

1 **Title: Big Data and Machine Learning to tackle Diabetes Management**

2 **Running Title: Machine Learning and Diabetes Management**

3  
4 **Ana F. Pina<sup>1,2</sup>, Maria João Meneses<sup>1,3,4</sup>, Inês Sousa-Lima<sup>1</sup>, Roberto Henriques<sup>5</sup>,**

5 **João F. Raposo<sup>1,3</sup>, M. Paula Macedo<sup>1,3,6</sup>**

6 1 iNOVA4Health, NOVA Medical School|Faculdade de Ciências Médicas, NMS|FCM,  
7 Universidade Nova de Lisboa; Lisboa, Portugal.

8 2 ProRegeM PhD Programme, NOVA Medical School|Faculdade de Ciências Médicas,  
9 NMS|FCM, Universidade Nova de Lisboa; Lisboa, Portugal.

10 3 Portuguese Diabetes Association - Education and Research Center (APDP-ERC), Lisboa,  
11 Portugal.

12 4 DECSIS II Iberia, Évora, Portugal.

13 5 NOVA Information Management School (NOVA IMS), Universidade NOVA de Lisboa,  
14 Lisboa, Portugal.

15 6 Department of Medical Sciences, University of Aveiro, Aveiro, Portugal.

16  
17 **Corresponding Author:**

18 M. Paula Macedo, e-mail: [paula.macedo@nms.unl.pt](mailto:paula.macedo@nms.unl.pt)

19 Faculdade de Ciências Médicas (FCM), Universidade Nova de Lisboa (UNL), 1169-056 Lisboa,  
20 Portugal

21  
22 **Acknowledgements**

23 This work was supported by “Fundação para a Ciência e a Tecnologia” – FCT to AP  
24 (PD/BD/136887/2018); MPM (PTDC/MEC-MET/29314/2017 and PTDC/BIM-  
25 MET/2115/2014); iNOVA4Health (UIDB/Multi/04462/2020), by the European Commission  
26 Marie Skłodowska-Curie Action H2020 (grant agreement n. 734719), and by the Sociedade  
27 Portuguesa de Diabetologia.

28 **Abstract**

29 Type 2 Diabetes (T2D) diagnosis is based solely on glycemia, even though it is an endpoint of  
30 numerous dysmetabolic pathways. T2D complexity is challenging in a real-world scenario, thus  
31 dissecting T2D heterogeneity is a priority. Cluster analysis, which identifies natural clusters  
32 within multidimensional data based on similarity measures, poses as a promising tool to unravel  
33 Diabetes complexity.

34 Herein, we aimed at scrutinizing and integrate the results obtained in most of the works up to date.  
35 We conclude that to correctly stratify subjects and to differentiate and individualize a preventive  
36 or therapeutic approach to Diabetes management, cluster analysis should be informed with more  
37 parameters than the traditional ones, such as etiological factors, pathophysiological mechanisms,  
38 other dysmetabolic co-morbidities, and biochemical factors i.e. the *milieu*. Ultimately the  
39 abovementioned factors may impact on Diabetes and its complications.

40 Lastly, we propose another theoretical model, which we named the Integrative Model. We  
41 differentiate three types of components: etiological factors, mechanisms, and milieu. Each  
42 component encompasses several factors to be projected in separate 2D planes allowing an holistic  
43 interpretation of the individual pathology.

44 Fully profiling the individuals, considering genomic and environmental factors, and exposure  
45 time, will allow the drive to precision medicine and prevention of complications.

46

47 **Keywords:** diabetes; machine learning; cluster analysis; big data

48

49 **Abbreviations**

- 50 BMI – body mass index
- 51 CAD – coronary artery disease
- 52 CKD – chronic kidney disease
- 53 CV – cardiovascular
- 54 DKD – diabetic kidney disease
- 55 eGFR – estimated glomerular filtration rate
- 56 GRS – genetic risk score
- 57 HOMA-B – Homeostatic Model Assessment for beta-cell function
- 58 HOMA-IR – Homeostatic Model Assessment for Insulin Resistance
- 59 MARD – mild-age related Diabetes
- 60 ML – machine learning
- 61 MOD – mild-obesity related Diabetes
- 62 MR – Mendelian Randomisation
- 63 MRI – magnetic resonance imaging
- 64 NAFLD – non-alcoholic fatty liver disease
- 65 OAD – oral antidiabetic drugs
- 66 OGTT – oral glucose tolerance test
- 67 PAM – partition around medoids
- 68 PD – Prediabetes
- 69 SAID – severe autoimmune Diabetes
- 70 SIDD – severe insulin-deficient Diabetes
- 71 SIRD – severe insulin-resistant Diabetes
- 72 SNPs – single nucleotide polymorphisms
- 73 SOM – self organizing maps
- 74 T1D – Type 1 Diabetes *mellitus*
- 75 T2D – Type 2 Diabetes *mellitus*
- 76 UACR – urine albumin creatinine ratio

## 77 1. Introduction

78 In Diabetes glucose metabolism is affected due to individual or simultaneous changes in insulin  
79 secretion, action or metabolism. Diabetes is diagnosed based on glycemia and cut-off values were  
80 defined based on the presence of microvascular complications, namely retinopathy.<sup>1</sup> However,  
81 dysglycemia, or the glucose altered metabolism, is not an all-or-nothing phenomenon on the  
82 contrary, it occurs continuously. Prediabetes (PD) is a less severe hyperglycemic state that depicts  
83 a higher risk of progression to Diabetes. Importantly, individuals with PD can develop Diabetes  
84 complications, whereas others with Diabetes may never develop them, showing the limitations of  
85 the current clinical classification.<sup>2</sup> Therefore, glycemic levels are not sufficient to inform about  
86 the onset and severity of the condition.

87 Notwithstanding all investment in Diabetes, specifically in Type 2 Diabetes *mellitus* (T2D), it is  
88 still one of the main non-communicable diseases, and its mortality increased 70% since 2000.<sup>3</sup>  
89 T2D is extremely heterogenous,<sup>4,5</sup> both in its initial presentation and complications' development,  
90 which is crucial to explain the sustained morbidity and increased mortality attributable to this  
91 condition.<sup>3,6</sup> The empirical individualisation of therapy in Diabetes dates back to 19<sup>th</sup> century,<sup>7,8</sup>  
92 and is still practised. The latest therapeutic guidelines for T2D include several recent drugs that  
93 are giving better results regarding cardiometabolic complications<sup>9</sup> and start to have an increased  
94 focus on the patient's co-morbidities.<sup>10</sup> The concept of *precision medicine* has been proposed,  
95 aiming at defining the most effective approach for a similar group of patients regarding genetic,  
96 environmental, lifestyle, clinical factors, amongst others.<sup>6</sup> However, further advances in the  
97 ability to define *precise* therapies for Diabetes also depend on the acquired knowledge regarding  
98 the heterogeneity of the condition.

