

# 1 **Type 1 Diabetes Risk Phenotypes Using Cluster Analysis**

2 Lu You<sup>1</sup>, Lauric A. Ferrat<sup>2</sup>, Richard A. Oram<sup>2</sup>, Hemang M. Parikh<sup>1</sup>, Andrea K. Steck<sup>3</sup>, Jeffrey

3 Krischer<sup>1</sup>, Maria J. Redondo<sup>4</sup> and the Type 1 Diabetes TrialNet Study Group\*

4 <sup>1</sup>Health Informatics Institute, Morsani College of Medicine, University of South Florida, Tampa,

5 FL, USA <sup>2</sup>University of Exeter, Exeter, UK, <sup>3</sup>Barbara Davis Center for Diabetes, University of

6 Colorado Anschutz Medical Campus, Aurora, CO, USA <sup>4</sup>Baylor College of Medicine, Texas

7 Children's Hospital, Houston, TX, USA

8 \*Listed in Supplementary Table S1

9 Corresponding author: [Lu.You@epi.usf.edu](mailto:Lu.You@epi.usf.edu)

10 Word count: 4432, Tables: 2, Figures: 4, Supplementary materials: 2.

## 1 **Abstract**

### 2 **Background**

3 Although statistical models for predicting type 1 diabetes risk have been developed, approaches  
4 that reveal clinically meaningful clusters in the at-risk population and allow for non-linear  
5 relationships between predictors are lacking. We aimed to identify and characterize clusters of  
6 islet autoantibody-positive individuals that share similar characteristics and type 1 diabetes risk.

### 7 **Methods**

8 We tested a novel outcome-guided clustering method in initially non-diabetic autoantibody-  
9 positive relatives of individuals with type 1 diabetes, using the TrialNet Pathway to Prevention  
10 (PTP) study data (n=1127). The outcome of the analysis was time to type 1 diabetes and  
11 variables in the model included demographics, genetics, metabolic factors and islet  
12 autoantibodies. An independent dataset (Diabetes Prevention Trial of Type 1 Diabetes, DPT-1  
13 study) (n=704) was used for validation.

### 14 **Findings**

15 The analysis revealed 8 clusters with varying type 1 diabetes risks, categorized into three groups.  
16 Group A had three clusters with high glucose levels and high risk. Group B included four  
17 clusters with elevated autoantibody titers. Group C had three lower-risk clusters with lower  
18 autoantibody titers and glucose levels. Within the groups, the clusters exhibit variations in  
19 characteristics such as glucose levels, C-peptide levels, age, and genetic risk. A decision rule for  
20 assigning individuals to clusters was developed. The validation dataset confirms that the clusters  
21 can identify individuals with similar characteristics.

1 **Interpretation**

2 Demographic, metabolic, immunological, and genetic markers can be used to identify clusters of  
3 distinctive characteristics and different risks of progression to type 1 diabetes among  
4 autoantibody-positive individuals with a family history of type 1 diabetes. The results also  
5 revealed the heterogeneity in the population and complex interactions between variables.

6 **Funding**

7 National Institute of Diabetes and Digestive and Kidney Diseases (R01DK121843), the National  
8 Institute of Allergy and Infectious Diseases, the Eunice Kennedy Shriver National Institute of  
9 Child Health and Human Development (cooperative agreements U01 DK061010, U01  
10 DK061034, U01 DK061042, U01 DK061058, U01 DK085461, U01 DK085465, U01  
11 DK085466, U01 DK085476, U01 DK085499, U01 DK085509, U01 DK103180, U01  
12 DK103153, U01 DK103266, U01 DK103282, U01 DK106984, U01 DK106994, U01  
13 DK107013, U01 DK107014, UC4 DK106993, UC4DK117009), and the Juvenile Diabetes  
14 Research Foundation.

## 1 **Research in Context**

### 2 **Evidence Before This Study**

3 We searched PubMed on June 12, 2023, using the keywords “cluster type 1 diabetes”,  
4 “heterogeneity type 1 diabetes”, and “endotypes type 1 diabetes” with no restrictions on  
5 publication date. Only articles published in English were considered. Existing literature suggests  
6 that individuals at risk of type 1 diabetes form a diverse population influenced by a combination  
7 of genetic and environmental factors. However, there is a scarcity of research focusing on  
8 defining clusters within this at-risk population; previous studies using clustering analysis to  
9 define diabetes subtypes have primarily utilized unsupervised clustering methods, which may not  
10 be as effective in capturing the critical variables that inform disease risks.

### 11 **Added Value of This Study**

12 We applied an outcome-guided clustering analysis to unravel the complexity and heterogeneity  
13 of type 1 diabetes. This study introduces a new method for identifying clusters of individuals  
14 based on key risk factors and the observed disease outcomes. Within each cluster, individuals  
15 will exhibit similar characteristics and diabetes risk, while the clusters themselves represent  
16 distinct levels of risk. Unlike previous approaches that employ regression models relying on  
17 linearity and additivity assumptions, this method provides advantages by uncovering underlying  
18 interactions and correlations among risk factors. Our results were validated in the DPT-1 study  
19 cohort.

### 20 **Implications of All the Available Evidence**

21 Our findings suggest that demographic, metabolic, immunological, and genetic markers can be  
22 used to identify distinct clusters with varying risks of progression to type 1 diabetes within  
23 autoantibody-positive individuals with a family history of the disease. The results demonstrate

- 1 how different combinations of these risk factors contribute to the risk and how clusters of similar
- 2 risks can exhibit unique characteristics with regard to the risk factors, which highlights the
- 3 heterogeneity in this population.

## 1 **Introduction**

2 Type 1 diabetes is a chronic disease characterized by immune-mediated destruction of insulin-  
3 producing beta cells in the pancreas<sup>1</sup>. Progression to type 1 diabetes is a complicated process that  
4 involves interactions between multiple genetic and environmental factors, leading to  
5 immunological dysfunction and metabolic abnormalities<sup>2</sup>. Heterogeneity in the causes and  
6 processes contributes to explaining the high degree of variation in the risk and rate of  
7 progression to clinical disease within individuals at risk for type 1 diabetes<sup>3</sup>. Although previous  
8 studies have identified major risk factors<sup>4,5</sup> and developed models for predicting type 1 diabetes  
9 risk<sup>6-10</sup>, there is a need to further understand the heterogeneity of the disease by exploring  
10 various levels and patterns of the risk factors in relation to one another and how they contribute  
11 to the risks associated with the disease. A majority of the prediction models in the literature used  
12 regression models to quantify risk for type 1 diabetes by a linear combination of several selected  
13 risk factors.<sup>6-10</sup> These models mainly focus on the prediction of outcomes and are usually  
14 constrained by model assumptions such as linearity and proportionality of hazard rates, that is,  
15 the supposition that predictors proportionally affect the risk along the entire spectrum of  
16 variation. Therefore, such models are usually not flexible enough to capture the complex  
17 interactions between risk factors and thus, may overlook the heterogeneity of the population.

18 The objective of this study is to use an outcome-guided clustering method to identify groups of  
19 individuals with comparable characteristics informed by their type 1 diabetes risks. The study  
20 will focus on autoantibody-positive relatives of people with type 1 diabetes. There is a lack of  
21 literature on defining clusters among individuals at risk of type 1 diabetes, and previous studies  
22 on clustering analysis of diabetes subtypes are limited to unsupervised clustering methods which  
23 may be less efficient in capturing the key variables that inform disease risks<sup>11</sup>. The proposed

1 method is a nonparametric approach (meaning that it does not rely on assumptions about the data  
2 distribution or linear relationships) for identifying clusters of individuals informed by several key  
3 risk factors ascertained from the data. Individuals within the same cluster will share similar  
4 characteristics and diabetes risk while the clusters correspond to different risks. This method  
5 differs from previous attempts to measure risk using regression models<sup>6,7</sup> in that it does not rely  
6 on stringent distributional assumptions to describe the relationship between predictors and  
7 outcomes. Instead, the proposed method is more advantageous to allow for relationships that  
8 change along the spectrum of the characteristics under study and reveal underlying interactions  
9 and correlations between risk factors to further define the etiopathogenesis of type 1 diabetes.

