

SynthA1c: Towards Clinically Interpretable Patient Representations for Diabetes Risk Stratification

Michael S. Yao^{†,1,2}[0000-0002-7008-6028], Allison Chae^{†,2}[0000-0001-7029-556X],
 Matthew T. MacLean³[0000-0002-0514-7218], Anurag
 Verma⁴[0000-0002-5063-9107], Jeffrey Duda³[0000-0002-5031-5735], James C.
 Gee³[0000-0002-2258-0187], Drew A. Torigian³[0000-0001-8999-9735], Daniel
 Rader⁴[0000-0002-9245-9876], Charles E. Kahn Jr.^{2,3}[0000-0002-6654-7434], Walter
 R. Witschey^{‡,2,3}[0000-0003-1669-2120], and Hersh
 Sagreiya^{‡,2,3,*}[0000-0002-2909-6793]

- ¹ Department of Bioengineering, University of Pennsylvania, Philadelphia PA 19104, USA
² Perelman School of Medicine, University of Pennsylvania, Philadelphia PA 19104, USA
³ Department of Radiology, University of Pennsylvania, Philadelphia PA 19104
⁴ Department of Medicine, University of Pennsylvania, Philadelphia PA 19104

Abstract. Early diagnosis of Type 2 Diabetes Mellitus (T2DM) is crucial to enable timely therapeutic interventions and lifestyle modifications. As the time available for clinical office visits shortens and medical imaging data become more widely available, patient image data could be used to opportunistically identify patients for additional T2DM diagnostic workup by physicians. We investigated whether image-derived phenotypic data could be leveraged in tabular learning classifier models to predict T2DM risk in an automated fashion to flag high-risk patients *without* the need for additional blood laboratory measurements. In contrast to traditional binary classifiers, we leverage neural networks and decision tree models to represent patient data as ‘SynthA1c’ latent variables, which mimic blood hemoglobin A1c empirical lab measurements, that achieve sensitivities as high as 87.6%. To evaluate how SynthA1c models may generalize to other patient populations, we introduce a novel generalizable metric that uses vanilla data augmentation techniques to predict model performance on input out-of-domain covariates. We show that image-derived phenotypes and physical examination data together can accurately predict diabetes risk as a means of opportunistic risk stratification enabled by artificial intelligence and medical imaging. Our code is available at <https://github.com/allisonjchae/DMT2RiskAssessment>.

Keywords: Disease Prediction · Representation Learning · Radiomics.

* Corresponding author: hersh.sagreiya@penncmedicine.upenn.edu

† Denotes equal contribution. ‡ Denotes equal contribution.

1 Introduction

Type 2 Diabetes Mellitus (T2DM) affects over 30 million patients in the United States, and is most commonly characterized by elevated serum hemoglobin A1c (HbA1c) levels measured through a blood sample [1,2]. Formally, a patient is considered diabetic if their HbA1c is greater than 6.5% A1c. While patients diagnosed with T2DM are at an increased risk of many comorbidities, early diagnosis and lifestyle interventions can improve patient outcomes [3].

However, delayed diagnosis of T2DM is frequent due to a low rate of screening. Up to a third of patients are not screened for T2DM as recommended by current national guidelines [4,5], and Porter et al. [6] estimate that it would take over 24 hours per day for primary care physicians to follow national screening recommendations for every adult visit. Furthermore, T2DM screening using patient bloodwork is not routinely performed in acute urgent care settings or emergency department (ED) visits. Given these obstacles, machine learning (ML) is a promising tool to predict patient risk of T2DM and other diseases [7].

Simultaneously, the usage of radiologic imaging in clinical medicine continues to increase every year [8,9]. Over 70 million computed tomography (CT) scans are performed annually and their utilization has become increasingly common in both primary care and ED visits [10]. Consequently, the wealth of CT radiographic data can potentially be used to estimate patient risk of T2DM as an incidental finding in these clinical settings. For example, T2DM risk factors include central adiposity and the buildup of excess fat in the liver that can be readily estimated from clinical CT scans. Liver fat excess can be estimated by calculating the spleen-hepatic attenuation difference (SHAD), which is the difference between liver and spleen CT attenuation [11]. These metrics are examples of **image-derived phenotypes** (IDPs) derived from patient CT scans and other imaging modalities. Other IDPs, such as volume estimation of subcutaneous fat and visceral fat, can also be used to quantify central adiposity. Using these metrics, a prediction model could report estimated T2DM risk as an incidental finding during an unrelated outpatient imaging study or ED visit workup as a means of opportunistic risk stratification from analysis of CT scans and patient information, with automated referral of high-risk patients for downstream screening without the need for intermediate physician intervention.

