

1 **Title**

2 A data-driven approach to identify clusters of HbA1c longitudinal trajectories and associated
3 outcomes in type 2 diabetes mellitus: a large population-based cohort study

4 **Author Names and Affiliations**

5 Adrian Martinez-De la Torre¹, M.Sc

6 Maria Luisa Faquetti¹, M.Sc

7 Fernando Perez-Cruz^{2,3}, Ph.D

8 Christian Meier⁴, MD

9 Stefan Weiler^{1,5}, MD, Ph.D

10 Andrea M. Burden¹, Ph.D

11

12 ¹ Institute of Pharmaceutical Sciences, Department of Chemistry and Applied Biosciences, ETH Zurich,
13 Zurich, Switzerland.

14 ² Swiss Data Science Center, ETH Zurich and EPFL, Switzerland

15 ³ Institute of Machine Learning in the Computer Science Department at ETH Zurich

16 ⁴ Department of Endocrinology, Diabetology and Metabolism, University Hospital Basel, Basel,
17 Switzerland.

18 ⁵ Clinical Pharmacology and Toxicology, Department of General Internal Medicine, Inselspital, Bern
19 University Hospital, University of Bern, Bern, Switzerland

20

21

22

23 **Corresponding Author**

24 Prof. Dr. Andrea Burden

25 Institute of Pharmaceutical Sciences, Department of Chemistry and Applied Biosciences

26 ETH Zurich

27 Vladimir-Prelog-Weg 1-5/10

28 8093 Zurich

29 Switzerland

30

31

32

33 **NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.**

34 Abstract

35 Title

36 A data-driven approach to identify clusters of HbA1c longitudinal trajectories and associated
37 outcomes in type 2 diabetes mellitus: a large population-based cohort study

38 Background

39 We aimed to identify and characterize common patterns of HbA1c progression among type 2
40 diabetes mellitus patients who initiate a non-insulin antidiabetic drug (NIAD).

41 Methods

42 The IQVIA Medical Research Data incorporating data from THIN, a Cegedim database of
43 anonymized electronic health records, was used to identify a cohort of patients with a first-ever
44 prescription for a NIAD between 2006 and 2019. Trajectory clusters were identified using an
45 Expectation-Maximization algorithm by iteratively fitting k thin-plate splines and reassigning
46 each patient to the nearest cluster. Cox proportional hazards models calculated the hazard ratios
47 (HR) and 95% confidence intervals (CI) for the estimated risk of microvascular (e.g.,
48 retinopathy, diabetic polyneuropathy [DPN]) and macrovascular events.

49 Findings

50 Among 116,251 new users of NIADs we found five distinct clusters of HbA1c progression,
51 which were characterized as: optimally controlled (**OC**), adequately controlled (**AC**), sub-
52 optimally controlled (**SOC**), poorly controlled (**PC**), and uncontrolled (**UC**). The UC and AC
53 clusters had similar index HbA1C (>9%) but the AC cluster achieved HbA1c control (HbA1C
54 <7.5%), while the UC cluster HbA1c remained >9.0%. Compared to the OC cluster, there was
55 a 21% (HR: 1.21, 95% CI: 1.14-1.28) and 30% (HR: 1.30, 95% CI: 1.21-1.40) elevated risk of
56 retinopathy in the AC and UC clusters, respectively. While the PC and UC clusters had a
57 significant 23% (HR 1.23, 95% CI 1.12 – 1.35) and 45% (HR 1.45, 95% CI: 1.27 – 1.64)
58 increased risk of DPN, respectively.

59 Interpretation

60 The five identified HbA1c trajectory clusters had different risk profiles. Despite achieving
61 diabetic control, patients categorized in the AC cluster had similar outcomes to the UC cluster,
62 suggesting baseline HbA1c is an important indicator of health outcomes.

63 Funding

64 The Swiss Data Science Centre

65 Introduction

66 Type 2 diabetes mellitus (T2DM) is a chronic disease arising from the body's inefficient use
67 of insulin or progressive inability to secrete insulin, which results in abnormal blood levels of
68 glucose.[1–3] T2DM is characterized by hyperglycemia, ultimately leading to microvascular
69 (e.g., retinopathy) and macrovascular complications (e.g., cardiovascular disease). Disease
70 management includes lifestyle modification and glucose control using pharmacotherapy (e.g.,
71 non-insulin antidiabetic drugs [NIADs]) to reduce diabetic complications and mortality risk.[4]

72 Glycated hemoglobin (HbA1c) provides a long-term trend of glucose levels in the blood over
73 the last two to three months, and it is often used as a clinical target for glycemic control.
74 Guidelines for disease management frequently consider HbA1c <7% (53mmol/mol) a general
75 target for glucose control.[5–7] However, T2DM has a high degree of heterogeneity in
76 individual patient characteristics leading to different treatment strategies and distinct treatment
77 responses. Therefore, guidelines for T2DM management, such as the National Institute for
78 Health and Care Excellence (NICE) in the UK, recommend a more individualized treatment
79 approach considering patient's preferences and individual characteristics (e.g., age,
80 comorbidities, and multiple medications).[5]

81 Previous studies have shown that hyperglycemia, and therefore elevated HbA1c levels, is
82 associated with the risk of diabetes complications and mortality.[8,9] The United Kingdom
83 Prospective Diabetes Study (UKPDS) identified that intensive glycemic control had important
84 long-term clinical implications on microvascular and macrovascular outcomes.[10] Using data
85 from the Action to Control Cardiovascular Risk in Diabetes (ACCORD) trial, Wang et al.
86 found worse cardiovascular outcomes among patients with persistent poor glycemic control,
87 when compared to patients with HbA1c around 7%. [11] Similarly, using observational data,
88 Liateerapong and colleagues found that 10-year glycemic control was associated with
89 improved microvascular outcomes. [12]

90 Thus, identifying groups of patients with specific HbA1c courses may help to develop more
91 personalized strategies for T2DM management. However, the evaluation of HbA1c
92 progression with time is challenging, particularly due to patients with an unequal length of
93 observations, unevenly spaced in time, and heterogeneous observation windows. Large data
94 and advanced statistical models are required to identify if patient trajectories can be grouped
95 into similar clusters, and if these trajectories affect microvascular and macrovascular outcomes.
96 Thus, in this study we applied an Expectation-Maximization (EM) algorithm using k -means

97 clustering and thin-plate splines to identify distinct HbA1c trajectories in new users of NIADs
98 using a large population-based electronic health record database. Additionally, we evaluate the
99 association between the trajectory clusters and subsequent microvascular and macrovascular
100 events.