## 1 **Materials and Methods**

### 2 **Population**

3 The TrialNet Pathway to Prevention (PTP) study is a prospective cohort study that follows  
4 participants who have a family member with type 1 diabetes and at least one positive islet  
5 autoantibody. Enrolled participants are routinely monitored for islet autoantibodies and  
6 metabolic status, including hemoglobin A1c (HbA1c) and oral glucose tolerance tests (OGTT).  
7 This study included 1127 individuals from the TrialNet PTP study without diabetes at baseline  
8 and who were genotyped on the ImmunoChip array, which is an Illumina Infinium genotyping  
9 chip designed to study variants in genes known to be associated with autoimmune diseases<sup>12</sup>  
10 (Supplementary Figure S1). The individuals were enrolled between 2004 to 2012 and followed  
11 for a median of 2.7 years. To validate the results from the analysis, we used a dataset of 704  
12 initially non-diabetic autoantibody positive participants from the Diabetes Prevention Trial of  
13 Type 1 Diabetes (DPT-1) study.

### 14 **Data Collection**

15 The primary endpoint in this analysis is the time from baseline OGTT visits to the diagnosis of  
16 clinical (stage 3) type 1 diabetes. The last follow-up date is defined as the date of type 1 diabetes  
17 diagnosis or, for those who did not progress to diabetes, the last OGTT test. Demographic  
18 variables in the TrialNet PTP data included gender, ethnicity, relationship to proband with type 1  
19 diabetes, age and body mass index (BMI) z-scores at baseline visit. BMI z-scores for participants  
20 under 20 were calculated using the Centers for Disease Control and Prevention (CDC) growth  
21 charts while for participants over 20 years of age the formula for 20-year-olds was used. Islet  
22 autoantibodies include glutamic acid decarboxylase autoantibodies (GADA), islet antigen-2  
23 autoantibodies (IA2A), autoantibodies to insulin (mIAA), and islet cell antibodies (ICA).<sup>13</sup>



1 GADA and IA2A titers from non-harmonized assays<sup>14</sup> were transformed to the scale of the  
2 harmonized ones by the constrained nonparametric B-spline regression method in the R package  
3 “cobs”<sup>15</sup>. Metabolic functions were assessed by OGTT and HbA1c tests. The metabolic variables  
4 considered are fasting glucose levels, fasting C-peptide levels, the area under the 2-hour glucose  
5 response curve (glucose AUC), the area under the 2-hour C-peptide response curve (C-peptide  
6 AUC), and HbA1c. Genotype information was collected based on ImmunoChip data<sup>12</sup>. Genetic  
7 risk factors in the analysis include the type 1 diabetes genetic risk score-2 (GRS2)<sup>16</sup>, 8 SNPs that  
8 were shown to be associated with islet autoantibody positivity by Törn et al.<sup>17</sup>, and the 5 most  
9 susceptible human leukocyte antigen (HLA) alleles discovered in the Type 1 Diabetes Genetics  
10 Consortium<sup>18</sup> (labeled S1 to S5) (see Supplementary Tables 1 and 2 for details). The validation  
11 dataset from the DPT-1 study contains demographic variables (gender, ethnicity, relationship to  
12 proband with type 1 diabetes, age, and BMI z-scores), islet autoantibody titers data (GADA,  
13 IA2A, mIAA, and ICA), and metabolic variables (fasting glucose, fasting C-peptide, glucose  
14 AUC, and C-peptide AUC) listed above, but no ImmunoChip genotype data is available in the  
15 DPT-1 population.

## 16 [Statistical Analysis](#)

17 The descriptive analysis summarizes the distribution of continuous variables by median and  
18 interquartile range (IQR), and categorical variables by counts and proportions.

19 A statistical machine learning method developed by the authors is applied to the clustering  
20 analysis of the data. It is an outcome-guided clustering method that can identify clusters of  
21 individuals with varying levels of type 1 diabetes risk while maintaining similarity in terms of  
22 their characteristics within each cluster. Details about the clustering method is provided in the  
23 Appendix with a graphic illustration in Supplementary Figure S2. This approach ensures that the

1 identified clusters are differentiated by their disease risk, while maintaining internal homogeneity  
2 within each cluster. As a result, the method can fulfill the goal of identifying clusters of  
3 individuals with similar risks but disparate disease risks across clusters. After the clusters are  
4 identified, data and variables were visualized in a heatmap organized by their cluster  
5 membership, and disease risk in each cluster is assessed by Kaplan-Meier curves. To examine  
6 the relationship and closeness of the clusters, we used the multi-dimensional scaling (MDS)  
7 method<sup>19</sup> to map the clusters to a two-dimensional space. The distributions of variables are  
8 summarized by medians, IQR and boxplots. To examine how the risks of clusters align with the  
9 established risk scores, we summarized the distributions of risk scores (e.g., Diabetes Prevention  
10 Trial Risk Score<sup>6</sup> [DPTRS], and Index60<sup>7</sup>) by clusters. Pseudo R-squared statistics are used to  
11 examine the variation in the survival outcomes explained by the clusters as well as the risk  
12 scores. We applied the classification tree algorithm in the R package “rpart” to find a simple  
13 tree-structured decision rule for assigning new individuals to specific clusters. Finally, we used  
14 the DPT-1 cohort to validate the decision rule. Boxplots and Kaplan-Meier curves are used to  
15 compare the distribution of variables and type 1 diabetes risk of the clusters in the training and  
16 testing datasets.

### 17 [Role of the Funding Source](#)

18 The funder of the study had no role in study design, data collection, data analysis, data  
19 interpretation, or writing of the manuscript.

## 1 **Results**

2 We conducted clustering analysis of 1127 initially non-diabetic TrialNet participants with  
3 available genetic data. A total of 173 individuals were excluded as their cluster membership  
4 cannot be determined due to missing data (Supplementary Figure S1). The baseline  
5 characteristics of the 954 individuals with cluster assignment are presented in Table 1 (second  
6 column). The results from hierarchical clustering can be visualized in a tree-structured  
7 dendrogram (Column 1 of Figure 1). The average Silhouette width<sup>20</sup> that measures the similarity  
8 of individuals within a cluster is maximized at three clusters, with another peak at eight clusters.  
9 Other measures to determine the number of clusters indicate that the optimal number of clusters  
10 is either three or eight (Supplementary Figure S4), suggesting the possibility of a higher-order  
11 structure to the clustering. We analyzed the results at two levels, with 8 clusters categorized into  
12 3 groups. We labeled the three groups by A (with cluster A1), B (with clusters B1, B2, B3 and  
13 B4), and C (with clusters C1, C2, and C3) (Figure 1, Column 6). Heatmaps were used to  
14 visualize the values of variables and their cluster membership (Figure 1, Columns 2-6). Risk of  
15 progression to type 1 diabetes is represented by the Kaplan-Meier curves in Column 7 of Figure  
16 1 and Figure 2 (outer). The estimated diabetes-free survival probability in 1-5 years for each  
17 cluster is given in Table 2. The spatial distribution of the 8 clusters on the MDS map is shown in  
18 Column 8 of Figure 1 and Figure 2 (center). In Figure 2, we used radar plots to display the  
19 distribution of the variables in each cluster. The medians and IQR of variables by cluster are  
20 summarized in Supplementary Table S4.