Existing machine learning methods for disease prediction have largely focused on developing classification models that output probability values for different physiologic states [12,13,14]. However, these metrics are difficult for clinicians to interpret at face value and cannot be intelligently integrated into existing clinician workflows, such as diagnostic pathways based on clinical lab findings [15].

In this study, we hypothesized that radiomic metrics derived from CT scans could be used in conjunction with physical examination data to predict patient T2DM risk using SynthA1c, a novel synthetic *in silico* measurement approximating patient blood hemoglobin A1c (HbA1c) (Fig. 1). To predict model generalizability, we also propose a generalizable data augmentation-based model smoothness metric that predicts SynthA1c accuracy on previously unseen out-of-domain patient datasets.

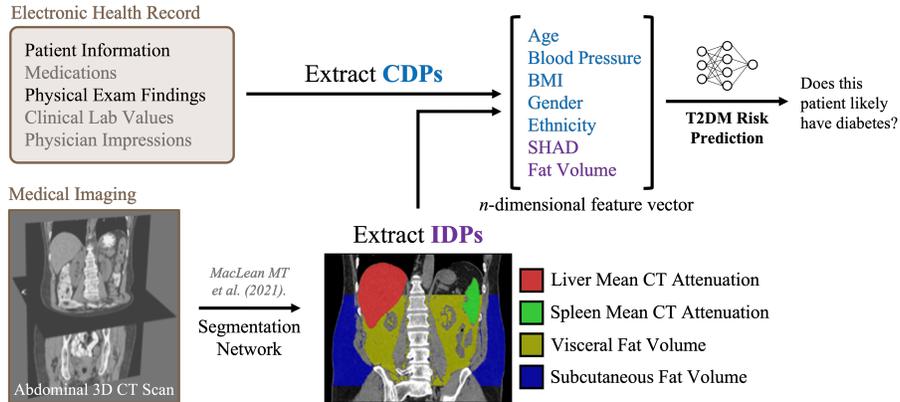


Fig. 1. Overview of our proposed work. Using the IDP extraction pipeline from MacLean et al. [11], we can estimate quantitative IDPs from abdominal CT scans associated with an increased risk for T2DM. IDPs and CDPs from corresponding patient electronic health records can then be used to train T2DM risk prediction models. *IDP*: image-derived phenotype; *CT*: computed tomography; *T2DM*: Type 2 Diabetes Mellitus; *CDP*: clinically derived phenotype.

2 Materials and Methods

2.1 Patient Cohort and Data Declaration

The data used for our retrospective study were made available by the Penn Medicine BioBank (PMBB), an academic biobank established by the University of Pennsylvania. All patients provided informed consent to utilization of de-identified patient data, which was approved by the Institutional Review Board of the University of Pennsylvania (IRB protocol 813913). From the PMBB outpatient dataset, we obtained patient ages, genders, ethnicities, heights, weights, blood pressures, abdominal CT scans, and blood HbA1c measurements. Notably, the only laboratory value used was HbA1c as a ground truth metric in model training and evaluation—no blood biomarkers were used as model inputs.

From the PMBB outpatient dataset, we obtained patient ages, genders, ethnicities, heights, weights, blood pressures, abdominal CT scans, and blood HbA1c measurements. Notably, the only clinical laboratory value extracted was HbA1c to be used as a ground truth—no blood biomarkers were used as model inputs. Patients with any missing features were excluded.

Using the pre-trained abdominal CT segmentation network trained and reported by MacLean et al. [11], we estimated four IDPs from any given CT of the abdomen and pelvis study (either with or without contrast) to be used as model inputs. Our four IDPs of interest were mean liver CT attenuation, mean spleen CT attenuation, and estimated volume of subcutaneous fat and visceral fat. Briefly, their segmentation network achieved mean Sørensen-Dice coefficients

of at least 98% for all IDP extraction tasks assessed (including our four IDPs of interest) and is detailed further in their work.

Any patient i has a set of measured values of any particular feature within the dataset. To construct a feature vector \mathbf{x} associated with an HbA1c measurement y_i , we selected the patient’s measurements that minimized the time between the date the feature was measured and the date y_i was measured.

2.2 Machine Learning Models: GBDT, NODE, and FT-Transformer

Current supervised methods for disease detection work with feature vectors derived from patient physical examinations and clinical laboratory values [7,12,13]. Our work builds on these prior advances by incorporating IDPs as additional input vector dimensions. Previously, Chen and Guestrin [16] introduced gradient-boosted decision trees (**GBDTs**) that incorporate scalable gradient boosting with forest classifiers for state-of-the-art prediction accuracy across tasks. A separate class of machine learning models is deep neural networks (DNNs). Recently, neural oblivious decision ensemble (**NODE**) DNNs achieved classification performance on par with decision tree models on certain tasks [17] and the Feature Tokenizer + Transformer (**FT-Transformer**) [18] effectively adopts transformer architectures to tabular data. Here, we assessed NODE, FT-Transformer, and GBDT architectures as backbones for our SynthA1c encoders.