99 As early as 1965, two major groups were acknowledged in Diabetes pathophysiology: insulin  
100 resistant and insulin deficient individuals.<sup>11</sup> The two pathophysiological mechanisms associated  
101 with these groups were assumed to be related with two main organs: insulin secretion impairment  
102 in the pancreas; and insulin resistance at the skeletal muscle. Since then, much more complexity  
103 was added to Diabetes pathophysiology, especially to T2D.<sup>12</sup> More recently, it has been shown  
104 that other organs and factors, such as the lung and microbiome, can impact on T2D onset and

105 progression.<sup>13-15</sup> Additionally, it is currently accepted that T2D etiology encompasses thousands  
106 of low impactful genes, as well as environmental and lifestyle factors, that interact with each  
107 other.<sup>16</sup>

108 Glucose metabolism is part of an intricate metabolic network where carbohydrates, lipids and  
109 other metabolic pathways should be considered as a whole and, when affected, result in  
110 dysmetabolism and/or hemodynamic alterations. Thus, depending on the affected mechanisms,  
111 Diabetes can appear in distinct dysmetabolic contexts. Interestingly, there are lipodystrophic  
112 phenotypes in which the inability of white adipose tissue to expand, despite diverse BMI values,  
113 causes ectopic fat deposition.<sup>17</sup> These subjects are exposed to atherogenic dyslipidemia<sup>18</sup> and, in  
114 the liver, development of fatty liver may progress to steatohepatitis<sup>19</sup> that can be further impacted  
115 by different adipose tissue amounts and function. Despite showing similar patterns regarding  
116 hyperglycemia and hyperlipidemia, subjects with lipodystrophy, might require a distinct  
117 treatment.<sup>20</sup> Another example relates to Diabetes and hypertension bidirectional association. Both  
118 conditions have several common pathophysiological mechanisms, namely hyperinsulinemia,  
119 increased sympathetic nervous activity, activation of renin-angiotensin-aldosterone system,  
120 endothelial dysfunction, etc.<sup>21</sup> The onset of hypertension in subjects with Type 1 Diabetes (T1D)  
121 has been related with the onset of kidney dysfunction; however, in subjects with T2D, it can  
122 appear before<sup>22</sup> and they can show a prehypertensive profile some years earlier.<sup>23</sup> The causal  
123 association of T2D in hypertension was depicted in a Mendelian Randomisation (MR) study, but  
124 does not explain the onset of T2D in hypertensive subjects.<sup>24</sup> However, a higher incidence of T2D  
125 in hypertensive subjects as compared with normotensive subjects is evident.<sup>24</sup> The above-  
126 described complexity, although easy to understand in concept, is very hard to demonstrate and  
127 tackle in clinical practice. Dissecting and understanding T2D heterogeneity is a priority to reverse  
128 the current scenario.<sup>25</sup>

129 To tackle the overly complex clinical challenges, involving multiple etiological factors, organs  
130 and mechanisms, classical statistical analyses are frankly insufficient. Recent progress in memory  
131 and computation power allowed for the development and implementation of more complex

132 algorithms, including a collection of tools that can learn from data, named machine learning (ML).  
133 Specifically cluster analysis, using unsupervised learning algorithms (algorithms that deal with  
134 observations that do not have a label to learn from<sup>26</sup>) are promising tools to unravel Diabetes  
135 complexity.

136 We will critically review distinct cluster analysis methodologies currently used to study Diabetes  
137 and integrate results from different studies. Since all analyses aimed at understanding  
138 Diabetes/T2D pathophysiology, we anticipate their conclusions to fit as pieces on a puzzle.  
139 Finally, we suggest a model that can be applied to Diabetes precision medicine and from a wider  
140 perspective to dysmetabolism overall.

141

## 142 **2. Advancement of Diabetes Management – travelling on the road to precision medicine**

143 The word Diabetes ("to go through" or siphon) is attributed to Apollonius of Memphis in Greece  
144 around 250BC. However, its clinical description and some complications date back to 3500 years  
145 ago in Egypt.<sup>27</sup> Interestingly, two types of Diabetes - congenital and late onset - and their  
146 relationship to heredity, obesity, sedentariness and diet, were already recognized in medical  
147 treatments in ancient India.<sup>8,28</sup> At the time Diabetes resulted in death and preventing it was the  
148 main goal. Additionally, complications of Diabetes, as peripheral neuropathy, gangrene and  
149 erectile dysfunction were described by an Arab doctor, Avicenna (AD 960-1037).<sup>27</sup> Centuries  
150 later Matthew Dobson (1732-1784) and Michel Chevreul (1786-1889), through the application of  
151 chemistry to diagnosis, identified glucose as the sugar that was increased both in urine and serum  
152 of these patients.<sup>8</sup> Arguing that glucose appeared in the urine because the body was unable to  
153 assimilate it, Dobson considered Diabetes a systemic disease rather than a kidney disease, as it  
154 was considered until then.<sup>28</sup> These findings led to the research on the metabolism of  
155 carbohydrates. However, insulin was not yet available and treatments were based on  
156 individualisation of diets, rest or other lifestyle changes,<sup>7</sup> unable to prevent death from acute  
157 complications. Neurological complications were also quite frequent, the association of  
158 neuropathy, vascular disease, plantar ulcers and gangrene with Diabetes was also described, rising  
159 the hypothesis that microvascular disease was the cause of some complications.<sup>28</sup>

160 In 1921-22 Banting and Best isolated insulin, one of the great discoveries in medicine, which has  
161 allowed most people with insulin-dependent Diabetes to be treated to this day. On the other hand  
162 it led to the distinction of T1D, in which people needed insulin, from T2D, in which insulin was  
163 present but ineffective.<sup>27</sup> Since the problem in question was hyperglycemia, other therapeutic  
164 strategies would be developed based on glycemic control.<sup>27</sup> In the 1950's the first sulfonylurea  
165 appeared - the first oral antidiabetic drug (OAD) for people with T2D.<sup>29</sup> Metformin, the most used  
166 OAD, appeared a few years later with its mechanism of action only recently fully understood.<sup>30</sup>  
167 Since then, other groups have been made available as the involvement of other organs and  
168 mechanisms is known.<sup>10,12,29</sup> In a paradigm of therapy which in the meantime has become  
169 evidence-based clinical guidelines began to be published, with the main therapeutic focus on  
170 glycemic control.<sup>31</sup> It was also recognized that the reduction of complications implied  
171 simultaneous treatment of other diseases that represent risk factors for the same complications,  
172 such as dyslipidemia and hypertension.<sup>31</sup>  
173 The etiologic classification of Diabetes recognizes several types besides Type 1, Type 2 and  
174 gestational Diabetes.<sup>1</sup> The recognition that there is still a high degree of heterogeneity leads to an  
175 effort to adapt the numerous drugs with distinct mechanisms to the patients who benefit most  
176 from them.<sup>32</sup> Weight control, hypertension and dyslipidemia, among others, have gained  
177 increasing relevance along with glycemic control.<sup>10</sup> Nowadays, these diseases are recognized as  
178 co-morbidities but treated as independent conditions.

179

### 180 3. Cluster Analysis

181 Cluster analysis is a ML methodology that uses a group of algorithms that can deal with non-  
182 labelled data, named unsupervised learning (Figure 1). Cluster analysis aims to stratify population  
183 observations' in natural groups/clusters without needing *a priori* categorization. Within each  
184 cluster observations' similarity are maximized whilst minimized between clusters.<sup>33</sup>  
185 Distinct clustering algorithms have advantages and drawbacks related to computation time, the  
186 need for an *a priori* knowledge regarding the number of groups, and cluster shape in a  
187 multidimensionality space that they can find (Table 1).<sup>26</sup> In (dys)glycemia, specifically in the

188 resolution of T2D heterogeneity, one should consider several parameters with distinct and specific  
189 characteristics (e.g. genes, environmental factors, biochemical analysis, omics, etc.). Therefore,  
190 it is natural that the best result is obtained using an ensemble of algorithms.  
191 Cluster analysis workflow implies taking several decisions (e.g. choosing the algorithm, variables  
192 to inform the cluster, similarity and distance measures, etc.). When algorithms are not able to find  
193 the best number of clusters (Table 2), there is the need to determine *a priori* a number of clusters.<sup>34</sup>  
194 Still, different measures can give a distinct optimal number of clusters and therefore should be  
195 carefully selected and interpreted. Of note, the found groups should be clinically relevant.  
196 Furthermore, aside from finding natural groups in data, cluster analysis is a powerful tool in data  
197 exploration and visualization. In the context of (dys)glycemia heterogeneity, by profiling the  
198 found groups, we can explore what characterizes them, posing a promising tool to explore and  
199 tackle (dys)glycemia complexity.