21 Group A consists of one cluster A1 with the highest risk (2-year diabetes-free survival rate 0.44  
22 [95% CI: 0.35,0.55]). From Figure 1 and Supplementary Table S4, we can observe that this  
23 cluster has the highest glucose AUC, fasting glucose, and HbA1c.

1 Group B consists of four clusters B1, B2, B3, and B4 (2-year diabetes-free survival rate, B1:  
2 0.90 [95% CI: 0.81,0.99], B2: 0.72 [95% CI: 0.60,0.87], B3: 0.84 [95% CI: 0.76,0.93], B4:  
3 0.57 [95% CI: 0.45,0.72]). These clusters are characterized by increased autoantibody titers,  
4 especially in IA2A titer and ICA titer, and numbers of positive autoantibodies, with additional  
5 distinctive characteristics that are unique to each of them. Cluster B4 represents individuals with  
6 high glucose AUC. Cluster B2 represents individuals with younger ages, lower glucose AUC,  
7 lower C-peptide AUC, lower fasting C-peptide. The remaining two clusters B1 and B3 have  
8 similar metabolic characteristics in terms of glucose AUC, fasting glucose, C-peptide AUC, and  
9 fasting C-peptide, but their immunological characteristics are quite different. Cluster B3  
10 represents individuals with increased IA2A titer, ICA titer, and number of autoantibodies.  
11 Cluster B1 is a lower risk cluster with older age and lower IA2A titer. (See Figure 1 and  
12 Supplementary Table S4.)

13 Group C consists of three clusters C1, C2, and C3 with lower glucose and autoantibody titers (2-  
14 year diabetes-free survival rate, C1: 0.98 [95% CI: 0.96-1.00], C2: 0.90 [95% CI: 0.84-0.97],  
15 C3: 0.76 [95% CI: 0.67-0.86]). Cluster C1 is the cluster with the lowest risk, characterized by  
16 older age, lower glucose AUC, higher C-peptide AUC, lower GRS2, and lower IA2A titer.  
17 Compared to Cluster C1, Cluster C2 had an increased risk characterized by younger age, lower  
18 C-peptide AUC, lower fasting C-peptide, slightly higher GRS2, and higher number of positive  
19 autoantibodies. Cluster C3 had the highest risk among the three clusters with higher glucose  
20 AUC and higher fasting glucose. (See Figure 1 and Supplementary Table S4.)

21 To compare our results with established risk scores in the literature, in Supplementary Figure S3,  
22 we present the distributions of DPTRS<sup>6</sup> and Index60<sup>7</sup> by cluster. DPTRS is a risk score derived  
23 from OGTT measures, BMI, and age. Index60 is another risk score that combines OGTT-derived

1 glucose and C-peptide values. Both DPTRS and Index60 have been found to correlate well with  
2 disease risks<sup>7,21-23</sup>. Supplementary Figure S3 shows that there is a reasonable alignment between  
3 DPTRS, Index60, and disease risks, where higher risk clusters generally have higher DPTRS and  
4 Index60. We also fitted three separate proportional hazards models with DPTRS, Index60, and  
5 the cluster membership being the covariates respectively, and compared the pseudo R-squared  
6 values<sup>24</sup> of the three model fits. The pseudo R-squared values are 0.30 (95% CI: 0.25-0.35),  
7 0.22 (95% CI: 0.17-0.27), and 0.28 (95% CI: 0.23-0.32) respectively, which suggests that the  
8 variability in the disease outcome explained by the cluster information is comparable to the  
9 DPTRS and Index60. We note that while DPTRS and Index60 are continuous risk scores used to  
10 predict disease outcomes, the cluster information offers a unique perspective on the  
11 heterogeneity of the disease by identifying different combinations and levels of risk factors that  
12 contribute to similar risk of developing T1D, resulting in distinct groups of individuals.

13 Last, we applied the classification tree algorithm to find a decision rule to assign individuals to  
14 one of the 8 identified clusters. As shown in Figure 3, glucose AUC, IA2A titers and C-peptide  
15 AUC are the primary variables for assigning individuals to the clusters. Individuals with glucose  
16 AUC higher than 162 mg/dl are assigned to Clusters A1. Individuals with glucose AUC between  
17 144 and 162 mg/dl are assigned to clusters B4 or C3 based on IA2A titers. Individuals with  
18 glucose AUC under 144 are classified based on IA2 titers and C-peptide AUC. The cutoff values  
19 were determined based on the points that most effectively distinguish the clusters.

20 As a validation of the results, we applied the decision rules to the 706 individuals in the DPT-1  
21 dataset. The demographics of the individuals included in the validation data are given in the third  
22 column of Table 1. The distribution of variables by cluster is summarized in Supplementary  
23 Table S5. The comparison of the diabetes-free survival rate in the training and the validation

1 dataset by cluster is shown in Figure 4. The two-year diabetes-free survival rate in the 8 clusters  
2 are A1: 0.54 [95% CI: 0.44-0.67], B1: 0.98 [95% CI: 0.94-1.00], B2: 0.86 [95% CI: 0.79-0.94],  
3 B3: 0.91 [95% CI: 0.85-0.98], B4: 0.65 [95% CI: 0.53-0.79], C1: 0.93 [95% CI: 0.89-0.97],  
4 C2: 0.95 [95% CI: 0.91-0.99], C3: 0.68 [95% CI: 0.56-0.82]. P-values from log-rank tests show  
5 that there were no significant differences in the survival rate in the two populations in the  
6 identified clusters except that the lowest risk cluster, Cluster A1, has a slightly lower diabetes-  
7 free survival rate in DPT-1 than in TrialNet PTP. This difference is likely attributed to the  
8 difference in study population, with DPT-1 having more participants with higher numbers of  
9 positive autoantibodies and higher autoantibody titers (see Table 1).

10 The comparison of variable distribution in the analysis dataset and validation dataset is presented  
11 in Supplementary Figure S5. Even though the two cohorts, TrialNet PTP and DPT-1 exhibit  
12 differences in variable distribution, there is an alignment in the distribution of variables by  
13 cluster, i.e., higher variable values in a cluster tend to align with higher values in the same cluster  
14 in the other cohort. For example, the DPT-1 population had higher BMI-z scores, higher GADA  
15 titers, lower mIAA titers, and higher ICA titers compared to the TrialNet PTP population. As a  
16 result, those variables are statistically different between the two cohorts for each of the clusters.  
17 However, similar to the relationships observed in TrialNet DPT, in DPT-1, cluster C1 still has a  
18 smaller proportion of individuals with positive autoantibodies, older ages and lower HbA1c  
19 levels. Cluster C2 identifies individuals with younger ages and a slightly higher rate of type 1  
20 diabetes. Clusters B1, B2, B3 and B4 identifies individuals with larger numbers of positive  
21 autoantibodies and higher autoantibody titers. Cluster B1 and B3 have more individuals with  
22 older ages, B2 has more individuals with younger ages, and Cluster B4 has individuals with

- 1 higher glucose levels. Clusters A1 consists of individuals with high glucose levels and higher
- 2 HbA1c.