We sought to compare our proposed SynthA1c models against a number of baselines. We looked at Ordinary Least Squares (OLS) encoders and traditional diabetes *binary classifier* models with the same three architectures as proposed above, in addition to a zero-rule classifier and a multi-rule questionnaire-based classifier currently recommended for clinical practice by the American Diabetes Association and Centers for Disease Control and Prevention [19].

2.3 Model Training and Evaluation Strategy

Our model inputs can be divided into two disjoint sets: clinically derived phenotypes (CDPs), which are derived from physical examination, and image-derived phenotypes (IDPs) that are estimated from abdominal CT scans herein. The specific CDPs and IDPs used depended on the model class—broadly, we explored two categories of models, which we refer to as r -type and p -type. r -type models were trained on ‘raw’ data types (CDPs: height, weight, race, gender, age, systolic blood pressure [SBP], diastolic blood pressure [DBP]; IDPs: liver CT attenuation, spleen CT attenuation, subcutaneous fat [SubQ Fat], visceral fat [Visc Fat]), while p -type models were trained on ‘processed’ data types (CDPs: BMI, race, gender, age, SBP, DBP; IDPs: SHAD, SubQ Fat, Visc Fat). Comparing the performance of r - and p -type models could help us better understand if using derivative processed metrics that are better clinically correlated with T2DM risk yields better model performance.

SynthA1c encoders were trained to minimize the L_2 distance from the ground truth HbA1c laboratory measurement, and evaluated using the root mean square error (RMSE) and Pearson correlation coefficient (PCC). We then compared

the predicted SynthA1c values with the traditional HbA1c $\geq 6.5\%$ A1c diabetes cutoff to assess the utility of SynthA1c outputs in diagnosing T2DM. A p value of $p < 0.05$ was used to indicate statistical significance.

2.4 Implementation Details

NODE models were trained with a batch size of 16 and a learning rate of $\eta = 0.03$, which decayed by half every 40 epochs for a total of 100 epochs. FT-Transformer models were trained with a batch size of 128 and a learning rate of $\eta = 0.001$, which decayed by half every 50 epochs for a total of 100 epochs. GBDT models were trained using 32 boosted trees with a maximum tree depth of 8 with a learning rate of $\eta = 0.1$.

2.5 Assessing Out-of-Domain Performance

An important consideration in high-stakes clinical applications of machine learning is the generalizability of T2DM classifiers to members of previously unseen patient groups. Generalizability is traditionally difficult to quantify and can be affected by training data heterogeneity and the geographic, environmental, and socioeconomic variables unique to the PMBB dataset.

Prior work has shown that model smoothness can be used to predict out-of-domain generalization of neural networks [20,21]. However, these works largely limit their analysis to classifier networks. To evaluate SynthA1c encoder robustness, we develop an estimation of model manifold smoothness \mathbb{M} for our encoder models. Under the mild assumption that our SynthA1c encoder function $y : \mathbb{R}^{|\mathbf{x}|} \rightarrow \mathbb{R}$ is Lipschitz continuous, we can define a local manifold smoothness metric μ at $\mathbf{x} = \tilde{\mathbf{x}}$ given by

$$\begin{aligned} \mu(\tilde{\mathbf{x}}) &= \mathbb{E}_{\mathcal{N}(\tilde{\mathbf{x}})} \left[\frac{\sigma_y^{-1} \|y(\mathbf{x}) - y(\tilde{\mathbf{x}})\|_1}{\|\delta\mathbf{x} \oslash \sigma_{\mathbf{x}}\|_2} \right] \\ &= \mathcal{V}[\mathcal{N}(\tilde{\mathbf{x}})]^{-1} \cdot \oint_{\mathcal{N}(\tilde{\mathbf{x}}) \in \mathcal{D}} d\mathbf{x} \frac{\sigma_y^{-1} |y(\mathbf{x}) - y(\tilde{\mathbf{x}})|}{\sqrt{(\delta\mathbf{x} \oslash \sigma_{\mathbf{x}})^T (\delta\mathbf{x} \oslash \sigma_{\mathbf{x}})}} \end{aligned} \quad (1)$$