101

102 **Methods**

103 *Data source*

104 We used the IQVIA Medical Research Data incorporating data from THIN, a Cegedim
105 database of anonymized electronic health records from general practitioners (GPs). The
106 database is comprised of over 18 million patients, from 800 general practices in the UK and
107 about 6% of the population. THIN provides detailed longitudinal information regarding patient
108 characteristics (e.g., sex, practice registration date, and ethnicity), medical conditions (e.g.,
109 diagnoses with dates, referrals to hospitals, and symptoms), medications (e.g., generic drug
110 name, dose, and prescription date), and additional health data (e.g., laboratory results including
111 HbA1c, creatinine and calcium blood levels, smoking status, height, weight, alcohol use, birth
112 and death dates).

113 The database contains information on drug prescriptions recorded by GPs. Medications are
114 recorded in the database using the British National Formulary (BNF) classification, and then
115 were mapped according to the international anatomical therapeutic codes (ATC) classification
116 system. All diagnoses are recorded using READ codes [13], a comprehensive coding system
117 with over 100,000 codes and comparable to the international classification of diseases (ICD)
118 system. The study protocol was approved by the THIN scientific research council (study
119 reference number: 20SR062).

120

121 *Study population*

122 We identified patients aged 18+ years between January 1st 2006 and December 31st 2019, with
123 new onset T2DM defined as a first-time NIAD prescription. The index date was defined as the
124 date of the first ever NIAD prescription after start of valid data collection. All patients were
125 required to have more than 1-year of eligible data collection. Patients with a diagnosis of
126 polycystic ovarian syndrome or gestational diabetes prior to index date were excluded, since
127 these conditions are often treated with NIADs. Additionally, we excluded patients with insulin

128 prescription prior to, or at, index date and patients with less than four records of HbA1c after
129 index date. For each the time-to-event analysis, we further excluded patients with a diagnosis
130 of the outcome of interest previous to index.

131

132 *Study variables*

133 We assessed variables of interest i.e., age, body mass index (BMI), smoking status, and alcohol
134 consumption at index date. Additionally, we included the following comorbidities, defined by
135 the presence of corresponding diagnostic or test Read codes: angina pectoris, anxiety and other
136 neurotic, stress related and somatoform disorders, arthropathy, atrial fibrillation, cancer,
137 chronic depression, chronic liver disease, congestive heart failure, high blood pressure,
138 hypercholesterolemia, hypothyroidism, irritable bowel syndrome, ischemic heart disease,
139 neuropathy, osteoarthritis, primary open-angle glaucoma, and senile cataract.

140 Moreover, we analyzed main laboratory results which are highly affected by T2DM, stratifying
141 by cluster i.e., estimated glomerular filtration rate (eGFR), bilirubin, vitamin B12, serum iron,
142 low-density lipoprotein (LDL), and triglycerides.

143

144 *Follow-up period*

145 For each patient, there were two exposure periods. The first one was defined as starting at the
146 date of a first-ever NIAD prescription and ending at the end of follow up (December 31st,
147 2019), death, or loss-to-follow up due to disenrollment from GP, whichever occurred first, and
148 it was used classify the type of HbA1c trajectory.

149 The second exposure time started at the date of a first-ever NIAD prescription until the
150 occurrence of any outcome of interest, ad it was used in the time-to-event analysis and the Cox-
151 proportional hazards model.

152

153 *Outcomes of interest*

154 The primary outcomes of interest were classified as microvascular conditions (retinopathy,
155 diabetic polyneuropathy [DPN], and erectile dysfunction [ED]), and macrovascular diseases
156 (acute myocardial infarction [AMI], coronary heart disease [CHD], peripheral arterial disease
157 [PAD]). All chronic disease conditions were identified based on read codes, while insulin use

158 was identified based on ATC codes. All included codes can be found in
159 https://github.com/adrianmartinez-ETH/hba1c_progression.

160

161 *Statistical analysis*

162 Analysis of longitudinal trajectories of HbA1c was conducted employing an Expectation-
163 Maximization (EM) algorithm using k -means clustering and thin-plate splines (TPSs).[14] TPS
164 are functions defined piece-wise by polynomials which are used to model relationships
165 between a predictor X and a variable Y . The functions are fitted using a generalized additive
166 model (GAM), as shown in **Equation 1**

$$167 \quad g(E(Y)) = \beta_0 + f(X) + \lambda \quad (1)$$

168 where β_0 is a constant, $f(X)$ a flexible function of X , and λ is the penalty term which
169 constrains the function to a certain degree of smoothness. A TPS depends on the m data points
170 with known coordinates and target values, and it can be described by $2(m + 3)$ parameters.
171 Therefore, the objective is to minimize the Residual Sum of Squares (RSS) where J is the
172 penalty function and λ controls the importance of this, as shown in Equation 2.

$$173 \quad RSS = \min_f \sum_{i=1}^N \{y_i - f(x_i)\}^2 + \lambda J[f] \quad (2)$$

174 In the extreme scenario of $\lambda = 0$ we would fit a spline that perfectly overfits the data points,
175 and with $\lambda \rightarrow \infty$ we would fit the polynomial base model fitted by ordinary least squares.

176 TPSs provide a very flexible framework for model fitting, since it does not require any prior
177 knowledge about the functional form, and there is no need to specify the number of nodes and
178 their location, which allows for optimal control of continuous confounders.

179 We implemented an EM algorithm which allows to estimate the latent structure of the data by
180 assuming that the observed data comes from a finite set of mixtures. We first fitted k different
181 random splines and assigned each patient to the nearest cluster based on the smaller mean
182 squared error (MSE). Then, we iteratively re-computed the k splines based on the clusters
183 formed (M-step) and re-assigned group membership (E-step) until convergence. In case
184 clusters become too small (i.e., with more observations than degrees of freedom) they may
185 merge, resulting in a small number of clusters.

186 In order to select a robust number of clusters, we used three different approaches. The first
187 approach consisted in computing a large model of 40 different clusters, followed by a

188 hierarchical clustering analysis. The latter was performed using complete linkage based on the
189 Euclidean distance between the fitted values of each cluster in order to visualize a potentially
190 optimal number of clusters.[15] The second approach consisted in performing a silhouette
191 analysis, which measured the separation between clusters, and therefore allowed for different
192 number of potential clusters.[16] The elbow method was used in the third approach. This
193 method was performed by computing the deviance of models of different numbers of clusters
194 in order to visually inspect the computed the deviance of models of different sizes. The location
195 of a bend in the plot was considered an indicator of the appropriate number of clusters.

