Predicting Future Type 1 Diabetes Onset Risk: A Machine Learning Approach Hood D¹, Walia K¹, Kshirsagar O¹, Poddar S¹, Mankin B¹, Lee KCS², Adamek L³, Rufino B³, Josleyn J⁴

¹Axtria, Berkley Heights, NJ, USA;

²Sanofi, Paris, France;

³Sanofi, Toronto Canada;

⁴Sanofi, Chicago, IL, USA

POSTER HIGHLIGHT: Novel BERT Transformer Model Predicts Symptomless T1D Earlier, Potentially Reducing DKA Risk and Enabling Early Intervention

INTRODUCTION

- Diabetes (T1D) is an autoimmune condition that affected Figure 1 details the cohort creation methodology from the US Optum Extreme Gradient Boosting (XGB) with a class weight of 200 for T1D and approximately 2 million people in the United States in 2021, including Clinformatics[®] Database⁸, comparing T1D and comparator groups. a prediction threshold of 0.8 effectively addressed class imbalance, around 304,000 children and adolescents (aged <20 years) and 1.7 million Feature selection reduced 256 one-hot encoded features to 50 grouped achieving precision: 2.1%, recall: 20%, F1 Score: 0.04, and Bayes adults (aged ≥ 20 years).¹ variables for efficiency, with LASSO regression retaining all 50 variables. Factor (BF): 4.67 (Figure 3). Figure 2 shows the input model features: including comorbidities, age Applying a BF of 4.67 to a baseline rate of 1:200 changed it to 1:43, and 200, but this risk increases to 1 in 20 among individuals with a family group, race, and LOINC lab results for T1D and Non-T1D cohort. history, representing a 15-fold increase compared to 1 in 300 for those for test data with a base rate of 0.9178 in 200, it became 1:47, equating Various ML models (random forest, decision tree, logistic regression, without.² to 448 true positives out of 21,362 predicted positives (Figure 3). XGBoost) were tested for predicting early diagnosis of pre-symptomatic Early identification of pre-symptomatic T1D reduces the risk of diabetic T1D. • Accuracy generally decreased as the time window moved from the index ketoacidosis (DKA) at diagnosis, with reported rates as low as 2.5%–5% in • Post hoc analysis using an optimized BERT model pre-trained completely date, with an "elbow" in the performance curve identified as the optimal screened populations, enables timely intervention, and gives patients and on Optum dataset (1 year of medical event data), following the Med-BERT⁹ time window (12-24 months before the index date) for early detection, families time to prepare for effective disease management, potentially architecture, and fine-tuned on three cohorts to enhance prediction improving long-term outcomes and reducing healthcare burden.³⁻⁵ balancing early prediction and accuracy (Figure 4). accuracy for diabetes-related endpoints. The most important variables used by the model to make a T1D This model was also applied to US claims data (T1D and non-T1D Figure 5: ROC curve for BERT model age, despite 80-90% of new cases having no family history. cohorts), utilizing key variables such as ICD, NDC, CPT codes, and prediction are age (17.79%), COM_HEALTH_EXAM (12.40%), demographics (age, race, gender, index date). comorbidities including: hyperlipidemia (7.62%) and musculoskeletal target more accurate predictions, enhancing the potential for early disorders (6.02%). Figure 2: Key input features for T1D and non-T1D detection and intervention.^{6,7} cohort post feature selection and engineering Figure 3: XGB model with W =200 and T =0.8 (12-24) months pre-index) Weight_loss_management BF = 4.67 Infectious And Parasitic Disease XGB 25% Heart failure Precision Recall 448 TP 0.2 -20% ■ F1 ROC curve (AUC = 0.85) 20,914 FP 15% -- Random selection Discases_oi_lite_cal_anu_i 0.0 0.2 Diseaes_of_respiratory_system 0.4 0.6 0.8 **US Optum Clinformatics® Database** False Positive Rate 10% FN 1,753 Cough/Fever/Headache/Fatique Coronary_heart_disease Inpatient Medical and Lab Results and Autoimmune Disorder atient Demographics 5% 458,695 ΤN **Provider Information** Data Pharmacy Cla 10 20 30 40 50 60 The BERT model achieved an overall F1 score of 81.82% across the **2.10%** Non-T1D (%) T1D (%) TP, true positive; FP, false excellent precision for the non-T1D cases with dataset, positive; FN, false negative; TN, Thyroid_Function_Tests (99.73%) and higher precision for T1D than XGB (2.59% vs 2.10%) true negative **T1D Cohort Comparator Cohor** Renal_Function_Tests (Table 1). Nutritional markers • **Comparator Cohort:** 1.5 million individuals randomly selected from the Optum database Liver_function_tests Figure 4: Sliding window results • Prevalence Basis: Sampling followed a 1:200 prevalence ratio • The model achieved high recall for T1D (83.42%) (Table 1). Lipid/Cardiovascular $(7,500 \times 200)$, aligned with real-world and Optum estimates Infectious_Disease • The ROC curve shows an AUC of 0.85, indicating good overall Hematologica 0.67% discriminatory ability between T1D and non-T1D cohort (Figure 5). HBA1C 0.662% 0.660% 0.659% **Diabetes Cohort** Glucose 0.656% 0.66% 0.656% Patients must have one of the Electroytes/Acid_Balance following three criteria: T1D cohort Autoimmune Antibody CONCLUSION • 2+ DM diagnosis codes (250.x, 0.65% **Inclusion criteria Exclusion** E11.x, E10.x) within 1 year 25 30 20 T2DM to T1DM ratio criteria 0.8
- Type 1 • In the general population, the lifetime risk of developing T1D is about 1 in • Traditional prediction methods rely on risk factors like family history and Machine learning (ML) uses real-world data and advanced algorithms to OBJECTIVE The objective of this research is to explore the potential of ML models in predicting the onset of T1D using US claims data and to identify the earliest timeframe at which T1D can be accurately predicted using ML techniques. **METHODS**