where we have a feature vector \mathbf{x} in domain \mathcal{D} and a neighborhood $\mathcal{N}(\tilde{\mathbf{x}}) \in \mathcal{D}$ around \mathbf{x} with an associated volume of $\mathcal{V}[\mathcal{N}(\tilde{\mathbf{x}})]$. We also define $\delta\mathbf{x} = \mathbf{x} - \tilde{\mathbf{x}}$, \oslash as the Hadamard division operator, and $\sigma_{\mathbf{x}}$ as the vector of the estimated standard deviations of each feature over \mathcal{D} . The exact expectation value over a given neighborhood $\mathcal{N}(\tilde{\mathbf{x}})$ is computationally intractable, but we can approximate it with a Monte Carlo integration through an empirical sampling of $Q \gg 1$ random feature points \mathbf{x}_k from $\mathcal{N}(\tilde{\mathbf{x}})$:

$$\mu(\tilde{\mathbf{x}}) = \frac{1}{Q} \sum_{k=1}^Q \frac{\sigma_y^{-1} |y(\mathbf{x}_k) - y(\tilde{\mathbf{x}})|}{\sqrt{(\delta\mathbf{x}_k \oslash \sigma_{\mathbf{x}})^T (\delta\mathbf{x}_k \oslash \sigma_{\mathbf{x}})}} \quad (2)$$

We can now define a metric \mathbb{M} for the global encoder manifold smoothness over a domain \mathcal{D} as the expectation value of $\mu(\tilde{\mathbf{x}})$ over \mathcal{D} , which can similarly be approximated by an empirical sampling of N feature vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \in \mathcal{D}$. We hypothesized that this global smoothness metric \mathbb{M} inversely correlates with model performance on out-of-domain datasets. To evaluate this experimentally, we assessed model performance on two previously unseen T2DM datasets: (1) the Iraqi Medical City Hospital dataset [22]; and (2) the PMBB inpatient dataset. The Iraqi dataset contains 1,000 sets of patient age, gender, BMI, and HbA1c measurements. Because of this limited feature set, we trained additional SynthA1c encoders (referred to as p -type models) on the PMBB outpatient dataset using only these features. The PMBB inpatient dataset consists of 2,066 measurements of the same datatypes as the outpatient dataset (Section 3.1).

3 Results

3.1 Summary Statistics

Our model-building dataset from the PMBB consisted of 2,077 unique HbA1c measurements (1,159 diabetic, 619 prediabetic, 299 nondiabetic) derived from 389 patients (Table 1). 208 (10%) samples were set aside as a holdout test set partition disjoint by patient identity. Each HbA1c measurement was used to construct an associated feature vector from that patient’s data collected closest in time to each HbA1c measurement. To quantify the temporal association between a given patient’s measurements, we defined the daterange of an observation vector \mathbf{x} as the maximum length of time between any two features/imaging studies. The median daterange in our dataset was 18 days.

3.2 SynthA1c Encoder Experimental Results

Our results suggest that the GBDT encoder predicted SynthA1c values closest to ground truth HbA1c values, followed by both the NODE and FT-Transformer DNN models (Table 2). All the learning-based architectures assessed outperformed the baseline OLS encoder. When comparing SynthA1c outputs against the clinical HbA1c cutoff of 6.5% A1c for the diagnosis of diabetes, the r -GBDT SynthA1c model demonstrated the highest sensitivity of the assessed models at 87.6% on par with the best-performing binary classifier model assessed. In terms of an opportunistic screening modality for T2DM, a high sensitivity ensures that a large proportion of patients with diabetes can be identified for additional lab-based diagnostic work-up with their primary care physicians. Although the accuracy of SynthA1c encoders was lower than the corresponding binary classifier models assessed, this may be partially explained by the fact that the latter’s threshold value for classification was empirically tuned to maximize the model’s accuracy. In contrast, our SynthA1c encoders used the fixed clinical HbA1c cutoff of 6.5% A1c for diabetes classification. When comparing r - and p -type SynthA1c models, we did not observe a consistently superior data representation strategy.

Table 1. PMBB outpatient dataset characteristics. To reduce the effects of selection bias, all patients presenting to the University of Pennsylvania Health System were given the opportunity to enroll in the PMBB so as to best capture the population of patients that seek medical care and avoid overrepresentation of healthy patients as in traditional office visit patient recruitment strategies. However, the PMBB is still affected by hesitations of patient sub-populations in study enrollment and the unique socioeconomic factors affecting different groups of patients. *HTN*: Hypertension.