200

#### 201 **4. Cluster Analysis Algorithm impact on Founded Clusters**

202 To perform a cluster analysis, impactful decisions must be made: inclusion and exclusion criteria,  
203 choice of variables, and the algorithm to perform the analysis, amongs others. Additionally,  
204 indexes that define the best number of clusters and distance metrics have to be selected.<sup>26</sup> Cluster  
205 analysis used to date to tackle T2D and dysmetabolism have a dissimilar methodology that must  
206 be considered when interpreting and integrating the results (Table 2).<sup>35-38</sup>

207 Hierarchical clustering and k-means are two of the most well-known clustering algorithms.  
208 Agglomerative hierarchical clustering<sup>26</sup> is a simple algorithm that hierarchically joins nested  
209 clusters in a bottom-up way, with its agglomerative process visualized in a dendrogram. This  
210 process does not need the pre-specification of the optimal number of clusters, though it requires  
211 an *a posteriori* cut-off to define them. Furthermore, data can be analyzed at different cut-off  
212 values, allowing us to understand how observations aggregate. However it can only find clusters  
213 with specific shapes, it gives distinct solutions depending on the chosen aggregation methodology  
214 to join the observations and has a high computation cost.<sup>26</sup> *K*-means is a simple and efficient  
215 algorithm. Besides not dealing well with categorical variables, the final solution is highly



216 impacted by its random initialization, requires an *a priori* specification of the number of clusters  
217 and, importantly, it is prone to find spherical clusters, even if this is not their natural shape.<sup>26</sup> The  
218 latter can limit its use. Partition around medoids (PAM) is a *k*-medoids algorithm, that is less  
219 sensitive to noise than *K*-means, but with a higher computational cost.<sup>26</sup>  
220 *K*-means, PAM and hierarchical clustering have been used mainly when few parameters are used  
221 to tackle T2D.<sup>39</sup> To perform more complex analyses, self-organising maps (SOMs) and  
222 topological based analysis have proven to be more efficient and able to find clusters that have  
223 non-spherical shapes.<sup>26,40</sup>  
224 Hierarchical SOMs, followed by hierarchical clustering,<sup>41</sup> have been used to solve multiple  
225 intricate problems, including clustering analysis of T1D complications.<sup>40</sup> SOM is a neural  
226 network-based algorithm, which maps observations to neurons in a grid that at the end will  
227 represent the cluster (cluster centroid).<sup>42</sup> In summary, the first algorithm allows data  
228 dimensionality reduction, whereas the second enables the stratification and understanding of how  
229 the units agglomerate together. Aside from dealing with large and complex data, SOMs can find  
230 different cluster formats. Nonetheless, it has drawbacks as requiring too many parameters to be  
231 set and optimised, its computational cost and the number of clusters must be set *a priori*.<sup>42</sup>  
232 Network analysis is a graph-based method that assesses subjects (nodes) in relation to each other  
233 (edges).<sup>36</sup>  
234 The abovementioned algorithms are classified as hard clustering algorithms, i.e. they group the  
235 population to assign one subject only to one cluster. Contrarily, soft clustering uses algorithms  
236 that define the probability of one observation belonging to distinct clusters;<sup>43,44</sup> thus, one subject  
237 can belong to multiple clusters at a given time. Despite computational cost and convergence  
238 drawbacks, soft clustering algorithms are extremely useful when an item can belong to more than  
239 one cluster, as is the case of clustering T2D related genes/SNP's and mechanisms.<sup>38</sup>

240

## 241 **5. Population and Parameter Set to resolve Type 2 Diabetes**

242 Clusters analyses to resolve T2D heterogeneity are also diverse regarding the analyzed  
243 population, set of parameters used to inform the cluster,<sup>40,43-45</sup> thus impacting on the groups

244 found. Methodological heterogeneity reveals the authors' distinct perspectives on Diabetes  
245 definition, where it stands within the wider dysmetabolism concept, and the number and type of  
246 parameters that allows a precision medicine approach to T2D.

247 Although T2D is classically considered an affection of glucose metabolism, glucose metabolism  
248 occurs integrated with other substrates'.<sup>45</sup> Glucose metabolism-related parameters though  
249 informing about groups with different conditions, do not give a broader perspective on  
250 metabolism nor account for the overall metabolic heterogeneity. T2D impact relies mainly on its  
251 complications' development that, in turn, relate to other factors.<sup>46</sup> Herein, we distinguish  
252 etiological factors (e.g. time, genes, environmental factors, lifestyle factors), pathophysiological  
253 mechanisms (e.g. overall or organ-specific insulin resistance, insulin secretion, overall or organ-  
254 specific insulin clearance), other dysmetabolic co-morbidities (e.g. hypertension, dyslipidemia),  
255 biochemical and other internal environment factors present in the organism or that the organism  
256 is exposed to, that is its *milieu* (e.g. glycemia, insulinemia, free fatty acids, blood pressure, body  
257 weight).

258 Ahlqvist *et al.* performed a cluster analysis on a population of individuals recently diagnosed with  
259 Diabetes.<sup>35</sup> The analysis considered six clinical parameters: the presence of GAD antibodies  
260 (GADA), age at diagnosis, HbA1c, BMI, HOMA-IR and HOMA-B. The analysis does not rely  
261 only on glycemia nor on insulin levels. Nevertheless, the population solely includes individuals  
262 that were diagnosed based on current criteria. The authors found five optimal clusters using the  
263 silhouette index and hierarchical clustering. One of these clusters, named severe autoimmune  
264 Diabetes (SAID), included GADA+ subjects. Afterwards, GADA+ subjects were excluded and  
265 the *k*-means algorithm was used to define the other 4 clusters: severe insulin-deficient Diabetes  
266 (SIDD); severe insulin-resistant Diabetes (SIRD); mild-obesity related Diabetes (MOD); and  
267 mild-age related Diabetes (MARD). These clusters were replicated in other northern European  
268 cohorts.<sup>35</sup> In brief, SAID subjects showed an early-onset condition, low BMI and poor metabolic  
269 control. Subjects in SIDD cluster were similar to SAID but GADA-; these subjects showed a  
270 higher risk of having diabetic retinopathy. A variant in human leukocyte antigen (HLA) locus  
271 (rs2854275) was found to be associated with SAID but not with SIDD.

272 Interestingly, Zaharia *et al.* showed that, in a German population, individuals that were GADA-  
273 at baseline could be GADA+ after five years, determining that for better classification of  
274 autoimmune Diabetes other antibodies should be used.<sup>47</sup> SIRD cluster included individuals with  
275 marked insulin resistance, high BMI and a high prevalence of non-alcoholic fatty liver disease  
276 (NAFLD). Of note, this cluster also revealed to have the highest  $\beta$ -cell function. Additionally,  
277 individuals in the SIRD cluster were at the highest risk of developing chronic kidney disease  
278 (CKD) and diabetic kidney disease (DKD, defined by persistent macroalbuminuria), despite  
279 proper glycemic control. Finally, subjects in MOD showed higher values of BMI but not insulin  
280 resistance, whereas MARD subjects were older, with only modest metabolic affection and were  
281 not associated with the evaluated Diabetes complications. These last two clusters included most  
282 of the population and still have a considerable proportion of subjects with Diabetes complications.  
283 Furthermore, not all Diabetes' complications were evaluated. In fact, it has been suggested that  
284 borderline Diabetes is associated with an increased risk of dementia and Alzheimer disease, which  
285 is potentiated when hypertension is present. Regarding gene loci, rs7903146 (a TCF7L2 SNP  
286 associated with T2D) was associated with SIDD, MOD and MARD; whereas rs10401969 (a  
287 TM6sf2 gene variant associated with NAFLD) was associated with SIRD but not with MOD.<sup>35</sup>  
288 The above-mentioned four subgroups of T2D have been overall replicated, using the same  
289 methodology as Ahlqvist *et al.*, in distinct geographical locations and ethnicities. This further  
290 confirms the already known association of Diabetes with younger subjects, with lower BMI and  
291 more insulin deficiency in Asian and Indian populations.<sup>48,49</sup> Moreover, 23% of subjects changed  
292 cluster in the five year follow-up.<sup>47</sup> Particularly, people in the insulin-deficient cluster (SIDD)  
293 were changed to clusters with better insulin secretion (MOD and MARD).

