

1 **AI-driven Integration of Multimodal Imaging Pixel Data and Genome-**
2 **wide Genotype Data Enhances Precision Health for Type 2 Diabetes:**
3 **Insights from a Large-scale Biobank Study**

4

5 Yi-Jia Huang¹, Chun-houh Chen¹, and Hsin-Chou Yang^{1,2,3,4,*}

6

7 ¹Institute of Statistical Science, Academia Sinica, Taipei, Taiwan

8 ²Biomedical Translation Research Center, Academia Sinica, Taipei, Taiwan

9 ³Department of Statistics, National Cheng Kung University, Tainan, Taiwan

10 ⁴Institute of Public Health, National Yang-Ming Chiao-Tung University, Taipei, Taiwan

11

12 *Corresponding author: Hsin-Chou Yang, Institute of Statistical Science, Academia

13 Sinica. No. 128, Sec. 2, Academia Road, Nankang 115, Taipei, Taiwan

14 (Fax) 886-2-27886833

15 (Tel) 886-2-27875686

16 (E-mail) hsinchou@stat.sinica.edu.tw

17

Abstract

18 The rising prevalence of Type 2 Diabetes (T2D) presents a critical global health
19 challenge. Effective risk assessment and prevention strategies not only improve
20 patient quality of life but also alleviate national healthcare expenditures. The
21 integration of medical imaging and genetic data from extensive biobanks, driven by
22 artificial intelligence (AI), is revolutionizing precision and smart health initiatives.

23 In this study, we applied these principles to T2D by analyzing medical images
24 (abdominal ultrasonography and bone density scans) alongside whole-genome single
25 nucleotide variations in 17,785 Han Chinese participants from the Taiwan Biobank.
26 Rigorous data cleaning and preprocessing procedures were applied. Imaging analysis
27 utilized densely connected convolutional neural networks, augmented by graph
28 neural networks to account for intra-individual image dependencies, while genetic
29 analysis employed Bayesian statistical learning to derive polygenic risk scores (PRS).
30 These modalities were integrated through eXtreme Gradient Boosting (XGBoost),
31 yielding several key findings.

32 First, pixel-based image analysis outperformed feature-centric image analysis in
33 accuracy, automation, and cost efficiency. Second, multi-modality analysis
34 significantly enhanced predictive accuracy compared to single-modality approaches.
35 Third, this comprehensive approach, combining medical imaging, genetic, and
36 demographic data, represents a promising frontier for fusion modeling, integrating AI
37 and statistical learning techniques in disease risk assessment. Our model achieved an
38 Area under the Receiver Operating Characteristic Curve (AUC) of 0.944, with an
39 accuracy of 0.875, sensitivity of 0.882, specificity of 0.875, and a Youden index of
40 0.754. Additionally, the analysis revealed significant positive correlations between
41 the multi-image risk score (MRS) and T2D, as well as between the PRS and T2D,
42 identifying high-risk subgroups within the cohort.

43 This study pioneers the integration of multimodal imaging pixels and genome-
44 wide genetic variation data for precise T2D risk assessment, advancing the
45 understanding of precision and smart health.

46

47 Introduction

48 Medical imaging has emerged as an indispensable auxiliary instrument for facilitating
49 clinical diagnostics, playing a crucial role in various applications, such as breast
50 cancer diagnosis using mammography (MMG) ¹, fatty liver diagnosis through
51 abdominal (ABD) ultrasonography ², stroke diagnosis utilizing magnetic resonance
52 imaging (MRI) and computed tomography (CT) ³, and carotid artery stenosis
53 screening using carotid artery ultrasonography (CAU) ⁴.

54 Recently, artificial intelligence (AI) has revolutionized imaging analysis and its
55 applications ^{5,6}. AI diagnostic models trained by deep learning neural networks, such
56 as convolutional neural network (CNN) ⁷ and generative adversarial network (GAN) ⁸,
57 based on different imaging-modality data, provide an automated and cost-benefit
58 way to assist medical doctors in disease diagnosis and lesion detection. AI deep
59 learning models have been developed for automated medical diagnostics for smart
60 health, including fatty liver diagnosis based on ABD images ^{9,10,11}, thyroid nodule
61 detection based on thyroid ultrasound (TU) images ¹², breast cancer diagnosis using
62 MMG images¹³, diabetic retinopathy diagnosis using color fundus (CF) images ^{14,15},
63 atrial fibrillation and normal sinus rhythm detection using electrocardiogram (ECG)
64 images ¹⁶, and osteoporosis diagnosis using bone mineral density radiography (BMD)
65 images ^{17,18}.

66 In clinical practice, Imaging-Derived Features (IDFs) are obtained either through
67 imaging technology utilizing automated measurement algorithms or through the
68 expertise of medical technologists who strategically enhance medical imaging with
69 manual or semi-manual annotations. These IDFs are crucial in helping medical
70 doctors optimize disease diagnosis and select the Region of Interest (ROI). Influential
71 IDFs can link medical imaging and diseases, constituting the primary imaging
72 biomarkers for disease diagnosis and classification.

73 **“Feature-Centric Analysis (FECA)”** and **“Pixel-Based Analysis (PIXA)”** present two
74 major imaging data analytical approaches. They employ distinct data and
75 methodologies for disease diagnosis, presenting two contrasting approaches, each
76 with its advantages and limitations. FECA leverages the IDFs derived from medical
77 doctors’ annotations and observations, incorporating expert insights from medical
78 technologists to enhance diagnostic accuracy. However, this approach significantly
79 escalates the workload and cost of extracting annotations from medical images.
80 Additionally, incorporating numerous IDFs may complicate data and variable
81 collection, diminishing clinical applicability.

82 In contrast, PIXA eliminates the need for manual imaging labeling, offering an
83 automated, low-labor, and cost-effective analysis. This approach is particularly

84 beneficial for large-scale data processing. However, due to certain technical
85 constraints, PIXA may overlook clinical expertise and contextual understanding¹⁹. This
86 raises a practical question regarding the information richness and clinical plausibility
87 of the two strategies in precision medicine: Can medical imaging inherently provide
88 all requisite information, eliminating the need for human annotation by domain
89 experts? In other words, is a hybrid approach, where data-driven methods are used
90 initially and expert consultation is sought for final decision-making, recommended?
91 Or does the knowledge of medical experts significantly contribute insights beyond
92 medical imaging for disease diagnostics and classification, necessitating their
93 inclusion at an early stage?

94 In addition to comparing PIXA vs. FECA, this study also compared multi-modality
95 analysis (MUMA) vs. single-modality analysis (SIMA) in disease risk evaluation.
96 MUMA allows for a more comprehensive assessment of the studied subject by
97 integrating structural, functional, anatomical, and molecular information from
98 multiple imaging modalities with increased sensitivity and specificity in disease
99 diagnosis, classification, and subtyping, such as MRI and PET imaging, were
100 combined to the classification of Alzheimer's disease^{20, 21}. Integration of medical
101 imaging and clinical features for breast cancer classification and subtyping^{22, 23} and
102 lung cancer classification and subtyping²⁴. Integrating these modalities provides a
103 more holistic understanding of the etiology of the studied diseases under
104 investigation. However, this approach may increase technical complexity,
105 computational demand, and data acquisition cost compared to SIMA.

106 The convergence of medical imaging and genetic data within large-scale
107 biobanks, driven by artificial intelligence and data sciences, marks a transformative
108 paradigm shift in precision health for T2D. Our previous research aimed to
109 consolidate IDFs from four distinct medical imaging modalities—abdominal
110 ultrasonography (ABD), carotid artery ultrasonography (CAU), bone density scan
111 (BMD), and electrocardiography (ECG)—alongside genome-wide single-nucleotide-
112 polymorphism (SNP) data to assess T2D risk²⁵. This innovative analysis resulted in a
113 high-accuracy risk evaluation model, polygenic risk score (PRS), and multi-imaging
114 risk score (MRS), facilitating the identification of high-risk subgroups. Moreover, the
115 model recommended eight crucial risk factors, including family history, age, fatty
116 liver, spine thickness, PRS, end-diastolic velocity in the right common carotid artery,
117 RR interval, and end-diastolic velocity in the left common carotid artery. The result
118 highlights the importance of genetics and medical imaging in a precision medicine
119 revolution.

120 Building on the previous work²⁵ based on a FECA, the current study
121 concentrates explicitly on PIXA of ABD and BMD images alongside whole-genome

122 SNPs for 17,785 Han Chinese participants from the Taiwan Biobank for T2D risk
123 evaluation. The Taiwan Biobank, a national data repository from a Han Chinese
124 population in Taiwan, aims to recruit 200 thousand participants with comprehensive
125 data, including medical imaging, whole-genome genotyping, questionnaires, and lab
126 tests ²⁶. Given the predominant focus of many biobanks on European populations,
127 the Taiwan Biobank stands out as a valuable resource for exploring medical imaging,
128 genetic data, and precision medicine in East Asian populations ²⁷.

129 This study uses CNN-based deep-learning models to analyze raw pixel data to
130 generate a convolutional activation representation. Concurrently, using Bayesian
131 statistical learning models, we analyze genome-wide SNP data to derive a PRS.
132 Subsequently, the eXtreme Gradient Boosting (XGBoost) machine learning approach
133 ²⁸ integrates the imaging activation vector, genetic PRS, and demographic variables as
134 classification and disease risk evaluation predictors. In imaging analysis, we also
135 employ a graph neural network (GNN) to account for correlations among images
136 within an individual and to integrate these multiple images into a unified
137 representation. These advancements in PIXA and MUMA, coupled with genetic and
138 demographic variable integration, present a promising avenue for developing fusion
139 models encompassing deep learning, machine learning, and statistical learning in
140 artificial intelligence and data sciences dedicated to disease risk evaluation. This
141 contributes significantly to enhancing our understanding of precision health for T2D.
142

143 **Study participants and materials**

144 **Participants**

145 The study included 17,785 Han Chinese participants from the Taiwan Biobank, each
146 of whom possessed both genetic and medical imaging data; the medical imaging
147 included abdominal (ABD) imaging data (multiple-organ images) and bone mineral
148 density (BMD) imaging data (spine, left hip, and right hip images). A participant was
149 classified as a Type 2 Diabetes (T2D) case if they self-reported T2D, and had
150 hemoglobin A1C (HbA1c) levels $\geq 6.5\%$ or fasting glucose (GLU-AC) levels ≥ 126 . A
151 control was a participant who self-reported non-T2D, with HbA1C levels $\leq 5.6\%$ and
152 GLU-AC levels < 100 . These criteria included 7,786 participants consisting of 1,118
153 T2D cases and 6,668 non-T2D controls (**Fig. 1A**).