Self-Reported Ethnicity	Count (%)
White	720 (34.7)
Hispanic	40 (1.9)
Black	1248 (60.1)
Asian	36 (1.7)
Pacific Islander	6 (0.3)
Native American	5 (0.2)
Other/Unknown	22 (1.1)
Self-Reported Gender	Count (%)
Male	880 (42.4)
Female	1197 (57.6)
Age Decade	Count (%)
20-29	31 (1.5)
30-39	89 (4.3)
40-49	362 (17.4)
50-59	593 (28.6)
60-69	680 (32.7)
70-79	299 (14.4)
80-89	23 (1.1)
Blood Pressure	Count (%)
Normal (SBP < 120 mmHg and DBP < 80 mmHg)	421 (20.2)
Elevated (120 ≤ SBP < 130 mmHg and DBP < 80 mmHg)	398 (19.2)
Stage 1 HTN (130 ≤ SBP < 140 mmHg or 80 ≤ DBP < 90 mmHg)	652 (31.4)
Stage 2 HTN (SBP ≥ 140 mmHg or DBP ≥ 90 mmHg)	606 (29.2)
BMI	Count (%)
Underweight or Healthy Weight (BMI < 25.0)	275 (13.2)
Overweight (25.0 ≤ BMI < 30.0)	443 (21.3)
Class 1 Obesity (30.0 ≤ BMI < 35.0)	556 (26.8)
Class 2 Obesity (35.0 ≤ BMI < 40.0)	389 (18.7)
Class 3 Obesity (BMI ≥ 40.0)	414 (20.0)
HbA1c	Count (%)
Not Diabetic (HbA1c < 6.5% A1c)	918 (44.2)
Diabetic (HbA1c ≥ 6.5% A1c)	1159 (55.8)
CT Abdomen and Pelvis Enhancement	Count (%)
With Contrast	1570 (75.6)
Without Contrast	507 (24.4)
Image Derived Phenotypes (IDPs) Statistics	Mean ± SD
Spleen CT Attenuation (HU)	36.2 ± 16.7
Liver CT Attenuation (HU)	42.8 ± 20.2
Subcutaneous Fat Area (cm ²)	321.3 ± 170.1
Visceral Fat Area (cm ²)	172.4 ± 104.9
Total Count	2077

Table 2. SynthA1c prediction results using different encoder models. *r*- (*p*-) prefixed models are fed raw (processed) inputs as outlined in Section 2.3. RMSE in units of % A1c. For the SynthA1c encoder models, recall, precision, specificity, and accuracy metrics are reported based on the traditional T2DM cutoff of 6.5% A1c. The Multi-Rule binary classifier is the current deterministic risk stratification tool recommended by American Diabetes Association [19].

SynthA1c Encoder	RMSE	PCC	Recall	Precision	Specificity	Accuracy
<i>r</i> -OLS	1.67	0.206	85.3	56.0	26.3	57.2
<i>p</i> -OLS	1.73	0.159	80.7	57.5	34.3	58.6
<i>r</i> -FT-Transformer	1.44	0.517	87.6	63.4	55.9	70.7
<i>p</i> -FT-Transformer	1.51	0.441	83.5	61.4	54.1	67.8
<i>r</i> -NODE	1.60	0.378	85.6	55.0	38.7	60.6
<i>p</i> -NODE	1.57	0.649	77.3	59.5	54.1	64.9
<i>r</i> -GBDT	1.36	0.567	87.2	66.4	51.5	70.2
<i>p</i> -GBDT	1.36	0.591	77.1	72.4	67.7	72.6

Binary Classifier	AUROC (%)	Recall	Precision	Specificity	Accuracy
Zero-Rule	—	100	52.4	0.0	52.4
Multi-Rule	56.3	67.0	54.9	39.4	53.8
<i>r</i> -FT-Transformer	82.1	85.3	73.8	66.7	76.4
<i>r</i> -NODE	83.5	82.6	76.9	72.7	77.9
<i>r</i> -GBDT	83.1	87.2	76.6	70.7	79.3

To further interrogate our SynthA1c encoders, we investigated whether model performance varied as a function of demographic features. Defining the difference between the model prediction and ground truth HbA1c values as a proxy for model performance, all SynthA1c encoders showed no statistically significant difference in performance when stratified by gender or BMI (Fig. 2).

3.3 Ablation Studies: Relative Importance of CDPs and IDPs

Until now, prior T2DM classifiers have used only blood lab measurements and physical examination data to predict T2DM. In contrast, our models presented herein are the first to incorporate IDPs as input model features for the task of diabetes risk stratification. To better understand the benefit and value-add of using IDPs in conjunction with CDPs, we evaluated classifier performance on models trained using either only CDPs or only IDPs and compared them to corresponding models trained using both input types.

Our results suggest that while classifier models trained only on CDPs generally outperform those trained only on IDPs, the best performance is achieved when combining CDPs and IDPs together (Table 3). This further validates the clinical utility of incorporating IDPs into patient diagnosis and disease risk stratification first proposed by MacLean et al. [11] and related work.

