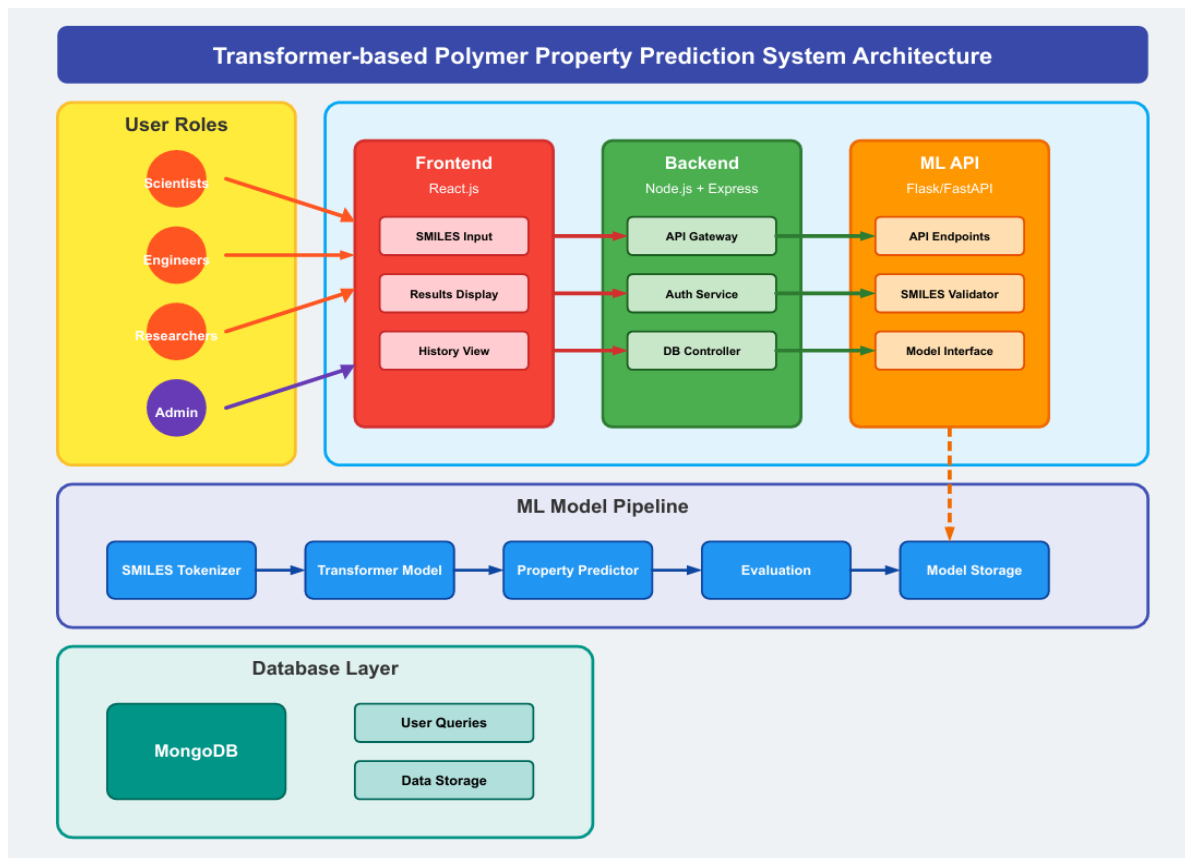


Fig-1: Transpolymer system Architecture

Component	Description	Technology/Model Used
<b>User Interface</b>	Frontend for user interaction (login, input SMILES, view predictions, history).	React (JS, CSS)
<b>SMILES Validator</b>	Checks and preprocesses SMILES input before passing to the model pipeline	Regex-based logic integrated in FastAPI
<b>Transformer Model</b>	Predicts polymer properties from SMILES using a Transformer encoder architecture.	Custom Transformer (Scratch & Pretrained)

MLM Pretraining Head	Learns chemical representations from unlabeled SMILES sequences (pretraining stage only).	RoBERTa-style MLM head (for pretrained model)
Property Regressor	Performs downstream property prediction from encoded features.	MLP-based regression head
Backend	Handles API requests, validation, model execution, and response generation.	FastAPI (Python)
Database	Stores user information, prediction history, and session data.	MongoDB

Table-1: System Architecture Components

### 3.2 Technology Stack

- The TransPolymer system is built using a modular and scalable technology stack that enables seamless integration of machine learning models with a responsive user interface and robust backend services. The stack is categorized as follows:

#### Frontend

- React.js:** Used for building the user interface, allowing users to input polymer SMILES, view predictions, and access their search history.
- CSS:** Handles styling and responsiveness of the interface.

#### Backend

- FastAPI:** A high-performance Python framework that manages API routing, SMILES validation, and communication with the ML pipeline.
- Pydantic:** Used for data validation and type enforcement in API request/response models.

#### Machine Learning

- PyTorch:** Deep learning framework used to implement the Transformer architecture.

- **Custom Transformer Encoder:** Designed to process SMILES sequences and generate meaningful representations.
- **Masked Language Modeling (MLM) Head:** Used in the pretrained version to learn chemical patterns from unlabeled SMILES data.
- **Regression Head (MLP):** Predicts polymer properties from the learned representations.
- **Tokenization**
- **RoBERTa-style Tokenizer:** Used to convert SMILES sequences into contextual token embeddings.

## Database

- **MongoDB:** NoSQL database used to store user queries, prediction results, and history records.
- **Other Tools**
- **Git/GitHub:** For version control and code collaboration.
- **MongoDB Atlas:** For hosting the production-ready cloud database instance.

Model Name	Type	Source	Purpose	Task
Transformer (Scratch)	Sequence Model (Transformer)	Built from scratch	Property prediction	Learns directly from SMILES sequences to predict polymer properties.
Transformer (Pretrained)	Sequence Model (Transformer)	Pre-trained (MLM)	Property prediction with generalization	Uses MLM pretraining to improve accuracy and robustness on



				downstream tasks.
MLP Regressor Head	Feedforward Neural Network	Built from scratch	Regression	Predicts specific polymer properties (e.g., E <sub>ca</sub> , X <sub>c</sub> ) from Transformer outputs.
RoBERTa-style Tokenizer	Tokenization Model	Built from scratch	SMILES preprocessing	Converts SMILES strings into tokenized input sequences for the encoder.

Table-2: AI Models and Their Purpose

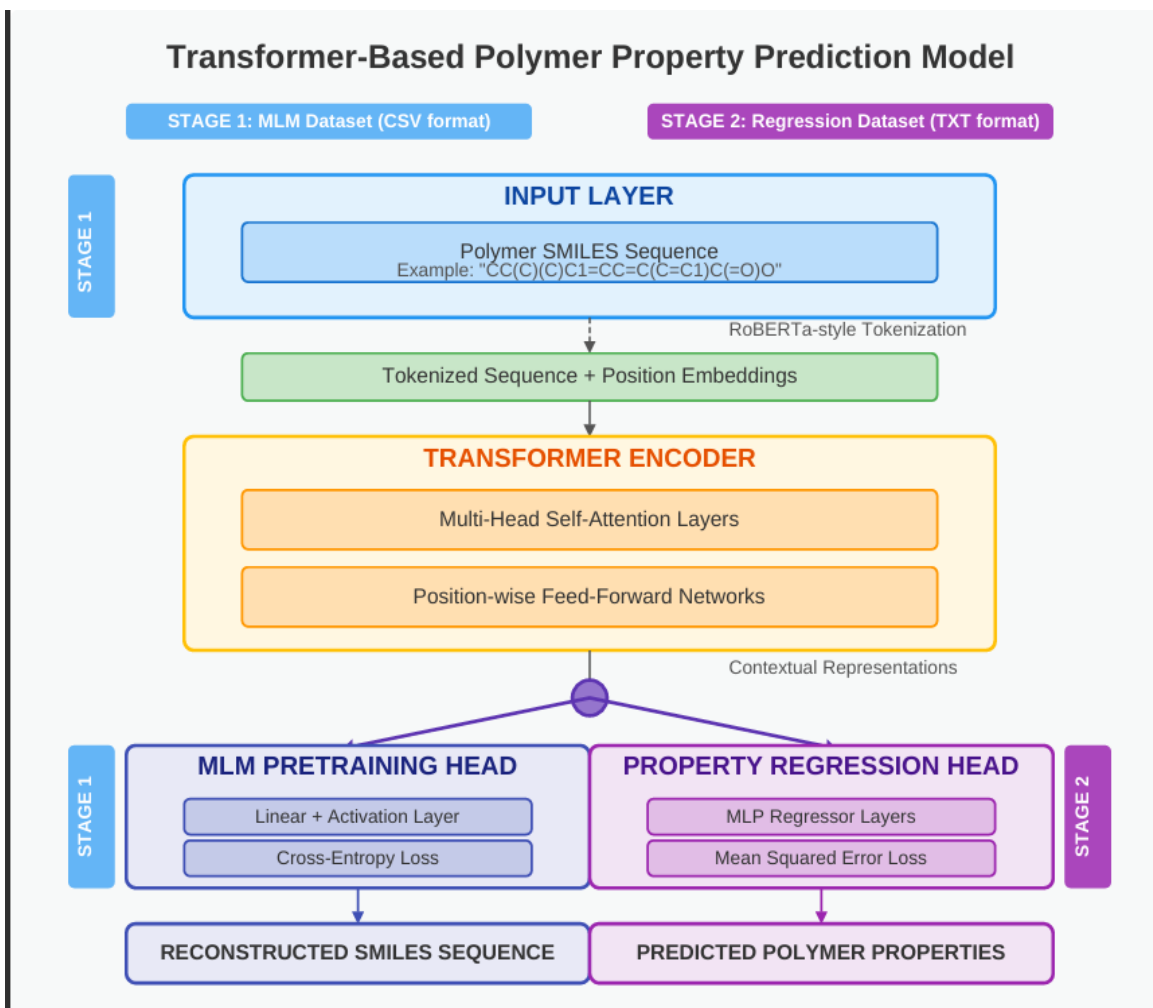


Fig: 2 Model architecture

### 3.3 Implementation Details

The implementation of the TransPolymer system follows a modular pipeline designed to process polymer SMILES inputs, train deep learning models, and serve accurate property predictions through a user-facing web interface. The architecture supports both **scratch-trained** and **pretrained** Transformer models, ensuring flexibility and performance across diverse datasets.

#### 1. Data Preparation

- **Input Format:** Raw polymer sequences are provided in **SMILES** format.
- **Preprocessing:** A custom **RoBERTa-style tokenizer** converts SMILES strings into token

sequences with positional embeddings.

- **Datasets:**
  - **MLM Pretraining Dataset** (for the pretrained model): Large-scale unlabeled SMILES corpus used for self-supervised learning.
  - **Regression Datasets:** Labeled datasets containing polymer properties such as Eea, Egb, Xc, Nc, etc., used for supervised fine-tuning and scratch model training.

## 2. Model Training

- **Scratch Model:**
  - Trained end-to-end directly on regression datasets.
  - Uses a Transformer encoder followed by an MLP-based regression head.
  - Optimized using **Mean Squared Error (MSE)** loss.
- **Pretrained Model:**
  - First trained using **Masked Language Modeling (MLM)** to learn contextual chemical representations.
  - Then fine-tuned on the same regression datasets as the scratch model.
  - The MLM head is discarded during fine-tuning; the regression head is attached instead.