154 Among 7,786 participants, the dataset of 7,342 participants was collected
155 earlier and initially divided into a training + validation set and testing set, named
156 “Testing Dataset 1,” at an 8:2 ratio. Subsequently, the training + validation set was
157 further divided into a training dataset and a validation dataset, named “Validation
158 Dataset,” using an 8:2 ratio. Furthermore, the training set was randomly partitioned

159 into two distinct subsets at a 5:5 ratio – the first subset, named “Training Dataset 1,”
160 was utilized for training the feature extraction model, while the second subset,
161 named “Training Dataset 2,” was employed to establish the classification model
162 independently. Finally, an additional 444 participants were recruited later, thus
163 regarded as a new cohort for the independent “Testing Dataset 2” (**Fig. 1B**).

164

165 **Demographic variables**

166 Demographic variables were collected through questionnaires. The study
167 incorporated age, sex, and family history of T2D, where the family history was
168 quantified by the count of T2D cases among the father, mother, brother, and sister
169 (ranging from 0 to 4). The epidemiological characteristics of these variables within
170 the study population were detailed (**Table S1**).

171

172 **Image pixel files and image-derived features**

173 For ABD medical imaging, each sample comprised multiple images depicting various
174 organs. The imaging data encompassed the raw image file (DICOM format) (**Fig. S1**)
175 and 28 image-derived features (IDFs) (**Table S1**). All the 28 IDFs were obtained
176 through medical experts’ assessment. For BMD medical imaging, each sample
177 included a single image for each type of BMD medical image, explicitly focusing on
178 the spine, left hip, and right hip (**Fig. S1**). The imaging data for BMD included the raw
179 image file (DICOM format) (**Fig. S1**) and 79 IDFs (**Table S1**). BMD machines
180 automatically generated all 79 IDFs. Further details about the medical imaging
181 protocol can be found on the Taiwan Biobank website
182 (<https://www.biobank.org.tw/english.php>).

183

184 **Single nucleotide polymorphisms and imputation**

185 All participants underwent genotyping at the National Center for Genomic Medicine
186 at Academia Sinica using the Axiom TWB1.0 and TWB2.0 SNP arrays, comprising 653
187 thousand and 750 thousand SNPs, respectively. Additional information on the SNP
188 annotation is available on the Taiwan Biobank website
189 (https://www.biobank.org.tw/about_value.php).

190 Sample and marker quality controls followed the procedures in our previous
191 study^{25, 29}. Pre-phasing imputation was performed for TWB1.0 and TWB2.0
192 individually using SHAPEIT2 and IMPUTE2 (v2.3.1). The imputation process yields a
193 probability distribution for each locus and each genotype of an individual. The PLINK
194 command “--hard-call-threshold 0.1” was used to convert probabilities into actual
195 genotypes, with interpretations made only when probabilities were greater than or
196 equal to 0.9. If all three genotype probabilities fell below 0.9, the locus for that

197 individual was considered missing. TWB1.0 and TWB2.0 imputation data were
198 merged, and loci with missing rates exceeding 5% and minor allele frequency (MAF)
199 less than 0.01% were subsequently removed. Finally, 9,814,944 loci remained in the
200 dataset.

201

202 **Methods**

203 **Inclusion and ethics declarations**

204 The TWB obtained written informed consent from all participants. The TWB
205 (TWBR10911-01 and TWBR11005-04) and the Institute Review Board at Academia
206 Sinica (AS-IRB01-17049 and AS-IRB01-21009) approved our data application and use.

207

208 **Image quality control and pre-processing**

209 *ABD images*

210 For ABD medical image quality control and pre-processing, we employed the
211 following steps for image quality control sequentially (**Fig. S2A**): (1) Removal of
212 images with inconsistent sizes compared to others ($m = 522$ images); (2) Exclusion of
213 images without pixel content ($m = 127$ images); (3) Elimination of B-mode images
214 with annotations ($m = 65,718$ images); (4) Removal of Doppler mode images ($m =$
215 $3,829$ images); (5) Removal of entirely black images ($m = 2$ images); (6) Removal of
216 images with multiple windows ($m = 2,090$ images); (7) Exclusion of images that are
217 not ABD ($m = 52$ images). After applying these exclusion criteria, it remains $m =$
218 $547,162$ images from $n = 22,062$ participants. Subsequently, we employed the
219 following step for image pre-processing sequentially (**Fig. S2A**): (1) Conversion of RGB
220 images into grayscale; (2) Cropping and selection of the Region of Interest (ROI) and
221 exclusion of the surrounding text. After cropping the ABD medical images, the
222 resulting image size was shortened from 614×816 pixels to 496×685 pixels (width \times
223 height) (**Fig. S2A**).

224

225 *BMD images*

226 For BMD medical image quality control and pre-processing, we employed the
227 following steps for image quality control procedures sequentially (**Fig. S2B**): (1)
228 Cropping and selecting the Region of Interest (ROI) and excluding the surrounding
229 text. (2) Remove the existing contours or bounding box from the images ($m =$
230 $22,071$). After applying these exclusion criteria, $m = 21,725$ spine images, $21,749$ left
231 hip images, and $21,749$ right hip images remained. For spine images, we removed
232 images with fewer than three vertebrae or image size height <163 pixels (**Fig. S2B**).
233 Finally, the detail of BMD image size after pre-processing was shown, and a padding

234 method was used to ensure consistent image sizes across different images
235 (**Supplementary Text 1**).

236

237 **Image activation vector extraction and classification for T2D**

238 Convolution-based DenseNet121 algorithm³⁰ was applied for activation vector
239 extraction based on the first training data (i.e., **Dataset 1** in **Fig. 1B**) and validation
240 data (i.e., **Dataset 3** in **Fig. 1B**). For ABD medical image, Average Activation Vector
241 (AAV) or Graph Neural Networks (GNN) was applied for integrative activation vectors
242 from the same sample's several ABD medical images. XGBoost algorithm²⁸ was
243 applied to classify T2D based on imaging activation vectors, genetic PRS, and
244 demographic variables. Classification was trained and validated based on the second
245 training data (i.e., **Dataset 2** in **Fig. 1B**) and validation data (i.e., **Dataset 3** in **Fig. 1B**).
246 All classification models were tested based on the first testing data (i.e., **Dataset 4** in
247 **Fig. 1B**). The best model was further validated based on the second independent
248 testing dataset (i.e., **Dataset 5** in **Fig. 1B**).

249 The DenseNet121 architecture was shown (**Fig. 2A**). The DenseNet121 models
250 were trained with the following settings: image size (64 x 64, 224 x 224, 256 x 256),
251 batch size (64, 128, 256), and pre-trained (w/ or w/o). All the DenseNet121 models
252 were trained using an initial learning rate of 0.05 for 500 epochs. The best
253 parameters set was used.

254 The GNN architecture was shown (**Fig. 2B**). For GNN, we obtained localized node
255 embedding using "GNNConv"³¹ applied from PyTorch Geometric (PyG)³² and used
256 mean or weighted mean pooling for a graph readout before applying classifier. The
257 GNN models were trained with the following parameters: edge weight calculated by
258 Euclidean distance or cosine similarity, number of features after convolution (64 or
259 1024), and graph readout according to node weight given equal weight (mean
260 pooling) or weighted by number of edge node. The best distance cutoff between two
261 nodes determines whether an edge exists between two nodes. All the GNN models
262 are trained using an initial learning rate 0.05 for 500 epochs. The best parameters set
263 was used.

264 The XGBoost models²⁸ were trained with the following default parameter
265 settings: maximum depth equal to 6, learning rate equal to 0.3, the value of the
266 regularization parameter alpha (L1) was set to 0, and lambda (L2) was set as 1, the
267 number of boosting stages was 100, and the early-stop parameter was set to 30.
268 Parameter tuning was conducted to establish the best model (**Supplemental Table**
269 **S2**).

270 The model's effectiveness was evaluated by computing the area under the
271 receiver operating curve (AUC). The performance of the created models was

272 assessed using accuracy, sensitivity, specificity, and Youden index metrics. The
273 optimal threshold value for the XGBoost model on the validation data was
274 determined using the Youden index³³.

275

276 **Multi-image Risk Score and Polygenic Risk Score**

277 *Multi-image Risk Score (MRS)*

278 The computation of the MRS involved a sophisticated process. Initially, imaging pixels
279 were utilized as input, and an activation vector was meticulously derived and
280 consolidated through a series of operations, including convolution, pooling,
281 transition, and dense block within the DenseNet121 architecture (refer to **Fig. 2A** for
282 an illustration). Subsequently, this activation vector, also known as the feature vector,
283 was constructed. The feature vector was the input variable for assessing the T2D
284 disease status employing the XGBoost classifier²⁸. Notably, the XGBoost feature
285 importance algorithm was applied to identify crucial features, and the MRS was
286 ultimately calculated as the likelihood of an individual being classified as a T2D case.