## 1 Discussion

2 Metabolic, immunological and genetic markers can be used to identify clusters of distinctive  
3 characteristics and different risks of disease progression among islet autoantibody positive  
4 individuals with a family history of type 1 diabetes. Many risk factors that are known to predict  
5 disease were revealed by our analysis as variables for categorizing clusters. The clustering results  
6 demonstrate how different combinations and levels of these risk factors can contribute to the risk  
7 of type 1 diabetes, illustrating the heterogeneity of the population. Our approach ensures that the  
8 identified clusters are differentiated by their disease risk, while maintaining internal homogeneity  
9 within each cluster. As a result, the method can fulfill the goal of identifying clusters of  
10 individuals with similar risks but disparate disease risks across clusters. There is a parallelism  
11 between the identified clusters and the staging system for type 1 diabetes that is commonly used  
12 in the literature and clinical studies<sup>25</sup>. Cluster A consists of individuals with high glucose and  
13 HbA1c with elevated autoantibody titers and thus, there is a parallelism with stage 2 type 1  
14 diabetes, i.e., dysglycemia with two or more positive islet autoantibodies, with a highly elevated  
15 risk of progression to stage 3 type 1 diabetes. Cluster B consists of participants with lower but  
16 still increased risk of type 1 diabetes, with elevated frequency of multiple autoantibodies (28·4%,  
17 45·9%, 20·6%, and 5·1%, respectively, with four, three, two and single positive autoantibodies)  
18 and autoantibody titers (especially IA2A) but, compared with Cluster A, lower glucose levels.  
19 Therefore, Cluster B has similarities with Stage 1 type 1 diabetes, i.e., two or more islet  
20 autoantibodies with normal glucose tolerance. Finally, in line with previous findings that genetic  
21 susceptibility is the major risk stratification factor for pre-stage 1 type 1 diabetes, we observed  
22 that the lowest risk cluster (C1) has the lowest type 1 diabetes GRS2.



1 While there are similarities between the clustering scheme and the stages of type 1 diabetes when  
2 categorizing individuals by their glucose and autoantibody characteristics, this new method  
3 incorporates additional factors that interact with the key characteristics. For example, Cluster B2  
4 is a high-risk cluster that cannot be identified by glycemic status alone as both glucose AUC and  
5 fasting glucose are low in this cluster. The elevated risk is primarily attributed to the presence of  
6 younger individuals with elevated autoantibody titers and genetic risk. Likewise, Cluster C2  
7 includes younger individuals with lower glucose AUC, but the risk of type 1 diabetes in this  
8 cluster is elevated due to low C-peptide, high GRS2 and high IA2A titers.

9 From Figure 3, we can see that the decision tree for assigning individuals to clusters mainly  
10 relies on the three variables glucose AUC, C-peptide AUC, and IA2A titers. However, the  
11 clustering results show that the combination of the three variables is able to explain the  
12 differences we see in other variables. For example, in Cluster B2 which we characterize by  
13 glucose AUC  $<144$ , C-peptide AUC  $<3.26$ , and IA2A titer  $\geq 137$ , the individuals have younger  
14 age, higher genetic risk, higher numbers of autoantibodies, and higher mIAA titers. In Cluster  
15 C1, which we characterize by glucose AUC  $<144$ , C-peptide AUC  $<3.56$ , and IA2A titer  $<137$ ,  
16 there tend to be more individuals with older age and lower autoantibody titers of GADA, IAA  
17 and ICA as well. Furthermore, we can observe that clusters within the same branch but  
18 distinguished by a single variable can also identify distinct populations. For example, Clusters  
19 B4 and C3 are in the same branch with glucose AUC between 144 and 162 but differentiated by  
20 IA2A titers  $\geq 38$  vs  $<38$ . Individuals in Cluster B4 have higher numbers of autoantibodies and  
21 titers of IAA and ICA while individuals in Cluster C3 have older age, higher C-peptide levels  
22 and higher HbA1c. Our results suggest that GRS2 tends to be higher in Clusters B1, B2, and B3,  
23 and lower in Cluster C1 (Supplementary Table S4). However, GRS2 and genetic variants do not

1 emerge as pivotal in defining the clusters. Additional research is required to investigate the  
2 connection between clusters and genetic factors.

3 In addition, this study shows that among individuals at similar risk of disease, there could be  
4 substantial differences in their metabolic and immunological characteristics. For example, we  
5 can compare Clusters B2 and C3 whose diabetes-free proportions are 0.87 vs 0.87 at 1 year,  
6 0.72 vs 0.76 at 2 years, and 0.54 vs 0.66 at 3 years. The two clusters are similar in type 1  
7 diabetes risk, but different in several aspects. Cluster B2 includes individuals with younger age,  
8 higher genetic risk, higher autoantibody titers, lower glucose levels, and lower C-peptide levels.  
9 In contrast, Cluster C3 includes individuals with older age, higher glucose levels, higher C-  
10 peptide levels, and higher HbA1c levels. Similarly, we can compare Clusters B1 and C2 whose  
11 diabetes-free proportions are 1.00 vs 0.96 at 1 year, 0.90 vs 0.90 at 2 years, and 0.74 vs 0.82 at 3  
12 years. We can see that Cluster B1 has higher genetic risk, higher numbers of autoantibodies  
13 while Cluster C2 has younger age, lower HbA1c and lower numbers of autoantibodies.

14 There have been previous attempts to dissect the heterogeneity among individuals with pre-  
15 clinical type 1 diabetes. Our group previously observed that, in autoantibody-positive relatives  
16 who had normal 2-hour glucose in OGTT, Index60 (which combines C-peptide and glucose  
17 measures) could stratify participants with significant differences in age, autoantibody positivity,  
18 HLA associations and risk of progression to clinical (stage 3) type 1 diabetes<sup>21</sup>. In another study,  
19 it was found that among autoantibody positive relatives, individuals with a low Index60,  
20 regardless of their glucose levels, were older and more obese, and had a lower frequency of  
21 multiple autoantibody positivity compared to those with a higher Index60.<sup>23</sup> Among patients  
22 with stage 3 type 1 diabetes, Taka et al.<sup>26</sup> found similar heterogeneity in metabolic and  
23 immunological characteristics, with the presence of HLA-associated risk associated with

1 multiple autoantibody positivity and younger age, while its absence was associated with diabetic  
2 ketoacidosis and older age. Overall, these results indicate that metabolic and immunological  
3 factors vary across the population and within groups of individuals with similar risk of  
4 progression, suggesting that there can possibly exist different disease pathways. The above  
5 observations also highlight another major difference between clustering analysis and prediction  
6 models. Prediction models assign everyone a risk score based on the individual's characteristics  
7 while the clustering analysis will also reveal subgroups of individuals with similar characteristics  
8 within each risk stratum. Defining the trajectory of individuals as they progress will help  
9 characterize different pathways in the development of type 1 diabetes. This analysis will require  
10 using longitudinal data on the natural course of preclinical type 1 diabetes.