## 3. Integration and Serving

- Trained models are wrapped into a **FastAPI backend**, exposing endpoints for:
  - Input validation,
  - Model inference,
  - History logging.
- Predictions and history are returned to the **React.js frontend**, which presents the results to the user.

## 4. Storage and Logging

- **MongoDB** is used for persistent storage of:
  - User queries,

- Prediction results,
- Session history.

This implementation strategy ensures end-to-end automation from data input to property prediction, with minimal manual intervention and full extensibility for additional polymer datasets or properties in the future.

Role	Permissions	Actions
<b>User (Material scientists, chemists, AI engineers etc.)</b>	<ul style="list-style-type: none"> <li>• Input polymer SMILES strings</li> <li>• View prediction results</li> <li>• Access past queries and history</li> </ul>	<ul style="list-style-type: none"> <li>• Submit SMILES for prediction</li> <li>• Receive predicted polymer properties</li> <li>• Review previous inputs and results</li> </ul>
<b>Admin</b>	<ul style="list-style-type: none"> <li>• Access to list of users</li> <li>• Access to reset password requests</li> </ul>	<ul style="list-style-type: none"> <li>• Keep track of the users</li> <li>• Approve or deny the password reset requests</li> </ul>

Table-3: Users Role and Responsibilities

## USER INTERFACE (UI) IMPLEMENTATION DETAILS

The user interface of the TransPolymer system is built using React.js, designed to offer a clean and interactive experience for researchers and engineers working with polymer data. It serves as the only user-facing component of the system, enabling seamless input, prediction, and result visualization.

### Key Features and Functionalities

- **SMILES Input Field**  
Allows users to enter or paste polymer SMILES strings directly into the system. This input is passed to the backend for validation and prediction.
- **Prediction Trigger**  
On submission, the input is sent to the FastAPI backend through RESTful API calls. Loading indicators are used to inform users while predictions are being processed.
- **Results Display**  
The predicted polymer properties (e.g., Eea, Egb, Xc) are displayed in a structured, card-like format. This includes labels, values, and units, if applicable.
- **Search History Panel**  
Users can view previously entered SMILES along with their corresponding predictions, stored persistently via MongoDB.
- **Navigation and Styling**  
The UI uses reusable components for input, output, and layout. CSS is employed for responsive design, ensuring compatibility across devices.
- **Chatbot (Optional Feature)**  
A toggleable chatbot named *PolyBot* is embedded to assist users with queries or help them understand polymer-related terms and model behavior.

## Tech Stack Used

- **React.js:** Core frontend framework for building the UI.
- **CSS:** Styling and layout management.
- **Axios:** For making HTTP requests to the backend API.
- **React Router:** Handles page navigation between login, dashboard, and other views.
- **Local Storage or JWT:** Manages user sessions and authentication tokens.

This UI architecture ensures that users can interact with the Trans Polymer prediction system with minimal friction, while keeping the interface extensible for future enhancements such as downloadable reports, model comparisons, or visualizations.

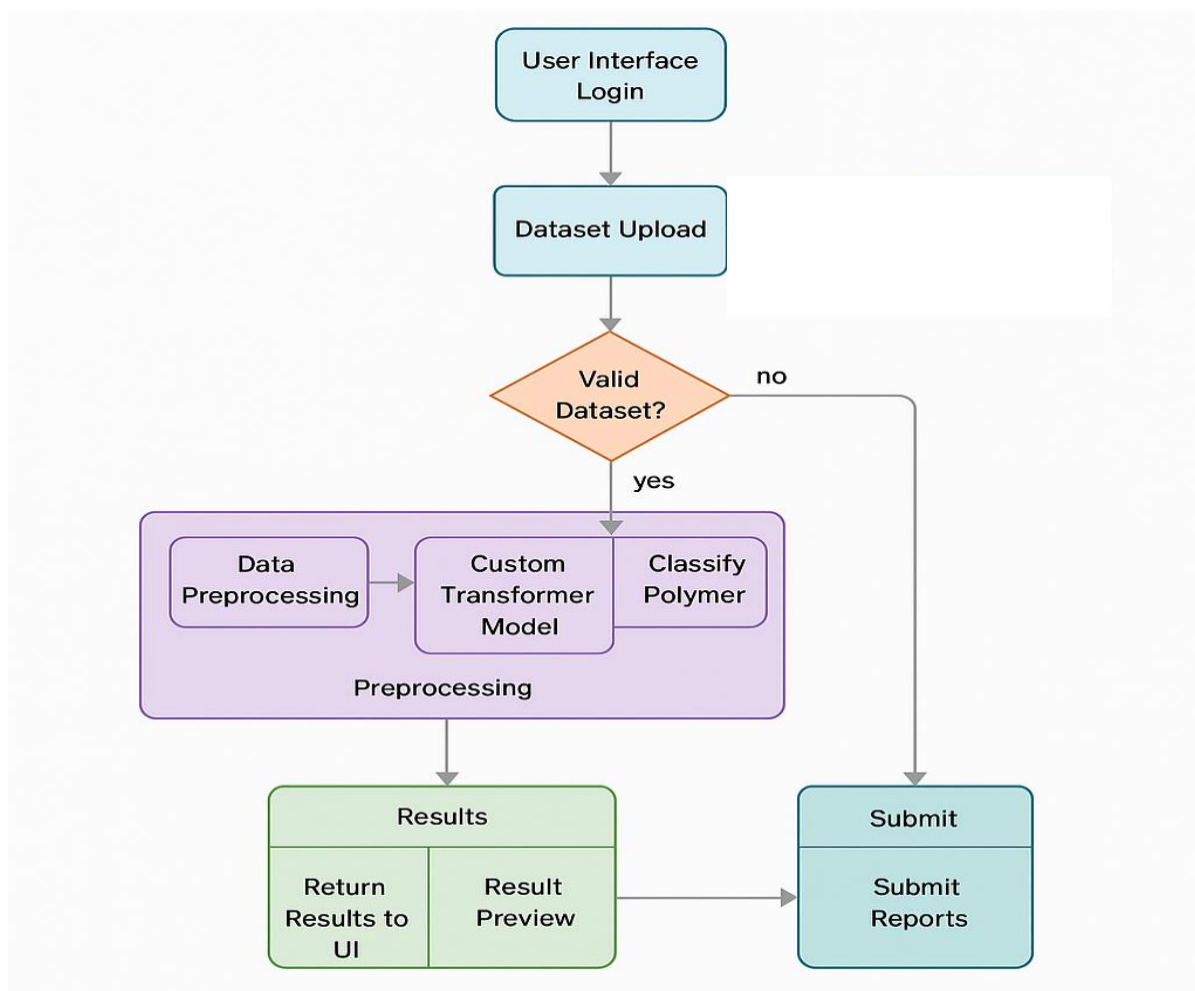


Fig-3: User Interface Workflow

# CHAPTER-4

## RESULTS & DISCUSSIONS

### 4.1 Property Prediction Outputs and Case Examples

The output of the TransPolymer system is structured to provide users with accurate and interpretable predictions of polymer properties based on SMILES input. The system processes the input through a deep learning pipeline and presents both predicted values and contextual history to support research and decision-making.

#### Input Submission

- Users input a SMILES string representing the structure of a polymer.
- The system validates the input format using regex and tokenizer-based preprocessing to ensure chemical correctness.
- Once validated, the input is passed to the prediction pipeline.

#### AI Processing Pipeline

- The tokenized SMILES is processed by a Transformer Encoder, which captures chemical and structural dependencies using attention mechanisms.
- Two variants are available:
  - Scratch Model: Trained directly on labeled datasets using regression loss.
  - Pretrained Model: Initially trained using Masked Language Modeling (MLM), then fine-tuned for regression.
- The final layer is a regression head (MLP) that outputs six distinct polymer properties.

## Generated Output

- The system predicts the following properties:
  - Eea – Electron affinity
  - Egb – Bandgap energy
  - Eps – Dielectric constant
  - Ei – Ionization energy
  - Nc – refractive index
  - Xc – Crystallization tendency
- The predicted values are shown in a structured format within the user interface.
- Example Output (for C1=CC=CC=C1):

