

Disease progression and clinical outcomes in latent osteoarthritis phenotypes: Data from the Osteoarthritis Initiative

Weihua Guo (✉ wguo@coh.org)

Beckman Research Institute at City of Hope <https://orcid.org/0000-0002-9580-513X>

ZeYu Huang

Duke University Hospital <https://orcid.org/0000-0002-0456-8379>

Zhao Zhang

West China Hospital, West China Medical School, Sichuan University

Mary Bucklin

Rush University <https://orcid.org/0000-0003-3893-2818>

John Martin

Duke University

Article

Keywords:

Posted Date: January 26th, 2024

DOI: <https://doi.org/10.21203/rs.3.rs-3855831/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: There is **NO** Competing Interest.

1 **Disease progression and clinical outcomes in latent osteoarthritis phenotypes: Data**
2 **from the Osteoarthritis Initiative**

3
4 **Authors:**

5 Zeyu Huang, M.D., PhD.¹ zey.huang@gmail.com

6 Mary A. Bucklin, PhD.² mary_a_bucklin@rush.edu

7 Weihua Guo, PhD.^{3**} wguo@coh.org

8 John T. Martin, PhD.^{2*} john_martin@rush.edu

9
10 ¹ Department of Orthopaedic Surgery, Orthopaedic Research Institute, West China Hospital,
11 West China Medical School, Sichuan University, Chengdu, Sichuan Province, People's
12 Republic of China

13 ² Department of Orthopedic Surgery, Rush University, Chicago, Illinois, USA

14 ³ Department of Immuno-oncology, City of Hope, National Medical Center, Duarte, California,
15 USA

16
17 **Correspondence:**

18 *John T. Martin, PhD.

19 Assistant Professor

20 Department of Orthopaedic Surgery

21 Rush University

22 Email: John_Martin@rush.edu

23
24 **Weihua Guo PhD.

25 Staff Scientist I

26 Department of Immuno-Oncology

27 Beckman Research Institute at City of Hope

28 Email: wguo@coh.org

29

30

31

32

33

34

35

36

37 **Abstract**

38 The prevalence of knee osteoarthritis (OA) is widespread and the heterogeneous patient
39 factors and clinical symptoms in OA patients impede developing personalized treatments for
40 OA patients. In this study, we used unsupervised and supervised machine learning to organize
41 the heterogeneity in knee OA patients and predict disease progression in individuals from the
42 Osteoarthritis Initiative (OAI) dataset. We identified four distinct knee OA phenotypes using
43 unsupervised learning that were defined by nutrition, disability, stiffness, and pain (knee and
44 back) and were strongly related to disease fate. Interestingly, the absence of supplemental
45 vitamins from an individual's diet was protective from disease progression. Moreover, we
46 established a phenotyping tool and prognostic model from 5 variables (WOMAC disability score
47 of the right knee, WOMAC total score of the right knee, WOMAC total score of the left knee,
48 supplemental vitamins and minerals frequency, and antioxidant combination multivitamins
49 frequency) that can be utilized in clinical practice to determine the risk of knee OA progression
50 in individual patients. We also developed a prognostic model to estimate the risk for total knee
51 replacement and provide suggestions for modifiable variables to improve long-term knee health.
52 This combination of unsupervised and supervised data-driven tools provides a framework to
53 identify knee OA phenotype in a clinical scenario and personalize treatment strategies.

54

55

56

57

58

59

60 **Introduction**

61 Osteoarthritis (OA) is the most common form of joint disease and a major cause of pain
62 and disability and is a heterogenous disease in which aging, obesity, trauma, and genetic
63 factors are implicated as drivers of pathogenesis¹. OA affects 9.6% of men and 18% of women
64 over 60 years of age² and 250 million people worldwide³. The United States Food and Drug
65 Administration (FDA), Centers for Disease Control (CDC), and National Institutes of Health (NIH)
66 all recognize the impact of OA and have guidelines and research agendas to reduce the
67 prevalence and burden. This public health issue is projected to worsen as life expectancy
68 increases and the US population skews towards older individuals⁴. Still, there are no disease-
69 modifying OA drugs (DMOADs) approved by the FDA or European Medicines Agency⁵ and as
70 a result, managing OA remains largely palliative.

71 One complicating factor is that OA phenotypes vary from patient to patient and there is
72 likely no “one size fits all” treatment⁶. It may be that the failure of numerous phase II/III OA
73 clinical trials, such as iNOS⁷, bisphosphonates⁸, and calcitonin^{9, 10}, has been due to the inability
74 to decipher the specific underlying drivers of OA at the individual patient level and therefore
75 DMOADs are not delivered to the most suitable subgroups. Thus, identifying OA phenotypes is
76 a critical task for the community. Machine learning (ML) is a computational tool that learns
77 complex non-linear patterns between many variables without precise instructions^{11, 12, 13}.
78 Classification ML models can identify novel, clinically significant features in patients^{14, 15}. These
79 methods have been used to determine disease phenotypes in many clinical populations¹⁶.
80 Furthermore, predictive ML models have been used to determine disease risk factors,
81 complications, and survival outcomes in clinical practice¹⁷. Our global hypothesis is that the

82 heterogeneity in knee OA phenotypes can be organized with unsupervised learning and that
83 supervised learning models can predict disease progression.

84 In the current study, we used unsupervised and supervised ML methods to identify knee
85 OA phenotypes and predict disease progression in the open access Osteoarthritis Initiative
86 (OAI) dataset (**Figure 1**). The OAI is a longitudinal, observational study of knee OA with 4,796
87 enrollees. It includes greater than 1,000 descriptive variables, including demographics, pain,
88 exercise habits, diet and nutrition, socioeconomic status, medical history, radiographic
89 evaluation, and psychological evaluation. We determined OA phenotypes by performing
90 unsupervised learning on enrollment data (k-means clustering) and visualized relationships
91 between phenotypes via dimensionality reduction. Then, we utilized data from multiple follow-
92 up time points over 8 years to develop supervised learning models that predicted long-term
93 disease progression, including the likelihood of total knee replacement (TKR).

94

95 **Materials and Methods**

96 *Data extraction and cleaning*

97 We included all 4,796 participants who enrolled in the OAI study with 1032 variables that
98 were measured at enrollment (variables: **Data S1**). We performed a data cleaning procedure
99 to remove individuals with incomplete data, remove variables that had missing values or low
100 variance, and remove variables that were highly correlated. All data was processed in either
101 Python or R as noted below.

102 First, we excluded 127 subjects with more than 595 variables (50% of total variables)
103 whose value were missing (**Figure S1A**). Next, we generated a correlation matrix for each

104 combination of numerical and categorical variables with the following calculations: 1) numerical
105 vs. numerical: Pearson's coefficient (pearsonr function from scipy.stats Python library, V1.10.1);
106 2) categorical vs. categorical: Cramers' V (customized function based on Python); 3) numerical
107 vs. categorical: R value from ordinary least squares linear regression (ols function from
108 statsmodels.formula.api Python library, V0.13.5). We performed hierarchical clustering
109 (Heatmap function from ComplexHeatmap R package, V2.14.0) to group variables in the
110 correlation matrix and found that variables with missing values were grouped together. We
111 screened different cutoffs (i.e., 25%, 50%, 80%) for the relative subject number of missing
112 values (number of subjects with missing value relative to total subject number) and found that
113 a 25% cutoff removed clustered variables with majority missing values (**Figure S1B, C, and D**)
114 (**Data S1**). Therefore, we removed 295 variables among which more than 25% data points were
115 missing (**Figure S1A**).

116

117 *Clustering and dimensionality reduction for identifying knee OA phenotypes.*

118 After data extraction and cleaning, we identified groups of similar individuals via
119 unsupervised learning and performed dimensionality reduction for data visualization. First, we
120 used the one-hot encoding method (get_dummies function from pandas Python library, V1.5.3)
121 to convert the categorical variables to numerical variables. We replaced missing values using
122 a k-Nearest Neighbors imputation (KNNImputer function from sklearn.impute Python library,
123 V1.2.2) with 2 neighboring samples and uniform weights. Imputed data was scaled and
124 normalized (StandardScaler function from sklearn.preprocessing Python library, V1.2.2) and
125 principal component reduction was performed (PCA function from sklearn.decomposition

126 Python library, V1.2.2). Based on the elbow method for variance thresholding (**Figure S2A**),
127 the top 16 principal components were selected for dimensionality reduction (Uniform Manifold
128 Approximation and Projection, UMAP, umap function from umap-learn Python library, V0.5.3)
129 and K-Means clustering (KMeans function from sklearn.cluster Python library, V1.2.2). We
130 calculated Silhouette scores (**Figure S2B**, silhouette_score function from sklearn.metrics
131 Python library, V1.2.2) for 2 to 21 clusters and identified that 4 was the optimal cluster number.

132 We performed statistical comparisons to identify variables that differentiated each cluster.
133 We used the Kruskal Wallis test for numerical variables (kruskal.test function from stats R
134 package, V4.2.3) and Fisher's exact test for categorical variables (fisher.test function from stats
135 R package, V4.2.3). P-values for both numerical and categorical variables were adjusted by
136 Benjamini & Hochberg method (**Data S2**, adjust_pvalue function from rstatix R package,
137 V0.7.2). We identified the top 10 variables that differentiated each cluster based on the following
138 criteria: numerical variables: maximum fold difference between means, categorical variables:
139 Chi-square statistic (chisq.test function from stats R package, V4.2.3). Cluster annotations were
140 determined by authors based on these cluster markers.