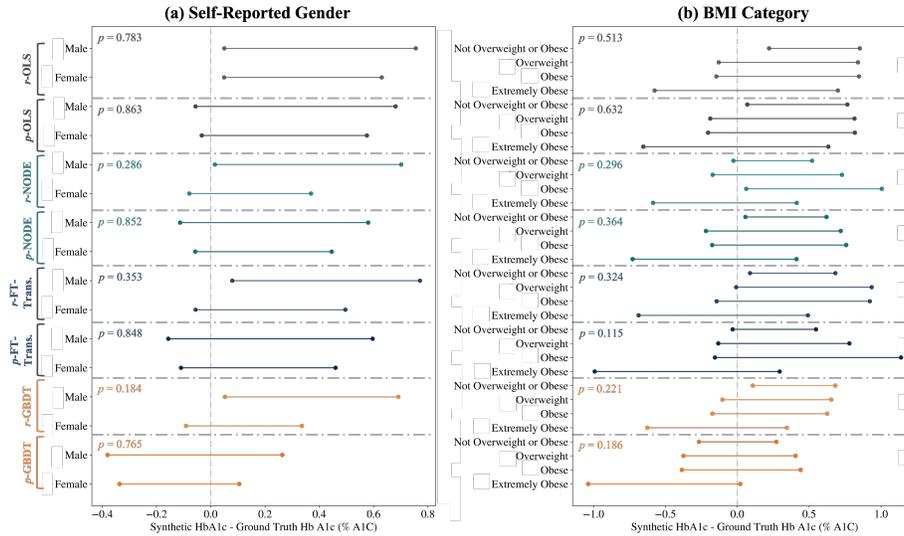


Fig. 2. Assessing for algorithmic bias in SynthA1c encoders. We plotted the 95% confidence interval of the mean difference between the SynthA1c model output and ground truth HbA1c as a function of self-reported (a) gender and (b) BMI category. p values comparing the differences in SynthA1c model performance when stratified by gender (two-sample T-test) and BMI category (one-way ANOVA) are shown.

3.4 Characterizing Out-of-Domain Model Performance

As our proposed global smoothness metric \mathbb{M} decrease across the three evaluated models, the RMSE in SynthA1c prediction decreases and the PCC increases, corresponding to better predictive performance on the out-of-domain Iraqi Medical Center Dataset (Table 4). This supports our initial hypothesis that smoother models may generalize better to unseen datasets. We also noted larger RMSE values using the Iraqi Medical Center Dataset when compared to the PMBB outpatient test dataset results (Table 2).

Interestingly, we found that this relationship did not ostensibly hold when considering the PMBB inpatient dataset; in fact, model predictive performance was *inversely* correlated with global smoothness. This suggested that the PMBB inpatient and outpatient dataset distributions were more similar than initially predicted. To validate this hypothesis, we computed the Kullback-Leibler (KL) divergence between each of the test dataset distributions and the training dataset distribution with respect to the features available in all datasets: ethnicity, gender, age, BMI, and HbA1c. We assumed the PMBB-derived outpatient training dataset was sampled from a distribution \mathcal{Q} and each of the PMBB outpatient test, PMBB inpatient, and Iraqi medical center datasets were sampled from $\mathcal{P}_{\text{Outpatient}}$, $\mathcal{P}_{\text{Inpatient}}$, and $\mathcal{P}_{\text{Iraqi}}$, respectively. The greatest KL divergence was between the Iraqi medical center and training dataset distributions, as expected ($D_{KL}[\mathcal{P}_{\text{Iraqi}}||\mathcal{Q}] = 31.2$). Despite the fact that our training set included outpa-

Table 3. Ablation study assessing model performance as a function of clinically derived phenotypes (CDPs) and/or image-derived phenotypes (IDPs).

<i>r</i> - NODE	Recall	Precision	Specificity	Accuracy
CDPs Only	77.1	73.7	69.7	73.5
IDPs Only	73.4	76.9	75.8	74.5
CDPs + IDPs	82.6	76.9	72.7	77.9
<i>r</i> - FT-Transformer	Recall	Precision	Specificity	Accuracy
CDPs Only	78.0	76.6	73.7	75.9
IDPs Only	71.6	60.5	48.5	60.6
CDPs + IDPs	85.3	73.8	66.7	76.4
<i>r</i> - GBDT	Recall	Precision	Specificity	Accuracy
CDPs Only	80.7	68.6	59.6	70.7
IDPs Only	73.4	75.5	73.7	73.6
CDPs + IDPs	87.2	76.6	70.7	79.3

Table 4. SynthA1c model sensitivity and out-of-domain generalization results. Global smoothness metric values \mathbb{M} were evaluated on the PMBB outpatient dataset. *r*-type models could not be evaluated on the Iraqi dataset because IDPs and medical imaging data were not available. RMSE in units of % A1c.