287

288 *Polygenic risk score (PRS)*

289 The PRS construction closely followed the methodology in our prior investigation²⁵.
290 PRS-CSx³⁴ was employed, utilizing meta-GWAS summary statistics for T2D across
291 diverse ancestral populations. Specifically, the data encompassed East Asian (EAS)
292 populations (56,268 cases and 227,155 controls from the DIAGRAM Consortium³⁵),
293 European (EUR) populations (80,154 cases and 853,816 controls from the DIAGRAM
294 Consortium³⁵), and South Asian (SAS) populations (16,540 cases and 32,952 controls
295 from the DIAGRAM Consortium³⁵). Additionally, Linkage Disequilibrium (LD)
296 references from the 1000 Genomes Project³⁶ for each of the three populations (EAS,
297 EUR, and SAS) were incorporated. Weights of 884,327, 880,098, and 900,047 SNPs
298 for EAS, EUR, and SAS were applied to our genotype data to calculate the population-
299 specific PRS using the PLINK (--score command) tool. Subsequently, we combined the
300 population-specific PRS with equal weights to derive the final PRS. The R language
301 was utilized to standardize the PRS, setting the mean to 0 and the standard deviation
302 to 1.

303

304 **Results**

305 **All T2D risk evaluation models**

306 We constructed 12 T2D risk evaluation models (**M1–M12** in **Table 1**) by various
307 combinations of conditions, including imaging type (ABD and BMD), image analysis
308 unit (sample-based and image-based), image analysis type (FECA vs. PIXA), and

309 analysis modality (MUMA vs. SIMA).

310

311 **Pixel-based analysis (PIXA) demonstrated competitiveness compared to feature-**
312 **centric analysis (FECA)**

313 We constructed T2D risk evaluation models utilizing raw pixel or IDF data of ABD and
314 BMD medical imaging as predictors. In the case of ABD, compared to FECA (Model
315 **M1** in **Table 1**), PIXA (Model **M2** in **Table 1**) outperformed across all performance
316 metrics (**Table 1** and **Fig. 3A**). In the case of BMD, compared to the FECA (Model **M5**
317 in **Table 1**), PIXA (Model **M6** in **Table 1**) demonstrated similar performance (see
318 **Table 1** and **Fig. 3B**). The results highlight that PIXA offered valuable information for
319 T2D risk evaluation, even in the absence of consideration of clinical data from
320 medical technologists.

321

322 **Multi-modality analysis (MUMA) outperformed single-modality analysis (SIMA)**

323 *Combining ABD and BMD imaging*

324 A multi-modality PIXA, which concurrently analyzed the raw pixel data of both ABD
325 and BMD imaging data (Model **M11** in **Table 1**), consistently outperformed the
326 individual-modality PIXA of ABD imaging (Model **M3** in **Table 1**) and BMD imaging
327 (Model **M6** in **Table 1**) across majority of performance measures, particularly in AUC,
328 ACC, SPEC, and Youden index (refer to **Fig. 3C**). The results underscore the enhanced
329 performance of a MUMA compared to a SIMA, although caution is advised due to
330 the associated higher data-collection cost in a MUMA. A similar finding applies to
331 FECA: multi-modality FECA (Model **M10** in **Table 1**) outperforms a single-modality
332 FECA: ABD imaging (Model **M1** in **Table 1**) and BMD imaging (Model **M5** in **Table 1**)
333 (**Fig. S3**).

334

335 *Combining spine, left hip, and right hip BMD imaging*

336 As to BMD imaging, which consists of the spine (Model **M7** in **Table 1**), left hip
337 (Model **M8** in **Table 1**), and right hip (Model **M9** in **Table 1**) medical imaging, the
338 three individual PIXAs exhibited a close performance in T2D classification (**Fig. 3D**). A
339 MUMA which integrated spine, left hip, and right hip medical imaging (Model **M6** in
340 **Table 1**) exhibited an improvement compared to the three individual SIMAs (Models
341 **M7–M9** in **Table 1**), particularly in increased AUC, SEN, and Youden index (**Fig. 3D**).
342 The results once again illustrated a better performance for a MUMA compared to a
343 SIMA.

344

345 **Robustness analysis**

346 *Robustness analysis of DenseNet121 parameters suggests the constructed models*

347 *are robust*

348 Our robustness analysis for ABD imaging examined three image sizes (64×64,
349 128×128, and 224×224), three batch sizes (64, 128, and 256), and usage of the pre-
350 trained model (yes or no) (**Fig. 4A**). All combinations of these settings demonstrated
351 similar performance in terms of AUC. The model attaining the highest AUC of 0.800 is
352 characterized by an image size of 224 x 224 pixels, batch size of 64, and no pre-
353 trained weights (**Fig. 4A**). These parameters were fixed in our subsequent analysis.
354 Other parameters in DenseNet121 were set as default values. The final parameters
355 were also applied to BMD imaging.

356

357 *Robustness analysis of Graph Neural Network parameters suggests the constructed*
358 *models are robust*

359 GNN was applied to account for the intra-individual dependency of multiple images
360 of ABD available for each individual. Our robustness analysis for GNN considered two
361 edge weights W_E (“Euclidean distance” and “Cosine similarity”), two node weights
362 W_N (“Equal-weight” and “Number of edges connected to a node”), and two numbers
363 of features after the graphical convolution (64 or 1,024) (**Fig. 4B**). All combinations of
364 these settings exhibited similar performance in AUC, suggesting that the constructed
365 GNN models are robust. The model attaining the highest AUC of 0.887 is
366 characterized by Euclidean-distance edge weight ($W_E = ED$), 1,024 features after
367 convolution ($F = 1,204$), and node weight proportional to the number of edge nodes
368 ($W_N = n_{edge}$) (**Fig. 4B**). Existence of an edge/link between two nodes was determined
369 by a threshold “ED”. The results show that the optimal ED cutoff, which attained the
370 highest AUC of 0.887, was $ED = 0.5$ (**Fig. S4**).

371

372 *Robustness analysis suggests the sample-based method outperforms the image-*
373 *based method and accounting for within-sample correlation further improves*
374 *performance*

375 The results suggest that the sample-based method (Models **M2** and **M3** in **Table 1**),
376 which integrates multiple images of each individual, outperforms the image-based
377 method (Model **M4** in **Table 1**) (**Fig. 4C**). Furthermore, a sample-based method using
378 GNN (i.e., the model with accounting for within-sample image correlation) (Model
379 **M3** in **Table 1**) is slightly better than the sample-based method using a direct cross-
380 image average (i.e., without accounting for image correlation) (Model **M2** in **Table 1**)
381 (**Fig. 4C**).

382

383 *Robustness analysis of classifiers (XGBoost and MLP) have similar performance,*
384 *suggesting the finding is robust*

385 To examine the robustness of our findings concerning classifiers, in addition to the
386 eXtreme Gradient Boosting (XGBoost) classifier, we also compared it with the multi-
387 layer perception (MLP) classifier based on the ABD imaging. The results show that
388 XGBoost and MLP exhibited similar results in terms of AUC, ACC, SEN, SPEC, and
389 Youden index (**Fig. 4D**), which were averaged across three PIXA variants – sample-
390 based PIXA with an average weight (Model **M2** in **Table 1**), sample-based PIXA with a
391 GNN weight (Model **M3** in **Table 1**), and image-based PIXA (Model **M4** in **Table 1**).
392 The detailed results of the three models based on ABD (**Fig. S5A**) and BMD (**Fig. S5B**)
393 are demonstrated.

394

395 **Clinical consideration**

396 *Integrative model*

397 In addition to the imaging data, other crucial features for T2D diagnosis including
398 genetic component (PRS) and demographic variables (sex, age, and family history of
399 T2D) were also considered in the ABD genetic–imaging integrative analysis (**Fig. 5A-1**)
400 and BMD genetic–imaging integrative analysis (**Fig. 5A-2**), and ABD+BMD integrative
401 analysis (**Fig. 5A-3**); additional results for the spine, left hip, and right hip are
402 presented (**Figs. S6A-1 – S6A-3**). Considering five performance metrics, the
403 combination of imaging features, demographic factors, and genetic PRS performed
404 the best. Imaging features performed exceptionally well, surpassing the performance
405 of demographic characteristics and genetic PRS when considered individually,
406 particularly in the ABD analysis.

407

408 *Best risk assessment model for T2D*

409 The best medical imaging model (**M11** in **Table 1**) is a sample-based PIXA and
410 MUMA, which combines ABD and BMD imaging (spine, left hip, and right hip
411 imaging). On top of the medical imaging in this model, demographic variables (D) and
412 genetic components (G) were further included. Finally, the best risk evaluation model
413 is “ $ABD_{GNN}^{PIX} + BMD_{COMBINE}^{PIX} + G + D$ ” (**M12** in **Table 1**), attaining AUC = 0.944, ACC
414 = 0.868, SEN = 0.889, SPE = 0.865, and Youden index = 0.754 in the first testing data
415 based on a classification threshold 0.01. The model was further validated well in the
416 second independent testing dataset with AUC = 0.954, ACC = 0.875, SEN = 0.882, SPE
417 = 0.875, and Youden index = 0.757. The XGBoost model parameter tuning and
418 performance evaluation for this model are shown (**Supplementary Table S2**).

419

420 *Multi-image risk score (MRS)*

421 ABD-based and BMD-based multi-image risk scores (MRSs) were calculated for each
422 participant. The odds ratio of T2D and its corresponding confidence interval revealed

423 a positive correlation between MRS and T2D (**Fig. 5B-1** for ABD and **Fig. 5B-2** for
424 BMD). This suggests that the risk of developing T2D increases with higher MRS values
425 than the reference group: 40%–60% decile group of MRS. These findings could
426 potentially lead to more effective treatments in the future. Similar results can also be
427 found at the SIMA of the spine (**Fig. S6B-1**), left hip (**Fig. S6B-2**), and right hip (**Fig.**
428 **S6B-3**).