141

142 *Long-term outcomes across clusters and cohorts*

143 For the four clusters identified in our study and for the three cohorts defined at OAI data
144 collection, we performed Kaplan-Meier (KM) survival analysis using data from enrollment and
145 each follow-up visit on the following 6 outcome variables: Kellgren-Lawrence (KL) grade, joint
146 space width (minimum joint space width in the medial compartment), Western Ontario and
147 McMaster Universities Osteoarthritis Index (WOMAC) disability score, WOMAC stiffness score,

148 WOMAC pain score, WOMAC total score (WOMTS). We defined a survival event as the change
149 of each outcome variable from the first visit to any follow-up visit above a defined threshold (KL
150 grade $\Delta \geq 1$; joint space width $\Delta \leq -25\%$; all WOMAC scores $\Delta \geq 25\%$). We used exact enrollment
151 and visit dates to account for variability in time between visits (**Data S3**). Once a progression
152 event was identified, all following visits were discarded. We also extracted whether an individual
153 received a total knee replacement (TKR) in either knee, where TKR was considered as the
154 survival event. With such converted survival information, we built KM curves for all outcome
155 variables for both knees (surv and survfit function from survival R package, V3.5.5). To quantify
156 the hazard ratios for each cluster, we built Cox regression models (coxph function from survival
157 R package, V3.5.5).

158 To further examine the prognostic values of our clusters, we implemented the same KM
159 survival analysis on all four clusters within the progression cohort and incidence cohort
160 separately. We built KM curves for all outcomes variables for both knees (surv and survfit
161 function from survival R package, V3.5.5). To quantify the hazard ratios for each cluster within
162 these two cohorts, we built Cox regression models (coxph function from survival R package,
163 V3.5.5).

164

165 *Development of a clinical tool to predict cluster assignment via supervised learning.*

166 With well-defined clusters and survival outcomes by cluster, we developed a clinical tool
167 that assigns individual patients to the appropriate cluster to determine their long-term knee
168 health. To do so, we benchmarked common supervised learning models to predict cluster
169 assignment. We evaluated logistic regression (LogisticRegression function from

170 sklearn.linear_model Python library, V1.2.2; solver: newton-cg solver, maximum iterations:
171 1000), random forest (RandomForestClassifier function from sklearn.ensemble Python library,
172 V1.2.2; trees: 100, entropy criterion), and support vector machine (SVC function from
173 sklearn.svm Python library, V1.2.2; kernel: sigmoid, probability estimation enabled).

174 As above, we utilized numerical data and one-hot encoded categorical data as input data.
175 We scaled each variable to its corresponding minimum and maximum range (MinMaxScaler
176 function from sklearn.preprocessing Python library, V1.2.2). To determine the optimal number
177 of input variables, we first ranked input variables based on the importance metrics calculated
178 by fitting a random forest classifier (RandomForestClassifier function from sklearn.ensemble
179 Python library, V1.2.2) to all the input variables with the cluster labels (**Data S4**) and then
180 screened the input variable number from 2 to 50 for all ML models. To obtain robust accuracies
181 of each input variable number, we utilized a random permutation cross-validator with 20 splits,
182 and within each split, 90% samples were considered as training data while the left 10% were
183 validation data (ShuffleSplit function from sklearn.model_selection Python library, V1.2.2). As a
184 multi-classification problem, we computed the accuracy classification score (accuracy_score
185 function from sklearn.metrics Python library, V1.2.2) and area under the receiver operating
186 characteristic curve (ROC AUC) using both one-vs-rest and one-vs-one approaches
187 (roc_auc_score function from sklearn.metrics Python library, V1.2.2) (**Data S5**). We averaged
188 the above metrics across all 20 test splits for each input variable number.

189 The most accurate model was exported and built on a web-based interface
190 (www.predictoaphenotpe.org). With free registration, users will be able to fill in required
191 information of the patient and the website will provide a prediction of the cluster (phenotype)

192 this patient could belong to.

193

194 *Supervised learning for predicting WOMTS and identifying key predictor variables.*

195 We benchmarked common supervising learning models to predict WOMTS at 4 and 8
196 years from enrollment data. To identify effective predictor variables, we computed the
197 correlations between input variables and WOMTS across all yearly visits for both knees. We
198 computed Spearman's correlation coefficients (spearmanr function from scipy.stats.stats
199 Python library, V1.10.1) or R from ordinary least squares regression (ols function
200 statsmodels.formula.api Python library, V0.13.5) to quantify the correlation between WOMTS
201 and numerical or categorical variables (**Data S6**). We visualized the top 10 highly correlated
202 variables based on their average correlation coefficients.

203 We directly predicted the WOMTS for both knees at 4th and 8th year visit. We evaluated
204 linear regression (LinearRegression function from sklearn.linear_model Python library, V1.2.2),
205 random forest (RandomForestRegressor function from sklearn.ensemble Python library, V1.2.2;
206 trees: 10, 20, 40, 60, 80, 100), support vector machine (SVR function from sklearn.svm Python
207 library, V1.2.2; kernels: linear, polynomial, rbf, sigmoid; regularization: 100, kernel coefficient:
208 reciprocal of variable number), and an artificial neural network (ANN). We followed the same
209 scaling, input variable selection, and cross-validation procedures used in predicting clusters.
210 As WOMTS is a continuous variable, all ML models were regression models and used to
211 compare the measured and predicted WOMTS we calculated root mean square error (RMSE,
212 mean_squared_error function from sklearn.metrics Python library, V1.2.2) and Pearson's
213 correlation coefficient (PCC, pearsonr function from scipy.stats Python library, V1.10.1) as

214 accuracy metrics (**Data S7**). Average accuracy metrics across all cross-validation tests were
215 calculated to select the optimal input variable number.

216 For the ANN, we built a sequential model (Sequential function from keras.models Python
217 library, V2.11.0) with one input layer, adaptive hidden layers, and one output layer (Dense
218 function from keras.layers Python library, V2.11.0). The node number of the input layer was
219 dependent on the number of input variables during the screening, and the output layer had one
220 node to represent the WOMTS. The hidden layers were adaptively designed based on the
221 number of input variables, where each hidden layer was 75% of its previous layer (including
222 input layer). All activation functions were linear functions, Adam optimization with 0.001 as the
223 learning rate (optimizer.Adam function from tensorflow.keras Python library, V2.11.0) was used
224 to train the model, and mean squared error was taken as the loss function. We trained the
225 model with 100 epochs and 10 as the batch size.

226 As the ANN achieved the most accurate and robust predictions, we utilized the ANN model
227 to identify the most effective predictor variables using a customized random search algorithm.
228 We firstly built the same sequential model (Sequential function from keras.models Python
229 library, V2.11.0) with one input layer, adaptive hidden layers, and one output layer, adaptive
230 hidden layers, and one output layer (Dense function from keras.layers Python library, V2.11.0)
231 as the above ANN model. The node number of the input layer was 25 based on the screening
232 results, and the output layer had one node to represent the WOMTS. Similarly, the hidden
233 layers were adaptively designed based on the number of input variables, where hidden layer
234 was 75% of its previous layer (including input layer). All activation functions were linear function,
235 Adam optimization with 0.001 was the learning rate (optimizer. Adam function from

236 tensorflow.keras Python library, V2.11.0) was used to train the model, and mean squared error
237 was taken as the loss function. We trained the model with 100 epochs and 10 as the batch size.
238 Here, we randomly selected 25 variables to train an ANN model based on the above design
239 principles. To reduce the number of potential combinations, we only selected variables from the
240 cluster markers identified from the unsupervised clustering (adjust p-value <0.05, **Data S2**).
241 Within each test, we also used the same random perturbation cross-validator with the same
242 parameters to obtain the accuracies. After 10,000 random selection tests, we ordered the test
243 based on their average prediction accuracy and selected the top 10 to 1,000 most accurate
244 tests to investigate the composition of their input variables. We quantified the popularity of each
245 variable by computing the relative occurrences of each variable within the most accurate tests
246 to the total 10,000 tests.

247

248 *Additional Statistics*

249 Graphs and statistics were performed using R (v4.2.3), and Python (v3.9.16) as described.
250 The Kruskal Wallis test, Fisher's exact test, and log-rank test were implemented to compare
251 numeric, categorical, and survival data across different phenotypes or cohorts. Pearson's
252 correlation coefficient, Spearman's correlation coefficients, and Cramer's V were calculated to
253 quantify the associations. Accuracy, one-vs-one, and one-vs-rest AUC were calculated from
254 multi-class prediction. Root-mean-square-error and correlation between prediction and
255 measurements were calculated for regression. Experiment specific detailed statistical methods
256 are described in corresponding figure legends and Methods sections. Calculated p values are
257 displayed as *, p<0.05; **, p<0.01; ***p<0.001; ****, p<0.0001.

258

259 **Data and code availability**

260 All scripts used in this publication are available in
261 https://github.com/weihuaguo/cluster_oai. All other data are available in the main text or the
262 supplementary materials.