196 For each cluster, we summarized the patient characteristics at index date. The history of
197 comorbidities was identified if a valid read code was present any time prior to index. We
198 compared patients' characteristics at index date using t-test and chi-square tests for continuous
199 and categorical variables, respectively, and computed the standardized mean differences
200 (SMDs) between groups to assess the magnitude of the difference. We defined a SMD >0.2 to
201 indicate significance.

202 Time-to-event analysis was conducted for microvascular conditions (retinopathy, DPN, and
203 ED), and for macrovascular diseases (AMI, CHD, and PAD], as well as the time to first insulin
204 use. Kaplan-Meier curves stratified by cluster were plotted for each outcome. Patients with a
205 diagnosis of outcomes under investigation previous to index date were excluded from the time-
206 to event analysis. Additionally, we fitted a Cox-proportional hazards model for each outcome
207 adjusted for sex, age, and BMI, smoking, and alcohol consumption, at index date. Moreover,
208 we computed the Tukey's range test to determine if there were statistically significant
209 differences between clusters.[17]

210 Finally, we explored annual changes in NIAD utilization within the first 5-years after index
211 date. At each annual time point, we identified the proportion of patients receiving different
212 classes of NIADs within the prior 3-month window to identify if relevant differences in NIAD
213 utilization were present between the clusters.

214 Results

215 *Study Population*

216 We identified 116,251 new users of NIADs who had at least four HbA1c measurements after
217 index date (**Figure 1**). In order to inspect for potential selection bias, we compared patient
218 characteristics of included and excluded patients in this study (**Supplementary Table S1**). No
219 substantial differences between included and excluded patients were observed.

220

221 *Longitudinal HbA1c trajectories*

222 The three different approaches (hierarchical clustering analysis, the silhouette analysis, and the
223 elbow method) used to select a robust number of clusters identified an optimal number of
224 clusters between four and six clusters. Given that four clusters were too general and did not
225 manage to capture and represent all the patterns, and six clusters provided little benefit at the
226 expense of overfitting, we opted for five clusters, **Supplementary Figure S1**.

227 **Figure 2** provides a visual depiction of the cluster trajectories. While all clusters showed an
228 initial drop in HbA1c levels following NIAD start, HbA1c trajectories varied greatly between
229 clusters. When characterizing the five clusters we identified the following patterns:

230

- 231 (1) Optimally controlled HbA1c (**OC**),
- 232 (2) Adequately controlled HbA1c (**AC**),
- 233 (3) Suboptimally controlled HbA1c (**SOC**)
- 234 (4) Poorly controlled HbA1c (**PC**)
- 235 (5) Uncontrolled HbA1c (**UC**)

236

237 The first OC cluster (28.8%; n=33,531) had a modest decrease in HbA1c from baseline and
238 HbA1c levels that remained mostly below 7.0%. The second AC cluster (14.6%; n=16,962)
239 showed a significant initial reduction in HbA1c after index and HbA1c levels remained below
240 the clinical target of 7.5% from year 4 onward. Although patients in the third SOC cluster
241 (32.1%; n=37,325) had an initial drop in HbA1c levels a gradual increased over time was noted,
242 and by year 10 was above 7.5% and on par with the baseline HbA1c. The fourth PC cluster
243 (17.1%; n=19,832) had a moderate initial HbA1c reduction, followed by a rapid increase, with
244 the 10-year HbA1c reaching 9.0% and above the baseline value. And finally, patients in the
245 fifth UC cluster (7.4%; n=8,601) never achieved HbA1c control, with values remaining above
246 9.0% and reaching above 10% by year 10.

247

248 Patient characteristics at index date, stratified by cluster, are shown in **Table 1**. The AC and
249 UC clusters started with the highest HbA1c values (>9%), yet the long-term trajectories were
250 strikingly different, thus, we provide the significance test between the AC and UC clusters in

251 Table 1. Patients in the UC cluster were significantly younger when compared to the AC cluster
252 (53.8 years vs. 60.5 years, SMD 0.53), had a higher mean BMI (34.1 vs. 32.6, SMD 0.21), and
253 were more likely to be current smokers (26.6% vs. 18.2%, SMD 0.22). The UC group had a
254 lower prevalence of high blood pressure (34.3% vs. 42.9%, SMD 0.18) compared to the AC
255 cluster. Conversely, the UC cluster had the highest proportion of patients with chronic liver
256 disease (3.1%) and anxiety and stress (21.3%). Nonetheless, although there were statistically
257 significant differences between both groups, the SMDs were not significant for comorbidities.
258 Among the other clusters, we note that patients in the OC cluster were more frequently older
259 (average age of 63.8 years old) and less likely to be smokers (14.5%) when compared to other
260 clusters, **Table 1**. On the other hand, this group had the highest proportion of patients with a
261 history of high blood pressure (50.2%), hypercholesterolemia (21.6%), and osteoarthritis
262 (25.4%). The SOC and PC clusters both showed increasing HbA1c trajectories over time, and
263 were relatively similar at baseline. The SOC cluster was slightly older (62.7 years vs. 58.2
264 years) than the PC cluster, and slightly less likely to be current smokers (15.0% vs. 19.8%).
265 The SOC cluster had the lowest average BMI (32.2 kg/m²) compared to other groups.

266

267 We observed that patients in the UC and PC clusters had overall good kidney (normal values:
268 eGFR \geq 60 mL/min/1.73 m²) and liver function (normal values: bilirubin \leq 21 μ mol/L), as
269 well as good vitamin B12 levels (normal values: \geq 200 and \leq 950 pg/mL). On the other hand,
270 these patients presented with the highest levels of LDL (normal values: \geq 2.6 and \leq 4.1) and
271 triglycerides (normal values: \geq 1.69 and \leq 2.25 mmol/mol), **Supplementary Figure S2**.