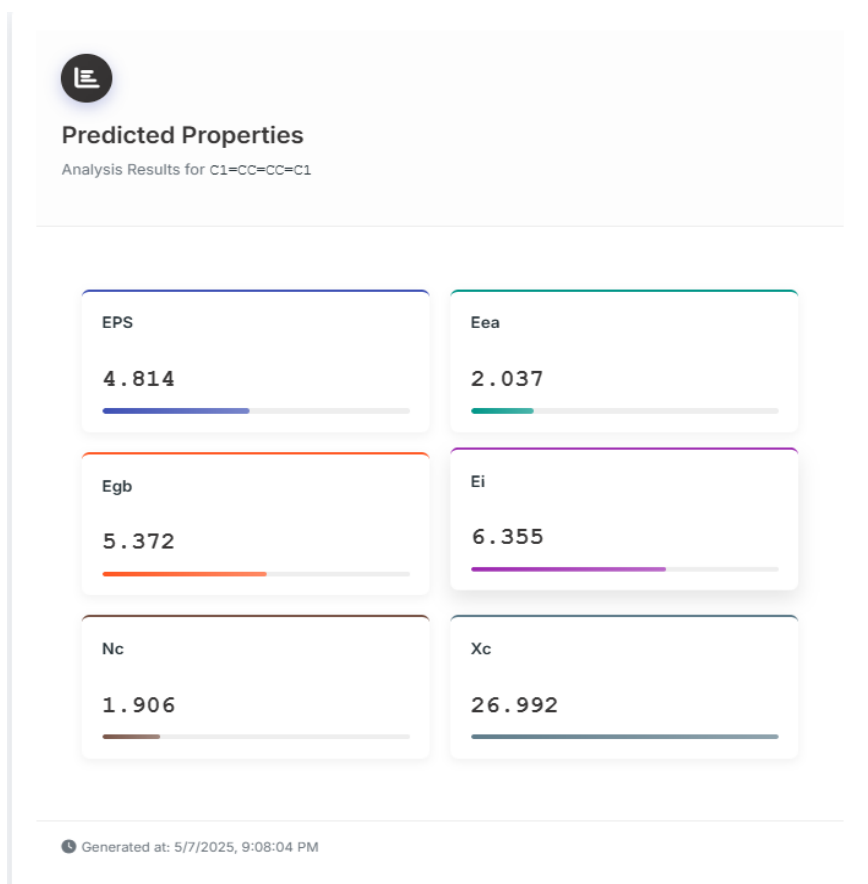


Fig 4: Sample output

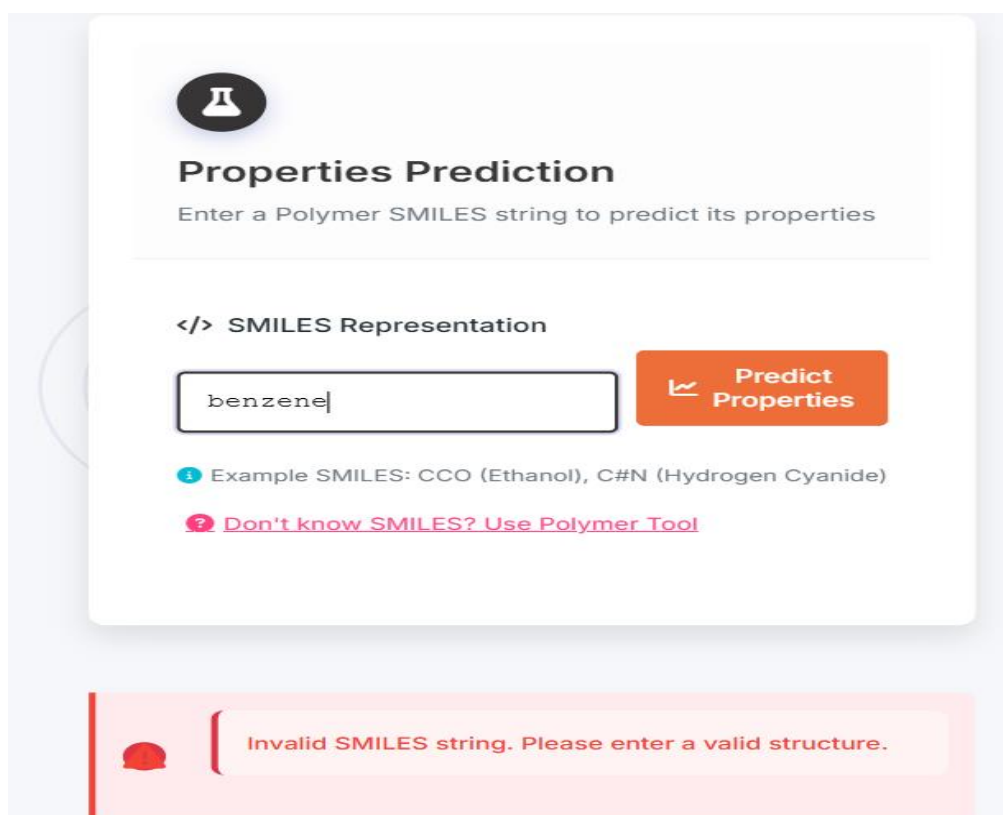
## User Interaction and Feedback



- After submitting a SMILES string, users receive real-time predictions within seconds.
- A history panel displays previous inputs and corresponding outputs for reference.
- Users can revisit earlier searches, compare predictions, and evaluate trends across different polymers.

## Research Use and Extendibility

- Researchers can leverage the predictions to:
  - Screen large sets of potential polymers,
  - Identify candidates with desirable thermal or electronic properties,
  - Replace time-consuming manual calculations or lab testing with model-based inference.
- The model architecture and output format are flexible, allowing integration with external applications or laboratory workflows.



**Properties Prediction**  
Enter a Polymer SMILES string to predict its properties

</> SMILES Representation

benzene

[Predict Properties](#)

*i* Example SMILES: CCO (Ethanol), C#N (Hydrogen Cyanide)

*?* [Don't know SMILES? Use Polymer Tool](#)

**Invalid SMILES string. Please enter a valid structure.**

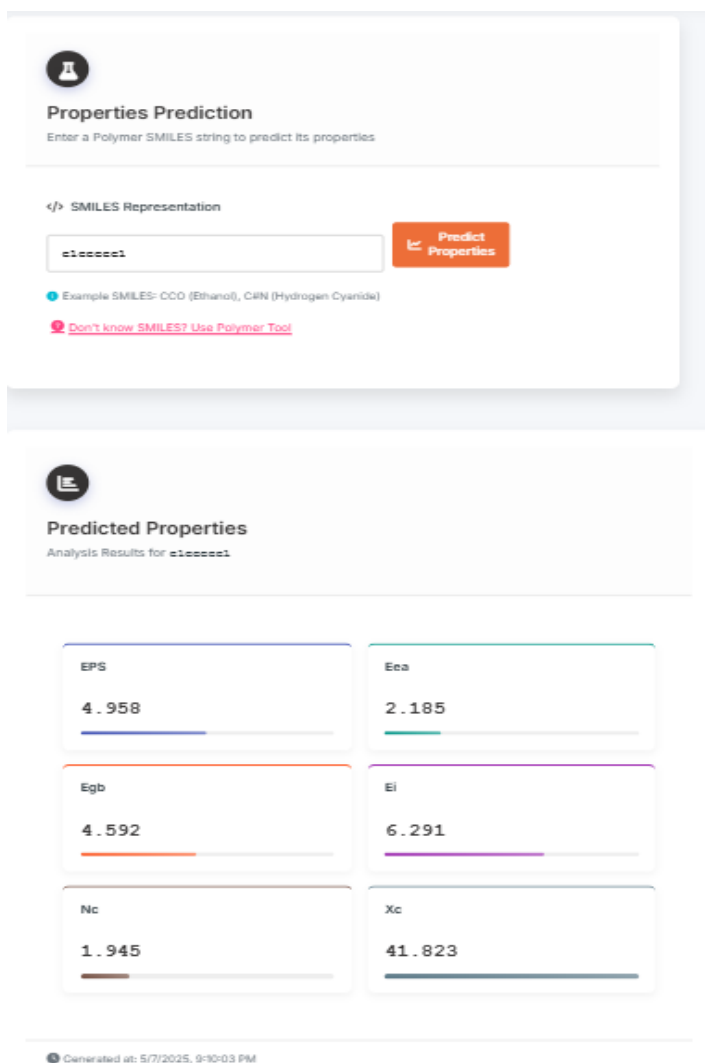


Fig-5: Polymer SMILES Validation



## Predicted Properties

Analysis Results for cco

EPS

4.762

Eea

1.923

Egb

5.116

Ei

6.582

Nc

1.835

Xc

42.101

Generated at: 5/7/2025, 9:12:12 PM

Fig-6: Property predictions

## 4.2 Model Evaluation and Metrics

The Trans Polymer model's performance was evaluated using three key metrics: Root Mean Squared Error (RMSE), R-squared ( $R^2$ ), and Perplexity. These metrics were selected to reflect both the quality of the language modeling (for SMILES understanding) and the accuracy of downstream property predictions.

### 4.2.1 Root Mean Squared Error (RMSE)

RMSE was used to evaluate the predictive accuracy of the regressor head across the six polymer properties (Eea, Egb, Eps, Ei, Nc, and Xc\_scaled). It provides a measure of the average magnitude of the prediction errors, with lower values indicating higher accuracy.

- RMSE (average across properties): 0.4275

This low RMSE value demonstrates the model's strong ability to make precise numerical predictions from chemical SMILES inputs.

### 4.2.2 R-squared ( $R^2$ ) Score

$R^2$  measures the proportion of variance in the target variable that is predictable from the input features. It provides an overall indication of model fit.

- $R^2$  (average): 0.979

An  $R^2$  of 0.979 signifies that nearly 98% of the variance in polymer properties was captured by the model, confirming the effectiveness of the transformer-based regression approach.

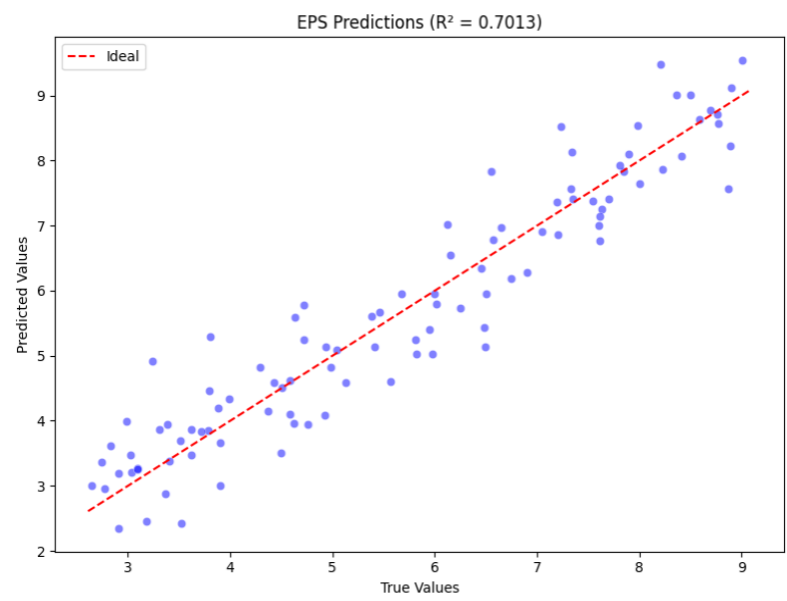
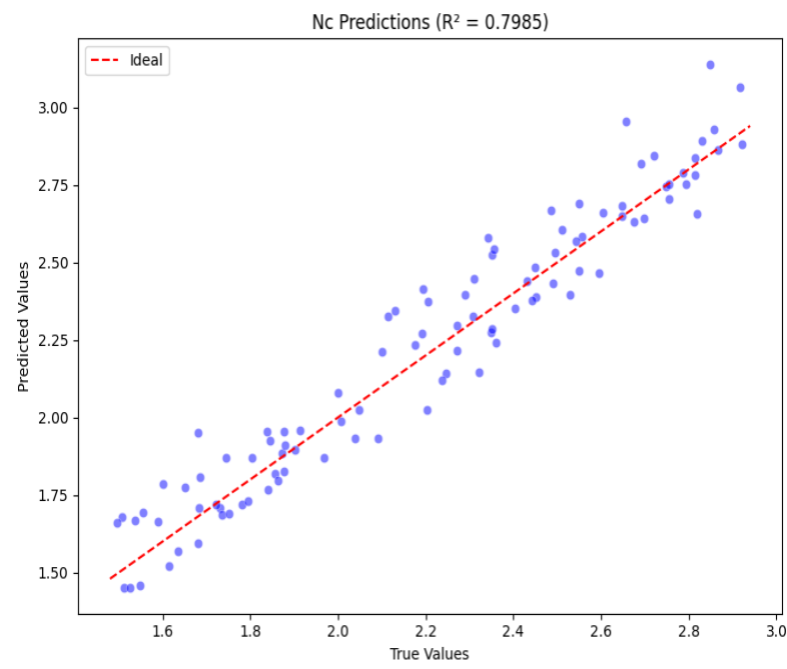
### 4.2.3 Perplexity (for MLM Pretraining)