263

264 **Results**

265 *Unsupervised learning identified four knee OA phenotypes in OAI.*

266 We identified four knee OA phenotypes by unsupervised learning: a group with low
267 supplemental and dietary vitamin intake ('Low Vitamin'), a group with poor knee health ('Poor
268 Knee'), a group with intermediate knee health ('Intermediate Knee'), and a group with good
269 knee health ('Good Knee') (**Figure 2A, Table 1**). These names are based on the most
270 significant and abundant variables between the groups (**Figure 2B&C, Figure S3A&B**).
271 Specifically, the Low Vitamin group was characterized by low frequency of vitamin
272 supplementation and low percentage of vitamins obtained from daily food intake, despite
273 demonstrating good knee health and daily function. The Poor Knee group was characterized
274 by poor knee health, in addition to low quality of life, poor general health, and poor daily function.
275 The Intermediate Knee group exhibited relatively poor knee health, intermediate quality of life,
276 and intermediate daily function. Lastly, the Good Knee group demonstrated good knee health,
277 along with good quality of life, good general health, good mental health, and good daily function,
278 (**Figure 2D-I**).

279

280 *Knee OA phenotypes are associated with disease progression.*

281 Survival analysis using KL grade (**Figure 3A**), joint space width (**Figure 3B**), WOMAC total
282 score (**Figure 3C**), WOMAC pain score (**Figure 3D**), WOMAC stiffness score (**Figure 3E**), and
283 WOMAC function score (**Figure 3F**) showed that the Good Knee group in both right and left
284 knees. KL grade matched these trends for both knees; however, joint space width matched this
285 trend in the right knee ($p=0.00027$) but not the left knee ($p=0.44$). By directly comparing these
286 outcomes through all visits, we found that Good Knee group always had the lowest KL grade
287 (**Figure S4A**), highest joint space width (**Figure S4B**), lowest WOMAC total score (**Figure S4C**),
288 lowest WOMAC pain score (**Figure S4D**), lowest WOMAC stiffness score (**Figure S4E**), and
289 lowest WOMAC function score (**Figure S4F**) on average. More importantly, survival analysis
290 with total knee replacement outcome showed that the Good Knee group had the highest
291 survival probability (**Figure 3G**).

292 Since the OAI has defined sub-cohorts (progression, incidence, and non-exposed control
293 group)¹⁸, we first examined the composition of these sub-cohorts within our knee OA
294 phenotypes (**Figure 4A**). We found that more than 85% of Good Knee subjects were from the
295 incidence cohort, more than 60% of Poor Knee subjects were from the progression cohort, and
296 more than 70% of Low Vitamin subjects were from incidence cohort. Since disease progression
297 in these sub-cohorts were clinically well-defined, we tested our definitions of disease
298 progression and survival by examining whether our cluster-based survival analysis results
299 using patient WOMAC total scores (**Figure S5A**), KL grade (**Figure S5B**) joint space width
300 (**Figure S5C**), and TKR (**Figure S5D**). As expected, our definition of disease progression and
301 survival analysis comprehensively captured the disease progression based on the pre-defined

302 sub-cohorts (i.e., non-exposed control group was the least progressed and progression cohort
303 was the most progressed). Additionally, we also investigated the prognostic values of our knee
304 OA phenotypes within incidence and progression sub-cohorts. The results showed that our
305 knee OA phenotypes remained partially significant in patient WOMAC total scores (**Figure 4B**),
306 KL grade (**Figure 4C**) and joint space width (**Figure 4D**), and TKR (**Figure 4E**) within the
307 incidence and progression cohort. Except for joint space width, our four OA phenotypes tended
308 to be associated with all the other clinical outcomes ($p < 0.10$).

309

310 *Supervised learning accurately predicts cluster assignment.*

311 To accurately predict knee OA phenotypes, we benchmarked commonly used supervised
312 learning models (four major types, i.e., logistic regression, random forest with six different tree
313 numbers, supporting vector classifier with four different kernels). Generally, all models reached
314 accuracy around 90%, above 0.975 AUC for ROC in both one-versus-one and one-versus-rest
315 analyses (**Figure 5A**). Furthermore, we found that a minimum of five variables were necessary
316 to achieve optimal predictive accuracy, namely: WOMAC disability score of the right knee,
317 WOMAC total score of the right knee, WOMAC total score of the left knee, supplemental
318 vitamins and minerals frequency, and antioxidant combination multivitamins frequency (**Figure**
319 **5B**).

320

321 *WOMAC total score predictive modeling*

322 Since WOMAC total score for both right and left knee is among the top variables for
323 constructing accurate group prediction model, we first used univariate analysis to identify

324 predictors at the screen phase (baseline) for the WOMAC total score of each visit (**Figure 6A**).
325 The results showed that variables that were positively correlated were baseline right knee
326 functional scores (difficulty in bathtub, standing, bending, car, shopping), baseline right knee
327 WOMAC pain and disability scores, and baseline right and left knee WOMAC total scores.
328 Variables that were negatively correlated include comorbidities and Knee Injury and
329 Osteoarthritis Outcome Score (KOOS) scores (left and right knee KOOS pain, right knee KOOS
330 quality of life, left and right knee KOOS symptom score).

331 We then developed ML prediction-based multivariate analysis to identify a set of key
332 variables related to WOMAC total scores (details in Methods). The principle of this analysis is
333 that the input variables, which are necessary to accurately predict WOMAC total scores through
334 ML models, are key variables. Based on this principle, we first benchmarked 12 multivariate
335 supervised ML models on their accuracies in predicting WOMAC total scores for either knee.
336 We found that the ML model built by ANN, linear regression, and an 80-tree random forest
337 showed the best predictive accuracy reflected by lower RMSE (**Figure 6B**) and higher PCC
338 between measurements and predictions (**Figure 6C**). Because ANN has the best robustness¹⁹,
339 we utilize ANN as the ML model and randomly selected 25 input variables to train the ANN
340 model and evaluate the corresponding prediction accuracy. With 10,000 random selections, we
341 analyzed the relative occurrences of the input variables within most accurate predictions. Based
342 on this analysis, we identified the top 5 variables with highest average occurrences from top 10
343 to 1,000 most accurate predictions of WOMAC total score at both 4-year and 8-year follow-ups.
344 The results showed that variables had the greatest relative occurrence were age, iron
345 supplement, knee difficulty – kneeling, difficulty with knees, B12 supplement, left knee WOMAC

346 disability score, left knee WOMAC pain score, and left knee WOMAC total score. Among them,
347 the baseline right knee WOMAC disability score had the greatest relative occurrence (**Figure**
348 **6D**).

349 **Discussion**

350 OA is a heterogeneous disease and modern multivariate solution is likely necessary to
351 identify disease phenotypes and progression patterns. In this study we identified four distinct
352 knee OA phenotypes using unsupervised learning in the 4,796 participants of the Osteoarthritis
353 Initiative. Phenotypes were primarily determined by nutrition and disability, stiffness, and pain
354 (knee and back) scores and were strongly related to disease fate. In addition, we established
355 a phenotyping tool from 5 variables that can be utilized in clinical practice to determine the risk
356 of knee OA progression in individual patients. We also developed a prognostic model that can
357 predict the risk of total knee replacement and provide suggestions for modifiable variables to
358 improve long-term knee health.

359 We utilized all available subjects and variables from 10 years of follow-up data in the OAI.
360 Our results show four distinct phenotypes that can be determined by simple questionnaires
361 related to general health, knee health, nutrition, and psychological evaluation. The groups
362 included a group with a hallmark of low supplemental and dietary vitamin intake ('Low vitamin'),
363 a group with hallmarks of poor knee health ('Poor Knee'), a group with hallmarks of intermediate
364 knee health ('Intermediate Knee'), and a group with hallmarks of good knee health ('Good Knee')
365 (**Figure 2A, Table 1**). The names of these groups are based on the most statistically significant
366 and prevalent variables between the groups (**Figure 2B&C**). Among them, the top variables
367 were related to the frequency of vitamins/minerals intake, the amount of the supplemental

368 Calcium, Beta-Carotene, Zinc, vitamin B6, B12, and D, WOMAC sub-scores, and WOMAC total
369 score. Previously, other studies have tried to identify knee OA phenotypes. For example, by
370 using biochemical markers data from IMI-APPROACH cohort, Angelini et al.²⁰ found that OA
371 patients could be divided into three phenotypes: low tissue turnover, structural damage, and
372 systemic inflammation. In addition, by using RNA sequencing data from knee OA patients tissue
373 (cartilage, subchondral bone, and synovium) Yuan et al.²¹ showed that OA patients could be
374 divided into four subtypes: metabolic disorder subtype, collagen metabolic disorder subtype,
375 activated sensory neuron subtypes, and inflammation subtype. In this work, we present a
376 concise and clinically applicable OA phenotyping method that does not require intra-articular
377 procedures, bloodwork, or sequencing that may be susceptible to error from environmental
378 factors²².