429

430 *Identification of high-risk subgroups using MRSs*

431 Furthermore, we identified specific high-risk subgroups within the study population.
432 In the analyses of ABD, BMD, and ABD+BMD, consistently high T2D risk (i.e., the ratio
433 of the number of cases vs. controls) was observed within the high MRS group (90–
434 100% decile group) for both women and men aged older than 62 years with a family
435 history of T2D, as follows: In the ABD analysis, T2D risks were 23 in the female group
436 (**Fig. 5C-1**) and 8 in the male group. In the BMD analysis, T2D risks were 5 in the male
437 group (**Fig. 5C-2**) and 4.4 in the female group. In the ABD+BMD analysis, T2D risks
438 were 21 in the female group (**Fig. 5C-3**) and 8 in the male group. These T2D risks in
439 the 90–100% MRS decile group were significantly higher than those in the lower-MRS
440 groups.

441

442 **Discussion**

443 **High-performance genetic-imaging integrated analysis of T2D risk evaluation and** 444 **diagnosis**

445 In our previous work²⁵, which considered an AI-enhanced integration of genetic and
446 medical imaging data for T2D risk assessment, a FECA analysis based on the IDFs of
447 four medical images (ABD, BMD, CAU, and ECG) in the Taiwan Biobank achieved an
448 AUC of 0.880, increasing to 0.945 after incorporating demographic (age and family
449 history of T2D) and genetic information (PRS). The current study, focusing on two
450 T2D-related medical images (ABD and BMD), demonstrates that PIXA outperforms
451 FECA for T2D risk evaluation. Despite analyzing fewer types of images than our
452 previous study²⁵, this investigation achieved an AUC of 0.902 based on ABD and
453 BMD pixel data. The AUC further increased to 0.953 after incorporating demographic
454 and genetic information. The result underscores the potential of an integrated multi-
455 modality study of genetic analysis and medical imaging PIXA for precision medicine.

456 For precision medicine, genetic information (genome-wide SNPs) and medical
457 imaging data (image-wide pixels) provide individualized information compared to
458 traditional biochemistry and body measurement indices, such as HbA1c and fasting
459 glucose have demonstrated limitations in disease risk prediction, as highlighted in

460 various studies^{37, 38, 39, 40, 41, 42, 43, 44}. These metrics often show limited sensitivity,
461 especially in specific populations, and lack accuracy in predicting pre-diabetes and
462 T2D. Genetic data offers stability for evaluating disease risk and diagnosis in
463 precision medicine, while medical imaging data provides detailed multi-modality
464 information for human organs, contributing to stable and insightful disease risk
465 evaluation, diagnosis, and classification.

466

467 **Comparison of pixel-based analysis (PIXA) and feature-centric analysis (FECA)**

468 Our investigation reveals that PIXA demonstrates competitive and, in some cases,
469 superior performance compared to FECA in T2D classification (**Figs. 3A and 3B**). In
470 the specific case examined in this paper, two potential explanations warrant careful
471 consideration. Firstly, it is plausible that IDFs did not entirely extract the complete
472 information embedded within the raw pixel data. Specifically, only 28 IDFs are
473 extracted in the case of ABD, and some of these features may lack direct association
474 with T2D. Secondly, IDFs defined by medical technologists might be suboptimal. For
475 instance, accurately labeling the exact fatty liver level (normal, mild, moderate, and
476 severe) poses challenges, particularly for closely related levels at the borderline.
477 Consequently, PIXA, without labor-intensive labeling and expensive annotation,
478 offers precise feature quantification with artificial intelligence, providing a high-
479 performance solution for disease classification and risk assessment. This approach
480 paves the way for effective and practical clinical applications.

481

482 **Comparison of MUMA and SIMA**

483 Our investigation reveals that a MUMA provides superior performance compared to
484 a SIMA (**Figs. 3C and 3D**) if there is non-overlapping information in the images of
485 different modalities. Our integrated analysis of ABD and BMD medical imaging
486 outperforms individual analyses of ABD and BMD (**Fig. 3C**). The enhanced
487 performance in the integrated analysis can be attributed to the non-overlapping
488 contributions of ABD and BMD to T2D. Conversely, our integrated classification
489 analysis of the spine, left hip, and right hip medical imaging performed similarly to
490 the three individual analyses (**Fig. 3D**). This suggests that these imaging modalities
491 provide highly correlated and redundant information for T2D despite representing
492 different body sections. These observations underscore the importance of carefully
493 selecting and integrating imaging modalities for disease classification and
494 considering each modality's unique contributions to enhance overall diagnostic
495 accuracy. Our findings align with previous studies demonstrating the superiority of
496 MUMA over SIMA in disease classification^{45, 46, 47, 48}.

497

498 **Robustness analysis for FECA**

499 We conducted a robustness analysis for the FECA, considering various factors such as
500 pre-trained models, image sizes, batch sizes, and classifiers. Firstly, employing a pre-
501 trained DenseNet121 model^{49, 50, 51, 52} did not enhance performance (**Fig. 4A**),
502 potentially due to differences in characteristics between ImageNet⁵³ and medical
503 imaging data and the limited number of images in ImageNet. Secondly, variations in
504 image sizes and batch sizes demonstrated minimal impact on AUCs (**Fig. 4A**). Thirdly,
505 variations in edge weight methods – “Euclidean distance (DS)” and “Cosine similarity
506 (CS),” node weight methods – “Equal weight (EW)” and “Unequal weight
507 proportional to the number of nodes connected to”), and several nodes after
508 convolution exhibited close AUCs (**Fig. 4B**). Thirdly, GNN which considers within-
509 individual image correlation performed slightly better than AVE which does not
510 consider the correlation (**Fig. 4C**); however, the difference in performance is limited.
511 Finally, the alternative classifier – multilayer perceptron (MLP), exhibited no
512 significant difference in performance compared to XGBoost (**Fig. 4D**). This
513 consistency across different classifiers underscores the robustness of our findings,
514 enhancing the credibility and generalizability of our proposed approach for T2D risk
515 evaluation and classification.

516

517 **Conclusion**

518 In conclusion, this study highlights the compelling findings that applying artificial
519 intelligence, comprising deep learning and machine learning, to integrated genetic
520 and medical imaging PIXA provides a fully automated, low-labor, cost-saving, and
521 high-accuracy analysis. Incorporating multi-modality data, encompassing diverse-
522 dimensional information, significantly enhances the performance compared to
523 single-modality data analysis. Notably, medical imaging PIXA emerges as a
524 competitive and, in many instances, superior performer compared to FECA.
525 Integrating genome-wide genetic data with multi-modality imaging marks a
526 revolutionary advancement in precision medicine and smart health for T2D. These
527 results provide crucial insights into the potential transformative impact of advanced
528 analytical methodologies on the future of T2D diagnosis and personalized
529 healthcare.

530

531 **Acknowledgments**

532 This work was supported by research grants from the Academia Sinica (AS-PH-109-01
533 and AS-SH-112-01). Data application and use were approved by the Taiwan Biobank
534 and the Institute Review Board (AS-IRB01-17049 and AS-IRB01-21009). We gratefully

535 acknowledge the Taiwan Biobank for providing the data used in this research. We
536 also extend our thanks to all the participants of the Taiwan Biobank for their
537 invaluable contributions. Technical support in genotyping from the National Center
538 for Genome Medicine of Taiwan is also acknowledged. We thank team members
539 Miss Chih-Ting Yang and Mr. Po-Wen Chen for imaging preprocessing of ABD and
540 BMD and Mr. Chia-Wei Chen and Dr. Shih-Kai Chu for genetic data quality control.
541

542 **Author Contributions Statement**

543 H.C.Y. conceptualized and supervised the study. Y.J.H. curated the data and applied
544 software. Y.J.H. & H.C.Y. conducted formal data analysis, visualized the results, and
545 wrote the paper. C.h.C. & H.C.Y. secured funding and provided resources. H.C.Y., Y.J.H.,
546 and C.h.C. validated the results.
547

548 **Competing Interests Statement**

549 The authors declare that they have no competing interests.
550

551 **Ethics declarations**

552 Competing interests

553 The authors declare that they have no competing interests.
554

555 **Data availability statement**

556 The data analyzed in this study were obtained from the Taiwan Biobank with proper
557 approval. The Taiwan Biobank retains ownership rights, so the data have not been
558 deposited in a public repository. Researchers interested in accessing the data must
559 apply through the Taiwan Biobank's formal process. Detailed instructions for data
560 access requests can be found on the Taiwan Biobank's official website
561 (<https://www.twbiobank.org.tw/index.php>). This paper provides Source data in the
562 Supplementary Information and Source Data files. Meta-GWAS summary statistics of
563 T2D in multiple populations from the DIAGRAM Consortium are available at
564 <https://diagram-consortium.org/downloads.html>. The linkage disequilibrium
565 reference from various populations of the 1000 Genomes Project can be downloaded
566 from <https://github.com/getian107/PRScsx>.
567

568 **Code availability statement**

569 We provide code at the repository at
570 https://github.com/yjhuang1119/Medical_Image_Risk_Assessment_Model for
571 medical image classification and risk assessment using a combination of Densely
572 Connected Convolutional Networks with 121 layers (DenseNet121) and eXtreme
573 Gradient Boosting (XGBoost). The pipeline includes establishing a disease risk
574 assessment model using DenseNet121, extracting feature maps, constructing a final
575 disease risk assessment model using XGBoost, and performance evaluation. The code
576 also computes performance metrics for model evaluation and feature importance
577 scores for model explainability. A README is provided.
578

579 **Supplemental information**

580 Supplemental information is available online.
581

582 **Figure legends**

583 **Figure 1. Flowchart of the study. (A) Data extraction and classification model**
584 **building.** In this dataset, 21,927 individuals possess ABD and BMD medical imaging
585 data, while 108,251 individuals have whole-genome genotyping and imputation data.
586 Moreover, 17,785 individuals possess both medical imaging and genetic profile data.
587 Finally, the final dataset comprises 7,786 individuals who meet the inclusion criteria
588 based on self-reported T2D status, HbA1C, and fasting glucose, consisting of 1,118
589 T2D cases and 6,668 normal controls. The complete dataset was initially divided into
590 training + validation and testing sets at 8:2. Subsequently, the training + validation
591 set was further separated into training and validation datasets with an 8:2 ratio.
592 The training set was divided into two independent subsets at a 5:5 ratio to mitigate
593 the winner's curse problem; a two-stage procedure was employed for feature
594 extraction and classification. In the first stage, which was dedicated to feature
595 extraction, the first training data was used to establish a DenseNet121 model based
596 on the initial training dataset (Training Dataset 1). Subsequently, a feature map
597 vector was obtained. The second stage was focused on sample classification. Utilizing
598 the data from the second training (Training Dataset 2) and validation datasets, a deep
599 learning model for T2D classification was developed, and the results were further
600 confirmed using the Testing Dataset 1. Ultimately, the best model's validity was
601 further confirmed using the second independent testing dataset (Testing Dataset 2)
602 (n = 444). Regarding the deep learning classification model, three methods—
603 multilayer perceptron (MLP), graph neural network (GNN), and eXtreme Gradient
604 Boosting (XGBoost)—were implemented. **(B) Sample Size.** Information on the total

605 sample size, number of cases, and number of controls is provided.