SynthA1c Encoder	\mathbb{M}	Iraqi Dataset		PMBB Inpatient	
		RMSE	PCC	RMSE	PCC
<i>p</i> '-/ <i>r</i> - NODE	1.43	3.62 / —	0.154 / —	1.76 / 1.23	0.512 / 0.795
<i>p</i> '-/ <i>r</i> - FT-Transformer	1.07	3.04 / —	0.246 / —	1.90 / 1.58	0.331 / 0.617
<i>p</i> '-/ <i>r</i> - GBDT	3.28	6.25 / —	0.021 / —	1.54 / 1.12	0.674 / 0.823

tient data alone, we found the KL divergence between the inpatient test and training datasets ($D_{KL}[\mathcal{P}_{\text{Inpatient}}||\mathcal{Q}] = 0.227$) was lower than that between the outpatient test and training dataset ($D_{KL}[\mathcal{P}_{\text{Outpatient}}||\mathcal{Q}] = 1.84$).

To further characterize the feature distributions within our datasets, we analyzed the pairwise relationships between BMI, age, and HbA1c. Individual feature distributions were statistically significant between either of the PMBB datasets and the Iraqi Medical Center dataset (two-sample Kolmogorov-Smirnov (KS) test; $p < 0.0001$ between [PMBB inpatient dataset, Iraqi Medical Center dataset] and [PMBB outpatient dataset, Iraqi Medical Center dataset] pairs for individual age, HbA1c, and BMI quantitative features), but not between the PMBB inpatient and outpatient datasets (two-sample KS test; age: $p = 0.315$, HbA1c: $p = 0.463$, BMI: $p = 0.345$). These results suggest that inpatients are a compact subset of outpatients within the PMBB with respect to T2DM risk assessment. This helps explain our initial findings regarding the relationship between \mathbb{M} and model generalization. Further work is warranted to validate the proposed metric \mathbb{M} across other tasks.

4 Conclusion

Our work highlights the value of using CT-derived IDPs and CDPs for opportunistic screening of T2DM. We show that tabular learning architectures can act as novel SynthA1c encoders to predict HbA1c measurements noninvasively. Furthermore, we demonstrate that model manifold smoothness may be correlated with prediction performance on previously unseen data sampled from out-of-domain patient populations, although additional validation studies on separate tasks are needed. Ultimately, we hope that our proposed work may be used in every outpatient and ED imaging study, regardless of chief complaint, for opportunistic screening of type 2 diabetes. Our proposed SynthA1c methodology will by no means replace existing diagnostic laboratory workups, but rather identify those at-risk patients who should consider consulting their physician for downstream clinical evaluation in an efficient and automated manner.

Acknowledgments MSY is supported by NIH T32 EB009384. AC is supported by the AQA Carolyn L. Kuckein Student Research Fellowship and the University of Pennsylvania Diagnostic Radiology Research Fellowship. WRW is supported by NIH R01 HL137984. MTM received funding from the Sarnoff Cardiovascular Research Foundation. HS received funding from the RSNA Scholar Grant.

References

1. Khan, M.A.B., Hashim, M.J., King, J.K., Govender, R.D., Mustafa, H., Al Kaabi, J.: Epidemiology of type 2 diabetes - Global burden of disease and forecasted trends. *J Epi Glob Health* **10**(1), 107-11 (2020). <https://doi.org/10.2991/jegh.k.191028.001>
2. Xu, G., Liu, B., Sun, Y., Du, Y., Snetselaar L.G., Hu F.B., Bao W.: Prevalence of diagnosed type 1 and type 2 diabetes among US adults in 2016 and 2017: Population based study. *BMJ* **362** (2018). <https://doi.org/10.1136/bmj.k1497>
3. Albarakat, M., Guzu, A.: Prevalence of type 2 diabetes and their complications among home health care patients at Al-Kharj military industries corporation hospital. *J Family Med Prim Care* **8**(10), 3303-12 (2019). <https://doi.org/10.4103/jfmpe.jfmpe.634.19>
4. Polubriaginof, F.C.G., Shang, N., Hripsak, G., Tatonetti, N.P., Vawdrey, D.K.: Low screening rates for diabetes mellitus among family members of affected relatives. *AMIA Annu Symp Proc*, 1471-7 (2019). PMID: 30815192
5. Kaul, P., Chu, L.M., Dover, D.C., Yeung, R.O., Eurich, D.T., Butalia, S.: Disparities in adherence to diabetes screening guidelines among males and females in a universal care setting: A population-based study of 1,380,697 adults. *Lancet Regional Health*, (2022). <https://doi.org/10.1016/j.lana.2022.100320>
6. Porter, J., Boyd, C., Skandari, M.R., Laiteerapong, N.: Revisiting the time needed to provide adult primary care. *J Gen Intern Med*, (2022). <https://doi.org/10.1007/s11606-022-07707-x>
7. Farran, B., Channanath, A.M., Behbehani, K., Thanaraj, T.A.: Predictive models to assess risk of type 2 diabetes, hypertension and comorbidity: Machine-learning algorithms and validation using national health data from Kuwait—A cohort study. *BMJ Open* **3**(5), (2013). <https://doi.org/10.1136/bmjopen-2012-002457>