379 Survival analysis revealed that the phenotypes defined by unsupervised learning were
380 associated with long-term knee symptom, structure, and clinical outcomes (WOMAC total score,
381 KL grade, TKR). More importantly, we developed phenotype prediction models and narrowed
382 the necessary parameters down to 5 variables (WOMAC disability score, right knee; WOMAC
383 total score, right knee; WOMAC total score, left knee; multivitamin frequency; antioxidant
384 multivitamin frequency) which can be conveniently deployed in daily clinical scenarios. In the
385 past, several studies tried to use ML methods to establish predictive models for TKR and
386 achieved good accuracy²³⁻²⁵. However, there are limitations prohibiting these models from wide
387 clinical use. Firstly, willingness to receive TKR is determined not only by medical related factors
388 but also by others such as socioeconomic status and culture. Secondly, not all models served
389 as a prognostic purpose. As OA is a chronic condition that is widespread, often ongoing, and

390 frequently marked by episodes of exacerbation, long-term management of the disease is crucial
391 for individualized treatment. Thus, the most important scientific question in this field is how to
392 identify the appropriate patient for the correct treatment. Previously, Driban et al.²⁶ utilized OAI
393 data and found that 80% of people with end-stage knee OA did not have progressive
394 radiographic severity, suggesting radiographic results alone are not an optimal variable for
395 disease stage definition. In addition, Pierson et al.²⁷ used an algorithmic approach and found
396 out that radiologist-based X-ray interpretation could only explain 9% of unexplained racial
397 disparities in pain, which makes determining the risk for TKR more difficult. We surmise that
398 our approach which incorporates a holistic view of knee health is well-suited to a clinical setting.

399 Interestingly, we identified a phenotype of Low Vitamin group with similar survival
400 probability to the Good Knee group. The signature variables associated with the population
401 from this group were the frequency of vitamin A and C intake. Antioxidant supplements such as
402 vitamin A and C have long been advocated for the treatment of OA²⁸. Although various
403 approaches have been employed to tackle this issue, there is still a dearth of substantial
404 evidence to support these treatments. In a systematic review, Canter et al.²⁹ summarized 9
405 RCTs results and found that no convincing evidence to support vitamin A and C in OA treatment.
406 Kraus et al. identified that Vitamin C can actually exacerbate OA in a guinea pig model³⁰. In
407 recent study, Qu et al.³¹ applied mendelian randomization to the data from UK Biobank and
408 failed to find the causal association between vitamin A and OA. Our findings were supported by
409 these data as a low supplemental vitamin intake did not worsen OA prognosis.

410 In the current study, other relevant factors like BMI, comorbidities, and depression
411 statistically differentiated phenotypes in addition to signature variables mentioned above. BMI

412 was one factor that contributed to the phenotypes and fate of the disease. This is consistent
413 with the literature that suggests that BMI has long been considered as a risk factor for OA³².
414 We also found that higher comorbidities were associated with worse knee OA phenotype and
415 disease progression. Gustafsson et al.³³ have shown that compared to matched references
416 from the general population, knee OA patients were more commonly associated with one or
417 more comorbidities, which was independent of socioeconomic status. We found the highest
418 depression score and worst prognostic results in the Poor Knee group, which implies the
419 importance of depression intervention in knee OA management. The association between
420 mental health, especially depression, and knee OA has long been established. In an OAI sub-
421 cohort, Rathbun et al.³⁴ has reported the association between depression and faster disease
422 progression and faster disease progression among individuals with radiographic knee OA.
423 Additionally, in an older OA cohort, Parmelee et al.³⁵ further confirmed depression as a
424 moderator between OA pain and negative affect.

425 Our study has several strengths. First, our phenotyping models are parsimonious and do
426 not rely on invasive or expensive genetic and biomarker outcomes. Secondly, all predictors can
427 be collected when a patient seeks clinical care using validated questionnaires. Thirdly, the
428 phenotyping we developed can predict long-term symptomatic and radiographic OA
429 progression using modifiable predictors. Thus, it could be used to assist clinicians for clinical
430 decisions to modify the risk factors and potentially lead to change of disease progression.
431 Finally, our models were designed to not only address end-stage knee OA patients, but also
432 individuals seeking clinical care due to recent knee pain, thus allowing for comprehensive
433 disease cycle management.

434 Several limitations of our study are worth noting. First, although the OAI dataset used for
435 our analysis enrolled a diverse patient group from sites across the USA, our findings need to
436 be validated in independent populations. Secondly, in the current study, it was not possible to
437 assess how using our phenotyping model as a decision aid would affect patient outcomes.
438 However, we have built our phenotyping model into an online platform which can be openly
439 accessed as validation step prior to its approval by regulatory bodies for clinical use.

440

441 **Conclusion**

442 In summary, we identified four distinct knee OA phenotypes using unsupervised ML
443 methods reflecting differences in knee symptoms and supplemental vitamin intake. Phenotypes
444 were strongly associated with long-term disease fate. Supervised ML results confirmed that this
445 phenotyping could be achieved with parsimonious, modifiable variables, and we propose this
446 strategy could improve clinical decisions.

447

448 **Acknowledgments**

449 Dr. Zeyu Huang wishes to acknowledge funding from the National Natural Science
450 Foundation of China (NSFC: 92049101; 81972097; 81702185) and Sichuan Science and
451 Technology Programs (No. 2022YFH0101, 2018HH0071). Dr. John T. Martin wishes to
452 acknowledge funding from the National Institute of Arthritis and Musculoskeletal and Skin
453 Diseases (K99 AR077685, R00 AR077685).

454

455 **Author Contributions**

456 John T. Martin, Zeyu Huang, and Weihua Guo designed the study. John T. Martin and
457 Weihua Guo analyzed data. John T. Martin, Zeyu Huang, Mary A. Bucklin, and Weihua Guo
458 interpreted the data. Zeyu Huang and Mary A. Bucklin drafted the manuscript. All authors
459 critically revised the manuscript.

460

461 **Conflict of Interest**

462 Zeyu Huang is a consultant for DePuy Synthes. Neither this company nor the funding
463 sources for this work contributed to the study design, data collection, data analysis, manuscript
464 preparation, or decision to submit this manuscript.

465

466

467

468

469

470

471

472

473

474

475

476

477 **Figure Legends**

478 **Figure 1. Overview of the experiment design.** Osteoarthritis Initiative (OAI) data was
479 organized and cleaned with 4,669 subjects (patients) and 737 variables. Unsupervised
480 clustering was used to stratify the patients into four clusters. The detailed characteristics of
481 each cluster were investigated with cluster annotation and survival analysis. A web-based
482 clinical tool was developed to predict the cluster new patient could belong to with required
483 information. Based on the most accurate WOMAC total score (WOMTS) prediction from an
484 artificial intelligence model, and OA care guideline was also provided for translational usage.

485

486 **Figure 2. Cluster characteristics of OAI.** (A) Four clusters on UMAP. (B) Top 10 numeric
487 variables of each cluster. Kruskal-Wallis test was used to determine the statistics between the
488 cluster of interests and all the other clusters together. Benjamini & Hochberg method was used
489 to adjust the p-value. The numerical variables with adjusted p-values <0.05 were ranked by the
490 log₂ fold changes (log₂FC) to select the top 10 of each cluster. (C) Top 10 categorical variables
491 of each cluster. Fisher's exact test was used to determine the statistics between the cluster of
492 interests and all the other clusters together. Benjamini & Hochberg method was used to adjust
493 the p-value. The categorical variables with adjusted p-values <0.05 were ranked by the
494 Pearson's chi-squared statistics to select the top 10 of each cluster. (D)~(I) Key variables
495 categorized into demographic (V00AGE, age; P02RACE, race; P02SEX, gender; P01BMI, BMI
496 at baseline), medical record (V00COMRB, Charlson Comorbidity Index; V00HSPSS, Short
497 Form 12 Physical Summary Score), pain evaluation (V00WOMKPL/R, WOMAC pain score of
498 left/right knee), diet & nutrition (V00VITCCV, Vitamin C single vitamin, how often taken in past

499 12 months; V00SUPVITC, average daily Vitamin C supplement, mg), psychological evaluation
500 (V00CESD, Center for Epidemiology Studies Depression Index; V00HSMSS, Short Form 12,
501 Mental Summary Score), socioeconomic status (V00INCOME, annual personal income;
502 V00EDCV, education level) on UMAP.

503

504 **Figure 3. Prognostic values of 4 OA phenotypes.** Kaplan-Meier plots (first and third from left)
505 and forest plots (second and fourth from left) considering good knee health cluster as reference
506 of KL grade (A), joint space width (B), WOMAC total score (C), WOMAC pain score (D),
507 WOMAC stiffness score (E), WOMAC function score (F), and total knee replacement (G) were
508 shown in a table format for both left (left two columns) and right (right two columns) knees. Log-
509 rank p-value was shown in the KM plots. A univariant cox regression model for each outcome
510 variable and each knee was built with the Good Knee group as the reference group and
511 visualized in the forest plots.

512

513 **Figure 4. Prognostic values of four OA phenotypes within baseline cohorts.** (A) Relative
514 distribution of 4 OA phenotypes within each baseline cohort. Kaplan-Meier plots for WOMAC
515 total score (B), KL grades (C), joint space width (D), total knee replacement (E) were shown in
516 a table format for both left (left first and third columns) and right (left second and fourth columns)
517 knees within incidence cohort (left two columns) and progression cohort (right two columns).
518 Log-rank p-value was shown in the KM plots.

519

520 **Figure 5. Prediction accuracies of cluster labels.** A) Screening the optimal number of input

521 variables (from 2 to 50) with different machine learning models (lr=linear regression, logistic
522 regression model; rf10/20/40/60/80/100tree = random forest model with 10/20/40/60/80/100
523 trees; svcllinear/poly/rbf/sigmoid = supporting vector classifier with linear/polynomial/radial
524 basis function/sigmoid kernels). As a multi-class prediction problem, three accuracy metrics
525 were used, i.e., accuracy (relative correct prediction numbers), roc_auc_ovo (area under curve
526 of receiver operating characteristic curve, one vs one), and roc_auc_ovr (area under cuve of
527 receiver operating characteristic curve, one vs rest). B) Detailed screening of the optimal
528 number of input variables (from 5 to 15). The dot represents the mean of corresponding metric
529 and the error bar represents the standard error of the mean from the cross-validation.

