

Multimodal Digital Twin for Type 2 Diabetes Patients



202255637 Qonysbekova Yenglik

202255631 Pak Elina

202255545 Bakhieva Aysuliu

Advisor: Professor Song Giltæ (sign)

Table of contents

1. Project Objectives	1
Background Statistics and Motivation	1
1.1 Introduction	2
1.2 Detailed Project Goals	2
2. Target Problem and Requirements Analysis	3
2.1 Analysis of Similar Systems	3
2.2 Functional Requirements	4
2.3 Non-Functional Requirements	4
3. Realistic Constraints and Countermeasures	5
3.1 Constraints	5
3.2 Countermeasures	5
4. Design Document and Implementation Plan	5
4.1 System Architecture Overview	5
4.2 Multimodal Data Pipeline	5
4.3 Feature Encoding and Model Pipeline	6
4.4 What-if Simulation Engine	6
4.5 Application Features	6
4.6 System Visualization Tools	6
4.7 Testing and Evaluation Plan	7
5. Development Schedule and Role Division	7
5.1 Development Schedule	7
5.2 Role Division	7
6. Datasets	8
8. Literature Review and References	10

1. Project Objectives

Background Statistics and Motivation

According to the International Diabetes Federation (IDF, 2024), as of 2021, over 537 million adults were living with diabetes, with the number projected to rise to 783 million by 2045. More than 90% of these cases are Type 2 Diabetes (T2D), and nearly half of all adults go undiagnosed with the disease. The estimated annual global cost of diabetes exceeds \$966 billion USD, leading to 6.7 million deaths per year due to complications such as stroke, cardiovascular disease, or kidney failure.

Traditional models of care are often fragmented, inconsistent, and lack personalization. According to the World Health Organization (WHO, 2024), 80% of T2D patients experience at least one comorbidity. Meanwhile, recent statistics from digital health platforms indicate that over 60% of users stop using diabetes apps after just one month, likely due to the absence of personalized insights or actionable feedback. The mission of our project is to address these gaps by developing an intelligent and predictive system — a Multimodal Digital Twin for T2D — leveraging AI, multimodal data, and clinical research to enhance the understanding, management, and personalization of chronic disease care.

1.1 Introduction

With more than 400 million affected, Type 2 Diabetes (T2D) stands as a grave global health crisis, further giving rise to serious complications, such as cardiovascular issues, kidney failure, neuropathy, and retinopathy. The management of this disease remains a tricky proposition due to high variability in patients as several genetic, behavioural, and environmental factors interplay.

Recent works, including Jiang et al. (2023), highlight the potential of Digital Twin technologies in simulating chronic disease states and real-time therapeutic planning (Jiang et al., 2023). Ma et al. (2024) add that multimodal data fusion combining EHR, sensor, and behavioral data drastically boosts predictive capability and customization (Ma et al., 2024).

Moreover, according to Choi et al. (2024), patient stratification by phenotypic clustering in diabetes helps guide more personalized and specific management strategies. This will further

promote the incorporation of ML methods that track individual variability such as TGNNs and multimodal feature fusion in customizing predictive insights (Choi et al., 2024).

In this regard, the study proposed to develop the Multimodal Digital Twin Application for T2D, functioning as an interactive mobile/web-based platform. This app will receive continuous multimodal patient data to simulate disease evolution via AI models and facilitate both patients and clinicians to explore various lifestyle and treatment interventions on a more user-friendly level.

1.2 Detailed Project Goals

Goal	Description
Mobile/Web App Interface	Develop a user-facing application to deliver real-time simulations, patient tracking, and visualization tools.
Multimodal Data Integration	Integrate EHR (from MIMIC-IV), behavioral and survey data (from NHANES), wearable sensor streams (Open mHealth), and publicly available health data (e.g., OhioT1DM, Synthea) into unified profiles.
Temporal Disease Simulation	Model T2D progression (e.g., HbA1c levels, insulin needs) over time using TGNNs and Transformer-based architectures.