8. Dowhanik, S.P.D., Schieda, N., Patlas, M.N., Salehi, F., van der Pol, C.B.: Doing more with less: CT and MRI utilization in Canada 2003-2019. *Canadian Association of Radiologists J* **73**(3), 592-4 (2022). <https://doi.org/10.1177/08465371211052012>
9. Hong, A.S., Levin, D., Parker, L., Rao, V.M., Ross-Degnan, D., Wharam, J.F.: Trends in diagnostic imaging utilization among Medicare and commercially insured adults from 2003 through 2016. *Radiology* **294**(2), 342-50 (2020). <https://doi.org/10.1148/radiol.2019191116>
10. Smith, J.W., Everhart, J.E., Dickson, W.C., Knowler, W.C., Johannes, R.S.: Using the ADAP learning algorithm to forecast the onset of diabetes mellitus. *Proc Sym Comp App and Med Care*, 261-5 (1988)
11. MacLean, M.T., Jehangir, Q., Vujkovic, M., Ko, Y., Litt, H., Borthakur, A., Sagreiya, H., Rosen, M., Mankoff, D.A., Schnall, M.D., Shou, H., Chirinos, J., Damrauer, S.M., Torigian, D.A., Carr, R., Rader, D.J., Witschey, W.R.: Quantification of abdominal fat from computed tomography using deep learning and its association with electronic health records in an academic biobank. *J Am Med Inform Assoc* **28**(6), 1178-87 (2021). <https://doi.org/10.1093/jamia/ocaa342>
12. Uddin, S., Khan, A., Hossain, M.E., Moni, M.A.: Comparing different supervised machine learning algorithms for disease prediction. *BMC Med Inform Decis Mak* **19**(281), (2019). <https://doi.org/10.1093/jamia/ocaa342>
13. Kopitar, L., Kocbek, P., Cilar, L., Sheikh, A., Stiglic, G.: Early detection of type 2 diabetes mellitus using machine learning-based prediction models. *Nat Sci Rep*, (2020). <https://doi.org/10.1038/s41598-020-68771-z>
14. Deberneh, H.M., Kim, I.: Prediction of type 2 diabetes based on machine learning algorithm. *Int J Environ Res Public Health* **18**(6), 3317 (2021). <https://doi.org/10.3390/ijerph18063317>
15. Sivaraman, V., Bukowski, L.A., Levin, J., Kahn, J.M., Perer, A.: Ignore, trust, or negotiate: Understanding clinician acceptance of AI-based treatment recommendations in health care. *arXiv*, (2023). <https://doi.org/10.48550/arxiv.2302.00096>
16. Chen, T., Guestrin, C.: XGBoost: A scalable tree boosting system. In: *Proc ACM SIGKDD Int Conf on Knowledge Discovery and Data Mining*, 785-94 (2016). <https://doi.org/10.1145/2939672.2939785>
17. Popov, S., Morozov, S., Babenko, A.: Neural oblivious decision ensembles for deep learning on tabular data. *arXiv*, (2019). <https://doi.org/10.48550/arxiv.1909.06312>
18. Gorishniy, Y., Rubachev, I., Khruklov, V., Babenko, A.: Revisiting deep learning models for tabular data. *arXiv*, (2021). <https://doi.org/10.48550/arxiv.2106.11959>
19. Bang, H., Edwards, A.M., Bomback, A.S., Ballantyne, C.M., Brillon, D., Callahan, M.A., Teutsch, S.M., Mushlin, A.I., Kern, L.M.: Development and validation of a patient self-assessment score for diabetes risk. *Ann Intern Med* **151**(11), 775-83 (2009). <https://doi.org/10.7326/0003-4819-151-11-200912010-00005>
20. Ng, N., Hulkund, N., Cho, K., Ghassemi, M.: Predicting out-of-domain generalization with local manifold smoothness. *arXiv*, (2022). <https://doi.org/10.48550/arxiv.2207.02093>
21. Jiang, Z., Zhou, J., Huang, H.: Relationship between manifold smoothness and adversarial vulnerability in deep learning with local errors. *Chinese Physics B* **30**(4), (2021). <https://doi.org/10.1088/1674-1056/abd68e>
22. Rashid, A.: Iraqi Diabetes Dataset. 2020. Available from: [data.mendeley.com/datasets/wj9rwkp9c2/1](https://doi.org/10.17632/wj9rwkp9c2/1). <https://doi.org/10.17632/wj9rwkp9c2.1>