606

607 **Figure 2. Architecture diagrams. (A) Densely Connected Convolutional Neural**
608 **Networks with 121 layers (DenseNet-121). (B) Graph Neural Network (GNN).**

609

610 **Figure 3. Model Comparison. (A) Comparison of PIXA (ABD_{AVE}^{PIX}) and FECA**
611 **(ABD^{IDF}) in ABD imaging analysis.** PIXA exhibited superior performance in T2D risk
612 evaluation compared to FECA. **(B) Comparison of PIXA ($BMD_{COMBINE}^{PIX}$) and FECA**
613 **(BMD^{IDF}) in BMD imaging analysis.** PIXA and FECA exhibited similar performance in
614 T2D risk evaluation. **(C) Comparison of multi-modality analysis (MUMA) and single-**
615 **modality analysis (MUMA) in ABD and BMD imaging analysis.** MUMA of ABD and
616 BMD outperforms SIMA of ABD and SIMA of BMD. **(D) Comparison of multi-modality**
617 **analysis (MUMA) and single-modality analysis (MUMA) in different BMD images,**
618 **including spine, left hip, and right hip.** MUMA of the three types of BMD images
619 outperforms SIMA of each of the three types of BMD images.

620

621 **Figure 4. Results of robustness analysis in ABD imaging. (A) The different parameter**
622 **settings of the DenseNet121 model for the ABD imaging analysis.** Testing AUC
623 results indicate that the configuration with an image size of 224 x 224, a batch size
624 64, and a pre-trained model set to False performed the best. **(B) The different**
625 **parameter settings of the GNN model for the ABD imaging analysis.** W_{eg} refers to
626 the weight of the edge used (Euclidean Distance - ED or Cosine Similarity - CS). W_n
627 represents the weight of nodes, either equal weight or weighted by the number of
628 edge nodes (n_{edge}). F indicates the number of features after graph convolution.
629 The testing AUC shows that the best performance is achieved with settings: $W_{eg} =$
630 ED, $F = 1024$, and $W_n = n_{edge}$. **(C) Comparative analysis between image-based**
631 **analysis, sample-based analysis with average weights, and sample-based analysis**
632 **with GNN weights in ABD imaging analysis.** The sample-based analysis with GNN
633 weights performs best. **(D) Comparison of two classifiers, Multilayer Perceptron**
634 **(MLP) and eXtreme Gradient Boosting (XGBoost).** MLP and XGBoost exhibited
635 similar performances in all measures: AUC, ACC, SEN, SPEC, and Youden index.

636

637 **Figure 5. Risk evaluation models for T2D. (A) Performance of various models**
638 **accounting for genetic PRS (G), demographic variables (D) – age, sex, and T2D**
639 **family history (D), and medical images of ABD (A-1), BMD (A-2), and ABD+BMD (A-**
640 **3).** The model comprising G, D, and medical images performs the best in T2D risk
641 evaluation. **(B) Positive correlation between MRS and T2D odds ratio for ABD (B-1),**
642 **BMD (B-2), and ABD+BMD (B-3).** In each decile of MRS based on the image, the odds

643 ratio of T2D risk and its 95% confidence interval were calculated based on an
644 unadjusted model (blue line) and model adjusted by age, sex, and T2D family history
645 (red line), where the MRS group in 40%–60% is set as the reference group. **(C)**
646 **Identification of high-risk subgroup based on MRS of ABD (C-1), BMD (C-2) and**
647 **ABD+BMD (C-3).** For ABD imaging and ABD+BMD imaging, the high-risk group was
648 females older than 62 with a T2D family history, and their MRS group was 90%–
649 100%. For BMD imaging, in addition to being identical to the group identified by
650 ABD-based MRS, another high-risk group was men older than 62 with a family history
651 of T2D.
652

653 **References**

- 654 1. Gøtzsche PC, Jørgensen KJ. Screening for breast cancer with mammography.
655 *Cochrane Database of Systematic Reviews*, (2013).
656
- 657 2. Ballestri S, *et al.* Nonalcoholic fatty liver disease is associated with an almost
658 twofold increased risk of incident type 2 diabetes and metabolic syndrome.
659 Evidence from a systematic review and meta-analysis. *Journal of*
660 *gastroenterology and hepatology* **31**, 936-944 (2016).
661
- 662 3. Al-Qazzaz NK, Ali SH, Ahmad SA, Islam S, Mohamad K. Cognitive impairment
663 and memory dysfunction after a stroke diagnosis: a post-stroke memory
664 assessment. *Neuropsychiatric Disease and Treatment* **10**, 1677-1691 (2014).
665
- 666 4. Takekawa H, Tsukui D, Kobayasi S, Suzuki K, Hamaguchi H. Ultrasound
667 diagnosis of carotid artery stenosis and occlusion. *J Med Ultrason (2001)* **49**,
668 675-687 (2022).
669
- 670 5. Alzubaidi L, *et al.* Review of deep learning: Concepts, CNN architectures,
671 challenges, applications, future directions. *Journal of big Data* **8**, 1-74 (2021).
672
- 673 6. Hatcher WG, Yu W. A survey of deep learning: Platforms, applications and
674 emerging research trends. *IEEE Access* **6**, 24411-24432 (2018).
675
- 676 7. LeCun Y, Bottou L, Bengio Y, Haffner P. Gradient-based learning applied to
677 document recognition. *Proceedings of the IEEE* **86**, 2278-2324 (1998).
678
- 679 8. Goodfellow I, *et al.* Generative adversarial nets. *Advances in neural*

- 680 *information processing systems* **27**, (2014).
681
- 682 9. Yang Y, *et al.* Nonalcoholic fatty liver disease (NAFLD) detection and deep
683 learning in a Chinese community-based population. *European Radiology* **33**,
684 5894-5906 (2023).
685
- 686 10. Zamanian H, Mostaar A, Azadeh P, Ahmadi M. Implementation of
687 Combinational Deep Learning Algorithm for Non-alcoholic Fatty Liver
688 Classification in Ultrasound Images. *J Biomed Phys Eng* **11**, 73-84 (2021).
689
- 690 11. Yen TJ, Yang CT, Lee YJ, Chen CH, Yang HC. Fatty liver classification via risk
691 controlled neural networks trained on grouped ultrasound image data.
692 *Scientific Reports* **14**, 13 (2024).
693
- 694 12. Ma J, Wu F, Zhu J, Xu D, Kong D. A pre-trained convolutional neural network
695 based method for thyroid nodule diagnosis. *Ultrasonics* **73**, 221-230 (2017).
696
- 697 13. Sun W, Tseng T-L, Zhang J, Qian W. Enhancing deep convolutional neural
698 network scheme for breast cancer diagnosis with unlabeled data.
699 *Computerized Medical Imaging and Graphics* **57**, 4-9 (2017).
700
- 701 14. Pratt H, Coenen F, Broadbent DM, Harding SP, Zheng Y. Convolutional Neural
702 Networks for Diabetic Retinopathy. *Procedia Computer Science* **90**, 200-205
703 (2016).
704
- 705 15. Gulshan V, *et al.* Development and validation of a deep learning algorithm for
706 detection of diabetic retinopathy in retinal fundus photographs. *Jama* **316**,
707 2402-2410 (2016).
708
- 709 16. Huang M-L, Wu Y-S. Classification of atrial fibrillation and normal sinus
710 rhythm based on convolutional neural network. *Biomedical Engineering*
711 *Letters* **10**, 183-193 (2020).
712
- 713 17. Jang R, Choi JH, Kim N, Chang JS, Yoon PW, Kim C-H. Prediction of
714 osteoporosis from simple hip radiography using deep learning algorithm.
715 *Scientific Reports* **11**, 19997 (2021).
716
- 717 18. Yamamoto N, *et al.* Deep Learning for Osteoporosis Classification Using Hip