What-if Scenario Forecasting	Enable patients and clinicians to simulate outcomes of lifestyle or medication changes through interactive scenario exploration.
Visual Progress Dashboard	Display disease trajectories and simulation results through intuitive graphs and charts in the app.
Patient-Specific Modeling	Ensure that all predictions and recommendations are tailored to each patient's personal data, with explainable AI outputs.

2. Target Problem and Requirements Analysis

2.1 Analysis of Similar Systems

Twin Health is a commercial platform to build digital twins catered to chronic conditions such as Type 2 Diabetes. It ingests clinical data, lifestyle habits, and data from wearable sensors to generate a digital profile for a patient, which can be monitored in real time and feedback given accordingly.

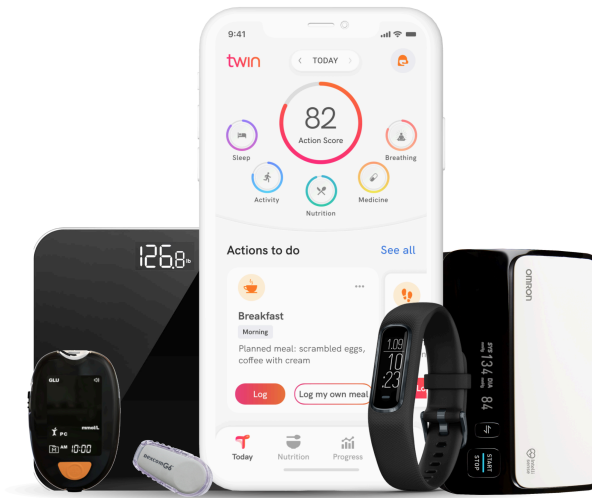
Key Advantages:

1. Combining health data across different domains for personalized care.
2. Allows continuous monitoring with interactive coaching.
3. Demonstrable effectiveness and user engagement.

Main Drawbacks:

1. Being proprietary, little transparency is offered.
2. Not of a great value for academic or open-source use.

-
3. Difficult to customize or extend on a stand-alone basis.



Omada Health is a mobile-based system focused on preventing diabetes and encouraging behavior change through remote care. It assists the user in managing obesity, activity, or diet using the structured digital tools.

Key Advantages:

1. Focuses on early intervention and behavior modification.
2. Systems of structured coaching with remote care support.
3. Clinically proved in bringing about a healthier life.

Major Drawbacks:

1. Lacks advanced AI functions such as prediction or simulation.
2. Very limited in terms of deep data integration and modeling.
3. Little to no customization outside of simple coaching needs.



2.2 Functional Requirements

1. Support login and secure access for patients and clinicians
2. Enable integration of EHR, wearables, and lifestyle data
3. Run disease simulation and visualize predicted outcomes
4. Provide interactive what-if simulations of lifestyle and treatment scenarios
5. Display historical trends and simulated projections
6. Ensure all predictions are explainable and patient-specific

2.3 Non-Functional Requirements

1. System should support mobile and web platforms (Flutter or React)
2. Ensure secure handling of sensitive health data (HIPAA-compliant)
3. Provide intuitive and accessible UI for patients and clinicians alike
4. Architecture must support scalable ML inference and data pipelines
5. Offline support for limited access scenarios with local caching

3. Realistic Constraints and Countermeasures

3.1 Constraints

- Limited access to real patient data due to privacy laws
- Difficulty in evaluating clinical accuracy without medical partners
- Synchronizing heterogeneous data (EHR, wearables) in real time

-
- Variability in patient behavior makes modeling complex

3.2 Countermeasures

- Use open-source datasets: MIMIC-IV (Johnson et al., 2021), NHANES (CDC, 2024), Open mHealth (2024), Synthea, OhioT1DM.
- Simulate and prototype using synthetic and de-identified data
- Start with modular integration pipelines for each modality
- Build personalization layers for patient-specific calibration

4. Design Document and Implementation Plan

4.1 System Architecture Overview

The system is structured into four main components:

1. **Frontend:** Built with Flutter (for mobile) or React (for web), providing a responsive and intuitive interface for both patients and clinicians.
2. **Backend API:** Developed using FastAPI (Python), exposing endpoints for data queries, simulation, user authentication, and interaction logging.
3. **AI/ML Engine:** Developed in PyTorch, using a modular architecture with Temporal Graph Neural Networks (TGNNs) and a Multimodal Transformer for data fusion.
4. **Database Layer:** PostgreSQL is used for user and scheduling data, while MongoDB stores time-series sensor data and intermediate model outputs.