530

531 **Figure 6. Prediction accuracy of WOMAC total score at fourth and eighth year. A)**

532 Correlation coefficients between WOMTS of all visiting years and baseline variables. The top
533 10 baseline variables were colored based on the average correlation coefficients crossing all
534 the visiting years. W. = WOMAC, K.=KOOS. B) Screening the optimal number of input variables
535 (from 2 to 50) with different machine learning models (linear, linear regression model;
536 rft10/20/40/60/80/100 = random forest regressor with 20/40/60/80/100 trees;
537 svrlinear/poly/rbf/sigmoid = supporting vector regressor with linear/polynomial/rbf/sigmoid
538 kernels, ann = artificial neural network). As a regression problem, two accuracy metrics were
539 used, i.e., RMSE (root mean squared error) and r (correlation coefficient between prediction
540 and measurements). D) Top 5 variables with highest occurrences from the top 10000 most
541 accurate prediction tests. We randomly selected 25 input variables from the cluster markers
542 and used these variables to train an ANN model with the same settings with cross-validation.

543 The above procedure was repeated 10,000 times. The top 1000 most accurate tests were
544 extracted and the relative occurrence of each variable to these 1000 tests was calculated. The
545 top 5 with highest relative occurrences for WOMTS of both left and right knees at fourth and
546 eighth year were selected to visualize here. The dot represents the relative occurrences. W. =
547 WOMAC, K.=KOOS
548

549 **References:**

- 550 1. Glyn-Jones, S. et al. Osteoarthritis. *Lancet* **386**, 376-387 (2015).
- 551 2. Woolf, A.D. & Pfleger, B. Burden of major musculoskeletal conditions. *Bull World Health*
552 *Organ* **81**, 646-656 (2003).
- 553 3. Hunter, D.J. & Bierma-Zeinstra, S. Osteoarthritis. *Lancet* **393**, 1745-1759 (2019).
- 554 4. Leifer, V.P., Katz, J.N. & Losina, E. The burden of OA-health services and economics.
555 *Osteoarthritis Cartilage* **30**, 10-16 (2022).
- 556 5. Kraus, V.B. et al. Predictive validity of biochemical biomarkers in knee osteoarthritis: data
557 from the FNIH OA Biomarkers Consortium. *Ann Rheum Dis* **76**, 186-195 (2017).
- 558 6. Karsdal, M.A. et al. Disease-modifying treatments for osteoarthritis (DMOADs) of the knee
559 and hip: lessons learned from failures and opportunities for the future. *Osteoarthritis*
560 *Cartilage* **24**, 2013-2021 (2016).
- 561 7. Hellio le Graverand, M.P. et al. A 2-year randomised, double-blind, placebo-controlled,
562 multicentre study of oral selective iNOS inhibitor, cindunistat (SD-6010), in patients with
563 symptomatic osteoarthritis of the knee. *Ann Rheum Dis* **72**, 187-195 (2013).
- 564 8. Bingham, C.O., 3rd et al. Risedronate decreases biochemical markers of cartilage
565 degradation but does not decrease symptoms or slow radiographic progression in
566 patients with medial compartment osteoarthritis of the knee: results of the two-year
567 multinational knee osteoarthritis structural arthritis study. *Arthritis Rheum* **54**, 3494-3507
568 (2006).
- 569 9. Karsdal, M.A. et al. The coupling of bone and cartilage turnover in osteoarthritis:
570 opportunities for bone antiresorptives and anabolics as potential treatments? *Ann Rheum*
571 *Dis* **73**, 336-348 (2014).
- 572 10. Karsdal, M.A. et al. The effect of oral salmon calcitonin delivered with 5-CNAC on bone
573 and cartilage degradation in osteoarthritic patients: a 14-day randomized study.
574 *Osteoarthritis Cartilage* **18**, 150-159 (2010).
- 575 11. Deo, R.C. Machine Learning in Medicine. *Circulation* **132**, 1920-1930 (2015).
- 576 12. Bibault, J.E., Giraud, P. & Burgun, A. Big Data and machine learning in radiation oncology:
577 State of the art and future prospects. *Cancer Lett* **382**, 110-117 (2016).
- 578 13. Obermeyer, Z. & Emanuel, E.J. Predicting the Future - Big Data, Machine Learning, and
579 Clinical Medicine. *N Engl J Med* **375**, 1216-1219 (2016).
- 580 14. Gaskin, G.L., Pershing, S., Cole, T.S. & Shah, N.H. Predictive modeling of risk factors and
581 complications of cataract surgery. *Eur J Ophthalmol* **26**, 328-337 (2016).
- 582 15. Huang, Z. et al. Analysis of a large data set to identify predictors of blood transfusion in
583 primary total hip and knee arthroplasty. *Transfusion* **58**, 1855-1862 (2018).
- 584 16. Huang, Z., Guo, W. & Martin, J.T. Socioeconomic status, mental health, and nutrition are
585 the principal traits for low back pain phenotyping: Data from the osteoarthritis initiative.
586 *JOR Spine* **6**, e1248 (2023).
- 587 17. Huang, Z. et al. Predicting postoperative transfusion in elective total HIP and knee
588 arthroplasty: Comparison of different machine learning models of a case-control study.
589 *Int J Surg* **96**, 106183 (2021).
- 590 18. Urish, K.L. et al. T2 texture index of cartilage can predict early symptomatic OA progression:
591 data from the osteoarthritis initiative. *Osteoarthritis Cartilage* **21**, 1550-1557 (2013).

- 592 19. Carlini, N. & Wagner, D. in 2017 IEEE Symposium on Security and Privacy (SP) 39-57 (IEEE, 593 2017).
- 594 20. Angelini, F. et al. Osteoarthritis endotype discovery via clustering of biochemical marker 595 data. *Ann Rheum Dis* **81**, 666-675 (2022).
- 596 21. Yuan, C. et al. Classification of four distinct osteoarthritis subtypes with a knee joint tissue 597 transcriptome atlas. *Bone Res* **8**, 38 (2020).
- 598 22. Enroth, S., Johansson, A., Enroth, S.B. & Gyllensten, U. Strong effects of genetic and 599 lifestyle factors on biomarker variation and use of personalized cutoffs. *Nat Commun* **5**, 600 4684 (2014).
- 601 23. Liu, Q. et al. Prediction models for the risk of total knee replacement: development and 602 validation using data from multicentre cohort studies. *Lancet Rheumatol* **4**, e125-e134 603 (2022).
- 604 24. Yu, D. et al. Development and validation of prediction models to estimate risk of primary 605 total hip and knee replacements using data from the UK: two prospective open cohorts 606 using the UK Clinical Practice Research Datalink. *Ann Rheum Dis* **78**, 91-99 (2019).
- 607 25. Leung, K. et al. Prediction of Total Knee Replacement and Diagnosis of Osteoarthritis by 608 Using Deep Learning on Knee Radiographs: Data from the Osteoarthritis Initiative. 609 *Radiology* **296**, 584-593 (2020).
- 610 26. Driban, J.B. et al. The natural history of end-stage knee osteoarthritis: Data from the 611 osteoarthritis initiative. *Semin Arthritis Rheum* **58**, 152148 (2023).
- 612 27. Pierson, E., Cutler, D.M., Leskovec, J., Mullainathan, S. & Obermeyer, Z. An algorithmic 613 approach to reducing unexplained pain disparities in underserved populations. *Nat Med* 614 **27**, 136-140 (2021).
- 615 28. Ovesen, L. Vitamin therapy in the absence of obvious deficiency. What is the evidence? 616 *Drugs* **27**, 148-170 (1984).
- 617 29. Canter, P.H., Wider, B. & Ernst, E. The antioxidant vitamins A, C, E and selenium in the 618 treatment of arthritis: a systematic review of randomized clinical trials. *Rheumatology 619 (Oxford)* **46**, 1223-1233 (2007).
- 620 30. Kraus, V.B. et al. Ascorbic acid increases the severity of spontaneous knee osteoarthritis in 621 a guinea pig model. *Arthritis Rheum* **50**, 1822-1831 (2004).
- 622 31. Qu, Z. et al. Causal relationship of serum nutritional factors with osteoarthritis: a 623 Mendelian randomization study. *Rheumatology (Oxford)* **60**, 2383-2390 (2021).
- 624 32. Wallace, I.J. et al. Knee osteoarthritis has doubled in prevalence since the mid-20th 625 century. *Proc Natl Acad Sci U S A* **114**, 9332-9336 (2017).
- 626 33. Gustafsson, K. et al. Health status of individuals referred to first-line intervention for hip 627 and knee osteoarthritis compared with the general population: an observational register- 628 based study. *BMJ Open* **11**, e049476 (2021).
- 629 34. Rathbun, A.M. et al. Association between disease progression and depression onset in 630 persons with radiographic knee osteoarthritis. *Rheumatology (Oxford)* **59**, 3390-3399 631 (2020).
- 632 35. Parmelee, P.A. et al. Momentary Associations of Osteoarthritis Pain and Affect: Depression 633 as Moderator. *J Gerontol B Psychol Sci Soc Sci* **77**, 1240-1249 (2022).