272

273 *Risk for microvascular and macrovascular outcomes*

274 The adjusted Cox-proportional hazards models for the microvascular and macrovascular
275 outcomes of interest are provided in **Table 3**, and the unadjusted survival curves and number
276 of events are provided in **Supplementary Figure S3** and **Supplementary Table S2**,
277 respectively. For retinopathy, only the AC and UC clusters presented statistically significant
278 hazard ratios (HRs) with a 21% (HR 1.21, 95% CI: 1.14 – 1.28) and 30% (HR 1.30, 95% CI:
279 1.21 – 1.40) increased risk with respect to the OC cluster, respectively. When assessing DPN,
280 the PC and UC clusters had a significant 23% (HR 1.23, 95% CI 1.12 – 1.35) and 45% (HR
281 1.45, 95% CI: 1.27 – 1.64) increased risk, respectively, when compared to the OC group. All
282 clusters showed an approximate 20% significant increased risk of ED, when compared to the

283 OC cluster. Between cluster differences are presented in **Supplementary Table S3**. There were
284 no statistically significant differences between the AC and UC clusters for retinopathy risk
285 ($p=0.312$), nor between the PC and UC clusters for DPN risk ($p=0.100$). Similarly, no between
286 group differences were noted for ED.

287

288 For macrovascular events, all clusters showed an increased risk of AMI, compared to the OC
289 group, with the highest risk observed among the UC cluster (HR 2.15, 95% CI: 1.84–2.52),
290 **Table 3**. When comparing the AC vs. the UC cluster with the Tukey’s range test we found
291 statistically significant differences between each other ($p<0.001$), **Supplementary Table S3**.
292 For CHD, all the clusters showed an elevated risk when compared to the OC cluster, with the
293 highest risk again being among the UC cluster (HR 1.64, 95% CI: 1.45–1.86). While for PAD,
294 the AC and UC clusters showed a 27% (HR 1.27, 95% CI: 1.11–1.45) and 62% (HR 1.62, 95%
295 CI: 1.36–1.93) increased risk, respectively, compared to the OC cluster, but no statistically
296 significant differences between each other.

297 *Prescription patterns*

298 Over 90% of the patients started treatment at index date with biguanides e.g., metformin,
299 followed by sulfonylureas, **Supplementary Figure S4**. Nonetheless, patients in the UC cluster
300 were the first ones to stop biguanides treatment the earliest and have them replaced the earliest
301 by other medications. This cluster was the one with the highest proportion of patients with
302 insulins, SGLT2 inhibitors, and GLP-1 at the five-years mark. The AC cluster, who had a
303 similar initial trajectory but a dramatic improvement with respect to the UC cluster, was the
304 group with the greatest number of users of thiazolidinediones (TZDs).

305

306 **Discussion**

307 In this population-based cohort study, we identified five distinct patterns of HbA1c progression
308 with different patient and clinical risk profiles. We found two clusters (UC and AC) with
309 similar high HbA1c values at the start of NIAD therapy, but with very different trajectories.
310 Nevertheless, similar microvascular and macrovascular risks were observed among both
311 groups, when compared to the OC cluster suggesting high baseline HbA1c may be an important
312 risk factor. Additionally, we found that the clusters with the highest risk of AMI and CHD were
313 observed among the clusters with elevated HbA1c during follow-up, the PC and UC clusters.

314 In addition to this, the UC cluster presented the highest use of insulin in the first five years after
315 NIAD start. Further work should aim to assess if patient-level predictors of cluster assignment
316 can be identified to aide in optimizing treatment and management strategies.

317 In the last twenty years there has been an increase in the prevalence of T2DM worldwide,
318 linked to a more sedentary lifestyle, diet, and an increasing aging population.[18] Moreover,
319 T2DM and its complications have contributed substantially to the global burden of mortality
320 and disability, as diabetes is one of the major causes of reduced life expectancy.[19] Thus,
321 improving our understanding the evolution patterns of HbA1c levels at T2DM diagnosis is
322 paramount to advancing tailored therapeutic management in order to achieve glycemic control.

323 While most of the studies that have aimed at modelling HbA1c progression have used latent
324 class growth modelling (LCGM) with smaller cohorts and shorter follow-up periods,[11,12,20]
325 we identified similar trajectory patterns. For example, we found that deteriorating HbA1c (i.e.,
326 PC) and extremely high baseline HbA1c (i.e., AC and UC) were associated with worse clinical
327 outcomes, when compared to the cluster with stable values (i.e., OC).

328 To date, only the work of Laiteerapong et al. used a LCGM in combination of a larger real-
329 world cohort of 28,016 individuals.[21] This study identified five distinct HbA1c trajectories,
330 and found an association between non-stable trajectories and greater risk of microvascular
331 events and mortality. While we looked at individual outcomes, rather than composite
332 endpoints, our results are comparable for microvascular events. The AC cluster was similar to
333 the “high decreasing early” cluster in the Laiteerapong study, which was associated with a 28%
334 increased risk of microvascular events.[12] In our analyses we could see differences between
335 the different microvascular events, which could not be observed in the Laiteerapong study. For
336 example, while the UC cluster had elevated risks for all outcomes, the AC cluster was only
337 significantly associated with diabetic retinopathy and ED, while the PC cluster was associated
338 with a significant increased risk of DPN and ED. As DPN typically develops 10-20 years after
339 the initial diabetes diagnosis, it is plausible that long-term diabetes control is an important
340 predictor. Conversely, retinopathy typical emerges within the first 5-years, and therefore, early
341 glycemic control is likely a key predictor.

342 While we did not cluster based on patient characteristics, we could identify distinct patient
343 profiles across the five clusters. For example, the OC cluster was older and included patients
344 with a history of age-related comorbidities at index, while patients in the worst performing
345 group (UC) were overall younger and more frequently obese. Additionally, some differences

346 between the AC and UC cluster were identified, namely that the AC cluster was older, had a
347 lower mean BMI, and were less frequently smokers. As we identified differences in the
348 likelihood of clinical outcomes, particularly for the AC and UC clusters, and two previous
349 studies found that data-driven clusters based on baseline characteristics (including HbA1c)
350 may be predictive of clinical outcomes.[22,23] Future work should further investigate the
351 patient-level factors that are predictive of the cluster orientation.

352 Finally, many guidelines for diabetes management are moving towards individualized
353 treatment to improve long-term glucose control and clinical outcomes. We could identify some
354 differences between the clusters regarding the treatment course. For example, the UC cluster
355 moved to second-line therapies earliest, especially insulins. Interestingly, the AC cluster
356 appeared to have the highest proportion of users switching to TZDs early in therapy. These
357 results are preliminary, but provide insight into the potential for tailored therapy in T2DM.