294 Li *et al.* performed topology-based cluster analysis of 2552 T2D subjects from several ethnicities,  
295 informed by 73 mixed features from electronic medical records derived clinical data.<sup>36</sup> These  
296 features included biochemical and clinical parameters besides glycemia, thus approaching T2D  
297 in a wider (dys)metabolic context. This was a landmark study and one of the first studies to show  
298 the ability to deal with a high number of variables to stratify subjects with T2D. However, the  
299 stratification results depend on the parameters selected to inform the cluster rather than the chosen

300 population. It is not clear if the authors have found three Diabetes subtypes or three subtypes of  
301 patients that have Diabetes, considering the highly mixed chosen parameters to inform the  
302 analysis that also included several diseases codes and medications. The chosen methodology  
303 renders it difficult to validate it in different populations.

304 To extend clusters' evaluation to subjects with normoglycemia and PD, we accounted for age as  
305 a surrogate of time exposition, anthropometry, and biochemical parameters (glycemia, insulin, c-  
306 peptide and free fatty acids) in three-time points of an OGTT.<sup>37</sup> In this study, we used a  
307 hierarchical SOM, followed by a hierarchical clustering algorithm. Subjects were then profiled  
308 concerning the abovementioned parameters and several mechanism's surrogate indexes, including  
309 overall and tissue-specific insulin resistance, insulin secretion, insulin clearance, NAFLD, and  
310 glomerular filtration rate (GFR). The sample had a limited number of subjects with non-treated  
311 T2D. Nonetheless, none of the subjects had Diabetes five years earlier. In this work, we found  
312 two main clusters: one that includes subjects with a median metabolic phenotype compared to the  
313 overall population; and the other with elevated insulin resistance and insulin secretion. However,  
314 these 2 clusters were highly heterogeneous when they were evaluated for a higher number of  
315 clusters. For example, despite the presence of a main insulin-resistant group, that comprised  
316 subjects with normoglycemia and dysglycemia it included subgroups that could be differentiated  
317 by their adipose tissue insulin resistance. Moreover, even though groups with lower estimated  
318 GFR (eGFR) were insulin resistant, not all insulin resistant groups showed this association.  
319 Additionally, we found that clusters including individuals with normo/dysglycemia and low  
320 eGFR could be further profiled and showed insulin resistance and NAFLD. Consistently, other  
321 groups have also shown that both high insulin resistance and NAFLD are related to kidney  
322 dysfunction in subjects with or without T2D.<sup>50</sup> In Ahlqvist *et al.* the group of individuals that had  
323 the highest risk of developing CKD/DKD, even considering proper glycemic control, was the  
324 most insulin resistant one.<sup>35</sup>

325 Furthermore, these subjects had the lowest GFR at baseline (when they had less than 12 months  
326 from diagnosis) on the German Diabetes Study cohort.<sup>47</sup> The impact of insulin resistance and  
327 NAFLD on GFR seems to be, at least partially, independent from glycemia. Importantly, both

328 conditions can be associated with hyperinsulinemia and insulin is a known trigger and a target of  
329 kidney (dys)function that may have a role in the pathophysiology of T2D.<sup>51</sup> Interestingly, the  
330 heterogeneity of affected mechanisms was not exclusive of people with T2D, including also  
331 subjects with PD and normoglycemia. Our work would benefit from being validated in other  
332 cohorts. Nevertheless, we highlight that T2D diagnosis should consider other parameters besides  
333 glycemia. In fact glucose levels impact is differently perceived by each individual. Therefore, it  
334 should include subjects with different ranges of glyceemic values together with other parameters.  
335 An interesting complementary approach to dissect T2D heterogeneity is the use of genetic  
336 markers. Reasoning that genetic variants remain constant despite disease progression and  
337 treatment, unlike clinical variables, thus being more likely to reveal T2D causal mechanisms, a  
338 cluster analysis including T2D gene-traits associations, including 94 genetic variants and 47 traits  
339 was performed.<sup>38</sup> Aside from genetic data the analysis was informed with clinical parameters,  
340 including surrogate indexes of insulin secretion and insulin resistance, as well as lipid parameters,  
341 that allowed for the identification of other insulin resistance-related groups. Importantly, in this  
342 work b-NMF, a soft clustering algorithm was used, allowing a SNP to be associated with more  
343 than one mechanism and one cluster. The authors identified five clusters of genetic loci – traits  
344 associations: two with variant-trait associations related to reduced  $\beta$ -cell function, distinct in pro-  
345 insulin levels; and three insulin resistance-related, namely obesity mediated, lipodystrophy-like  
346 fat distribution and disrupted liver lipid metabolism. Of note, this is also a potentially complex  
347 approach. As more than 100 loci were already found to be associated with T2D, each one with a  
348 very slight impact on the increased risk of the disease and in dysmetabolism etiology, we should  
349 consider, along with genetic factors, their interactions with environmental and lifestyle factors.  
350 Interestingly, Udler *et al.* evaluated the Genetic Risk Score (GRS) association with relevant  
351 outcomes in each cluster. Coronary artery disease (CAD) was mostly associated with the  
352 lipodystrophy and Beta-cell clusters. Beta-cell cluster was also associated with ischemic stroke.  
353 Increased blood pressure was only associated with lipodystrophy cluster, which also showed an  
354 association with higher urine albumin-creatinine ratio (UACR). Liver/Lipid cluster was

355 associated with decreased renal function and diminished UACR. GRS outcomes were validated  
356 in T2D cohorts by profiling subjects' characteristics in top quantile GRS's subjects.<sup>38</sup>  
357 More recently, Wagner *et al.* focused on a German population considered at risk of developing  
358 Diabetes based on BMI, previous history and family history (TUEF/TULIP cohort).<sup>52</sup> Besides  
359 OGTT-based measures reflecting blood glucose, insulin resistance and insulin secretion, liver,  
360 subcutaneous and visceral fat values measured by MRI, and HDL levels, polygenic risk score for  
361 Diabetes were also included. The defined six clusters were then evaluated in a larger cohort  
362 (Whitehall II). However to assign the latter individuals to the clusters, the authors used less  
363 profiling variables, still based on OGTT measurements. The authors reported a relocation rate of  
364 only 60% in the original cohort, which suggests that MRI fat measurements do not appear to be  
365 superior to measurements such as BMI and waist circumference.<sup>37</sup> Importantly, progression to  
366 Diabetes, CKD, CV events and all cause mortality were assessed.<sup>52</sup> Consistent with our findings,<sup>37</sup>  
367 Wagner *et al.* demonstrated that pathophysiological affection is already present before Diabetes  
368 diagnosis.<sup>52</sup> Within the six defined clusters, three that were older and/or more obese showed  
369 higher glycemia (clusters 3, 5 and 6); one related with insulin deficiency and raised genetic risk  
370 (cluster 3); and two with insulin resistance (clusters 5 and 6). Cluster 6 showed a dissociation of  
371 both risks of progression to Diabetes and CKD in Whitehall II cohort. However, considering that  
372 GFR is not depicted in TULIP/TULIF and CKD progression models in Whitehall II were not  
373 adjusted to GFR at the baseline these results should be carefully interpreted. Cluster 4 is consistent  
374 with a metabolically health obese profile that includes younger subjects than the most  
375 dysmetabolic groups and did not show a protected profile overtime, namely regarding CV events.  
376 In fact, although clusters in TULIP/TULIF cohort differ in intima-media thickness, in the  
377 Whitehall II cohort, the clusters did not differ in CV outcomes risk, after adjustment for BMI and  
378 age, except for Cluster 2 that had a protected profile. Considering the relevance of CV events in  
379 Diabetes this highlights the importance of an enriched milieu to a better stratification.<sup>37</sup>