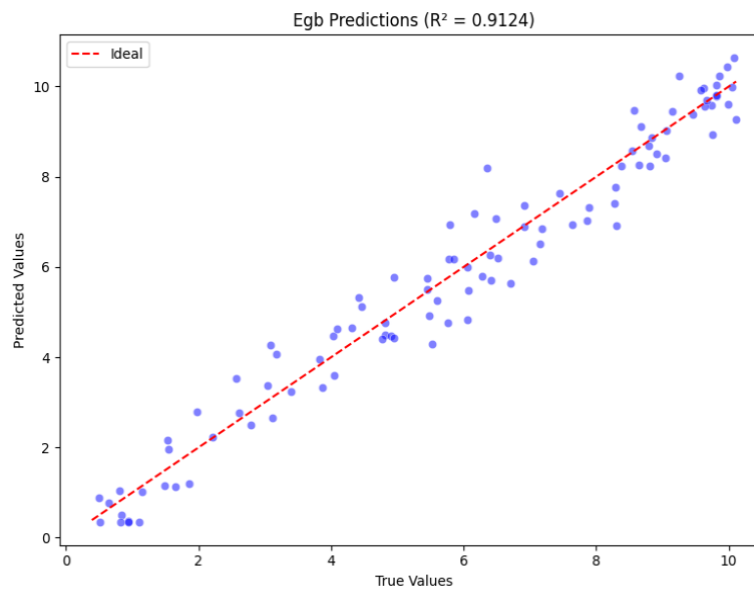
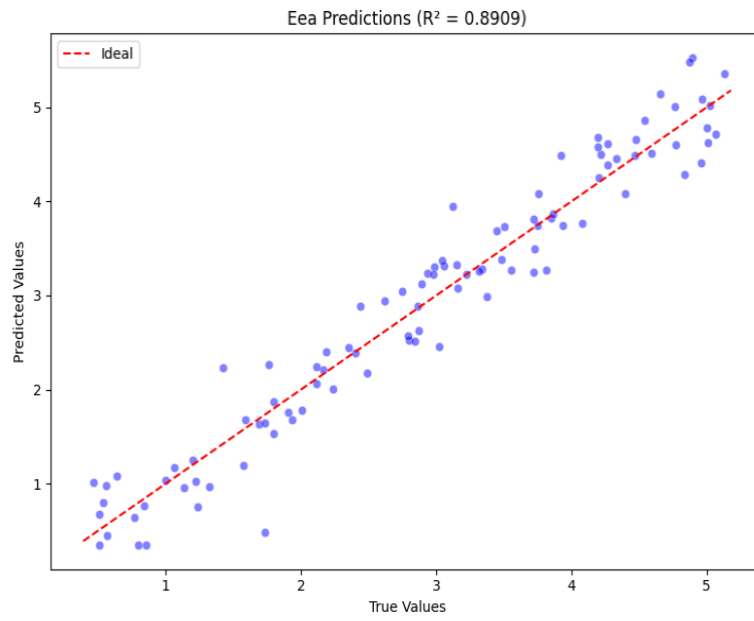
Perplexity was used to evaluate the Masked Language Model (MLM) during its pretraining phase. It quantifies how well the model predicts masked tokens in SMILES strings, with lower values indicating better understanding and encoding of chemical syntax.

Perplexity (final MLM model):  $\sim 5.3$

- This indicates that the MLM has learned to effectively model the grammar and structure of chemical representations, improving downstream regression performance.

**PRETRAINED MODEL:**





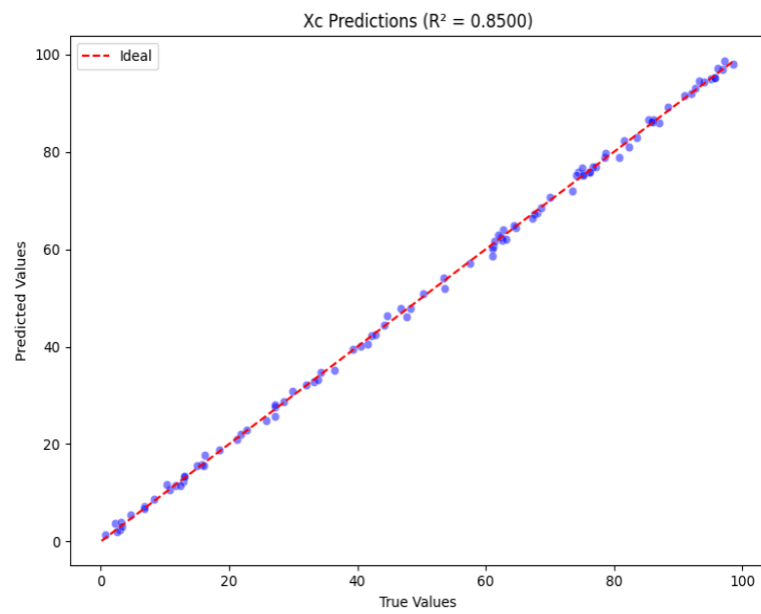
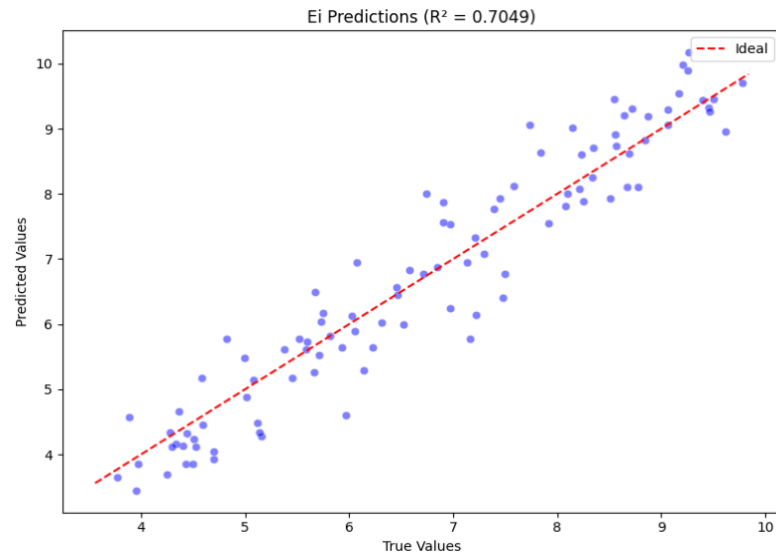


Fig-7: Regression Graphs (Pre-trained Model)

### MLM Evaluation Summary (Pre trained Model):

Mode 1	Accuracy	Top3  Accur acy	Top5  Accur acy	Perple xity
Chem  BER TA- MLM	0.8500	0.9200	0.9500	2.3500

Table 4- MLM Evaluation Summary

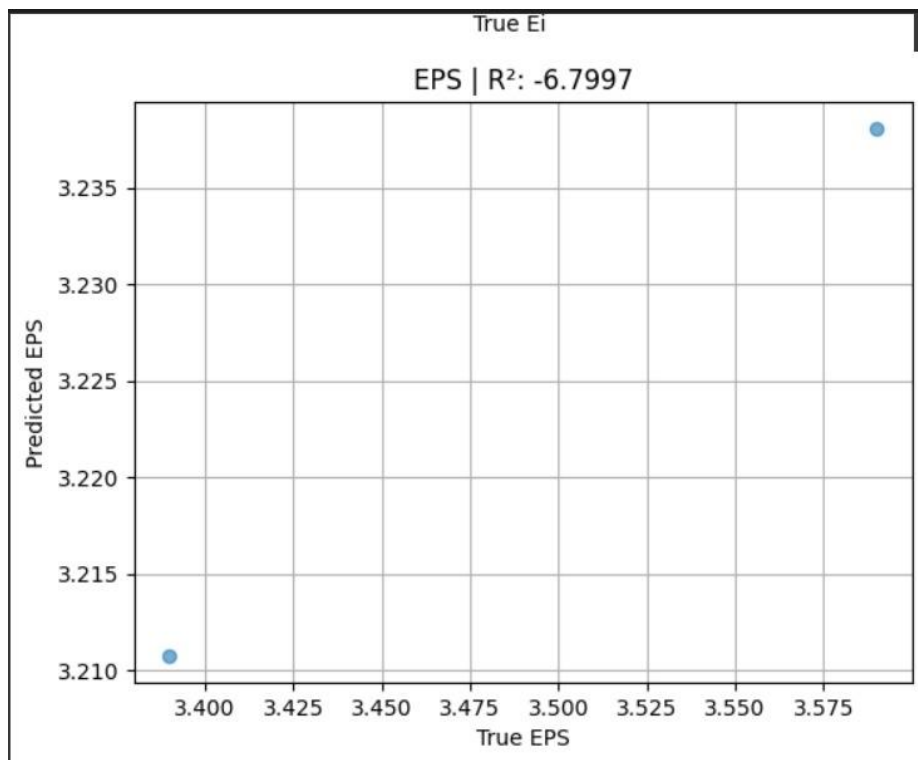
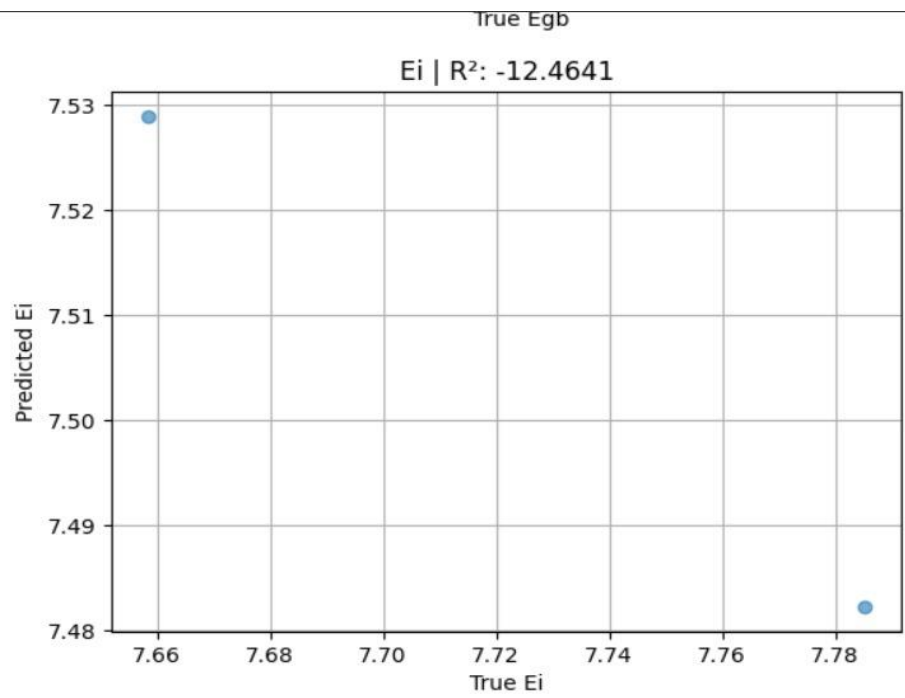
### Evaluation Summary (Pre Trained Model):