358 When interpreting our results, we have to acknowledge several limitations. First of all, we only
359 looked at patients after the start of their first NIAD prescription, thus we might have missed
360 patterns of patients that were first instructed to change their diet and physical activity levels in
361 order to attain glycemic control. Additionally, we were restricted to data from the UK where
362 intrinsic social factors such as lifestyle, diet, or physical exercise might have a relevant impact
363 in glycemic control. In addition, treatment guidelines have changed several times during our
364 follow-up period i.e., 2009 and 2015. Therefore, the relevant thresholds considered to attain a
365 controlled glycemic status have been revised and updated, as well as the medications used.
366 Although metformin still remains the first-line therapy, several types of drugs have been
367 developed, approved, and marketed during our study window. For instance, SGLT2 inhibitors
368 such as dapaglifozin or canaglifozin were approved in 2012 and 2013 respectively. The fact
369 that we excluded patients with less than four HbA1c measurements after index date might could
370 have introduced a selection bias in favor of either healthier or more careless patients.
371 Nonetheless, we did not find relevant differences between included and excluded patients
372 overall, as shown in **Supplementary Table S1**. A potential limitation in the methodological
373 approach we used of splines in combination with *k*-means clustering is the fact that we had to
374 empirically select the number of trajectories. Although we minimized the potential source of
375 bias by performing three different methods i.e., elbow method, hierarchical clustering, and
376 silhouette analysis, we did not find clear cutoff points for the number of clusters. Finally, since
377 we included patients with a first NIAD prescription between January 2003 to December 2019,

378 not many patients have a follow-up period of more than ten years. This led to a higher variance
379 and volatility in the last years of the time-to-event analyses.

380 In conclusion, we found that clusters with worse baseline and long-term glycemic control were
381 associated with higher degree long term microvascular and macrovascular complications.
382 Moreover, as we identified that high baseline HbA1c, as seen in the UC and AC clusters, maybe
383 a strong indicator of retinopathy risk, while long-term HbA1c control is associated with DPN.
384 Further studies should aim at understanding the differences between clusters of similar profile
385 but diverging trajectories, and investigate if more tailored therapy can help improve long-term
386 glycemic trajectories in patients with T2DM.

387

388

389

390 **Author Contributions:** Study Conception: AMB, FPC; data acquisition: AMB; data analysis:
391 AMDIT; data integrity and validity: AMDIT, AMB; data interpretation: AMDIT, FPC, MLF,
392 CM, SW, AMB; manuscript preparation: AMDIT, MLF, AMB; critical revisions: AMDIT,
393 MLF, CM, FPC, SW, AMB.

394

395 **Funding:** This research was funded by a Swiss Data Science Centre Collaboration Grant (C19-
396 09).

397

398 **Conflicts of Interest:** SW is a member of the Human Medicines Expert Committee (HMEC)
399 board of Swissmedic. The views expressed in this article are the personal views of the authors
400 and may not be understood or quoted as being made on behalf of or reflecting the position of
401 Swissmedic or one of its committees or working parties. The professorship of AMB was partly
402 endowed by the National Association of Pharmacists (PharmaSuisse) and the ETH Foundation,
403 but funds are not provided for research and the current project was not funded. AMDIT, MLF,
404 CM, and FPC have no conflicts of interest to declare regarding this research.

405

406 **Ethics statement:** The protocol for this project was approved by the THIN scientific research
407 council (reference number: 20SR062).

408

409 **Consent to participate and for publication:** Consent to participate and for publication was
410 granted by the THIN scientific research council (reference number: 20SR062).

411

412 **Availability of data and material:** The IQVIA Medical Research Data (IMRD) were
413 obtained from IQVIA, a Cegedim database of anonymized electronic health records. For
414 further information on access to the database, please contact IQVIA (contact details can be
415 found at [https://www.iqvia.com/locations/united-kingdom/information-for-members-of-the-](https://www.iqvia.com/locations/united-kingdom/information-for-members-of-the-public/medical-research-data)
416 [public/medical-research-data](https://www.iqvia.com/locations/united-kingdom/information-for-members-of-the-public/medical-research-data)).

417 References

- 418 1. Stumvoll, M.; Goldstein, B.J.; van Haefen, T.W. Type 2 Diabetes: Principles of
419 Pathogenesis and Therapy. *The Lancet* **2005**, *365*, 1333–1346, doi:10.1016/S0140-
420 6736(05)61032-X.
- 421 2. *Global Report on Diabetes*; Roglic, G., World Health Organization, Eds.; World Health
422 Organization: Geneva, Switzerland, 2016; ISBN 978-92-4-156525-7.
- 423 3. Zheng, Y.; Ley, S.H.; Hu, F.B. Global Aetiology and Epidemiology of Type 2 Diabetes
424 Mellitus and Its Complications. *Nat Rev Endocrinol* **2018**, *14*, 88–98,
425 doi:10.1038/nrendo.2017.151.
- 426 4. Reusch, J.E.B.; Manson, J.E. Management of Type 2 Diabetes in 2017: Getting to Goal.
427 *JAMA* **2017**, *317*, 1015, doi:10.1001/jama.2017.0241.
- 428 5. Overview | Type 2 Diabetes in Adults: Management | Guidance | NICE Available online:
429 <https://www.nice.org.uk/guidance/ng28> (accessed on 24 January 2020).
- 430 6. *Recommendations For Managing Type 2 Diabetes In Primary Care*; International
431 Diabetes Federation, 2017;
- 432 7. American Diabetes Association Professional Practice Committee 6. Glycemic Targets:
433 Standards of Medical Care in Diabetes—2022. *Diabetes Care* **2021**, *45*, S83–S96,
434 doi:10.2337/dc22-S006.
- 435 8. de Vegt, F.; Dekker, J.M.; Ruhé, H.G.; Stehouwer, C.D.A.; Nijpels, G.; Bouter, L.M.;
436 Heine, R.J. Hyperglycaemia Is Associated with All-Cause and Cardiovascular Mortality
437 in the Hoorn Population: The Hoorn Study. *Diabetologia* **1999**, *42*, 926–931,
438 doi:10.1007/s001250051249.
- 439 9. Stratton, I.M. Association of Glycaemia with Macrovascular and Microvascular
440 Complications of Type 2 Diabetes (UKPDS 35): Prospective Observational Study. *BMJ*
441 **2000**, *321*, 405–412, doi:10.1136/bmj.321.7258.405.
- 442 10. Holman, R.R.; Paul, S.K.; Bethel, M.A.; Matthews, D.R.; Neil, H.A.W. 10-Year Follow-
443 up of Intensive Glucose Control in Type 2 Diabetes. *N Engl J Med* **2008**, *359*, 1577–
444 1589, doi:10.1056/NEJMoa0806470.
- 445 11. Wang, J.-S.; Liu, W.-J.; Lee, C.-L. HbA1c Trajectory and Cardiovascular Outcomes: An
446 Analysis of Data from the Action to Control Cardiovascular Risk in Diabetes (ACCORD)
447 Study. *Therapeutic Advances in Chronic Disease* **2021**, *12*, 20406223211026390,
448 doi:10.1177/20406223211026391.
- 449 12. Laiteerapong, N.; Ham, S.A.; Gao, Y.; Moffet, H.H.; Liu, J.Y.; Huang, E.S.; Karter, A.J.
450 The Legacy Effect in Type 2 Diabetes: Impact of Early Glycemic Control on Future
451 Complications (The Diabetes & Aging Study). *Diabetes Care* **2018**, *42*, 416–426,
452 doi:10.2337/dc17-1144.
- 453 13. Chisholm, J. The Read Clinical Classification. *BMJ* **1990**, *300*, 1092–1092,
454 doi:10.1136/bmj.300.6732.1092.
- 455 14. Abraham, C.; Cornillon, P.A.; Matzner-Løber, E.; Molinari, N. Unsupervised Curve
456 Clustering Using B-Splines. *Scandinavian Journal of Statistics* **2003**, *30*, 581–595,
457 doi:10.1111/1467-9469.00350.
- 458 15. Salvador, S.; Chan, P. Determining the Number of Clusters/Segments in Hierarchical
459 Clustering/Segmentation Algorithms. In Proceedings of the 16th IEEE International
460 Conference on Tools with Artificial Intelligence; November 2004; pp. 576–584.
- 461 16. Rousseeuw, P.J. Silhouettes: A Graphical Aid to the Interpretation and Validation of
462 Cluster Analysis. *Journal of Computational and Applied Mathematics* **1987**, *20*, 53–65,
463 doi:10.1016/0377-0427(87)90125-7.