380

381 **6. What can we learn from cluster analysis?**

382 Insulin secretion and resistance have been included in parameters informing cluster analyses.  
383 However there can be different mechanisms that lead to insulin deficiency and resistance.<sup>37</sup> It has  
384 been suggested that insulin resistance can be considered a defensive mechanism against elevated  
385 insulin secretion due to a highly nutritional load in a sensitive  $\beta$ -cell.<sup>53</sup> In distinct cluster analysis,  
386 most of the groups found to be insulin resistant were the ones with the highest insulin  
387 secretion.<sup>35,37</sup> Nevertheless, the amount of circulating insulin depends not only on the cells'  
388 secretion capacity but on overall insulin metabolism and on insulin clearance.<sup>54</sup> Changes in insulin  
389 clearance have also been linked to hyperinsulinemia.<sup>37,54</sup> Insulin resistance has been associated  
390 with age and BMI. Interestingly, in work by Alqvist *et al.*, MARD and MOD groups differ from  
391 the SIRD in that they are less fat or younger, respectively, showing better metabolic control.<sup>35</sup>  
392 Several questions remain to be clarified concerning the mechanisms leading to insulin resistance.  
393 One concerns the mechanisms through which age and BMI impact on insulin resistance and  
394 whether this implies a different therapeutic approach. Secondly, in the setting of insulin  
395 resistance, it is known the association between liver and adipose tissue but whether insulin  
396 resistance develops through distinct pathways, implying distinct therapeutic approaches, remains  
397 elusive. Thirdly, when it comes to Diabetes complications, the majority of the results were  
398 obtained using patients undergoing treatments, which may, in its turn, promote complication's  
399 onset.<sup>55</sup> Finally, cluster analysis showed the association between GFR and albuminuria with  
400 insulin-resistant states;<sup>40,45,52,56</sup> however, the presence of an association does not necessarily imply  
401 homogeneity between clusters, when it comes to kidney function, making this an etiological factor  
402 of the uttermost importance in Diabetes stratification.

403 Udler using SNP and traits, in addition to HOMA-IR and HOMA-B, namely lipid profile, found  
404 three groups of insulin-resistant subjects that showed involvement of different mechanisms and  
405 organs.<sup>38</sup> We and others have shown that distinct insulin resistance patterns can be present in  
406 subjects with normoglycemia and PD.<sup>37,52</sup>

407 Altogether these support the view that, in order to stratify subjects to differentiate a preventive or  
408 therapeutic approach to Diabetes, one should inform the cluster analysis with more parameters  
409 reflecting other mechanisms metabolites and factors (e.g., lipids, blood pressure, insulin).

410 Additionally, Diabetes pathophysiology occurs continuously and people without Diabetes can  
411 already have Diabetes's complications, hinting to different susceptibilities to glycemic levels. This  
412 may be due to concomitant exposure to other factors such as hypertension or dyslipidemia, or due  
413 to the common underlying pathophysiologic mechanisms.

414

## 415 **7. New models for an approach to Diabetes in precision medicine**

416 Cluster analysis is contributing to uncover the heterogeneity of Diabetes.<sup>35,37,38,47</sup> However, its  
417 superiority over simple predictive models (e.g., predicting complications such as renal  
418 dysfunction) is being questioned.<sup>56</sup>

419 McCarthy proposed the palette model to resolve T2D heterogeneity.<sup>57</sup> The model defined  
420 component planes, such as mechanisms, etiological factors and others, that can be affected,  
421 comparing them to a palette hue. The characterization of subjects by their component planes  
422 places them in a bidimensional plane where the path from normoglycemia to Diabetes can be  
423 assessed for each individual. Importantly this model includes subjects with normoglycemia and  
424 dysglycemia, which have different affected mechanisms. Ahlqvist *et al.* suggested a model based  
425 on the assumption that there is a dominant pathway that gives at least to the majority of patients  
426 with Diabetes a well-defined "palette colour".<sup>58</sup> Additionally, few clinical parameters render  
427 larger groups.

428 In our view, a precision medicine model to approach Diabetes must consider glycemia and  
429 glucose metabolism, as well as other substrates and factors, that impact on dysglycemia and/or  
430 Diabetes complications onset and progression. Diabetes complications occur for different values  
431 of glycemia, impacted by the metabolic context of the individual. In fact, dysmetabolic factors  
432 interaction might potentiate the risk for specific conditions, as is the case of glycemia and blood  
433 pressure interaction in the development of Alzheimer's disease.<sup>55</sup> Finally, the model must be  
434 holistic and applicable to different ethnicities. There are ethnicities that show a higher risk for the  
435 onset of T2D at younger ages and for lower BMI.<sup>48</sup> Interestingly, subjects with an Asian genetic  
436 background seem to have diminished insulin secretory capacity, but one cannot exclude the  
437 environmental and culture-related factors.



438 We propose to paint another picture, the Integrative model (Figure 2). We consider that the  
439 approach can only be attained by being detailed in the metabolic characterization of the  
440 individuals, and by placing it in a wider context of dysmetabolism. Thus, we consider the path  
441 from normometabolism to dysmetabolism, in which dysglycemia is one axis among other factors  
442 that can impact on complications onset/progression and organ dysfunction. Therefore, the  
443 metabolic condition of each subject is approached in an integrated way. Also, we differentiate  
444 three types of components: etiological factors, mechanisms and milieu. Each encompasses  
445 several factors or axis that are projected in separate 2D planes. We postulate that, by deeply  
446 profiling a subject for one type of component, we can place him in the corresponding plan.  
447 Furthermore, we postulate that we can predict where the individual is in one plan by knowing the  
448 others. Ultimately it will allow placing each individual in a last plan where his metabolic state is  
449 known. It is natural that there are groups in the data. However, given the possible combinations  
450 of affected mechanisms and organs, it is clear that their number is too high for human  
451 understanding.

452 This model differs from McCarthy's palette model in two main points: 1) it considers the path to  
453 dysmetabolism and not to hyperglycemia; 2) it separates the different etiological factors from the  
454 affected mechanisms and from the internal environment to which the person is exposed on  
455 different levels (Figure 3). The different planes are thus projected among themselves, giving us  
456 the possibility to know one when we fully evaluate the others. This differentiation can be relevant  
457 to prioritize the clinical approach to the individual and to delineate distinct integrated therapeutic  
458 and preventive strategies to be adopted in the different planes that nonetheless should be validated  
459 in clinical studies.

460 Currently, in therapeutic individualization, therapy is first prescribed to hyperglycemia and then  
461 adapted according to the individual characteristics of each patient. In contrast, in precision  
462 medicine, the therapeutic approach is chosen after assigning the patient to a group that already  
463 considers the individual specificities. For example, in the individualized treatment of T2D, a  
464 subject without atherosclerotic disease or CKD but with hypertension and poorly controlled  
465 glycemia, when on metformin, can be medicated with one of five drugs (DPP4, GLP-1, SGLT-2,

466 thiazolidinediones, sulfonylureas). This will be chosen by each doctor considering some  
467 characteristics of the patient, such as weight. In addition, an antihypertensive is associated. In real  
468 life, situations are not so clear as in guidelines. For instance, what to do with a patient with T2D  
469 on metformin, with good glyceic control (average HbA1c 6.8%) but with evidence of early  
470 DKD and without other metabolic risk factors? How intensive and with which agents should he  
471 be treated to have the best health outcome? Is it better to use a SGLT-2 inhibitor or/and start an  
472 ACE2 inhibitor? Is this the best treatment for all the patients in this condition? Or what to do with  
473 another patient with 15 years of T2D, mostly with poor control (HbA1c >8.5%) under different  
474 antihyperglycemic medication, without other risk factors or evidence of Diabetes complications?  
475 Should we keep trying to put him in a good track of glyceic control? For what purpose? In a  
476 precision medicine approach, he would first be assigned into a group of people sharing common  
477 features of the overall metabolic condition, already accounting with all his specificities (including  
478 *milieu*, mechanisms and etiological factors) for which the optimal treatment of that group would  
479 be already tested, defined and can then be prescribed for that individual.  
480 In order to train and validate this theoretical model, datasets that consider the overall metabolism  
481 and deep phenotyping subjects in the distinct proposed planes are needed. Ultimately this model  
482 may be implemented in a decision support system that predicts where people are in their overall  
483 metabolism. This would assign the individual to a homogeneous group, eventually unravelling  
484 his metabolic footprint.