- 718 Radiographs and Patient Clinical Covariates. *Biomolecules* **10**, 1534 (2020).
719
- 720 19. Razzak MI, Naz S, Zaib A. Deep learning for medical image processing:
721 Overview, challenges and the future. *Classification in BioApps: Automation of*
722 *Decision Making*, 323-350 (2018).
723
- 724 20. Liu M, Cheng D, Wang K, Wang Y, Initiative AsDN. Multi-modality cascaded
725 convolutional neural networks for Alzheimer's disease diagnosis.
726 *Neuroinformatics* **16**, 295-308 (2018).
727
- 728 21. Zhang D, Wang Y, Zhou L, Yuan H, Shen D. Multimodal classification of
729 Alzheimer's disease and mild cognitive impairment. *Neuroimage* **55**, 856-867
730 (2011).
731
- 732 22. Holste G, Partridge SC, Rahbar H, Biswas D, Lee CI, Alessio AM. End-to-end
733 learning of fused image and non-image features for improved breast cancer
734 classification from mri. In: *Proceedings of the IEEE/CVF International*
735 *Conference on Computer Vision* (2021).
736
- 737 23. Yan R, *et al.* Richer fusion network for breast cancer classification based on
738 multimodal data. *BMC Medical Informatics and Decision Making* **21**, 1-15
739 (2021).
740
- 741 24. Sousa JV, Matos P, Silva F, Freitas P, Oliveira HP, Pereira T. Single Modality vs.
742 Multimodality: What Works Best for Lung Cancer Screening? *Sensors (Basel)*
743 **23**, (2023).
744
- 745 25. Huang Y-J, Chen C-h, Yang H-C. AI-Enhanced Integration of Genetic and
746 Medical Imaging Data for Risk Assessment of Type 2 Diabetes. *medRxiv*,
747 2023.2008.2014.23294093 (2023).
748
- 749 26. Fan CT, Lin JC, Lee CH. Taiwan Biobank: a project aiming to aid Taiwan's
750 transition into a biomedical island. *Pharmacogenomics* **9**, 235-246 (2008).
751
- 752 27. Lin JC, Hsiao WW, Fan CT. Managing "incidental findings" in biobank research:
753 Recommendations of the Taiwan biobank. *Comput Struct Biotechnol J* **17**,
754 1135-1142 (2019).
755

- 756 28. Chen T, Guestrin C. Xgboost: A scalable tree boosting system. In: *Proceedings*
757 *of the 22nd acm sigkdd international conference on knowledge discovery and*
758 *data mining*) (2016).
759
- 760 29. Yang HC, *et al.* Genome-Wide Pharmacogenomic Study on Methadone
761 Maintenance Treatment Identifies SNP rs17180299 and Multiple Haplotypes
762 on CYP2B6, SPON1, and GSG1L Associated with Plasma Concentrations of
763 Methadone R- and S-enantiomers in Heroin-Dependent Patients. *PLoS Genet*
764 **12**, e1005910 (2016).
765
- 766 30. Huang G, Liu Z, Van Der Maaten L, Weinberger KQ. Densely connected
767 convolutional networks. In: *Proceedings of the IEEE conference on computer*
768 *vision and pattern recognition*) (2017).
769
- 770 31. Kipf TN, Welling M. Semi-supervised classification with graph convolutional
771 networks. *arXiv preprint arXiv:160902907*, (2016).
772
- 773 32. Fey M, Lenssen JE. Fast graph representation learning with PyTorch
774 Geometric. *arXiv preprint arXiv:190302428*, (2019).
775
- 776 33. Youden WJ. Index for rating diagnostic tests. *Cancer* **3**, 32-35 (1950).
777
- 778 34. Ruan Y, *et al.* Improving polygenic prediction in ancestrally diverse
779 populations. *Nature Genetics* **54**, 573-580 (2022).
780
- 781 35. Mahajan A, *et al.* Multi-ancestry genetic study of type 2 diabetes highlights
782 the power of diverse populations for discovery and translation. *Nat Genet* **54**,
783 560-572 (2022).
784
- 785 36. Auton A, *et al.* A global reference for human genetic variation. *Nature* **526**,
786 68-74 (2015).
787
- 788 37. Li G, *et al.* Evaluation of ADA HbA1c criteria in the diagnosis of pre-diabetes
789 and diabetes in a population of Chinese adolescents and young adults at high
790 risk for diabetes: a cross-sectional study. *BMJ open* **8**, e020665 (2018).
791
- 792 38. Greenhalgh T, *et al.* New models of self-management education for minority
793 ethnic groups: pilot randomized trial of a story-sharing intervention. *Journal*

- 794 *of Health Services Research & Policy* **16**, 28-36 (2011).
795
- 796 39. Nowicka P, *et al.* Utility of hemoglobin A1c for diagnosing prediabetes and
797 diabetes in obese children and adolescents. *Diabetes care* **34**, 1306-1311
798 (2011).
799
- 800 40. Eehalt S, *et al.* Diabetes screening in overweight and obese children and
801 adolescents: choosing the right test. *European journal of pediatrics* **176**, 89-
802 97 (2017).
803
- 804 41. Spiller S, Blüher M, Hoffmann R. Plasma levels of free fatty acids correlate
805 with type 2 diabetes mellitus. *Diabetes, Obesity and Metabolism* **20**, 2661-
806 2669 (2018).
807
- 808 42. Panwar H, Rashmi HM, Batish VK, Grover S. Probiotics as potential
809 biotherapeutics in the management of type 2 diabetes—prospects and
810 perspectives. *Diabetes/metabolism research and reviews* **29**, 103-112 (2013).
811
- 812 43. John RM. The Well Pediatric Primary Care Visit and Screening Laboratory
813 Tests. In: *Pediatric Diagnostic Labs for Primary Care: An Evidence-based*
814 *Approach*). Springer (2022).
815
- 816 44. Buchanan G, John J, Whiteside A, Moisey R, Malik M, Beer S. Admission
817 glucose is poor predictor of an abnormal glucose tolerance in acute coronary
818 syndrome but abnormal oral glucose tolerance test predicts mortality.). BMJ
819 Publishing Group Ltd and British Cardiovascular Society (2009).
820
- 821 45. Fang X, Liu Z, Xu M. Ensemble of deep convolutional neural networks based
822 multi-modality images for Alzheimer's disease diagnosis. *IET Image*
823 *Processing* **14**, 318-326 (2020).
824
- 825 46. Song J, Zheng J, Li P, Lu X, Zhu G, Shen P. An effective multimodal image fusion
826 method using MRI and PET for Alzheimer's disease diagnosis. *Frontiers in*
827 *digital health* **3**, 637386 (2021).
828
- 829 47. Wei L, Osman S, Hatt M, El Naqa I. Machine learning for radiomics-based
830 multimodality and multiparametric modeling. *Q J Nucl Med Mol Imaging* **63**,
831 323-338 (2019).

832

833 48. Guo Y, Wang Q, Guo Y, Zhang Y, Fu Y, Zhang H. Preoperative prediction of
834 perineural invasion with multi-modality radiomics in rectal cancer. *Scientific*
835 *Reports* **11**, 9429 (2021).

836

837 49. Almezghwi K, Serte S. Improved Classification of White Blood Cells with the
838 Generative Adversarial Network and Deep Convolutional Neural Network.
839 *Computational Intelligence and Neuroscience* **2020**, 6490479 (2020).

840

841 50. Larochelle H, Bengio Y, Louradour J, Lamblin P. Exploring strategies for training
842 deep neural networks. *Journal of machine learning research* **10**, (2009).

843

844 51. Hinton GE, Osindero S, Teh Y-W. A fast learning algorithm for deep belief nets.
845 *Neural computation* **18**, 1527-1554 (2006).

846

847 52. Kim J, Calhoun VD, Shim E, Lee J-H. Deep neural network with weight sparsity
848 control and pre-training extracts hierarchical features and enhances
849 classification performance: Evidence from whole-brain resting-state
850 functional connectivity patterns of schizophrenia. *Neuroimage* **124**, 127-146
851 (2016).

852

853 53. Russakovsky O, *et al.* Imagenet large scale visual recognition challenge.
854 *International journal of computer vision* **115**, 211-252 (2015).

855

856

Table 1. Models and performance

No.	Model description	Image type	Image analysis unit	Image data analysis type	Analysis modality	AUC	Accuracy	Sensitivity	Specificity	Youden
M1	ABD^{IDF}	ABD	sample-based	FECA	SIMA	0.719	0.703	0.643	0.714	0.357
M2	ABD_{AVE}^{PIX}	ABD	sample-based	PIXA	SIMA	0.881	0.847	0.746	0.864	0.610
M3	ABD_{GNN}^{PIX}	ABD	sample-based	PIXA	SIMA	0.887	0.842	0.765	0.855	0.620
M4	ABD_{IMAGE}^{PIX}	ABD	image-based	PIXA	SIMA	0.808	0.735	0.728	0.736	0.464
M5	BMD^{IDF}	BMD	sample-based	FECA	SIMA	0.847	0.763	0.783	0.759	0.542
M6	$BMD_{COMBINE}^{PIX}$	BMD	sample-based	PIXA	MUMA	0.824	0.774	0.697	0.788	0.485
M7	BMD_{SPINE}^{PIX}	BMD	sample-based	PIXA	SIMA	0.765	0.725	0.653	0.737	0.390
M8	$BMD_{L.HIP}^{PIX}$	BMD	sample-based	PIXA	SIMA	0.763	0.736	0.649	0.751	0.400
M9	$BMD_{R.HIP}^{PIX}$	BMD	sample-based	PIXA	SIMA	0.757	0.786	0.528	0.828	0.356
M10	$ABD^{IDF} + BMD^{IDF}$	ABD+BMD	sample-based	FECA	MUMA	0.871	0.830	0.745	0.845	0.590
M11	$ABD_{GNN}^{PIX} + BMD_{COMBINE}^{PIX}$	ABD+BMD	sample-based	PIXA	MUMA	0.904	0.880	0.706	0.910	0.616
M12	$ABD_{GNN}^{PIX} + BMD_{COMBINE}^{PIX} + G + D$	ABD+BMD	sample-based	PIXA	MUMA	0.944	0.868	0.889	0.865	0.754

Abbreviation list. ABD: abdominal ultrasonography; BMD: bone density scan; FECA: feature-centric analysis; PIXA: pixel-based analysis; SIMA: single-modality analysis;

MUMA: multi-modality analysis; G: genetic predictor (PRS); D: demographic variables (age, sex, and family history of T2D).