- 464 17. Tukey, J.W. Comparing Individual Means in the Analysis of Variance. *Biometrics* **1949**,
465 5, 99–114.
- 466 18. Zghebi, S.S.; Steinke, D.T.; Carr, M.J.; Rutter, M.K.; Emsley, R.A.; Ashcroft, D.M.
467 Examining Trends in Type 2 Diabetes Incidence, Prevalence and Mortality in the UK
468 between 2004 and 2014. *Diabetes, Obesity and Metabolism* **2017**, *19*, 1537–1545,
469 doi:10.1111/dom.12964.
- 470 19. Saeedi, P.; Salpea, P.; Karuranga, S.; Petersohn, I.; Malanda, B.; Gregg, E.W.; Unwin,
471 N.; Wild, S.H.; Williams, R. Mortality Attributable to Diabetes in 20–79 Years Old
472 Adults, 2019 Estimates: Results from the International Diabetes Federation Diabetes
473 Atlas, 9th Edition. *Diabetes Research and Clinical Practice* **2020**, *162*, 108086,
474 doi:10.1016/j.diabres.2020.108086.
- 475 20. Luo, M.; Lim, W.Y.; Tan, C.S.; Ning, Y.; Chia, K.S.; van Dam, R.M.; Tang, W.E.; Tan,
476 N.C.; Chen, R.; Tai, E.S.; et al. Longitudinal Trends in HbA1c and Associations with
477 Comorbidity and All-Cause Mortality in Asian Patients with Type 2 Diabetes: A Cohort
478 Study. *Diabetes Res. Clin. Pract.* **2017**, *133*, 69–77, doi:10.1016/j.diabres.2017.08.013.
- 479 21. Laiteerapong, N.; Karter, A.J.; Moffet, H.H.; Cooper, J.M.; Gibbons, R.D.; Liu, J.Y.;
480 Gao, Y.; Huang, E.S. Ten-Year Hemoglobin A1c Trajectories and Outcomes in Type 2
481 Diabetes Mellitus: The Diabetes & Aging Study. *Journal of Diabetes and its*
482 *Complications* **2017**, *31*, 94–100, doi:10.1016/j.jdiacomp.2016.07.023.
- 483 22. Ahlqvist, E.; Storm, P.; Käräjämäki, A.; Martinell, M.; Dorkhan, M.; Carlsson, A.;
484 Vikman, P.; Prasad, R.B.; Aly, D.M.; Almgren, P.; et al. Novel Subgroups of Adult-
485 Onset Diabetes and Their Association with Outcomes: A Data-Driven Cluster Analysis of
486 Six Variables. *Lancet Diabetes Endocrinol* **2018**, *6*, 361–369, doi:10.1016/S2213-
487 8587(18)30051-2.
- 488 23. Dennis, J.M.; Shields, B.M.; Henley, W.E.; Jones, A.G.; Hattersley, A.T. Disease
489 Progression and Treatment Response in Data-Driven Subgroups of Type 2 Diabetes
490 Compared with Models Based on Simple Clinical Features: An Analysis Using Clinical
491 Trial Data. *Lancet Diabetes Endocrinol* **2019**, *7*, 442–451, doi:10.1016/S2213-
492 8587(19)30087-7.
493

Figure 1. Flowchart of included patients.