634

Figures

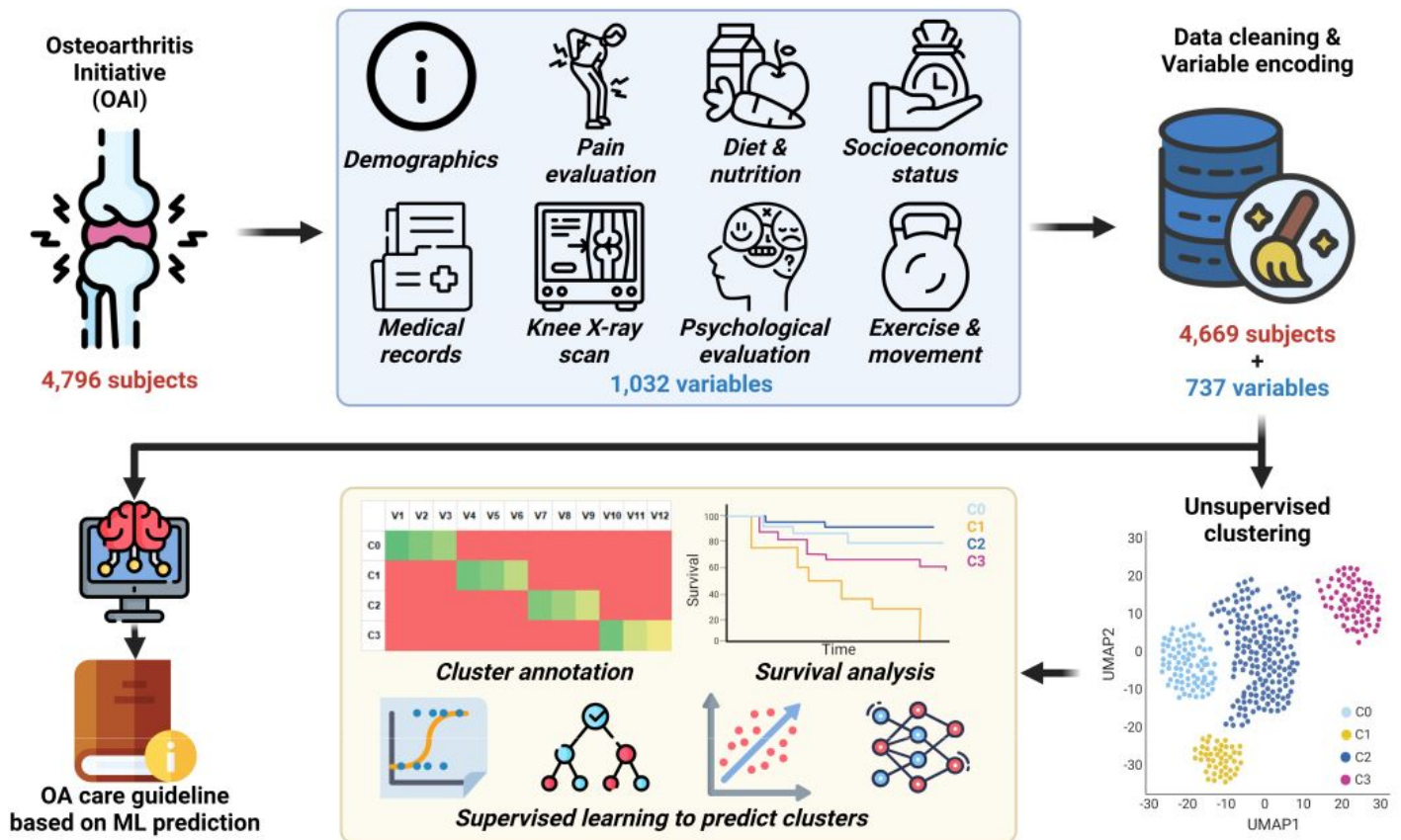


Figure 1

Overview of the experiment design. Osteoarthritis Initiative (OAI) data was organized and cleaned with 4,669 subjects (patients) and 737 variables. Unsupervised clustering was used to stratify the patients into four clusters. The detailed characteristics of each cluster were investigated with cluster annotation and survival analysis. A web-based clinical tool was developed to predict the cluster new patient could belong to with required information. Based on the most accurate WOMAC total score (WOMTS) prediction from an artificial intelligence model, and OA care guideline was also provided for translational usage.

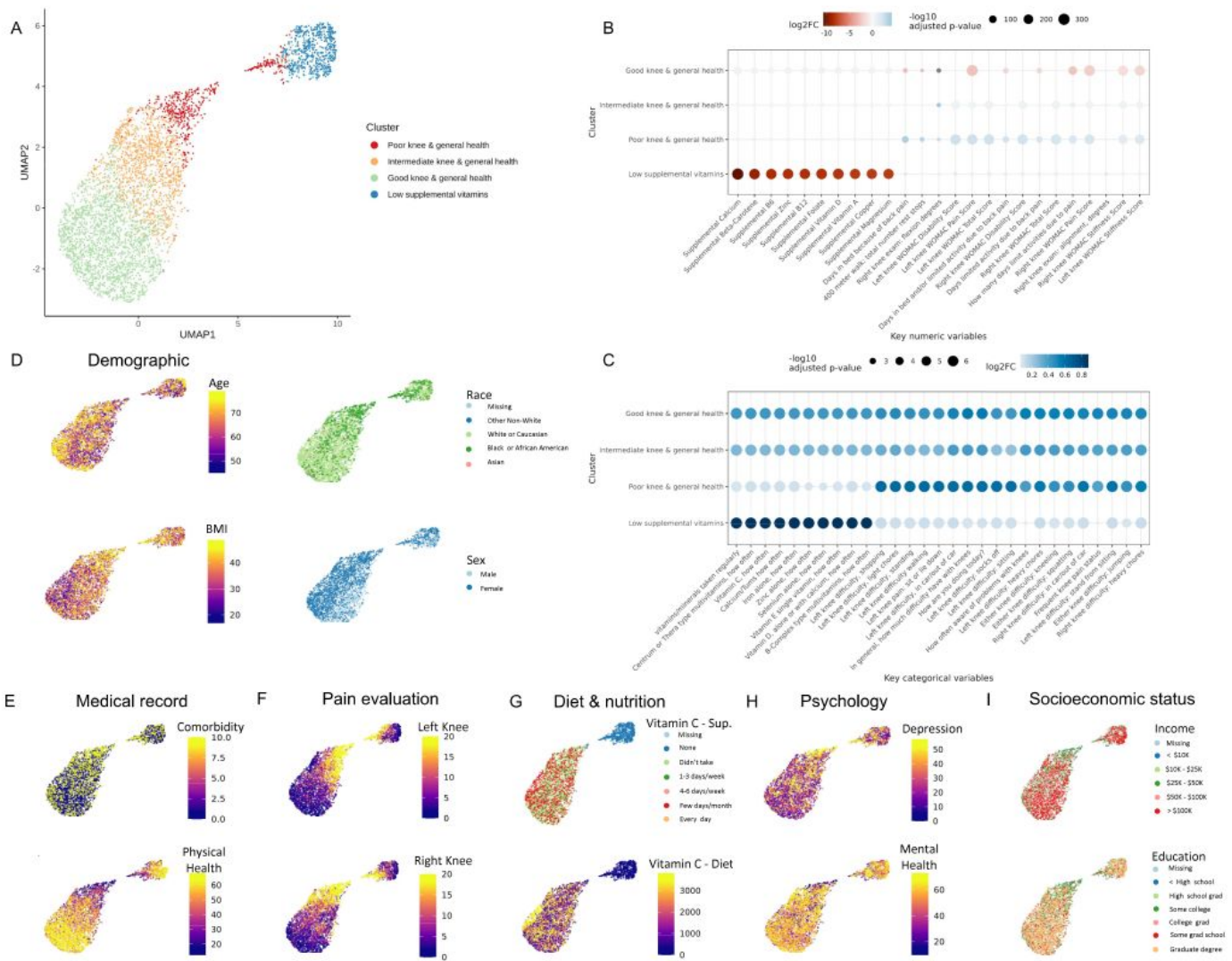


Figure 2

Cluster characteristics of OAI. (A) Four clusters on UMAP. (B) Top 10 numeric variables of each cluster. Kruskal-Wallis test was used to determine the statistics between the cluster of interests and all the other clusters together. Benjamini & Hochberg method was used to adjust the p-value. The numerical variables with adjusted p-values <0.05 were ranked by the log2 fold changes (log2FC) to select the top 10 of each cluster. (C) Top 10 categorical variables of each cluster. Fisher's exact test was used to determine the statistics between the cluster of interests and all the other clusters together. Benjamini & Hochberg method was used to adjust the p-value. The categorical variables with adjusted p-values <0.05 were ranked by the Pearson's chi-squared statistics to select the top 10 of each cluster. (D)~(I) Key variables categorized into demographic (V00AGE, age; P02RACE, race; P02SEX, gender; P01BMI, BMI at baseline), medical record (V00COMRB, Charlson Comorbidity Index; V00HSPSS, Short Form 12 Physical Summary Score), pain evaluation (V00WOMKPL/R, WOMAC pain score of left/right knee), diet & nutrition (V00VITCCV, Vitamin C single vitamin, how often taken in past

12 months; V00SUPVITC, average daily Vitamin C supplement, mg), psychological evaluation (V00CESD, Center for Epidemiology Studies Depression Index; V00HSMSS, Short Form 12, Mental Summary Score), socioeconomic status (V00INCOME, annual personal income; V00EDCV, education level) on UMAP.