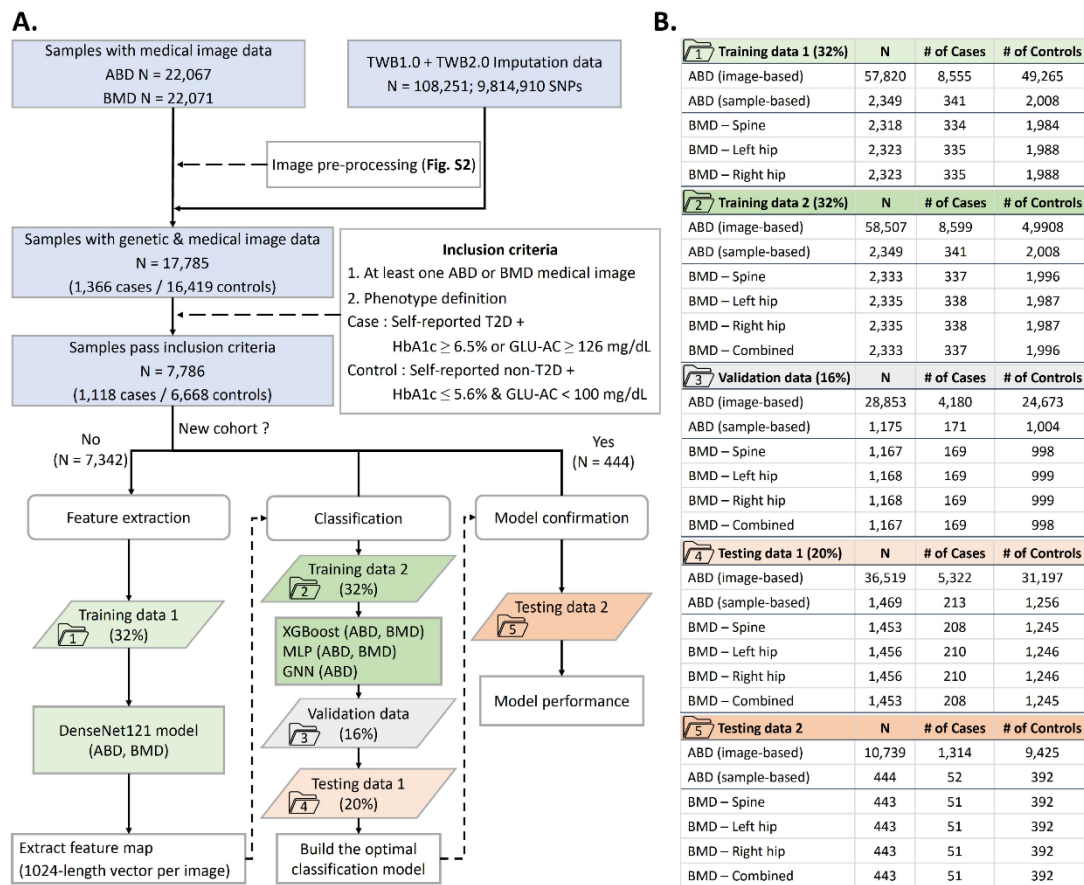


Figure 1. Flowchart of the study. (A) Data extraction and classification model building. In this dataset, 21,927 individuals possess ABD and BMD medical imaging data, while 108,251 individuals have whole-genome genotyping and imputation data. Moreover, 17,785 individuals possess both medical imaging and genetic profile data. Finally, the final dataset comprises 7,786 individuals who meet the inclusion criteria based on self-reported T2D status, HbA1C, and fasting glucose, consisting of 1,118 T2D cases and 6,668 normal controls. The complete dataset was initially divided into a training + validation set and a testing set at 8:2. Subsequently, the training + validation set was further separated into training and validation datasets with an 8:2 ratio. The training set was divided into two independent subsets at a 5:5 ratio to mitigate the winner’s curse problem; a two-stage procedure was employed for feature extraction and classification. In the first stage, which was dedicated to feature extraction, the first training data was used to establish a DenseNet121 model based on the initial training dataset (Training Dataset 1). Subsequently, a feature map vector was obtained. The second stage was focused on sample classification. Utilizing the data from the second training (Training Dataset 2) and validation datasets, a deep learning model for T2D classification was developed, and the results were further confirmed using the Testing Dataset 1. Ultimately, the best model’s validity was

further confirmed using the second independent testing dataset (Testing Dataset 2) (n = 444). Regarding the deep learning classification model, three methods—multilayer perceptron (MLP), graph neural network (GNN), and eXtreme Gradient Boosting (XGBoost)—were implemented. **(B) Sample Size.** Information on the total sample size, number of cases, and number of controls is provided.

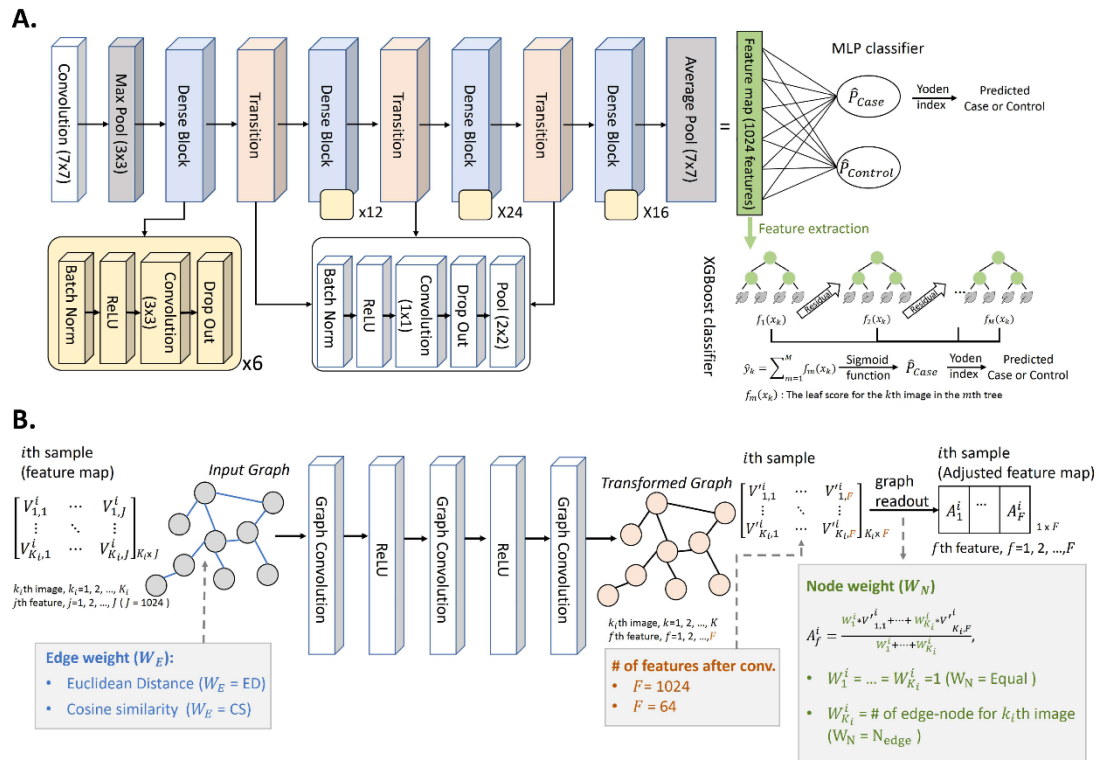


Figure 2. Architecture diagrams. (A) Densely Connected Convolutional Neural Networks with 121 layers (DenseNet-121). (B) Graph Neural Network (GNN).

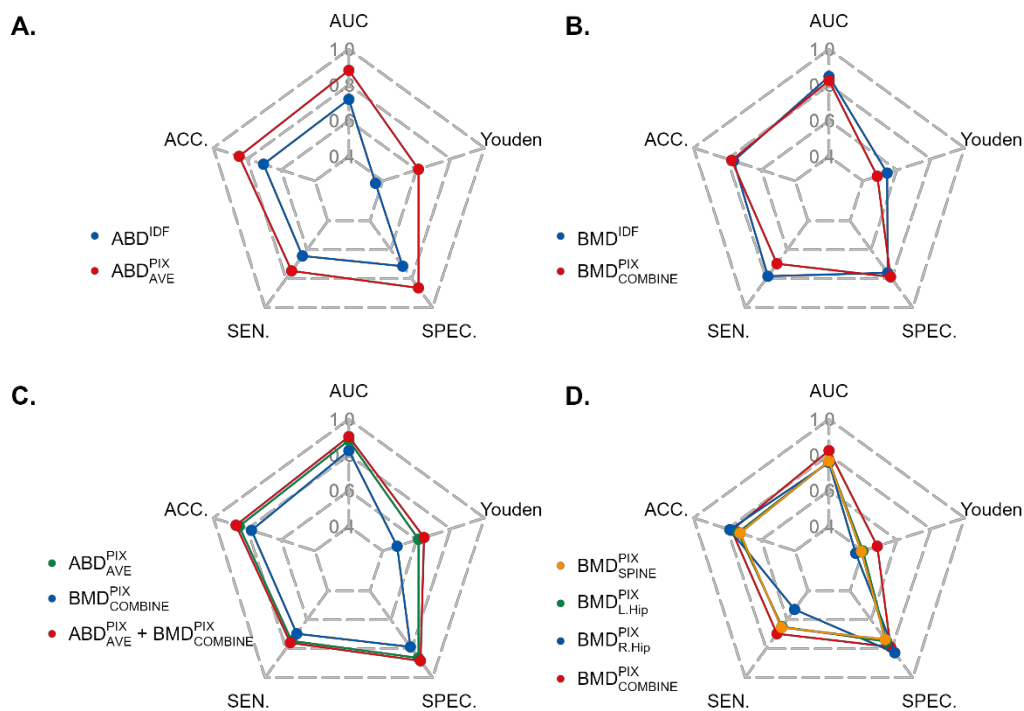


Figure 3. Model Comparison. (A) Comparison of PIXA ($ABD^{PIX_{AVE}}$) and FECA (ABD^{IDF}) in ABD imaging analysis. PIXA exhibited superior performance in T2D risk evaluation compared to FECA. (B) Comparison of PIXA ($BMD^{PIX_{COMBINE}}$) and FECA (BMD^{IDF}) in BMD imaging analysis. PIXA and FECA exhibited similar performance in T2D risk evaluation. (C) Comparison of multi-modality analysis (MUMA) and single-modality analysis (SIMA) in ABD and BMD imaging analysis. MUMA of ABD and BMD outperforms SIMA of ABD and SIMA of BMD. (D) Comparison of multi-modality analysis (MUMA) and single-modality analysis (SIMA) in different BMD images, including spine, left hip, and right hip. MUMA of the three types of BMD images outperforms SIMA of each of the three types of BMD images.