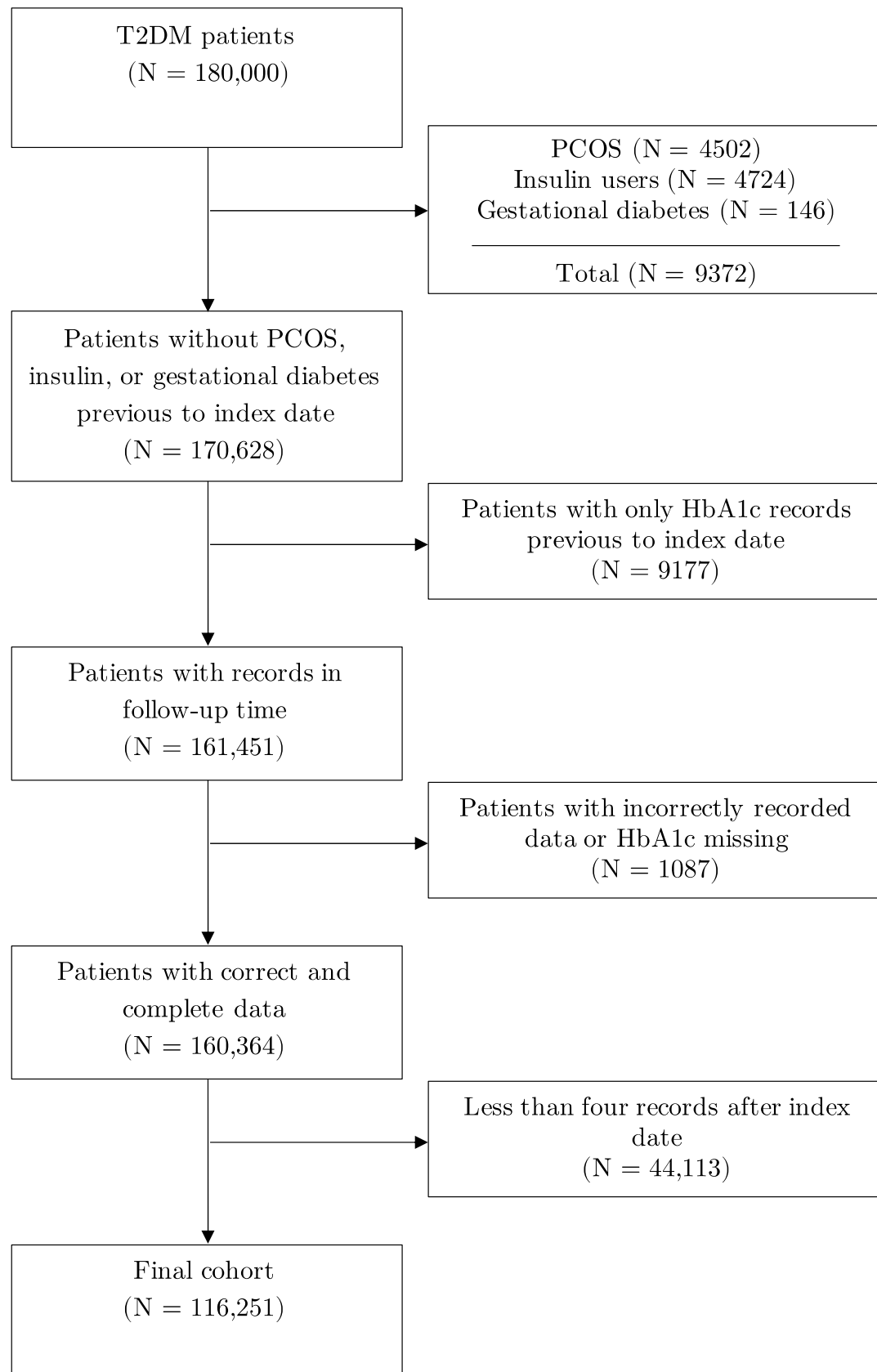
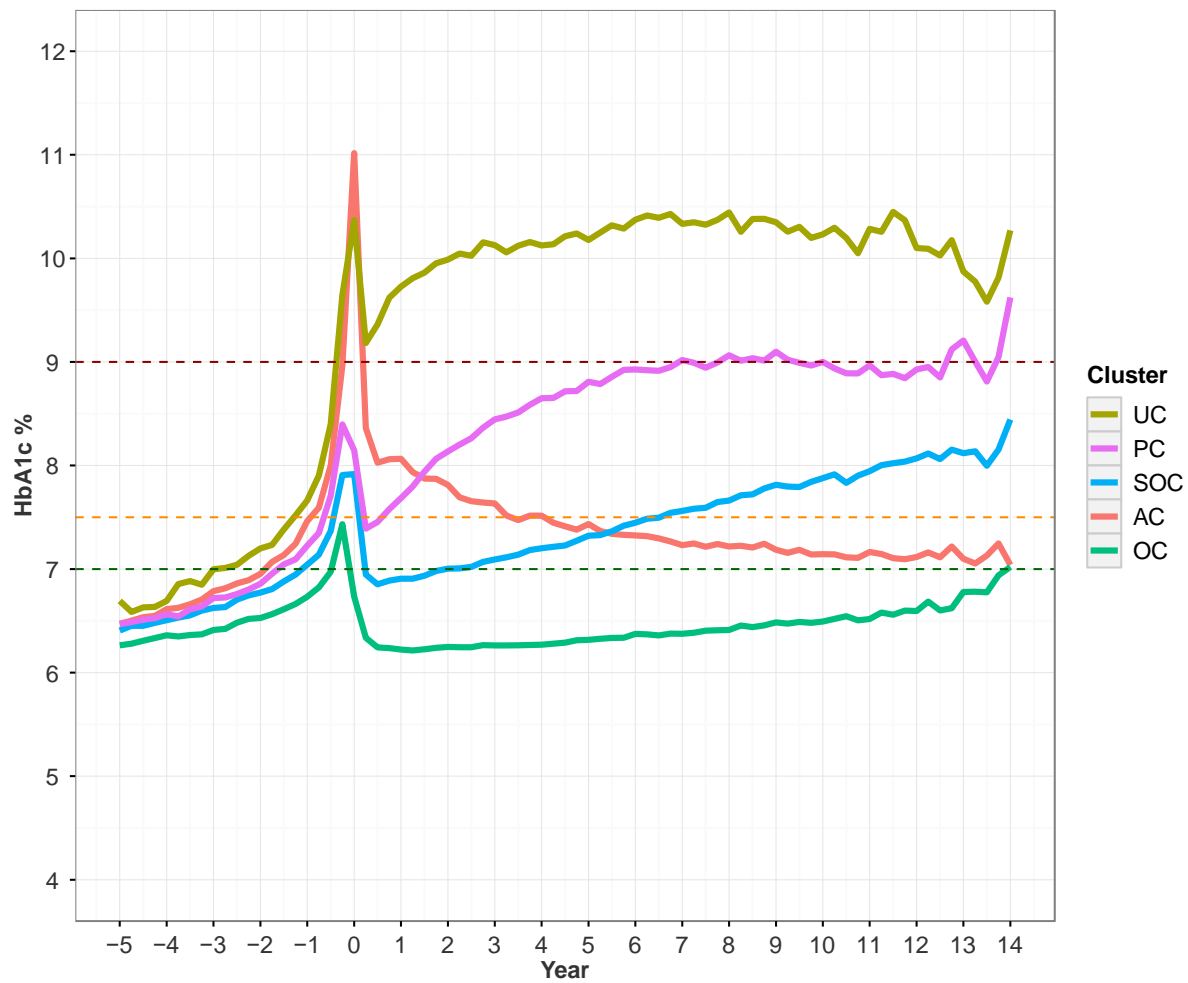


Figure 2. Evolution of cluster centroids of HbA1c level in percentage.