11 In this study, we also compared this method to previous prediction models based on multiple risk  
12 factors related to metabolic, immunological, and demographic characteristics. Most of the  
13 previous literature has used regression models to identify linear combinations of risk factors for  
14 disease prediction. For example, the DPTRS<sup>6</sup> and the subsequent DPTRS60<sup>7</sup> both use the  
15 proportional hazards model to define risk scores as linear combinations of age, BMI, and glucose  
16 and C-peptide measures from OGTT; DPTRS60 is a simplified version using only the first hour  
17 of OGTT results. Similarly, Index60 is a linear combination of OGTT-derived measure<sup>7</sup>.  
18 Bediaga et al.<sup>9</sup> additionally included gender, IA2A titers, and HbA1c as risk factors in the model.  
19 Others have further considered expanding the list of predictors to incorporate genetic,  
20 demographic and other information.<sup>8,27</sup> The risk scores derived from the regression models have  
21 been used to define risk groups by certain cutoffs. For example, a DPTRS threshold of 7.0 can  
22 reliably identify normoglycemic individuals at high risk<sup>28</sup>, and an Index60 threshold of 2.0 can  
23 result in an earlier diagnosis of type 1 diabetes compared to dysglycemia<sup>7</sup>. A combination of 2-

1 hour glucose and Index60 can identify a subgroup of individuals at high risk of type 1 diabetes  
2 with younger ages, lower C-peptide, and multiple autoantibodies<sup>23</sup>. These prior efforts aimed to  
3 identify more homogenous subgroups of individuals for better prediction. However, prior  
4 methods using regression models assume that the risk increases linearly with the risk factors and  
5 neglect the interactions between risk factors that usually exist in a heterogeneous population. In  
6 contrast, our method is flexible to capture nonlinear structures when complex interactions exist.<sup>13</sup>  
7 Although the regression methods can be improved to capture the interactions by adding  
8 interaction terms or subgroup analysis, the cutpoints for defining strata or subgroups are usually  
9 specified empirically by researchers or by evenly spaced quantile points<sup>6,28,29</sup>, while others also  
10 considered model-based approaches to search for an optimal cutpoint<sup>27</sup>. In contrast, the proposed  
11 tree-based clustering algorithm determines clusters and cutpoints by an automated algorithm.  
12 Additionally, Supplementary Figure S3 shows that overlaps between distributions of DPTRS and  
13 Index60 in different clusters exist and hence our results produce different separations of risk  
14 groups. For instance, Clusters C3 and C2 have similar DPTRS means, but different distributions  
15 of glucose AUC, C-peptide AUC, and age. Likewise, Clusters B2 and B4 have similar Index60  
16 distributions, but Cluster B4 has a higher risk and higher glucose AUC compared Cluster B3  
17 (Supplementary Figure S3).

18 There are several reasons why we chose the tree-based outcome-guided clustering method over  
19 other clustering and regression methods. As opposed to other methods, ours is outcome-guided,  
20 that is, it automatically selects variables that are predictive of progression to stage 3 type 1  
21 diabetes to identify clusters. In addition, this method can easily handle datasets with a mixture of  
22 continuous, ordinal, and categorical variables. Compared to other clustering analyses in diabetes-  
23 related research, our proposed analysis is more robust and stable and optimizes the information

1 gained from the dataset. Compared to the unsupervised clustering methods that identify clusters  
2 without reference to disease outcomes, the outcome-guided method we use defines clusters with  
3 respect to disease outcomes and tries to capture the clinically meaningful features in the data.

4 The clustering method that we present has limitations too. Our clustering method has limitations;  
5 for variables not in the decision tree, it is unclear how they affect disease risk simply based on  
6 the clustering results. The correlations between variables can also be obscured without further  
7 investigation. Also, we can observe that there still exists some unexplained heterogeneity within  
8 each cluster. For example, risk factors such as the type 1 diabetes GRS2 and the number of  
9 positive autoantibodies may additionally stratify Cluster C1; however, a larger sample size with  
10 more observed events can increase the number of clusters that can be robustly observed. In  
11 addition, the method assigns each individual to a discrete cluster, which can make it difficult to  
12 determine the impact of certain predictors on disease risk in the interpretation of results. The  
13 generalizability of the results is restricted by the population and variables that we selected.

14 Finally, the analysis used participant characteristics only at baseline and thus cannot capture the  
15 longitudinal change and evolution of markers.

16 In sum, our analysis is unique in the literature as it describes differing features and their  
17 associated risks among autoantibody-positive individuals with familial predisposition for type 1  
18 diabetes. The information can be integrated into prediction models and be used to identify the  
19 most beneficial subgroups for enrollment into prevention trials. Our analysis suggests that  
20 different pathways and mechanisms of disease progressions may exist. Assessment of the  
21 longitudinal pattern of risk factors would give more insights into disease staging and pathways.

## 1 **Software Sharing Plan**

- 2 The proposed method can be implemented using the R package “SurvivalClusteringTree”
- 3 (<https://cran.r-project.org/web/packages/SurvivalClusteringTree/index.html>) developed by the
- 4 authors.

## 1 **Acknowledgments**

2 This research was funded by the National Institutes of Health (NIH) through the National  
3 Institute of Diabetes and Digestive and Kidney Diseases (R01DK121843). The Type 1 Diabetes  
4 TrialNet Study Group is a clinical trials network currently funded by the National Institutes of  
5 Health (NIH) through the National Institute of Diabetes and Digestive and Kidney Diseases, the  
6 National Institute of Allergy and Infectious Diseases, and The Eunice Kennedy Shriver National  
7 Institute of Child Health and Human Development, through the cooperative agreements U01  
8 DK061010, U01 DK061034, U01 DK061042, U01 DK061058, U01 DK085461, U01  
9 DK085465, U01 DK085466, U01 DK085476, U01 DK085499, U01 DK085509, U01  
10 DK103180, U01 DK103153, U01 DK103266, U01 DK103282, U01 DK106984, U01  
11 DK106994, U01 DK107013, U01 DK107014, UC4 DK106993, UC4DK117009, and the JDRF.  
12 The authors would like to acknowledge TrialNet participants, their families, research staff, and  
13 investigators who directly or indirectly contribute to the studies.

## 1 **References**

- 2 1. Atkinson MA, Eisenbarth GS, Michels AW. Type 1 diabetes. *The Lancet*.  
3 2014;383(9911):69-82. doi:10.1016/S0140-6736(13)60591-7
- 4 2. Bluestone JA, Herold K, Eisenbarth G. Genetics, pathogenesis and clinical interventions  
5 in type 1 diabetes. *Nature*. 2010;464(7293):1293-1300. doi:10.1038/nature08933
- 6 3. Redondo MJ, Hagopian WA, Oram R, et al. The clinical consequences of heterogeneity  
7 within and between different diabetes types. *Diabetologia*. 2020;63(10):2040-2048.  
8 doi:10.1007/s00125-020-05211-7
- 9 4. Krischer JP, Liu X, Vehik K, et al. Predicting Islet Cell Autoimmunity and Type 1  
10 Diabetes: An 8-Year TEDDY Study Progress Report. *Diabetes Care*. 2019;42(6):1051-  
11 1060. doi:10.2337/dc18-2282
- 12 5. Krischer JP, Liu X, Lernmark Å, et al. Predictors of the Initiation of Islet Autoimmunity  
13 and Progression to Multiple Autoantibodies and Clinical Diabetes: The TEDDY Study.  
14 *Diabetes Care*. 2022;45(10):2271-2281. doi:10.2337/dc21-2612
- 15 6. Sosenko JM, Krischer JP, Palmer JP, et al. A Risk Score for Type 1 Diabetes Derived  
16 From Autoantibody-Positive Participants in the Diabetes Prevention Trial–Type 1.  
17 *Diabetes Care*. 2008;31(3):528-533. doi:10.2337/dc07-1459
- 18 7. Sosenko JM, Skyler JS, DiMeglio LA, et al. A New Approach for Diagnosing Type 1  
19 Diabetes in Autoantibody-Positive Individuals Based on Prediction and Natural History.  
20 *Diabetes Care*. 2015;38(2):271-276. doi:10.2337/dc14-1813