REFERENCES

- 1. Centers for Disease Control and Prevention (CDC). (2024, May 15). National Diabetes Statistics Report: Diabetes. Retrieved from
- https://www.cdc.gov/diabetes/php/data-research/?CDC_AAref_Val=https://www.cdc.gov/diabetes/data/statistics-report/index.html 2. TrialNet. (2025, April 15). Type 1 Diabetes Facts. Retrieved from https://www.trialnet.org/t1d-facts 3. Simmons, K. M., & Sims, E. K. (2023). Screening and prevention of type 1 diabetes: Where are we? The Journal of Clinical Endocrinology &
- Metabolism, 108(12), 3067–3079. 4. Schneider, J., Gemulla, G., Kiess, W., Berner, R., & Hommel, A. (2023). Presymptomatic type 1 diabetes and disease severity at onset.
- Diabetologia, 66(12), 2387–2388. 5. Hekkala, A. M., Ilonen, J., Toppari, J., Knip, M., & Veijola, R. (2018). Ketoacidosis at diagnosis of type 1 diabetes: Effect of prospective studies with newborn genetic screening and follow-up of risk children. Pediatric Diabetes, 19(2), 314–319.

METHODS (cont'd)



- 6. Ghalwash, M., Anand, V., Ng, K., Dunne, J. L., Lou, O., Lundgren, M., et al. (2024). Data-driven phenotyping of presymptomatic type 1 diabetes using longitudinal autoantibody profiles. Diabetes Care, 47(8), 1424–1431. 7. Meng, W., Qin, J., Wang, T., & Zhao, R. (2025). Application of data science in management of type 1 diabetes. In Type 1 Diabetes - Causes,
- Treatments and Management. IntechOpen. 8. Optum Clinformatics® Database. (Accessed 2024). Available from: https://www.optum.com/business/solutions/life-sciences/real-world-
- 9. Rasmy, L., Xiang, Y., Xie, Z., Tao, C., & Zhi, D. (2021). Med-BERT: Pretrained contextualized embeddings on large-scale structured electronic health records for disease prediction. NPJ Digital Medicine, 4(1), 86.

RESULTS



0.61%

0.60%

15-27

18-30

FUNDING

9-21

Sliding Window

12-24

The study was sponsored by Sanofi.

6-18

ACKNOWLEDGMENT was funded by Sanofi.

0-12

3-15

____E-Poster



Copies of this poster obtained through Quick Response (QR) Code are for personal use only

RESULTS (cont'd)

Table 1: Performance metrics for BERT model			
Class/Metric	Precision	Recall	F1 Score
Non-T1D	99.78%	70.38%	82.54%
T1D	2.59%	83.42%	5.03%
Accuracy			70.51%
Macro Average	51.18%	76.90%	43.78%
Weighted Average	98.87%	70.51%	81.82%



- The study successfully demonstrated that significant predictive signals exist within the data, allowing for the development of effective and reliable models for early T1D detection.
- Future work will further explore OSCAR, a BERT derivative, incorporating new capabilities and broader data sources including lab tests, clinical procedures, imaging, genomics, wearables, and SDoHrelated data.
- Exploring cost/benefit assumptions can refine model thresholds and weightings, ensuring accuracy and alignment with real-world clinical and economic contexts.

Editorial assistance was provided by **Navneet Kumar, PhD**, Shweta Shaw and Jatinder Kumar of Axtria India Pvt Ltd. and

DISCLOSURES

Lee KCS, Adamek L, Rufino B, and Josleyn J are employees of Sanofi and may hold stocks /shares in Sanofi. Hood D, Walia K, Poddar S, and Mankin B are employees of Axtria, Inc. and provided consulting services to Sanofi Kshirsagar O was an employee of Axtria at the time of the study