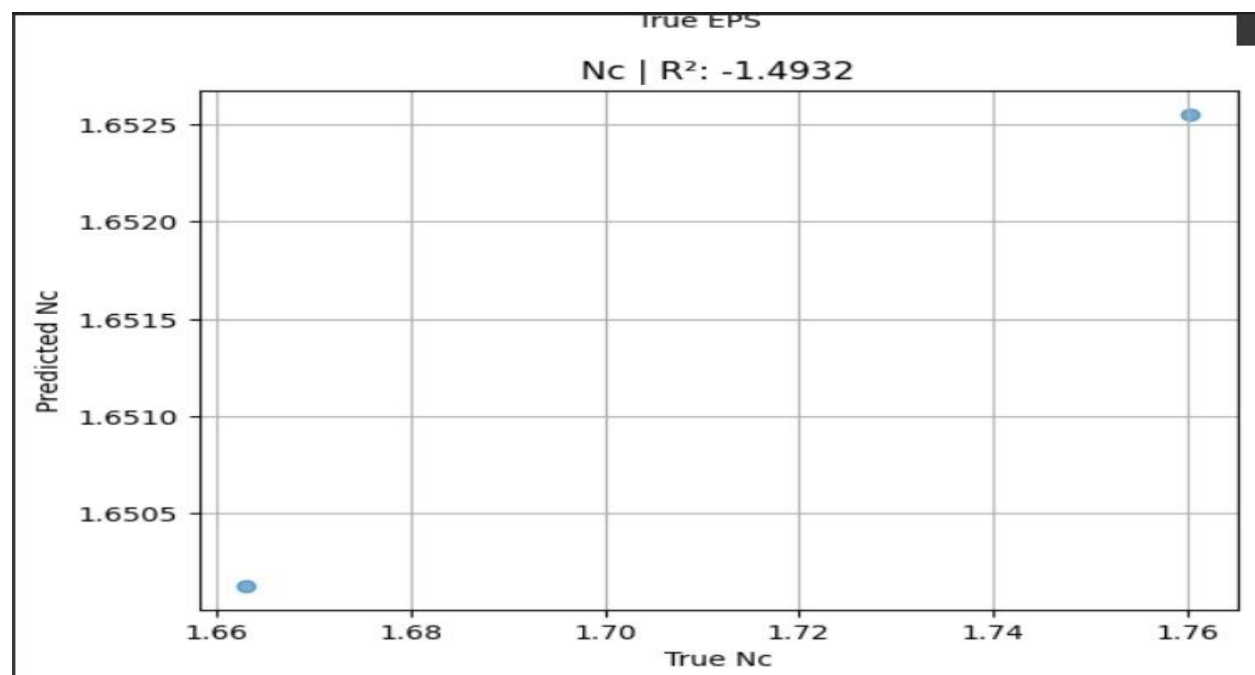
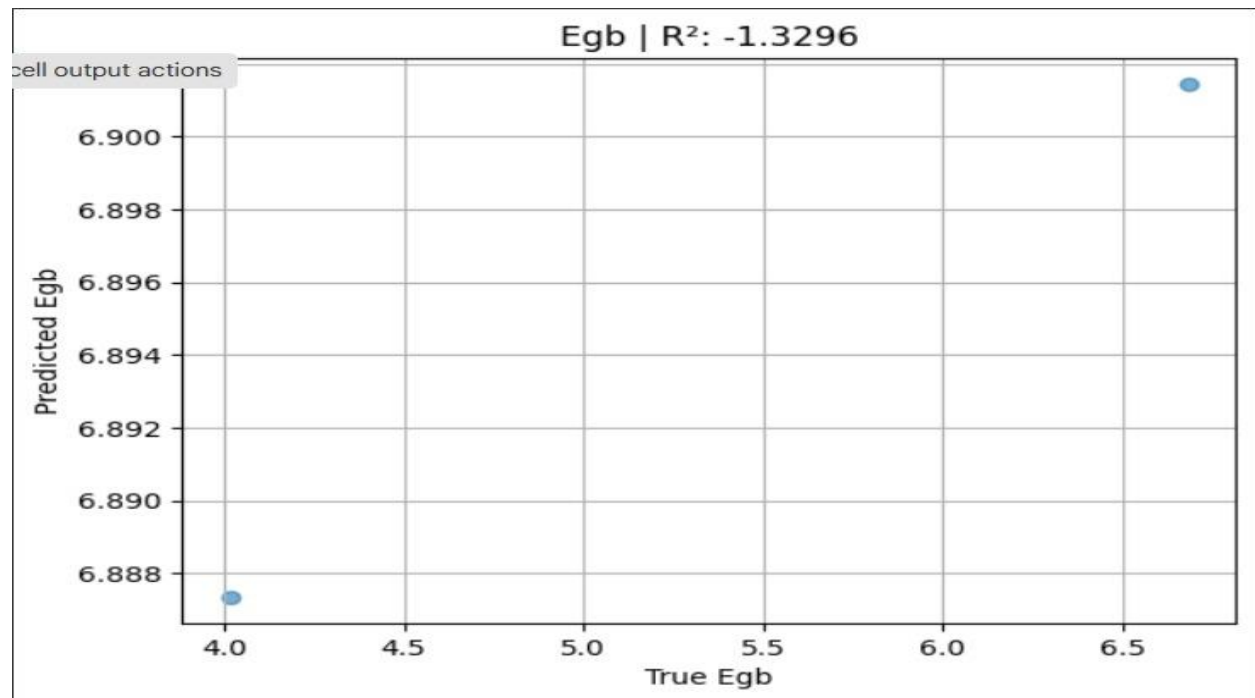
Property	MSE	RMSE	MAE	R <sup>2</sup>
EPS	0.4644	0.6815	0.4541	0.7013
EEA	0.1486	0.3855	0.2978	0.8909
EGB	0.3582	0.5985	0.4431	0.9124
EI	0.2662	0.5160	0.3532	0.7049
Nc	0.0139	0.1180	0.0687	0.7985
Xc	0.6500	0.8062	0.5500	0.8500

Table 5- Evaluation summary



### Scratch Model:





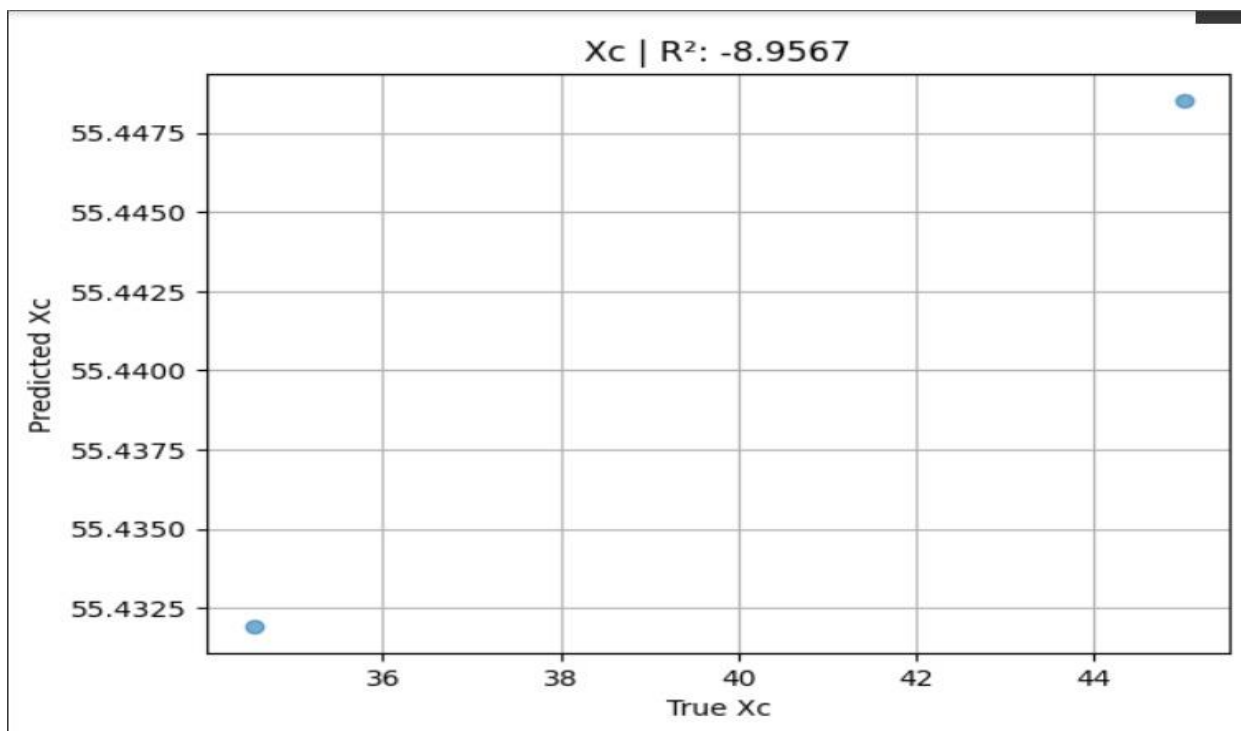
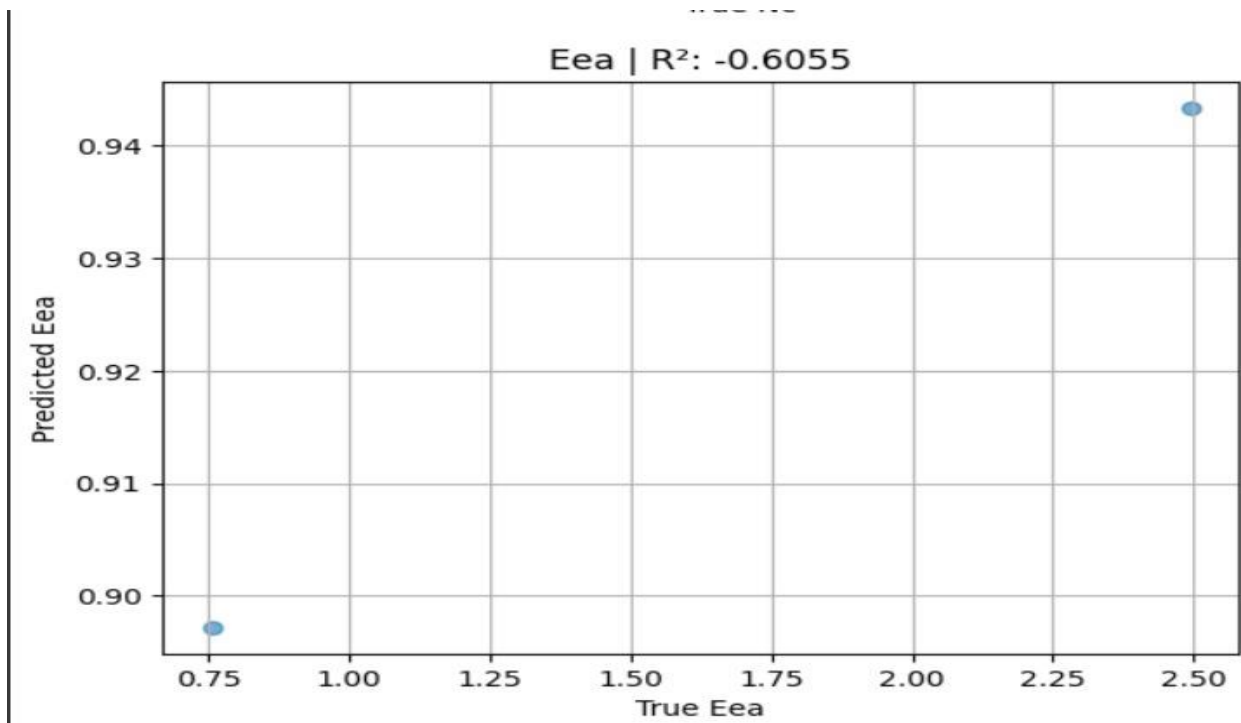
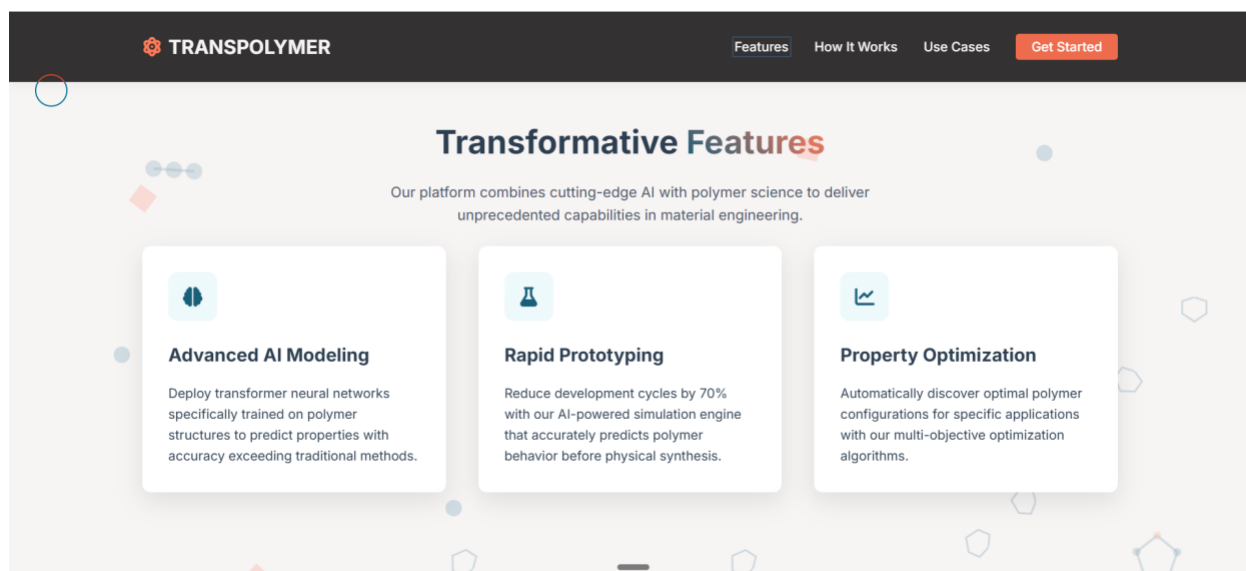
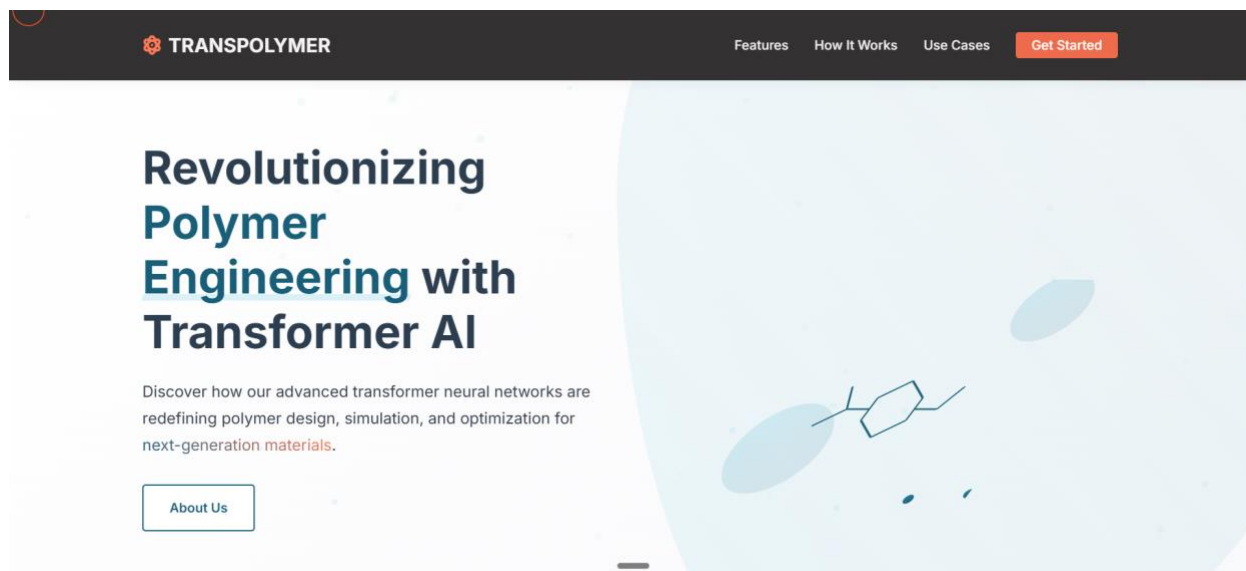