- 1 8. Ferrat LA, Vehik K, Sharp SA, et al. A combined risk score enhances prediction of type 1  
2 diabetes among susceptible children. *Nat Med*. 2020;26(8):1247-1255.
- 3 9. Bediaga NG, Li-Wai-Suen CSN, Haller MJ, et al. Simplifying prediction of disease  
4 progression in pre-symptomatic type 1 diabetes using a single blood sample.  
5 *Diabetologia*. 2021;64(11):2432-2444. doi:10.1007/s00125-021-05523-2
- 6 10. Jacobsen LM, Larsson HE, Tamura RN, et al. Predicting progression to type 1 diabetes  
7 from ages 3 to 6 in islet autoantibody positive TEDDY children. *Pediatr Diabetes*.  
8 2019;20(3):263-270. doi:10.1111/pedi.12812
- 9 11. Meng L, Avram D, Tseng G, Huo Z. Outcome-guided Sparse K-means for Disease  
10 Subtype Discovery via Integrating Phenotypic Data with High-dimensional  
11 Transcriptomic Data. *Journal of the Royal Statistical Society: Series C*. Published online  
12 March 17, 2022. <http://arxiv.org/abs/2103.09974>
- 13 12. Cortes A, Brown MA. Promise and pitfalls of the Immunochip. *Arthritis Res Ther*.  
14 2010;13(1):101. doi:10.1186/ar3204
- 15 13. Vehik K, Beam CA, Mahon JL, et al. Development of Autoantibodies in the TrialNet  
16 Natural History Study. *Diabetes Care*. 2011;34(9):1897-1901. doi:10.2337/dc11-0560
- 17 14. Bonifacio E, Yu L, Williams AK, et al. Harmonization of Glutamic Acid Decarboxylase  
18 and Islet Antigen-2 Autoantibody Assays for National Institute of Diabetes and Digestive  
19 and Kidney Diseases Consortia. *J Clin Endocrinol Metab*. 2010;95(7):3360-3367.  
20 doi:10.1210/jc.2010-0293

- 1 15. Ng P, Maechler M. A fast and efficient implementation of qualitatively constrained  
2 quantile smoothing splines. *Stat Modelling*. 2007;7(4):315-328.  
3 doi:10.1177/1471082X0700700403
- 4 16. Sharp SA, Rich SS, Wood AR, et al. Development and standardization of an improved  
5 type 1 diabetes genetic risk score for use in newborn screening and incident diagnosis.  
6 *Diabetes Care*. 2019;42(2):200-207.
- 7 17. Törn C, Hadley D, Lee HS, et al. Role of type 1 diabetes-associated SNPs on risk of  
8 autoantibody positivity in the TEDDY study. *Diabetes*. 2015;64(5):1818-1829.
- 9 18. Erlich H, Valdes AM, Noble J, et al. HLA DR-DQ haplotypes and genotypes and type 1  
10 diabetes risk: analysis of the type 1 diabetes genetics consortium families. *Diabetes*.  
11 2008;57(4):1084-1092.
- 12 19. Cox MAA, Cox TF. Multidimensional Scaling. In: *Handbook of Data Visualization*.  
13 Springer Berlin Heidelberg; 2008:315-347. doi:10.1007/978-3-540-33037-0\_14
- 14 20. Rousseeuw PJ. Silhouettes: A graphical aid to the interpretation and validation of cluster  
15 analysis. *J Comput Appl Math*. 1987;20:53-65. doi:10.1016/0377-0427(87)90125-7
- 16 21. Nathan BM, Redondo MJ, Ismail H, et al. Index60 Identifies Individuals at Appreciable  
17 Risk for Stage 3 Among an Autoantibody-Positive Population With Normal 2-Hour  
18 Glucose Levels: Implications for Current Staging Criteria of Type 1 Diabetes. *Diabetes*  
19 *Care*. 2022;45(2):311-318. doi:10.2337/dc21-0944

- 1 22. Nathan BM, Boulware D, Geyer S, et al. Dysglycemia and Index60 as Prediagnostic End  
2 Points for Type 1 Diabetes Prevention Trials. *Diabetes Care*. 2017;40(11):1494-1499.  
3 doi:10.2337/dc17-0916
- 4 23. Redondo MJ, Nathan BM, Jacobsen LM, et al. Index60 as an additional diagnostic  
5 criterion for type 1 diabetes. *Diabetologia*. 2021;64(4):836-844. doi:10.1007/s00125-020-  
6 05365-4
- 7 24. Royston P. Explained Variation for Survival Models. *The Stata Journal: Promoting*  
8 *communications on statistics and Stata*. 2006;6(1):83-96.  
9 doi:10.1177/1536867X0600600105
- 10 25. Insel RA, Dunne JL, Atkinson MA, et al. Staging Presymptomatic Type 1 Diabetes: A  
11 Scientific Statement of JDRF, the Endocrine Society, and the American Diabetes  
12 Association. *Diabetes Care*. 2015;38(10):1964-1974. doi:10.2337/dc15-1419
- 13 26. Taka A, Härkönen T, Vähäsalo P, et al. Heterogeneity in the presentation of clinical type 1  
14 diabetes defined by the level of risk conferred by human leukocyte antigen class II  
15 genotypes. *Pediatr Diabetes*. Published online December 23, 2021.  
16 doi:10.1111/pedi.13300
- 17 27. Redondo MJ, Geyer S, Steck AK, et al. A Type 1 Diabetes Genetic Risk Score Predicts  
18 Progression of Islet Autoimmunity and Development of Type 1 Diabetes in Individuals at  
19 Risk. *Diabetes Care*. 2018;41(9):1887-1894. doi:10.2337/dc18-0087

- 1 28. Sosenko JM, Skyler JS, Mahon J, et al. Use of the Diabetes Prevention Trial-Type 1 Risk  
2 Score (DPTRS) for Improving the Accuracy of the Risk Classification of Type 1 Diabetes.  
3 *Diabetes Care*. 2014;37(4):979-984. doi:10.2337/dc13-2359
  
- 4 29. Sosenko JM, Skyler JS, Mahon J, et al. Validation of the Diabetes Prevention Trial–Type  
5 1 Risk Score in the TrialNet Natural History Study. *Diabetes Care*. 2011;34(8):1785-  
6 1787. doi:10.2337/dc11-0641
  
- 7

Table 1. Baseline characteristics of TrialNet PTP participants and DPT-1 participants included in the analysis. Numbers in the table are count (frequency) for categorical variables and median (IQR) for continuous variables.

Variable	TrialNet PTP	DPT-1
N	954 (100.0%)	706 (100.0%)
Sex	..	..
Female	503 (52.9%)	311 (44.1%)
Male	448 (47.1%)	395 (55.9%)
Ethnicity	..	..
Hispanic or Latino	77 (8.3%)	30 (4.2%)
Not Hispanic or Latino	846 (91.7%)	676 (95.8%)
Relatives with Type 1 Diabetes	..	..
Child	151 (16.0%)	184 (26.1%)
Parent	183 (19.4%)	31 (4.4%)
Sibling	488 (51.8%)	385 (54.5%)
Two or More FDR†	48 (5.1%)	47 (6.7%)
SDR*	72 (7.6%)	59 (8.4%)
BMI-z Score	0.54 (-0.22,1.21)	0.89 (0.24,1.54)
Age (Years)	11 (7,21)	11 (8,16)
Number of Positive Autoantibodies	2 (1,3)	3 (2,3)
GADA Titer (DK units/ml)	282 (30,575)	475 (39,695)
IA2A Titer (DK units/ml)	1.6 (0.1,198.0)	9 (0,262)
mIAA Titer (index)	0.005 (0.001,0.022)	0.001 (0.000,0.013)
ICA Titer (JDF units)	0 (0,80)	80 (40,320)
GADA Positivity	716 (75.1%)	518 (73.4%)
IA2A Positivity	418 (43.8%)	372 (52.7%)
mIAA Positivity	360 (37.7%)	174 (24.6%)
ICA Positivity	444 (46.5%)	703 (99.6%)
Glucose AUC (mg/dl)	129 (115,148)	128 (113,147)
Fasting Glucose (mg/dl)	89 (84,95)	86 (80,92)
C-peptide AUC (pmol/ml)	4.9 (3.7,6.6)	3.7 (2.7,4.7)
Fasting C-peptide (pmol/ml)	1.4 (1.0,1.9)	0.9 (0.6,1.3)
HbA1c (%)	5.1 (4.9,5.3)	5.3 (5.1,5.6)
Type 1 Diabetes GRS2	12.5 (11.0,13.8)	Not Available

\*SDR: second-degree relative(s)

†FDR: first-degree relative(s)

Table 2. The diabetes-free survival probability of the 8 identified clusters. Values in the parentheses are the corresponding 95% CI.