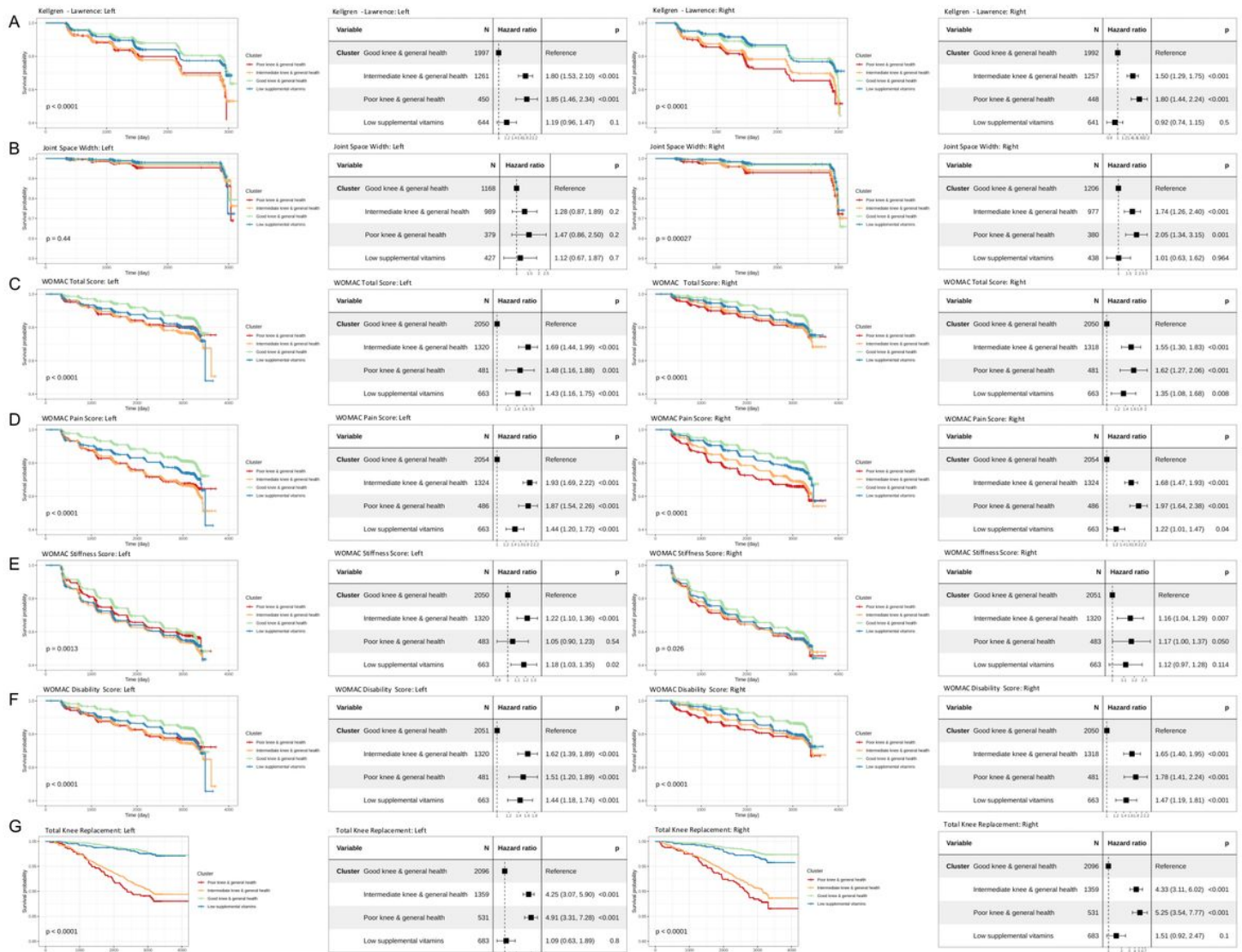


Figure 3

Prognostic values of 4 OA phenotypes. Kaplan-Meier plots (first and third from left) and forest plots (second and fourth from left) considering good knee health cluster as reference of KL knee grade (A), joint space width (B), WOMAC total score (C), WOMAC pain score (D), WOMAC stiffness score (E), WOMAC function score (F), and total knee replacement (G) were shown in a table format for both left (left two columns) and right (right two columns) knees. Log-rank p-value was shown in the KM plots. A univariate cox regression model for each outcome variable and each knee was built with the Good Knee group as the reference group and visualized in the forest plots.

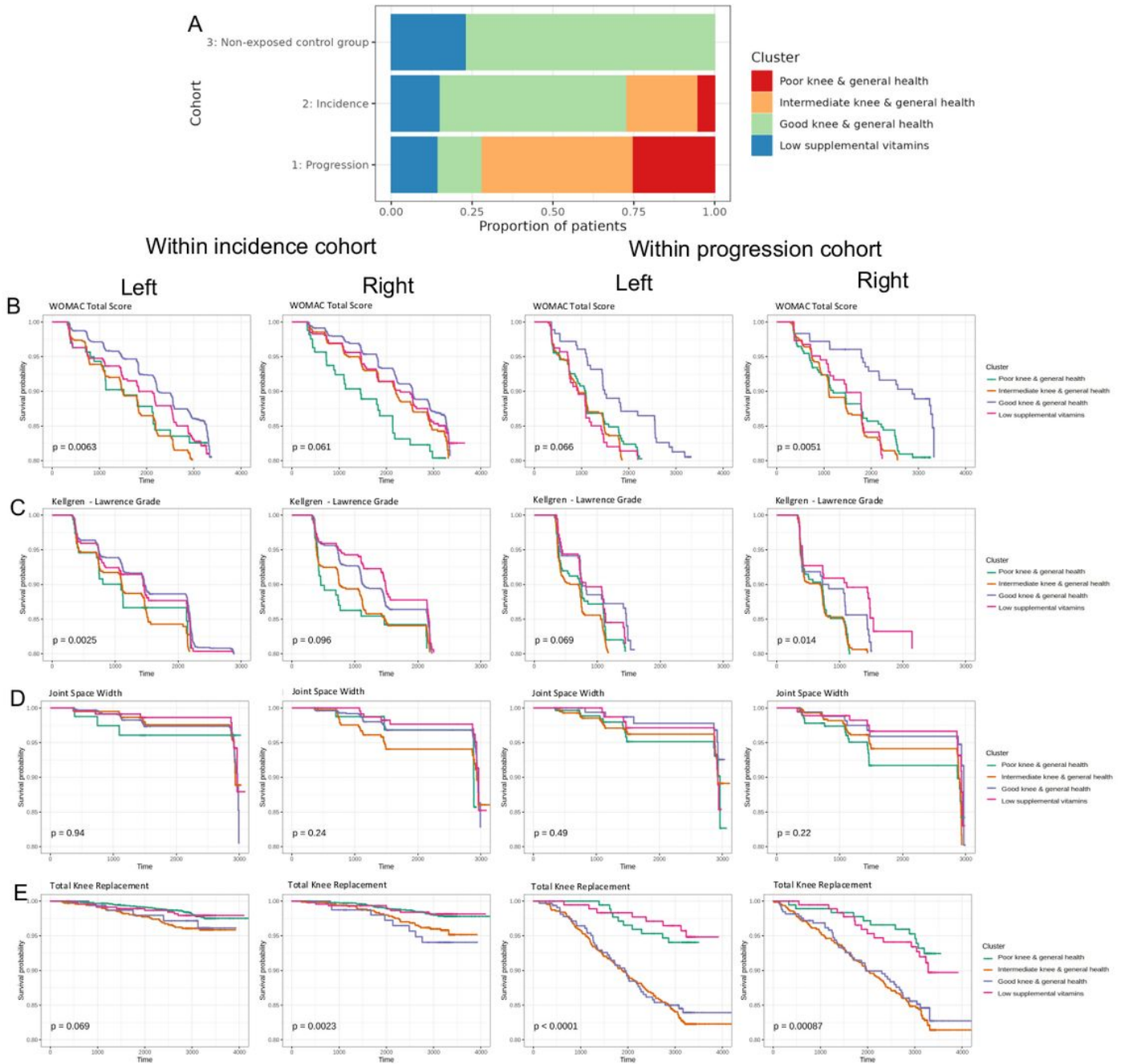


Figure 4

Prognostic values of four OA phenotypes within baseline cohorts. (A) Relative distribution of 4 OA phenotypes within each baseline cohort. Kaplan-Meier plots for WOMAC total score (B), KL grades (C), joint space width (D), total knee replacement (E) were shown in a table format for both left (left first and third columns) and right (left second and fourth columns) knees within incidence cohort (left two columns) and progression cohort (right two columns). Log-rank p-value was shown in the KM plots.