485

## 486 **8. Conclusion**

487 Precision medicine allows tailoring an approach or treatment to different individuals. In other  
488 words, a population is stratified into similar groups, considering relevant characteristics to the  
489 condition (e.g. T2D). Doing so for each group an appropriate therapeutic approach is defined.  
490 Although precision medicine approaches can make use of genetic data, they can also be based on  
491 many other types of clinical data. Observed complexity is solved with the help of mathematical  
492 algorithms that stratifies individuals into groups by similarity.

493 In the era of omics and digital health, in which we can extract and deal with thousands of features  
494 and use them to tailor care to Diabetes, it is not prudent to limit cluster analysis to a few already  
495 preestablished common mechanisms. Furthermore, these new strategies allow us to deal with  
496 blood glucose levels as a continuum, together with the overall milieu, surpassing the artificial  
497 glycemia-based cut-off approach. By fully profiling subjects regarding genomics, environmental  
498 factors and time exposition, we will be able to know which mechanism(s) is(are) affected and  
499 is(are) responsible for a dysmetabolic condition. This enables the use of drugs in a precise manner  
500 and the discovery of new ones. Additionally, prevention of complications, such as cardiovascular  
501 events, may be earlier and more effective. The great big challenge will be identifying which  
502 features are relevant to consider precise care and gather the data to perform these analyses. In a  
503 global village such as our world, we should gather robust clinical data working in a worldwide  
504 consortium.

505

#### 506 **Conflict of Interest**

507 The authors declare that the research was conducted in the absence of any commercial or financial  
508 relationships that could be construed as a potential conflict of interest.

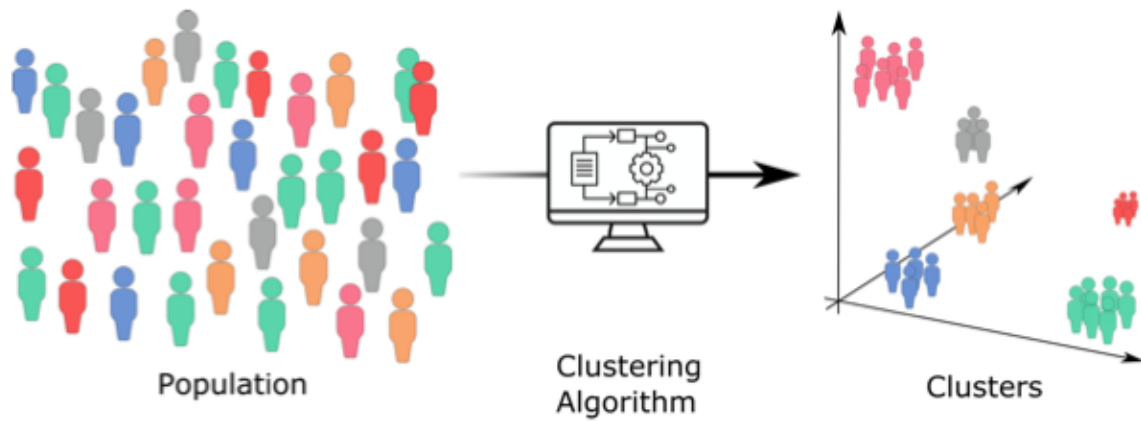
509 **References**

- 510 1. Soh SB, Topliss D. Classification and laboratory diagnosis of diabetes mellitus. *Diabetes Care*  
511 2014;27:S5–10.
- 512 2. Kumar R, Nandhini LP, Kamalanathan S, et al. Evidence for current diagnostic criteria of  
513 diabetes mellitus. *World J Diabetes* 2016;7:396.
- 514 3. Hunter WB. Diabetes as a public health problem. *N C Med J* 1950;11:289–92.
- 515 4. Fajans SS, Cloutier MC, Crowther RL. Clinical and Etiologic Heterogeneity of Idiopathic  
516 Diabetes Mellitus. *Diabetes* 1978;27:1112 LP – 1125.
- 517 5. Tuomi T, Santoro N, Caprio S, et al. The many faces of diabetes: A disease with increasing  
518 heterogeneity. *Lancet* 2014;383:1084–94.
- 519 6. Fradkin JE, Hanlon MC, Rodgers GP. NIH precision medicine initiative: Implications for  
520 diabetes research. *Diabetes Care* 2016;39:1080–4.
- 521 7. Joslin EP. Apollinaire Bouchardat 1806-1886. *Diabetes* 1952;1:490–1.
- 522 8. Schneider T. Diabetes through the ages: a salute to insulin. *South African Med J* 1972;46:1394–  
523 400.
- 524 9. Sattar N. Advances in the clinical management of type 2 diabetes: A brief history of the past 15  
525 years and challenges for the future. *BMC Med* 2019;17:2–5.
- 526 10. Care D, Suppl SS. Pharmacologic approaches to glycemic treatment: Standards of medical care in  
527 diabetesd2021. *Diabetes Care* 2021;44:S111–24.
- 528 11. World Health Organization. *Diabetes mellitus. Report of a WHO expert committee.*; 1965.
- 529 12. DeFronzo RA. From the triumvirate to the ominous octet: A new paradigm for the treatment of  
530 type 2 diabetes mellitus. *Diabetes* 2009;58:773–95.
- 531 13. Yang BY, Qian Z (Min), Li S, et al. Ambient air pollution in relation to diabetes and glucose-  
532 homoeostasis markers in China: a cross-sectional study with findings from the 33 Communities  
533 Chinese Health Study. *Lancet Planet Heal* 2018;2:e64–73.
- 534 14. Teeter JG, Riese RJ. Cross-sectional and prospective study of lung function in adults with type 2  
535 diabetes: The atherosclerosis risk in communities (ARIC) study. *Diabetes Care* 2008;31.
- 536 15. Gurung M, Li Z, You H, et al. Role of gut microbiota in type 2 diabetes pathophysiology.  
537 *EBioMedicine* 2020;51:102590.
- 538 16. Cuschieri S. The genetic side of type 2 diabetes – A review. *Diabetes Metab Syndr Clin Res Rev*  
539 2019;13:2503–6.
- 540 17. Rother KI, Brown RJ. Novel Forms of Lipodystrophy. *Diabetes Care* 2013;36:2142 LP – 2145.
- 541 18. Kinzer AB, Shamburek RD, Lightbourne M, et al. Advanced Lipoprotein Analysis Shows  
542 Atherogenic Lipid Profile That Improves after Metreleptin in Patients with Lipodystrophy. *J*  
543 *Endocr Soc* 2019;3:1503–17.
- 544 19. Polyzos SA, Perakakis N, Mantzoros CS. Fatty liver in lipodystrophy: A review with a focus on  
545 therapeutic perspectives of adiponectin and/or leptin replacement. *Metabolism* 2019;96:66–82.
- 546 20. Stefan N. Causes, consequences, and treatment of metabolically unhealthy fat distribution. *Lancet*  
547 *Diabetes Endocrinol* 2020;8:616–27.
- 548 21. Colussi GL, Da Porto A, Cavarape A. Hypertension and type 2 diabetes: lights and shadows  
549 about causality. *J Hum Hypertens* 2020;34:91–3.
- 550 22. Van Buren PN, Toto R. Hypertension in Diabetic Nephropathy: Epidemiology, Mechanisms, and  
551 Management. *Adv Chronic Kidney Dis* 2011;18:28–41.
- 552 23. Tsimihodimos V, Gonzalez-villalpando C, Meigs JB, et al. Coprediction and Time Trajectories.  
553 <https://doi.org/10.1161/HYPERTENSIONAHA.117.10546>.
- 554 24. Sun D, Zhou T, Heianza Y, et al. Type 2 Diabetes and Hypertension: A Study on Bidirectional  
555 Causality. *Circ Res* 2019;124:930–7.