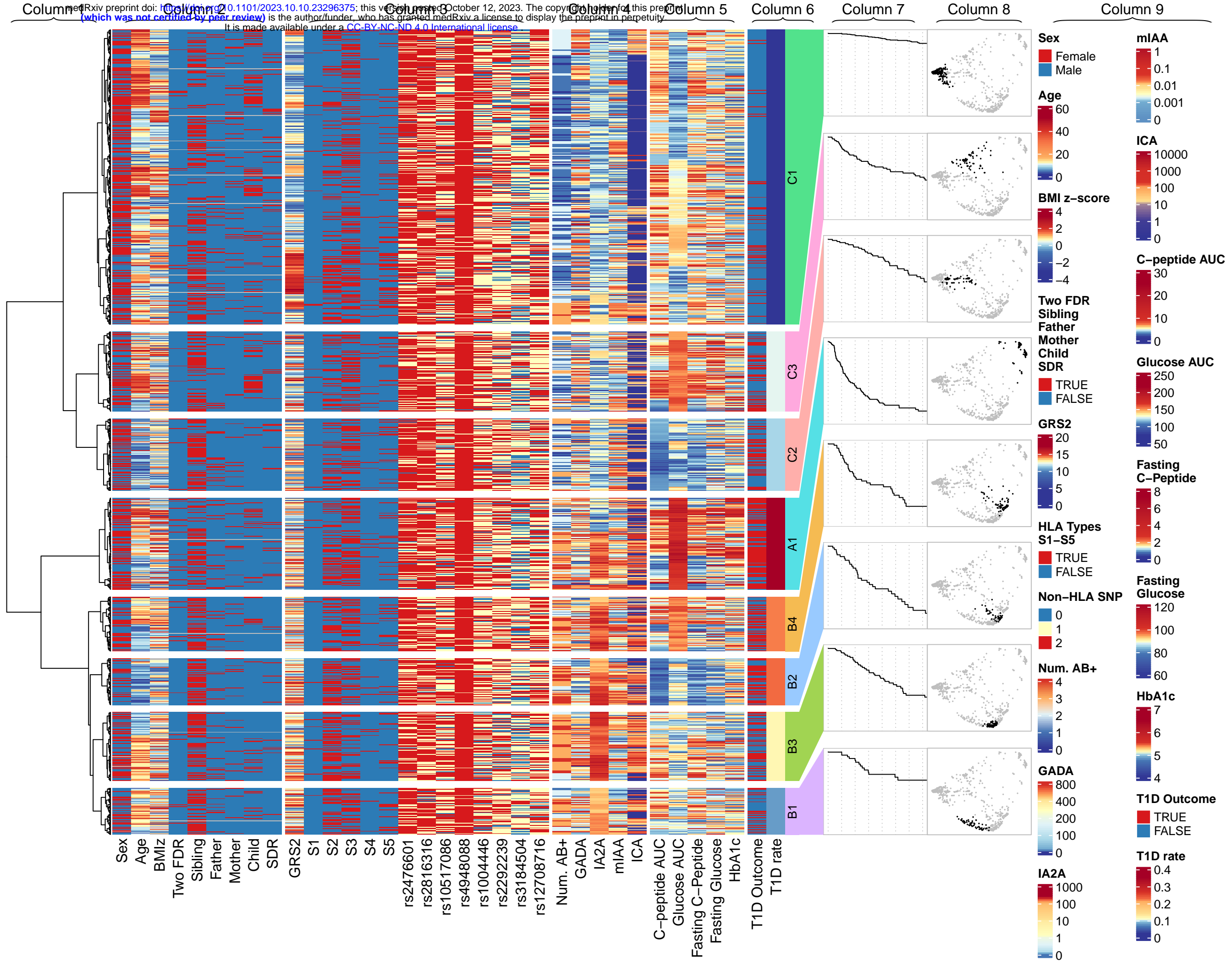
Cluster	N	1-Year Survival	2-Year Survival	3-Year Survival	4-Year Survival	5-Year Survival
A1	116	0.69 (0.60,0.78)	0.44 (0.35,0.55)	0.29 (0.21,0.40)	0.22 (0.15,0.32)	0.19 (0.13,0.30)
B1	59	1.00 (1.00,1.00)	0.90 (0.81,0.99)	0.74 (0.62,0.89)	0.72 (0.59,0.87)	0.68 (0.55,0.85)
B2	59	0.87 (0.78,0.97)	0.72 (0.60,0.87)	0.54 (0.41,0.72)	0.46 (0.33,0.64)	0.36 (0.24,0.56)
B3	87	0.94 (0.89,1.00)	0.84 (0.76,0.93)	0.73 (0.62,0.84)	0.62 (0.51,0.76)	0.54 (0.43,0.69)
B4	69	0.85 (0.76,0.95)	0.57 (0.45,0.72)	0.51 (0.39,0.66)	0.48 (0.36,0.64)	0.35 (0.23,0.52)
C1	372	1.00 (1.00,1.00)	0.98 (0.96,1.00)	0.96 (0.93,0.98)	0.94 (0.91,0.97)	0.92 (0.88,0.95)
C2	91	0.96 (0.92,1.00)	0.90 (0.84,0.97)	0.82 (0.74,0.92)	0.77 (0.68,0.88)	0.67 (0.56,0.80)
C3	101	0.87 (0.80,0.95)	0.76 (0.67,0.86)	0.66 (0.57,0.78)	0.62 (0.52,0.74)	0.56 (0.46,0.69)

Figure 1. Dendrogram and heatmap resulted from the clustering analysis. From left to right, the displayed information is the complete-linkage dendrogram (Column 1), the heatmap of demographic information (Column 2), the heatmap of genetic information (Column 3), the heatmap of immunological markers (Column 4), the heatmap of metabolic markers (Column 5), the heatmap of T1D outcomes and rates (Column 6), diabetes-free Kaplan-Meier survival curves of clusters (Column 7), distribution of clusters on the multidimensional scaling map (Column 8), and, in the foremost right side, the legends denoting the color scheme for the variables (Column 9).