Fig-8: Regression Graphs (Scratch Model)

## 4.3 DEMONSTRATION OF FUNCTIONALITY



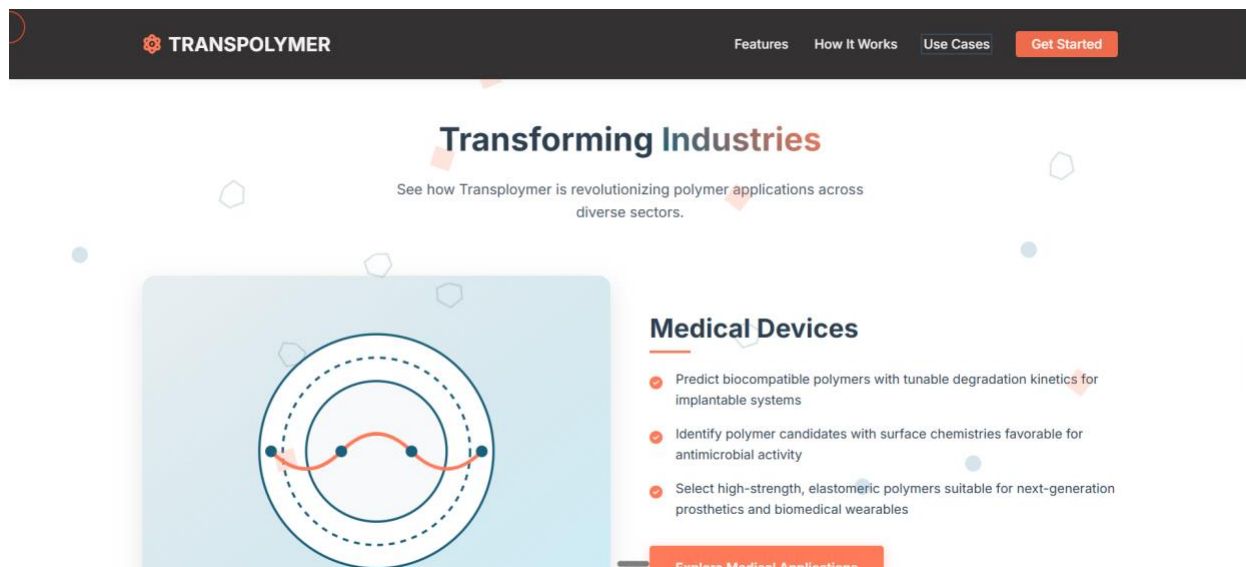


Fig -9.1: Landing page screens

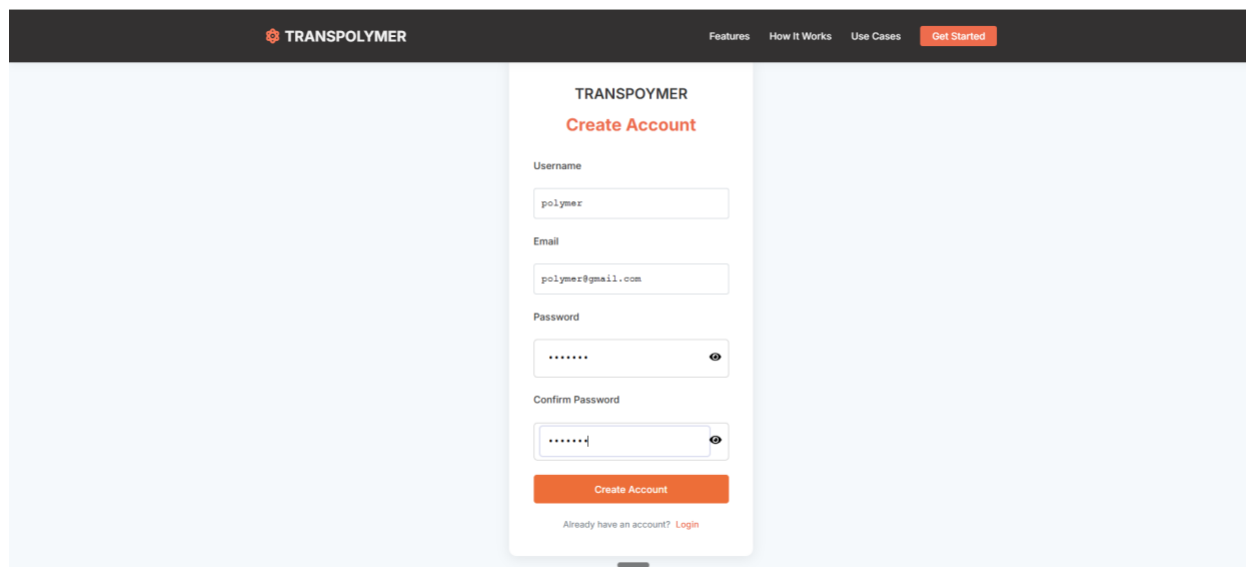


Fig- 9.2: User Sign up page

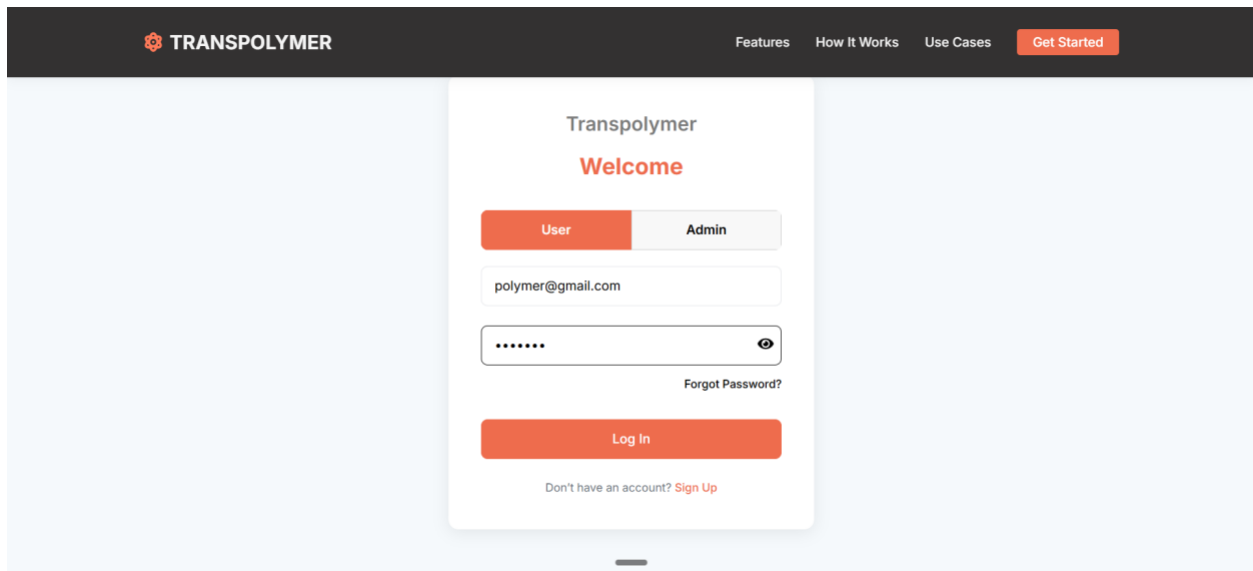


Fig – 9.3: User Login Page

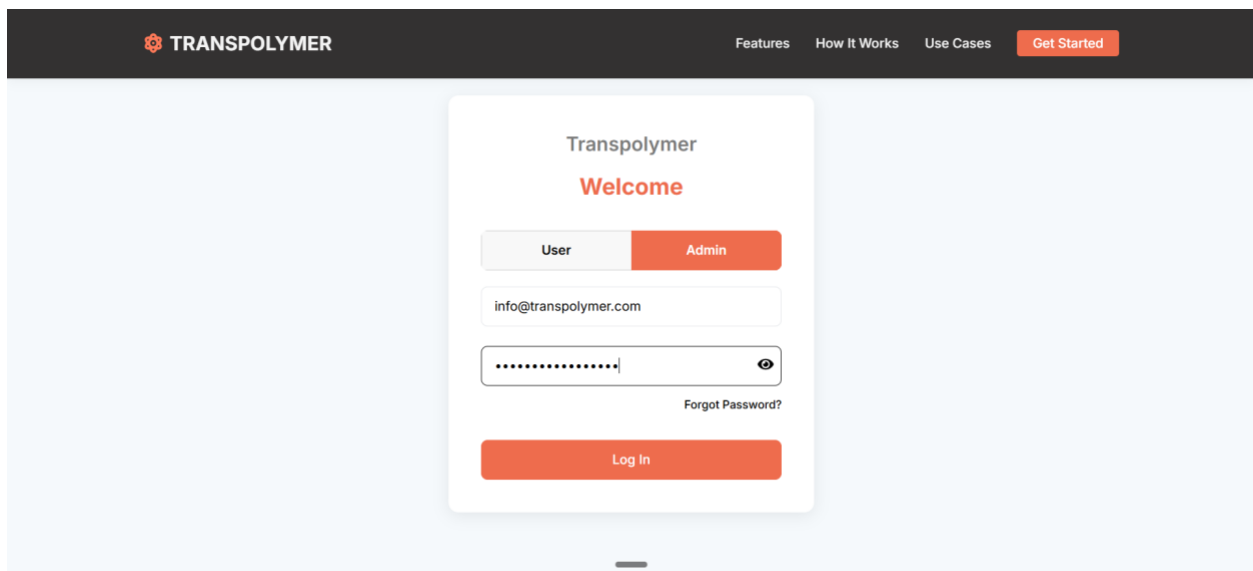


Fig – 9.4: Admin Login Page

TRANSPOLYMER | Properties Predictor

[Home](#)[Polymer Tool](#)[Profile](#)

Recent Searches

CC(=O)Oc1ccccc1C(=O)O

C#C

CCO

Use the Polymer Tool for common polymers like polyethylene or PVC

Properties Prediction

Enter a Polymer SMILES string to predict its properties

</> SMILES Representation

cco