- 556 25. Chung WK, Erion K, Florez JC, et al. Precision Medicine in Diabetes: A Consensus Report from  
557 the American Diabetes Association (ADA) and the European Association for the Study of  
558 Diabetes (EASD). *Diabetes Care* 2020;43:1617–35.
- 559 26. Pina A, Macedo MP, Henriques R. Clustering Clinical Data in R. In: Matthiesen R, ed. *Mass*  
560 *Spectrometry Data Analysis in Proteomics*. New York, NY: Springer New York; 2020:309–43.
- 561 27. Lasker SP, Mclachlan CS, Wang L, et al. Discovery , treatment and management of diabetes. *J*  
562 *Diabetol* 2010;1:1–8.
- 563 28. Ahmed AM. History of diabetes mellitus. *Saudi Med J* 2002;23:373–8.
- 564 29. White JR. A brief history of the development of diabetes medications. *Diabetes Spectr*  
565 2014;27:82–6.
- 566 30. Rena G, Pearson ER, Sakamoto K. Molecular mechanism of action of metformin: Old or new  
567 insights? *Diabetologia* 2013;56:1898–906.
- 568 31. Of S, Care diabetes M. Updates to the Standards of Medical Care in Diabetes-2018. *Diabetes Care*  
569 2018;41:2045–7.
- 570 32. Meneses MJ, Silva BM, Sousa M, et al. Antidiabetic drugs: Mechanisms of action and potential  
571 outcomes on cellular metabolism. *Curr Pharm Des* 2015;21.
- 572 33. Jain AK, Murty MN, Flynn PJ. Data clustering: a review. *ACM Comput Surv* 1999;31:264–323.
- 573 34. Xu S, Qiao X, Zhu L, et al. Reviews on determining the number of clusters. *Appl Math Inf Sci*  
574 2016;10:1493–512.
- 575 35. Ahlqvist E, Storm P, Käräjämäki A, et al. Novel subgroups of adult-onset diabetes and their  
576 association with outcomes: a data-driven cluster analysis of six variables. *Lancet Diabetes*  
577 *Endocrinol* 2018;6:361–9.
- 578 36. Li L, Cheng WY, Glicksberg BS, et al. Identification of type 2 diabetes subgroups through  
579 topological analysis of patient similarity. *Sci Transl Med* 2015;7.
- 580 37. Pina AF, Patarrão RS, Ribeiro RT, et al. Metabolic Footprint, towards Understanding Type 2  
581 Diabetes beyond Glycemia. *J Clin Med* 2020;9:2588.
- 582 38. Udler MS, Kim J, von Grotthuss M, et al. Type 2 diabetes genetic loci informed by multi-trait  
583 associations point to disease mechanisms and subtypes: A soft clustering analysis. *PLoS Med*  
584 2018;15:e1002654.
- 585 39. Xing L, Peng F, Liang Q, et al. Clinical Characteristics and Risk of Diabetic Complications in  
586 Data-Driven Clusters Among Type 2 Diabetes. *Front Endocrinol (Lausanne)* 2021;12:617628.
- 587 40. Toppila I. Identifying novel phenotype profiles of diabetic complications and their genetic  
588 components using machine learning approaches. 2016. [Epub ahead of print].
- 589 41. Vesanto J, Alhoniemi E. Clustering of the self-organizing map. *IEEE Trans Neural Networks*  
590 2000;11:586–600.
- 591 42. Kohonen T. The self-organizing map. *Proc IEEE* 1990;78:1464–80.
- 592 43. Dunn JC. A fuzzy relative of the ISODATA process and its use in detecting compact well-  
593 separated clusters. *J Cybern* 1973;3:32–57.
- 594 44. Devarajan K. Nonnegative matrix factorization: An analytical and interpretive tool in  
595 computational biology. *PLoS Comput Biol* 2008;4.
- 596 45. Cohain AT, Barrington WT, Jordan DM, et al. An integrative multiomic network model links  
597 lipid metabolism to glucose regulation in coronary artery disease. *Nat Commun* 2021;12.
- 598 46. Saeedi P, Salpea P, Karuranga S, et al. Mortality attributable to diabetes in 20–79 years old  
599 adults, 2019 estimates: Results from the International Diabetes Federation Diabetes Atlas, 9th  
600 edition. *Diabetes Res Clin Pract* 2020;162:108086.
- 601 47. Zaharia OP, Strassburger K, Strom A, et al. Risk of diabetes-associated diseases in subgroups of  
602 patients with recent-onset diabetes: a 5-year follow-up study. *Lancet Diabetes Endocrinol*  
603 2019;7:684–94.

- 604 48. Zou X, Zhou X, Zhu Z, et al. Novel subgroups of patients with adult-onset diabetes in Chinese  
605 and US populations. *Lancet Diabetes Endocrinol* 2019;7:9–11.
- 606 49. Anjana RM, Baskar V, Nair ATN, et al. Novel subgroups of type 2 diabetes and their association  
607 with microvascular outcomes in an Asian Indian population: A data-driven cluster analysis: The  
608 INSPIRED study. *BMJ Open Diabetes Res Care* 2020;8.
- 609 50. Byrne CD, Targher G. NAFLD as a driver of chronic kidney disease. *J Hepatol* 2020;72:785–  
610 801.
- 611 51. Pina AF, Borges DO, Meneses MJ, et al. Insulin: Trigger and Target of Renal Functions. *Front*  
612 *cell Dev Biol* 2020;8:519.
- 613 52. Wagner R, Heni M, Tabak AG, et al. Pathophysiology-based subphenotyping of individuals at  
614 elevated risk for type 2 diabetes. *Nat Med* 2021;27:49–57.
- 615 53. Nolan CJ, Prentki M. Insulin resistance and insulin hypersecretion in the metabolic syndrome and  
616 type 2 diabetes: Time for a conceptual framework shift. *Diabetes Vasc Dis Res* 2019;16:118–27.
- 617 54. Borges DO, Patarrão RS, Ribeiro RT, et al. Loss of postprandial insulin clearance control by  
618 Insulin-Degrading Enzyme drives dysmetabolism traits. *Metabolism*  
619 <https://doi.org/https://doi.org/10.1016/j.metabol.2021.154735>.
- 620 55. Xu W, Qiu C, Winblad B, et al. The effect of borderline diabetes on the risk of dementia and  
621 Alzheimer’s disease. *Diabetes* 2007;56:211–6.
- 622 56. Dennis JM, Shields BM, Henley WE, et al. Disease progression and treatment response in data-  
623 driven subgroups of type 2 diabetes compared with models based on simple clinical features: an  
624 analysis using clinical trial data. *Lancet Diabetes Endocrinol* 2019;8587:1–10.
- 625 57. McCarthy MI. Painting a new picture of personalised medicine for diabetes. *Diabetologia*  
626 2017;60:793–9.
- 627 58. Ahlqvist E, Prasad RB, Groop L. Subtypes of type 2 diabetes determined from clinical  
628 parameters. *Diabetes* 2020;69:2086–93.
- 629 59. Udler MS, Kim J, von Grotthuss M, et al. Clustering of Type 2 diabetes genetic loci by multi-trait  
630 associations identifies disease mechanisms and subtypes. *bioRxiv* <https://doi.org/10.1101/319509>.
- 631

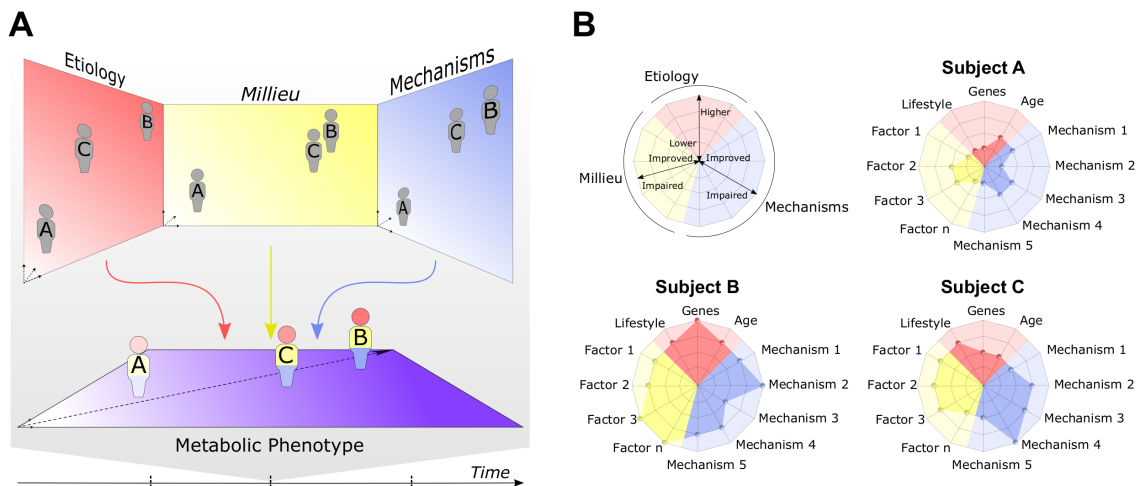
632 **Figures**



633

634 **Figure 1** – Cluster analysis scheme. An heterogenous population regarding characteristics of  
 635 interest is stratified by a chosen algorithm, that places them in a hyperplane, differentiating natural  
 636 homogenous groups.