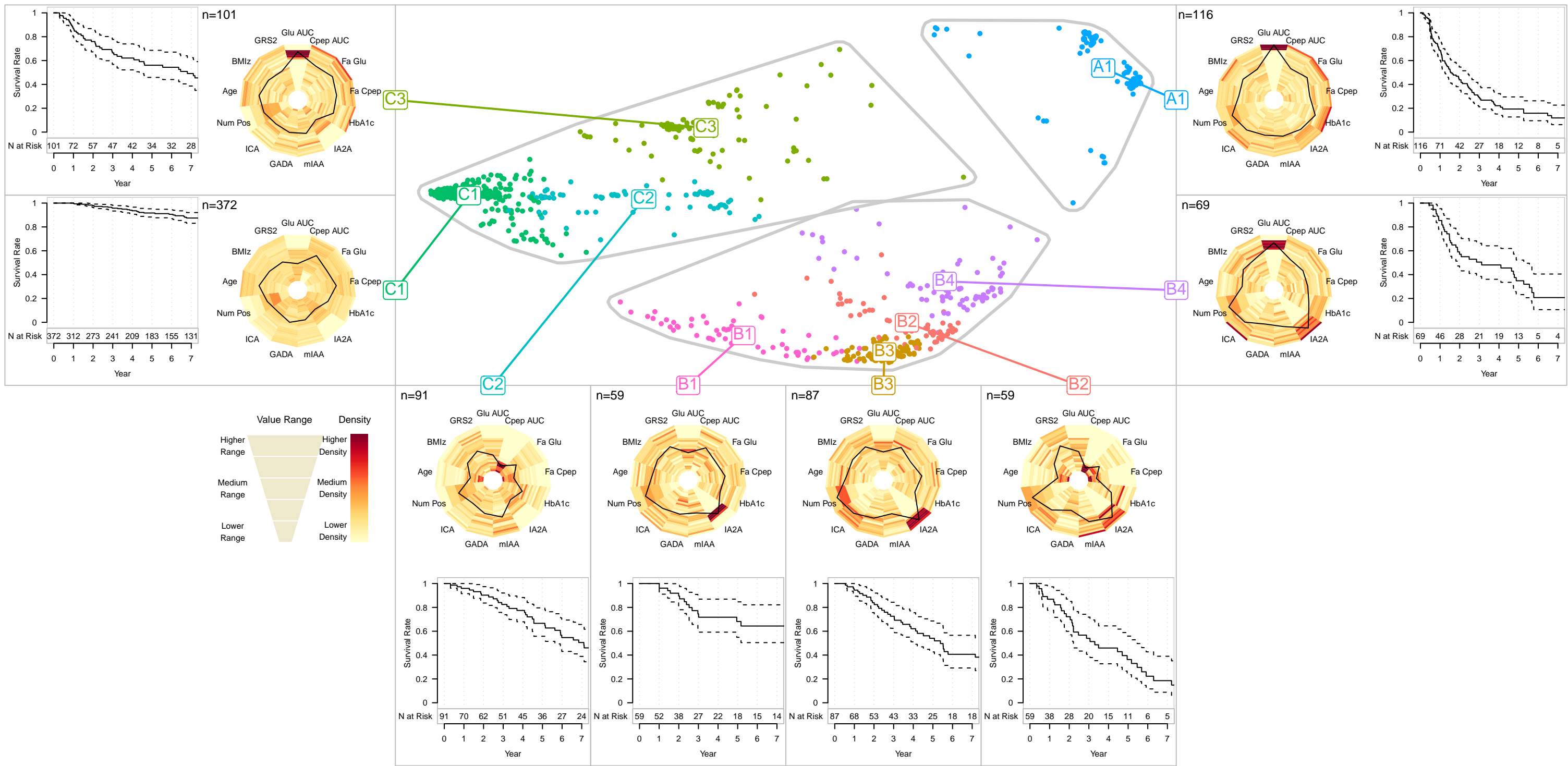
Figure 2. MDS map of clusters (center), flanked by the radar charts of each cluster (surrounding the MDS maps), and the Kaplan-Meier survival curves by cluster with their corresponding 95% confidence intervals (outer).

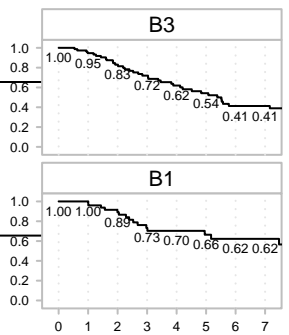
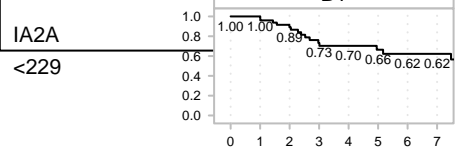
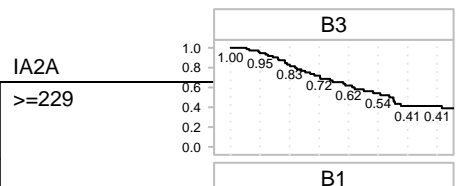
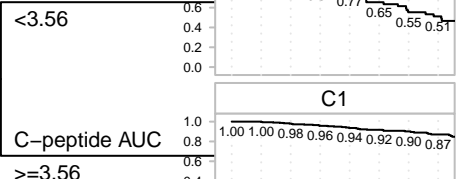
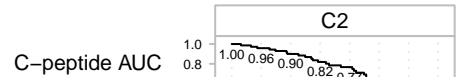
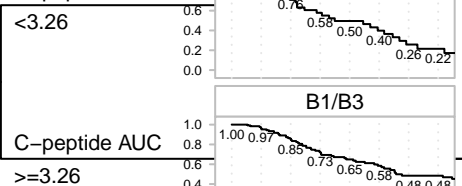
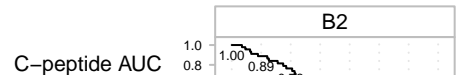
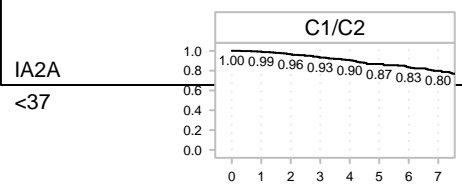
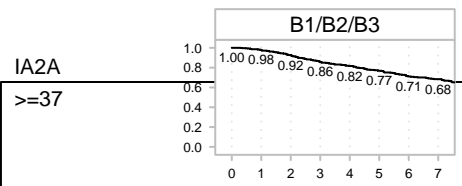
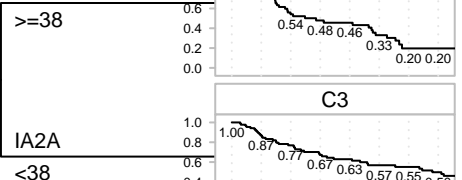
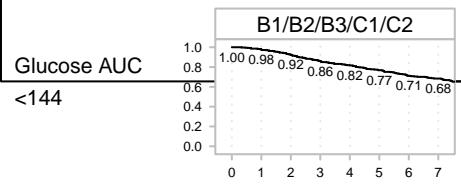
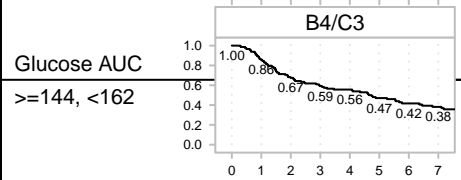
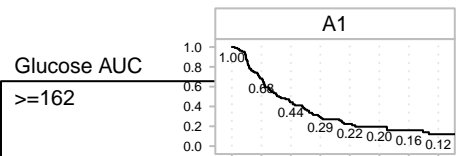
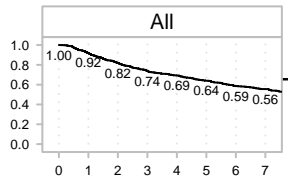
Figure 3. Decision rules to assign individuals to clusters.

Figure 4. Comparison of diabetes-free Kaplan-Meier survival curves of clusters in the analysis dataset (TrialNet PTP; black solid lines) and the validation dataset (DPT-1; blue dashed lines). The numbers under the curves are the numbers of subjects at risk in the cluster at each time point P-values on the upper right corners are based on log-rank tests.









T1D-Free Proportion