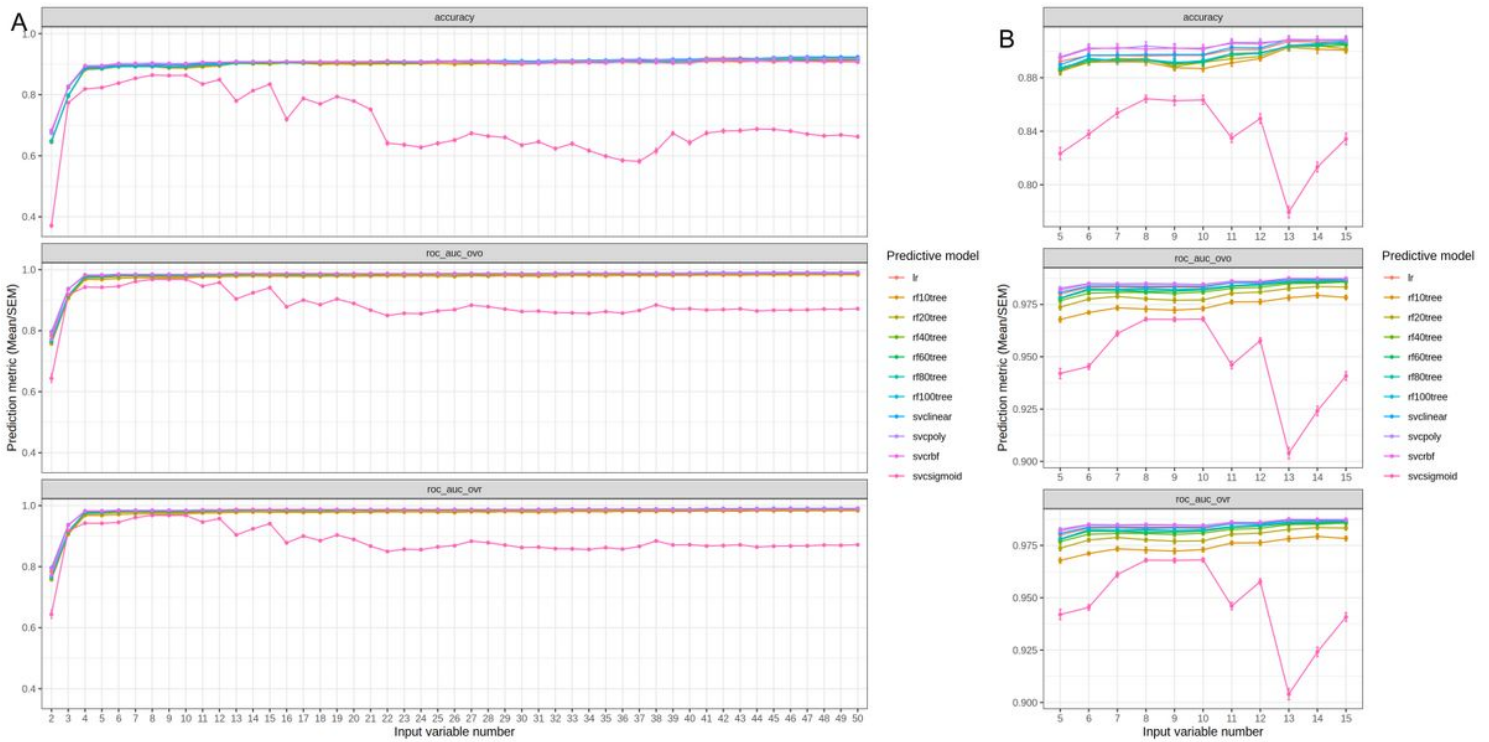


Figure 5

Prediction accuracies of cluster labels. A) Screening the optimal number of input

24

variables (from 2 to 50) with different machine learning models (lr=linear regression, logistic regression model; rf10/20/40/60/80/100tree = random forest model with 10/20/40/60/80/100 trees; svlinear/poly/rbf/sigmoid = supporting vector classifier with linear/polynomial/radial basis function/sigmoid kernels). As a multi-class prediction problem, three accuracy metrics were used, i.e., accuracy (relative correct prediction numbers), roc_auc_ovo (area under curve of receiver operating characteristic curve, one vs one), and roc_auc_ovr (area under curve of receiver operating characteristic curve, one vs rest). B) Detailed screening of the optimal number of input variables (from 5 to 15). The dot represents the mean of corresponding metric and the error bar represents the standard error of the mean from the cross-validation.

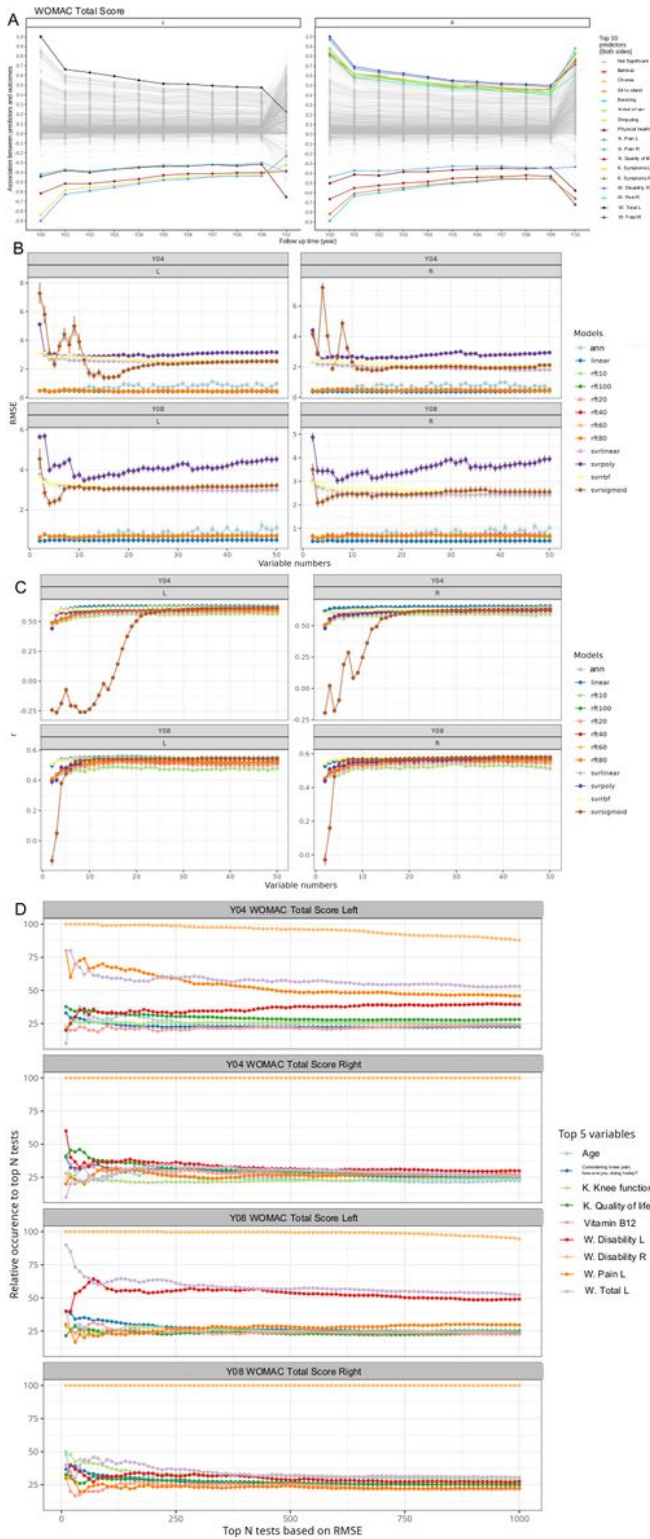


Figure 6

Prediction accuracy of WOMAC total score at fourth and eighth year. A) Correlation coefficients between WOMTS of all visiting years and baseline variables. The top 10 baseline variables were colored based on the average correlation coefficients crossing all the visiting years. W. = WOMAC, K.=KOOS. B) Screening the optimal number of input variables (from 2 to 50) with different machine learning models (linear, linear regression model; rft10/20/40/60/80/100 = random forest regressor with 20/40/60/80/100 trees;

svrlinear/poly/rbf/sigmoid = supporting vector regressor with linear/polynomial/rbf/sigmoid kernels, ann = artificial neural network). As a regression problem, two accuracy metrics were used, i.e., RMSE (root mean squared error) and r (correlation coefficient between prediction and measurements). D) Top 5 variables with highest occurrences from the top 10000 most accurate prediction tests. We randomly selected 25 input variables from the cluster markers and used these variables to train an ANN model with the same settings with cross-validation.

25

The above procedure was repeated 10,000 times. The top 1000 most accurate tests were extracted and the relative occurrence of each variable to these 1000 tests was calculated. The top 5 with highest relative occurrences for WOMTS of both left and right knees at fourth and eighth year were selected to visualize here. The dot represents the relative occurrences. W. = WOMAC, K.=KOOS

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [OAIKneeOASupplementaryInfov3.pdf](#)
- [Table1231208.pdf](#)
- [SD1usedinputvariables.csv](#)
- [SD2cleanv25finalcluster4kmeansdirectknn2imp12112020datacleanmarkerdflargeB.csv](#)
- [SD3outcomerealdateconversionyear.csv](#)
- [SD4directpredictcluster230109importancedataframe.xlsx](#)
- [SD5directpredictcluster230109mergescoredataframe.csv](#)
- [SD6predictorvis230213gathmerge.csv](#)
- [SD7mergescoredfptestmax9999230202.csv](#)
- [FigS1231207.pdf](#)
- [FigS2231207.pdf](#)
- [FigS3231207.pdf](#)
- [FigS4231207.pdf](#)
- [FigS5231207.pdf](#)
- [OAIKneeOASupplementaryInfov3.pdf](#)