637

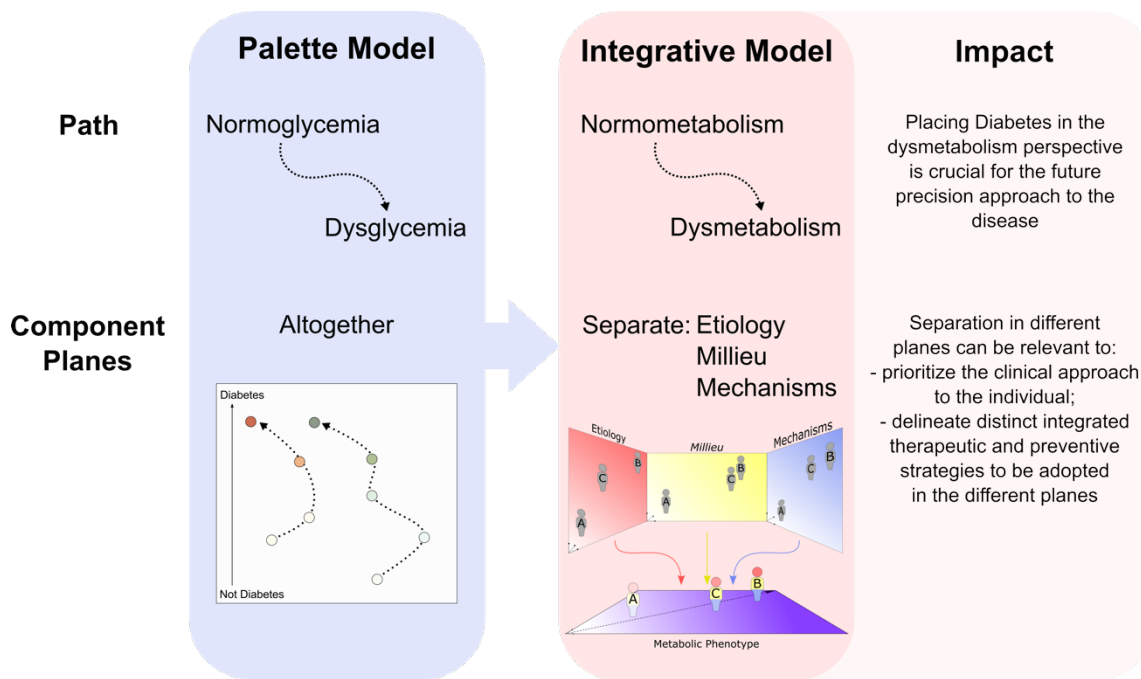


638

639 **Figure 2** – Integrative model of Diabetes. A) Subjects are deeply characterized regarding  
 640 etiological factors (including genes, lifestyle and environmental factors), underlying  
 641 physiopathological mechanisms and metabolic and hemodynamic factors that they are exposed  
 642 to. They are placed correspondingly onto the Etiology, Mechanisms and Millieu plan. The  
 643 location of a subject in each plan can be predicted by knowing their position in the others.  
 644 Ultimately, etiology, mechanisms and milieu projects the subject onto the Metabolic Phenotype

645 plan where its health condition is assessed also considering Diabetes complications as  
 646 nephropathy, retinopathy as well as cardiovascular complications. Each subject path through time  
 647 in the Metabolic Phenotype plan can be analyzed but also predicted, leveraging therapeutic and  
 648 preventive strategies. B) Etiology, Mechanisms and Millieu for each subject can be summarized  
 649 and more easily visible on a radarplot.

650



651

652 **Figure 3** – From the Palette Model to the proposed Integrative Model. The Integrative model that  
 653 we propose was based on the McCarthys' Palette Model,<sup>57</sup> but differs essentially in the path and  
 654 in the component planes of the model.



655 **Tables**

656 **Table 1** - Clustering algorithms used in Diabetes studies.

<i>Hierarchical</i>	<i>Partitioning</i>
<ul style="list-style-type: none"><li>• Agglomerative <sup>26,41</sup></li></ul>	<ul style="list-style-type: none"><li>• Hard clustering<ul style="list-style-type: none"><li>- k-means <sup>41</sup></li><li>- k-medoids (Partition around Medoids - PAM) <sup>52</sup></li><li>- Self-organising Maps (SOM) <sup>37,41</sup></li></ul></li><li>• Soft Clustering<ul style="list-style-type: none"><li>- Fuzzy c-mean <sup>59</sup></li></ul></li></ul>

657

658 **Table 2** – Advantages and drawbacks of clustering algorithms (Adapted from <sup>26</sup>).

<i>Clustering Algorithm</i>	<i>Advantages</i>	<i>Disadvantages</i>
<i>Hierarchical</i>	<ul style="list-style-type: none"> <li>• Does not need pre-specification of the number of clusters</li> <li>• Accepts any kind of distance function</li> <li>• Visualisation of number of clusters</li> <li>• Agglomerative good at identifying small clusters, divisive better identifying large clusters</li> </ul>	<ul style="list-style-type: none"> <li>• High computational cost, it does not scale properly</li> <li>• Difficult to alter once the analysis starts</li> <li>• Different clusters form according to the linkage function</li> <li>• More prone to identify spherical and convex clusters</li> <li>• Need to define the cophenetic distance cut-off</li> <li>• Sensitive to outliers</li> </ul>
<i>k-means</i>	<ul style="list-style-type: none"> <li>• Simple to implement and understand</li> <li>• Fast and efficient for large datasets</li> </ul>	<ul style="list-style-type: none"> <li>• Require specification of the number of clusters</li> <li>• Sensitive to the randomly chosen seeds</li> <li>• Some implementations use only</li> <li>• More prone to identify spherical and convex clusters</li> </ul>
<i>PAM</i>	<ul style="list-style-type: none"> <li>• Simple to understand and implement</li> <li>• Less sensitive to noise and outliers than k-means</li> <li>• Allows using general dissimilarities of objects</li> </ul>	<ul style="list-style-type: none"> <li>• Require specification of number of clusters</li> <li>• Sensitive to random initialization of medoids</li> <li>• Higher computational cost than k-means</li> <li>• More prone to identify spherical and convex clusters</li> <li>• Does not scale well for large datasets</li> </ul>
<i>SOM</i>	<ul style="list-style-type: none"> <li>• Easy to understand and interpret</li> <li>• Deals with large and complex data sets</li> <li>• Finds different clusters formats</li> </ul>	<ul style="list-style-type: none"> <li>• Many parameters to be set and optimised</li> <li>• Computational expensive</li> <li>• When initialized randomly, it is sensitive to the initial seeds</li> <li>• The number of clusters must be previously defined</li> </ul>
<i>b-NMF</i>	<ul style="list-style-type: none"> <li>• Best results for an overlapped data set</li> <li>• Datapoint may belong to more than one cluster.</li> </ul>	<ul style="list-style-type: none"> <li>• Require specification of the number of clusters</li> <li>• Computational cost</li> </ul>

659 PAM: partition around medoids; SOM: Self-Organizing Maps.

660