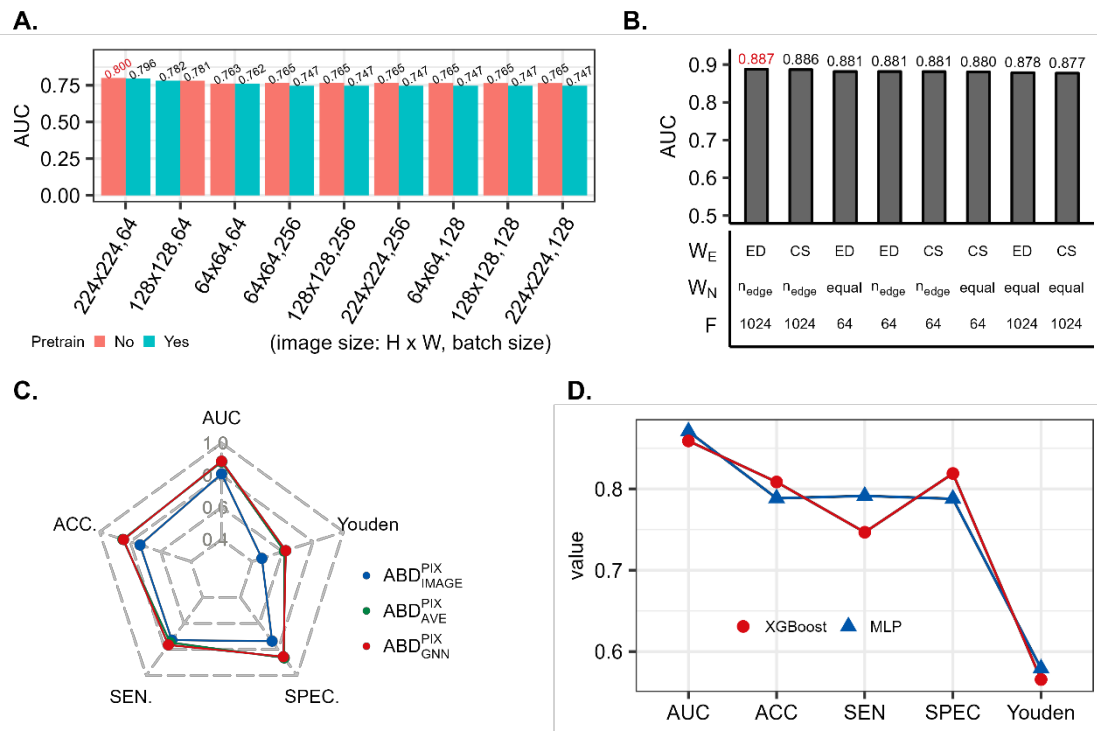


Figure 4. Results of robustness analysis in ABD imaging. (A) The different parameter settings of the DenseNet121 model for the ABD imaging analysis. Testing AUC results indicate that the configuration with an image size of 224 x 224, a batch size 64, and a pre-trained model set to False performed the best. **(B) The different parameter settings of the GNN model for the ABD imaging analysis.** W_{eg} refers to the weight of the edge used (Euclidean Distance - ED or Cosine Similarity - CS). W_n represents the weight of nodes, either equal weight or weighted by the number of edge nodes (n_{edge}). F indicates the number of features after graph convolution. The testing AUC shows that the best performance is achieved with settings: $W_{eg} = ED$, $F = 1024$, and $W_n = n_{edge}$. **(C) Comparative analysis between image-based analysis, sample-based analysis with average weights, and sample-based analysis with GNN weights in ABD imaging analysis.** The sample-based analysis with GNN weights performs best. **(D) Comparison of two classifiers, Multilayer Perceptron (MLP) and eXtreme Gradient Boosting (XGBoost).** MLP and XGBoost exhibited similar performances in all measures: AUC, ACC, SEN, SPEC, and Youden index.

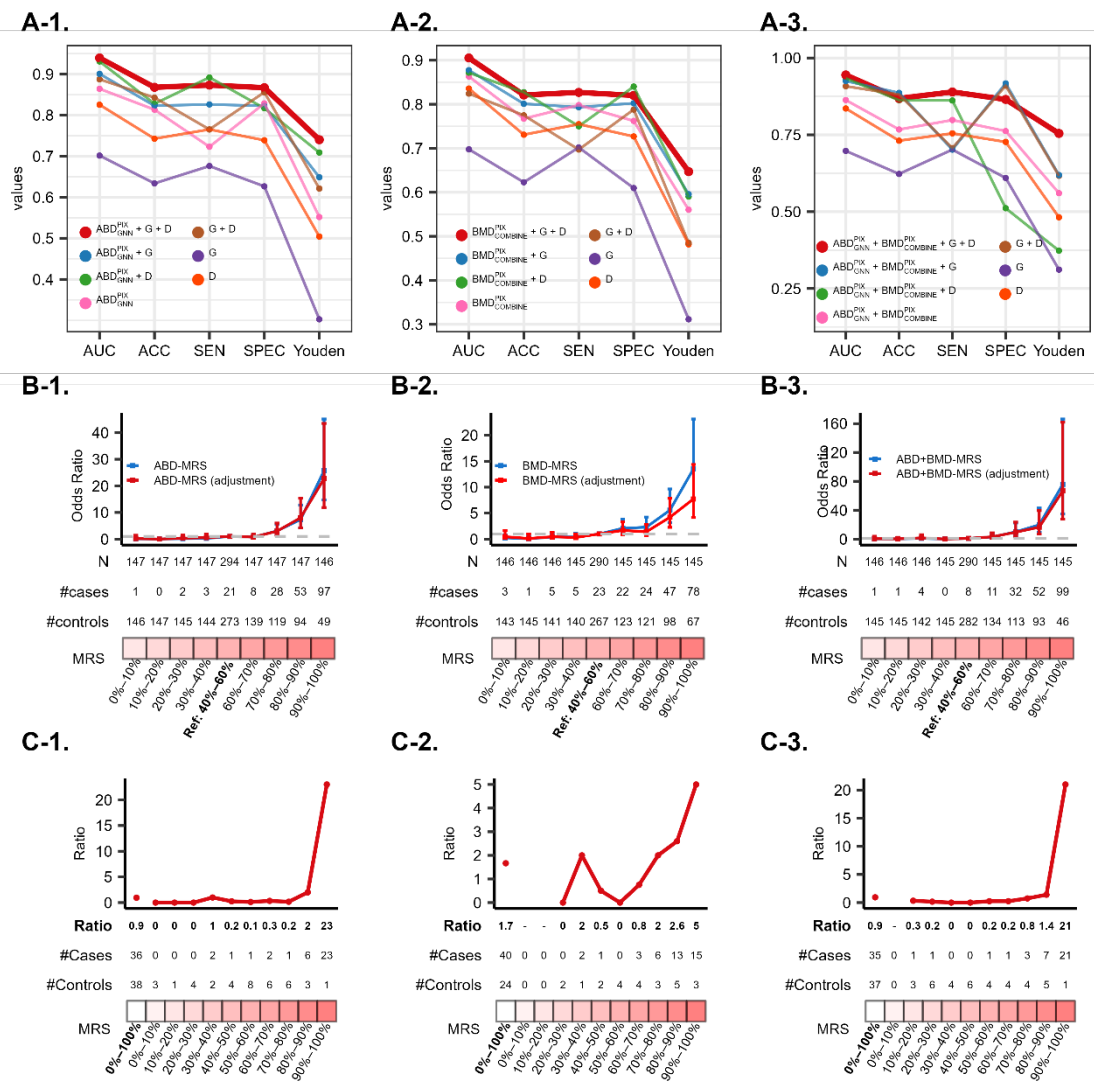


Figure 5. Risk evaluation models for T2D. (A) Performance of various models accounting for genetic PRS (G), demographic variables (D) – age, sex, and T2D family history (D), and medical images of ABD (A-1), BMD (A-2), and ABD+BMD (A-3). The model comprising G, D, and medical images performs the best in T2D risk evaluation. (B) Positive correlation between MRS and T2D odds ratio for ABD (B-1), BMD (B-2), and ABD+BMD (B-3). In each decile of MRS based on the image, the odds ratio of T2D risk and its 95% confidence interval were calculated based on an unadjusted model (blue line) and model adjusted by age, sex, and T2D family history (red line), where the MRS group in 40%–60% is set as the reference group. (C) Identification of high-risk subgroup based on MRS of ABD (C-1), BMD (C-2), and ABD+BMD (C-3). For ABD imaging and ABD+BMD imaging, the high-risk group was females older than 62 with a T2D family history, and their MRS group was 90%–100%. For BMD imaging, in addition to being identical to the group identified by ABD-based MRS, another high-risk group was men older than 62 with a family history of T2D.