4.2 Multimodal Data Pipeline

The pipeline supports ingestion, preprocessing, and transformation of four types of data:

- **EHR Data** (MIMIC-IV (Johnson et al., 2021)): Extract diagnosis, lab values, medications.
- **Survey & Lifestyle** (NHANES (CDC, 2024)): Clean and transform self-reported lifestyle variables.
- **Wearables** (Open mHealth (2024)): Time-align step count, heart rate, and sleep.

-
- **Synthetic EHRs** (Synthea, OhioT1DM): Prototype pipeline integration, glucose prediction modeling.

Steps:

5. Standardize column names and timestamps.
6. Impute missing data using time-aware or statistical imputation.
7. Normalize numerical features.
8. One-hot or embedding for categorical values.
9. Store preprocessed vectors in a unified feature store.

4.3 Feature Encoding and Model Pipeline

Each modality has a dedicated encoder:

- EHR → MLP
- Lifestyle/Survey → Fully connected layers
- Wearables → 1D CNN or LSTM

These features are fused in a **Multimodal Transformer**, then passed to a **TGNN** to simulate temporal dynamics and predict outcomes (e.g., glucose trajectory, HbA1c projection).

4.4 What-if Simulation Engine

Patients or clinicians can simulate lifestyle changes (e.g., increase exercise, improve diet).

- Simulation modifies feature vectors
- Model re-runs and forecasts new outcomes
- Display projected trends on dashboard for comparison

4.5 Application Features

- Secure Login: Google/email/Kakao integration
- Role-based Interface: Doctor, patient, researcher views
- Dashboard: Timeline of real + predicted metrics (glucose, HbA1c)
- Simulation Interface: Dropdown and sliders for hypothetical input changes
- Progress Reports: Generated PDFs summarizing patient evolution

4.6 System Visualization Tools

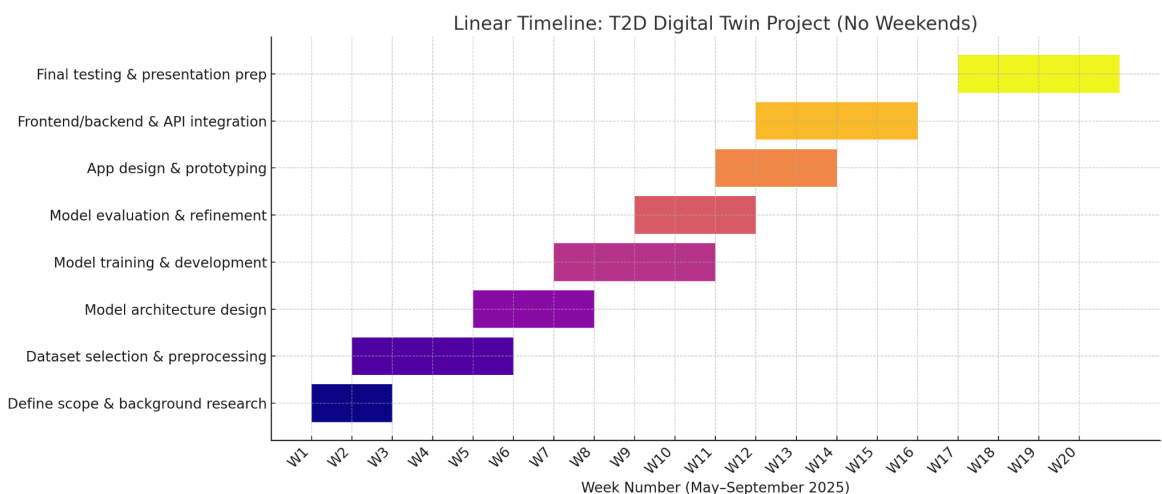
- Use matplotlib or plotly for dynamic graphing
- Use Chart.js for frontend charting
- Display clusters, predictions, and trends per patient visually