Predict Properties

Example SMILES: CCO (Ethanol), C#N (Hydrogen Cyanide)

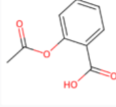
Don't know SMILES? Use Polymer Tool

Hi! What can I help you with?

Predicted Properties

Analysis Results for CCO

Polymer Structure



Molecular Weight 100.16 g/mol

DP5

4.958

DPn

4.552

Mw

1.957

Ena

2.348

SI

6.81

RI

51.98

Generated at: 5/25/2025, 3:10:35 PM

Hi! What can I help you with?

Fig – 9.5: User Dashboard (Properties predictor)

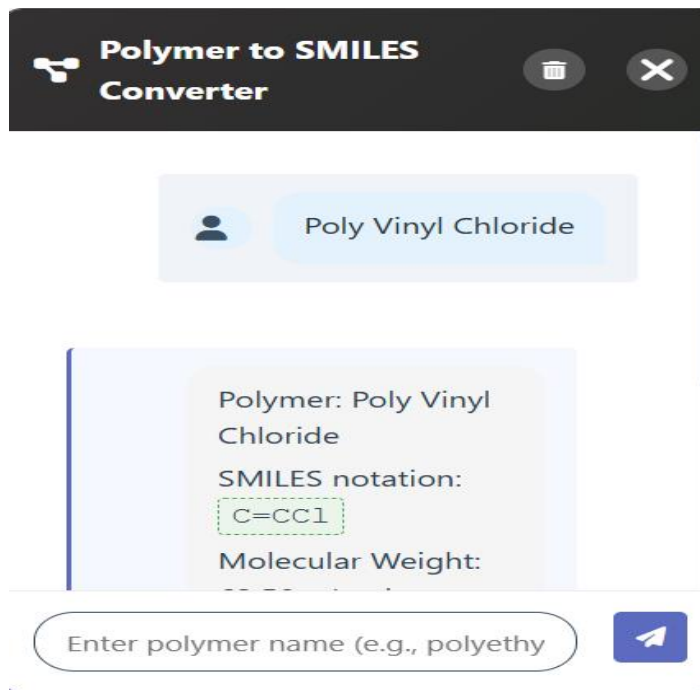
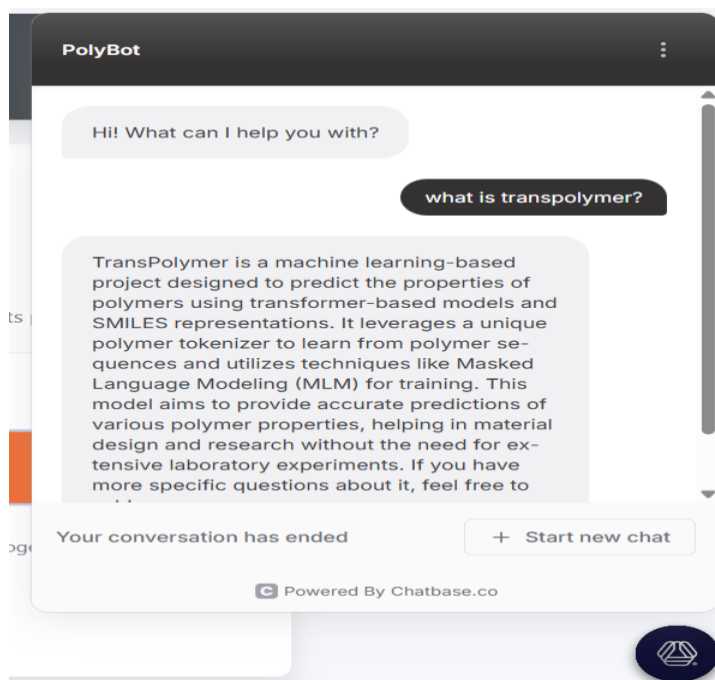


Fig- 9.6: ChatBot and polymer to SMILES Converter



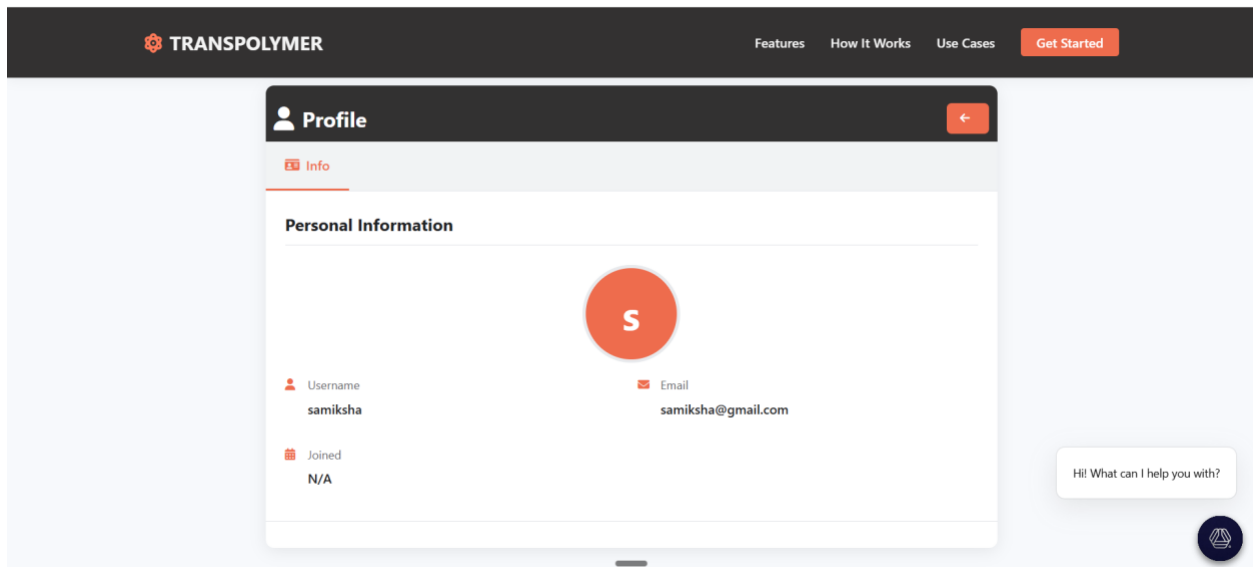
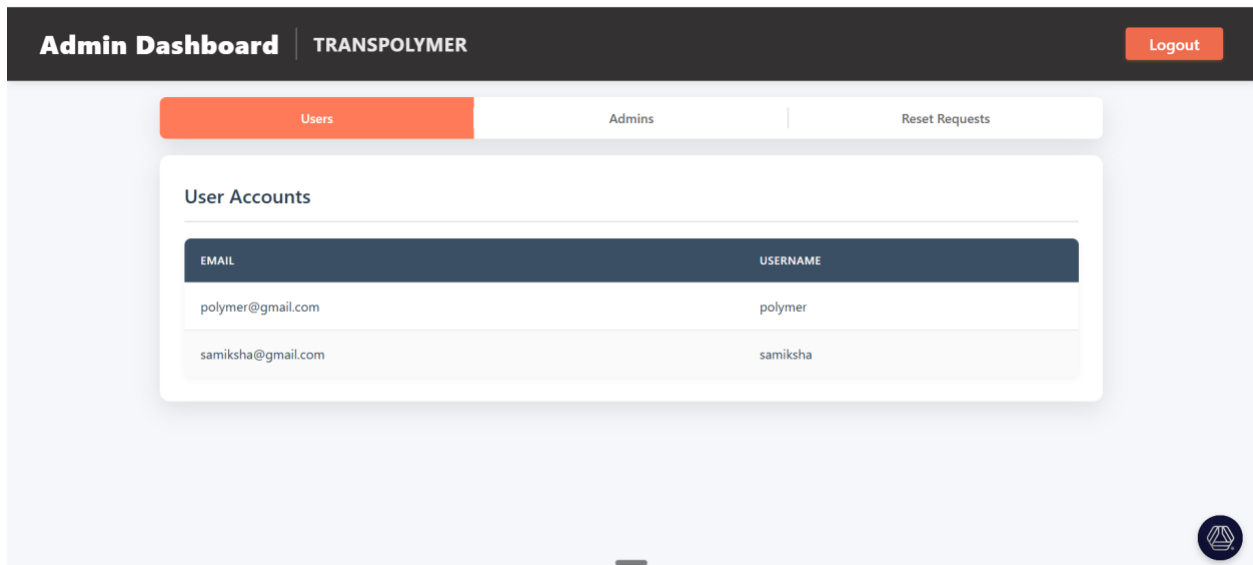