Abbreviations: UC, uncontrolled HbA1c; PC, poor HbA1c control; SOC, suboptimal HbA1c control; AC, adequate HbA1c response; OC, optimal HbA1c control. The gray area indicates the start of NIAD therapy (index date), the green, orange, and red horizontal dashed lines indicate an optimal control ($\leq 7\%$), adequate HbA1c level after treatment intensification ($\leq 7.5\%$), and level where insulin should be considered ($\geq 9\%$), respectively.

Table 1. Patient characteristics at index date, stratified by cluster.

	OC (N=33531)	AC (N=16962)	SOC (N=37325)	PC (N=19832)	UC (N=8601)	AC vs UC P value	AC vs UC SMD
Gender = Male (%)	18091 (54.0)	10376 (61.2)	21463 (57.5)	11756 (59.3)	5091 (59.2)	0.002	0.04
Age (mean (standard deviation))	63.8 (12.4)	60.5 (12.0)	62.7 (12.0)	58.2 (12.8)	53.8 (13.0)	<0.001	0.53
BMI (mean (standard deviation))	32.5 (6.8)	32.6 (6.9)	32.2 (6.4)	33.4 (6.9)	34.1 (7.6)	<0.001	0.21
Alcohol - Current use (%)	23061 (72.9)	11502 (73.7)	25742 (73.4)	13467 (72.9)	5465 (70.0)	<0.001	0.08
Smoking - Current use (%)	4845 (14.5)	3076 (18.2)	5597 (15.0)	3904 (19.8)	2269 (26.6)	<0.001	0.22
Conditions							
Angina pectoris	2802 (8.8)	1048 (6.8)	2920 (8.3)	1188 (6.4)	339 (4.4)	<0.001	0.10
Anxiety & other*	5828 (18.4)	2596 (16.7)	6126 (17.4)	3574 (19.3)	1644 (21.3)	<0.001	0.12
Arthropathy	2385 (7.5)	903 (5.8)	2338 (6.6)	1094 (5.9)	365 (4.7)	<0.001	0.05
Atrial fibrillation	2685 (8.5)	1183 (7.6)	2715 (7.7)	1374 (7.4)	468 (6.1)	<0.001	0.06
Cancer	9330 (29.5)	3977 (25.6)	9968 (28.3)	4731 (25.6)	1730 (22.4)	<0.001	0.02
Chronic depression	199 (0.6)	92 (0.6)	210 (0.6)	132 (0.7)	67 (0.9)	<0.001	0.03
Chronic liver disease	852 (2.7)	407 (2.6)	915 (2.6)	551 (3.0)	243 (3.1)	<0.001	0.03
Congestive heart failure	1016 (3.2)	524 (3.4)	1076 (3.1)	543 (2.9)	228 (2.9)	<0.001	0.02
High blood pressure	15889 (50.2)	6651 (42.9)	16414 (46.6)	7662 (41.4)	2652 (34.3)	<0.001	0.18
Hypercholesterolaemia	6831 (21.6)	2726 (17.6)	7187 (20.4)	3125 (16.9)	1061 (13.7)	<0.001	0.11
Hypothyroidism	3215 (10.2)	1303 (8.4)	3163 (9.0)	1556 (8.4)	670 (8.7)	<0.001	0.01
Intermittent claudication	1262 (4.0)	526 (3.4)	1293 (3.7)	581 (3.1)	184 (2.4)	<0.001	0.06
Irritable bowel syndrome	3541 (11.2)	1403 (9.0)	3735 (10.6)	1951 (10.5)	713 (9.2)	<0.001	0.01
Ischaemic heart disease	560 (1.8)	244 (1.6)	634 (1.8)	310 (1.7)	105 (1.4)	<0.001	0.02
Neuropathy	321 (1.0)	115 (0.7)	303 (0.9)	135 (0.7)	47 (0.6)	<0.001	0.02
Osteoarthritis	8035 (25.4)	3247 (20.9)	8287 (23.5)	3646 (19.7)	1204 (15.6)	<0.001	0.00
Primary open-angle glaucoma	1479 (4.7)	513 (3.3)	1422 (4.0)	568 (3.1)	165 (2.1)	<0.001	0.07
Senile cataract	1998 (6.3)	696 (4.5)	1888 (5.4)	786 (4.2)	217 (2.8)	<0.001	0.03

Abbreviations: UC, uncontrolled HbA1; PC, poor HbA1c control; SOC, suboptimal HbA1c control; AC, adequate HbA1c response; OC, optimal HbA1c control; BMI, body mass index; Anxiety & other*, anxiety and other neurotic, stress related, and somatoform disorders.

Table 2. Adjusted hazard ratio of the different outcomes.

	OC		AC		SOC		PC		UC	
	aHR*	95% CI	aHR*	95% CI	aHR*	95% CI	aHR*	95% CI	aHR*	95% CI
Micro-vascular Outcomes										
Retinopathy	Reference		1.21	1.14, 1.28	1.02	0.97, 1.07	1.05	0.99, 1.11	1.30	1.21, 1.40
Diabetic Polyneuropathy	Reference		1.11	1.00, 1.22	0.99	0.91, 1.08	1.23	1.12, 1.35	1.45	1.27, 1.64
Erectile Dysfunction**	Reference		1.22	1.13, 1.32	1.17	1.09, 1.25	1.20	1.11, 1.30	1.24	1.12, 1.36
Macro-vascular Outcomes										
Acute Myocardial Infarction	Reference		1.30	1.13, 1.48	1.18	1.06, 1.32	1.65	1.46, 1.87	2.15	1.84, 2.52
Coronary Heart Disease	Reference		1.11	1.01, 1.23	1.14	1.05, 1.23	1.38	1.25, 1.51	1.64	1.45, 1.86
Peripheral arterial disease	Reference		1.27	1.11, 1.45	1.11	0.99, 1.23	1.17	1.02, 1.33	1.62	1.36, 1.93

* The models were adjusted for the following variables at baseline: gender, age, body mass index, alcohol consumption, and smoking status.

**Erectile dysfunction outcome assessed along among males

Abbreviations: aHR, adjusted Hazard Ratio; CI, Confidence Interval; OC, optimal HbA1c control; AC, adequate HbA1c control; SOC, suboptimal HbA1c control; PC, poor HbA1c control; UC, uncontrolled HbA1c. aHR in bold are statistically significant.