4.7 Testing and Evaluation Plan

- Unit Tests: API endpoints, model inputs/outputs
- Evaluation Metrics: MAE, RMSE, accuracy, AUC
- Simulation Validity: Ensure plausible outcomes from behavior changes
- User Testing: Interface clarity and usefulness via feedback forms

5. Development Schedule and Role Division

5.1 Development Schedule



5.2 Role Division

Responsibility Area	Task Leader	Support Contributors
Project ideation & direction	All team members	—
Literature review & dataset exploration	All team members	—
Data collection &	Bakhieva Aisuliu	Qonysbekova Yenglik, Pak

preprocessing		Elina
Model development & evaluation	Pak Elina	Bakhieva Aisuliu, Qonysbekova Yenglik
Prototype application development	Qonysbekova Yenglik	Pak Elina, Bakhieva Aisuliu
Testing, documentation, final report & presentation	All team members	—

6. Datasets

AI-READI

The AI-READI dataset is a large-scale, multi-modal health dataset developed by NIH as part of the Bridge2AI initiative. It integrates electronic health records, biospecimens, imaging, and patient-reported data with a focus on enabling responsible AI development for health diagnostics and treatment.

Shanghai T2DM

The Shanghai T2DM dataset consists of real-world clinical records collected from patients with Type 2 Diabetes Mellitus (T2DM) in Shanghai. It includes demographic, laboratory, and lifestyle data and is used for research on disease progression, risk prediction, and intervention strategies in diabetic populations.

MIMIC-IV

MIMIC-IV (Medical Information Mart for Intensive Care IV) is a comprehensive, de-identified dataset comprising health data from ICU patients at the Beth Israel Deaconess Medical Center. It includes clinical notes, vitals, medications, and laboratory results, supporting research in critical care and AI in medicine.

NHANES

The National Health and Nutrition Examination Survey (NHANES) is a U.S. population-level dataset that combines interviews, physical examinations, and laboratory tests to assess the health and nutritional status of adults and children, often used in epidemiological

and public health research.

Open mHealth

Open mHealth is an open-source data standard initiative and platform providing structured mobile health data, such as step count, heart rate, and mood logs. It enables the integration and interoperability of data across mHealth devices and applications for personalized health analytics.

Synthea

Synthea is a synthetic patient data generator that produces realistic but fictional health records. It simulates medical histories using standard healthcare models and enables testing and development of health IT systems without privacy concerns.

OhioT1DM

The OhioT1DM dataset is a collection of data from individuals with Type 1 Diabetes Mellitus, including continuous glucose monitoring, insulin dosage, and meal intake. It supports research in glucose prediction, insulin management, and AI-based decision support for diabetes care.

8. Literature Review and References

1. Jiang, Y. et al. Digital twins for type 2 diabetes: from modeling to application. npj Digital Medicine (2023). <https://www.nature.com/articles/s41746-023-00933-5>
2. Ma, Y., et al. (2024). *Multimodal data fusion for disease prediction*. *Scientific Reports*. <https://www.nature.com/articles/s41598-024-71020-2>
3. Alamo, T., et al. (2020). *Explainable Artificial Intelligence in health: A study on T2D patients*. *Scientific Reports*. <https://www.nature.com/articles/s41598-020-68771-z>
4. ACM Digital Library. (2024). *Patient-centered digital twin approaches for chronic care*. <https://dl.acm.org/doi/abs/10.1145/3643479.3662049>
5. Choi, Y., et al. (2024). *Phenotypic clustering of type 2 diabetes to inform patient stratification and treatment*. *Frontiers in Endocrinology*. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10853398/>
6. International Diabetes Federation. (2024). *IDF Diabetes Atlas*. <https://diabetesatlas.org>
7. World Health Organization (WHO). (2024). *Diabetes Fact Sheet*. <https://www.who.int/news-room/fact-sheets/detail/diabetes>
8. Statista & ACM Digital Library. (2024). *Analysis Reports on Digital Health and Chronic Care*. <https://dl.acm.org/doi/abs/10.1145/3643479.3662049>