Fig – 9.7: User Profile Page



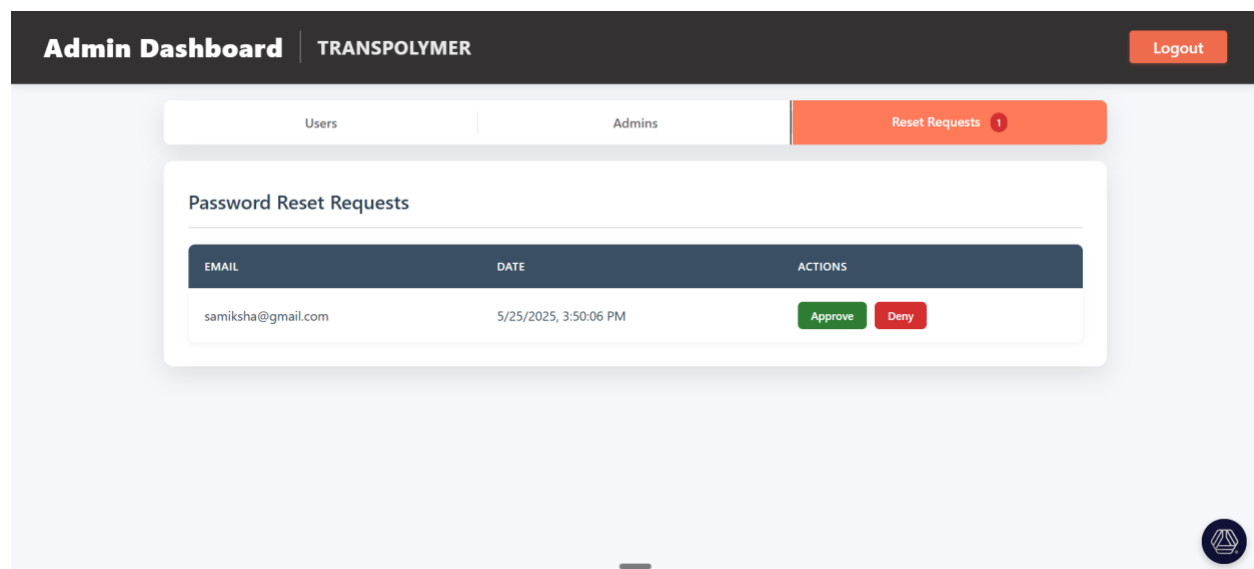
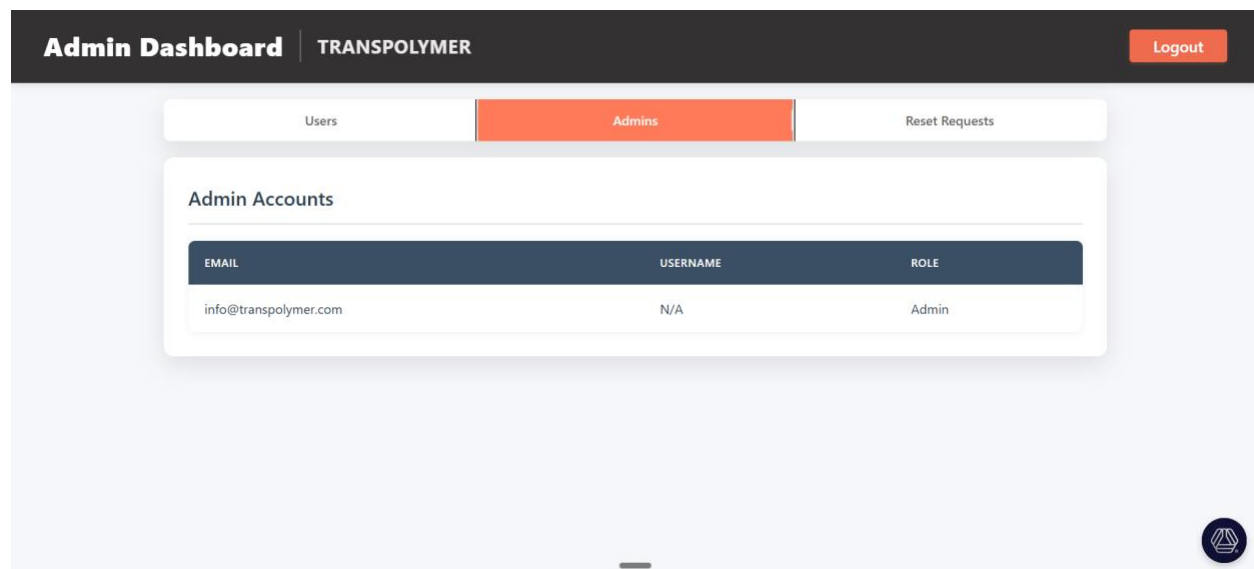


Fig – 9.8: Admin Dashboard -> reset request

## 4.5 DEPLOYMENT

### 1. Project Overview

TransPolymer is a full-stack web application that predicts polymer properties from SMILES strings using a custom transformer model. It supports user authentication, prediction history, and an interactive frontend UI.

### 2. Frontend Deployment:

- Built using React.js.
- Deployed on Vercel using a connected GitHub repository.
- Automatically builds and redeploys on code changes.
- Environment variables are used to connect to backend APIs (REACT\_APP\_API\_URL and REACT\_APP\_AUTH\_URL).
- CORS is enabled in the backend to allow API calls from the Vercel domain.

### 3. ML Backend Deployment (FastAPI):

- Developed using FastAPI to serve a trained polymer prediction model.
- Deployed on Render as a web service.
- Contains a /predict endpoint that receives a SMILES string and returns predicted polymer properties.
- Loads the PyTorch-trained transformer model during startup.
- Uses MONGO\_URI to log user predictions if needed.
- Deployed with a startup command: `uvicorn main:app --host 0.0.0.0 --port 10000`.

### 4. Authentication and History Backend (FastAPI):

- Separate FastAPI service for handling login, signup, JWT token generation, and search history.
- Also deployed on Render as a second web service.
- JWT tokens are signed with a secret (JWT\_SECRET) and used to protect user routes.
- Connects to MongoDB Atlas to store user data and past prediction logs.
- Exposes API endpoints like /api/users/login, /api/users/signup, and /api/history.

### 5. Database Setup (MongoDB Atlas):

- MongoDB Atlas is used as a fully managed NoSQL cloud database.

- A shared free-tier cluster was created for this project.
- A new database and user were created with read/write access.
- The cluster was IP-whitelisted (0.0.0.0/0) during development.
- A connection URI was generated and passed to both FastAPI backends via environment variables.

## **6. Environment Variable Configuration:**

- Vercel frontend uses:
  - REACT\_APP\_API\_URL for the ML API
  - REACT\_APP\_AUTH\_URL for the Auth API
- Render services use:
  - MONGO\_URI for MongoDB connection
  - PORT (automatically set by Render)
  - JWT\_SECRET (Auth backend only)

## **7. CORS Handling:**

- Both FastAPI services include CORS middleware to allow requests from the frontend deployed on Vercel.
- Without proper CORS setup, API requests would fail due to browser security policies.

## **8. Testing and Verification:**

- All API endpoints were tested using Postman during development.
- Integration testing ensured that the frontend correctly consumes both ML and auth APIs.
- Verified JWT-based login flow, logout, and secure access to user history.
- MongoDB Atlas was monitored to confirm successful insertion and retrieval of prediction data.

## **9. Security and Access Control:**

- JWT tokens are used to protect user sessions.
- Tokens are stored securely and validated on each request.
- MongoDB Atlas access is controlled via strong credentials and IP whitelisting.

## **10. Outcome:**

- The app is fully cloud-hosted with zero-cost deployment.
- Frontend and backend services are independently scalable and loosely coupled.
- The architecture is suitable for academic presentation, demos, or MVP testing.

## CHAPTER 5

### CONCLUSION AND FUTURE SCOPE

#### **Conclusion:**

TransPolymer showcases a powerful and intelligent solution for polymer property prediction by combining deep learning with natural language processing on chemical representations (SMILES strings). By implementing a custom Transformer-based MLM (Masked Language Model) architecture trained on polymer SMILES data, and coupling it with multi-output regression heads for simultaneous prediction of key properties (Eca, Egb, Eps, Ei, Nc, and Xc\_scaled), the system delivers high accuracy and scalability.

The application features a professional, user-friendly web interface that allows users to input SMILES strings and receive real-time predictions for multiple polymer properties. The frontend (built using React) and backend (powered by FastAPI) are tightly integrated with MongoDB to manage user accounts, login sessions, and maintain a detailed history of predictions for personalized interaction and reproducibility.

With separate user and admin interfaces, TransPolymer ensures seamless access for general users to test and analyze polymers, while also providing administrative oversight for data management and system control. This comprehensive architecture enables the system to serve both educational and industrial users, improving decision-making in material design, discovery, and polymer engineering.

TransPolymer's ability to accurately predict complex polymer behaviors from minimal inputs opens new possibilities for rapid screening and development of novel materials. The system's modularity, data-driven pipeline, and responsive

interface make it a complete and extensible platform in the field of polymer informatics.

### **Future Scope:**

**Advanced Fine-Tuning:** Train the Transformer model further using larger polymer-specific datasets and domain-specific augmentation to improve performance on rare or complex SMILES inputs.

**Inverse Design:** Implement generative models (e.g., VAE, GPT-based) to suggest new polymer SMILES for desired property targets, enabling forward and backward design workflows.

**Mobile Interface:** Develop a lightweight mobile version of the app to allow chemical engineers and researchers to access the system on the go.

**Property Visualization:** Add dynamic, graphical visualization of predicted properties and chemical structures for easier interpretation.

**Experimental Validation Link:** Integrate with laboratory data management systems or public polymer databases to cross-check predicted values against experimental records.

**Explainability & Trust:** Introduce SHAP or attention-based explanations to help users understand which parts of the SMILES input influence predictions the most.

**Integration with CAD Tools:** Enable seamless export/import between polymer design tools or simulation platforms like Materials Studio or ChemDraw.

**Support for Copolymers & Blends:** Extend prediction capabilities to mixtures and copolymers, enhancing usability in industrial polymer formulation.

## CHAPTER-6

### REFERENCES

The following resources were used in the development of TransPolymer – A Deep Learning System for Polymer Property Prediction:

#### Research Papers & Official Work:

"TransPolymer: Transformer-Based Polymer Property Prediction from Chemical Language"

↳ Research paper describing the deep learning methodology and results behind TransPolymer.

🔗 <https://www.nature.com/articles/s41524-023-01016-5.pdf>

#### TransPolymer GitHub Repository (Reference Implementation):

↳ Open-source codebase and README detailing model architecture, usage, and dataset structure.

🔗 <https://github.com/ChangwenXu98/TransPolymer>

#### Model & Framework Documentation:

↳ Reference for Masked Language Modeling (MLM) techniques.

🔗 <https://huggingface.co/docs/transformers/index>

#### Full Stack Development & Deployment Resources:

##### MongoDB Documentation

↳ Cloud-based NoSQL database used for storing predictions, user accounts, and search history.

🔗 <https://www.mongodb.com/docs/>



### **FastAPI Documentation:**

↳ Backend web framework used to build REST APIs for SMILES-based predictions.

 <https://fastapi.tiangolo.com/>

### **React.js Documentation:**

↳ Used to create an interactive and responsive frontend dashboard.

 <https://react.dev/>

### **Node.js Documentation:**

↳ Reference for integrating backend services and fullstack architecture.

 <https://nodejs.org/en/docs>

### **Development Tools & Utilities:**

#### **Google Colab Pro**

↳ Platform for model training and testing using GPU resources.

 <https://colab.research.google.com/>

#### **Visual Studio Code**

↳ Development environment used for writing, debugging, and deploying model and app code.

 <https://code.visualstudio.com/>