

# 1 Integrating Multidimensional Data Analytics for Precision Diagnosis 2 of Chronic Low Back Pain

3 Sam Vickery<sup>1\*</sup>, Frederick Junker<sup>1</sup>, Rebekka Döding<sup>1</sup>, Daniel L Belavy<sup>1</sup>, Maia Angelova<sup>2,3</sup>,  
4 Chandan Karmakar<sup>3</sup>, Luis Becker<sup>4,5</sup>, Nima Taheri<sup>4,5</sup>, Matthias Pumberger<sup>4</sup>, Sandra Reitmaier<sup>5</sup>,  
5 Hendrik Schmidt<sup>5</sup>

6  
7 <sup>1</sup> Fachbereich Pflege-, Hebammen- und Therapiewissenschaften (PHT) | Hochschule Bochum  
8 (University of Applied Sciences) | Bochum | Germany

9 <sup>2</sup> Aston Digital Futures Institute | Aston University | Birmingham | United Kingdom

10 <sup>3</sup> School of Information Technology | Deakin University | Geelong | Australia

11 <sup>4</sup> Center for Musculoskeletal Surgery, Charité – Universitätsmedizin Berlin, Berlin, Germany

12 <sup>5</sup> Julius Wolff Institut, Berlin Institute of Health - Charité at Universitätsmedizin Berlin |  
13 Berlin | Germany

14 **Short title:** Multimodal classification of chronic low back pain

15

## 16 **Abstract**

17 Low back pain (LBP) is a leading cause of disability worldwide, with up to 25% of cases  
18 become chronic (cLBP). Optimal diagnostic tools for cLBP remains unclear. Here we leveraged  
19 a comprehensive multi-dimensional data-set and machine learning-based variable importance  
20 selection to identify the most effective diagnostic tools for cLBP patient stratification. The  
21 dataset included questionnaire data, clinical and functional assessments, and spino-pelvic  
22 magnetic resonance imaging (MRI), encompassing a total of 144 parameters from 1,161 adults  
23 with (n=512) and without cLBP (n=649). Boruta and random forest were utilised for variable  
24 importance selection and cLBP classification respectively. Boruta variable selection led to  
25 pronounced variable reduction (median of all 15 datasets: 63.3%), while performing  
26 comparable to using all variables across all modality datasets. Boruta selected key variables  
27 from questionnaire, clinical, and MRI data were the most effective in distinguishing cLBP  
28 patients from controls. The most robust variables (n=9) across the whole dataset identified were  
29 psychosocial factors, neck and hip mobility, as well as lower lumbar disc herniation and  
30 degeneration. These critical variables outperformed all parameters in an unseen holdout dataset,  
31 demonstrating superior patient delineation. Paving the way for targeted diagnosis and  
32 personalized treatment strategies, ultimately enhancing clinical outcomes for cLBP patients.

33

34 **Key words**

35 Chronic low back pain; classification; data-driven; feature selection; multi-modality;  
36 psychosocial; MRI

37

## 38 Introduction

39

40 In recent years, chronic low back pain (cLBP) has become one of the most prevalent and  
41 challenging conditions in clinical practice, affecting a significant portion of the global  
42 population <sup>1,2</sup>. Despite its widespread occurrence, diagnosing and assessing cLBP remains  
43 difficult due to the complex and multifactorial nature of the disease <sup>3,4</sup>, which encompasses  
44 physical, psychological, and social dimensions <sup>5</sup>. Traditional diagnostic approaches often rely  
45 on self-reported symptoms, which can be subjective and prone to variability. Furthermore,  
46 surgical and non-surgical treatment outcomes are still inconsistent, reflected in a high rate of  
47 treatment-refractory in cLBP patients <sup>6</sup>. As such, there is an increasing interest in exploring more  
48 objective and comprehensive methods that incorporate multimodal data, including medical  
49 imaging, physical assessments, clinical evaluations, and patient-reported outcomes.

50

51 Machine learning algorithms provide models for identifying distinct subgroups that can help  
52 explain the occurrence and characteristics of a disease <sup>7,8</sup>. A previous systematic review by our  
53 team using machine learning applications in LBP <sup>9</sup> highlighted that a narrow range of  
54 mechanistic domains have been assessed, and sample sizes in these studies were consistently  
55 small, ranging up to only 171 participants. Consequently, using limited data and modalities  
56 limits the robustness and applicability of such models. Through gathering many data points  
57 across multiple modalities one can ascertain which variables and modalities are the most  
58 informative at distinguishing cLBP patients from asymptomatic controls. Reducing the number  
59 of variables to those that are most informative has been previously employed in predictive and  
60 classification modelling to improve accuracy <sup>10,11</sup>. This approach can be applied as the main  
61 outcome and not only in model preprocessing, in order to obtain a data-driven decision on the  
62 most important variables in multi-dimensional clinical data. Such a systematic data-informed  
63 investigation of back pain diagnosis in a large multi-modality sample is lacking to help inform

64 future studies in selecting which data to acquire and for clinicians in which tests to conduct.  
65 Therefore, the aim of this study was to identify and compare the most informative domains and  
66 variables in delineating patients with and without cLBP utilising a large multi-modality dataset.  
67

## 68 Results

### 69 Study sample

70 The prospective cross-sectional study draws its data from the ongoing “Berliner Rückenstudie”  
71 (“Berlin Back Study”; [https://spine.charite.de/en/spine\\_study/](https://spine.charite.de/en/spine_study/); running time: 01/01/2022 to  
72 31/12/2025), which was registered at the German Clinical Trial Register (DRKS-ID:  
73 [DRKS00027907](#)). Recruitment procedures vary from local promotion (i.e., postal flyers, notice  
74 boards, internet approaches, and social media) at the Charité-Universitätsmedizin Berlin, in  
75 the general public (i.e., newspapers, magazines, podcasts) to cooperation with local companies,  
76 administrative authorities, and word-of-mouth. The protocol is in accordance with the Helsinki  
77 Declaration of ethical principles<sup>12</sup> and has been approved by the Ethics Committee of the  
78 Charité – Universitätsmedizin Berlin (registry numbers: EA4/011/10, EA1/162/13). Written  
79 informed consent was obtained from all participants. The STROBE guideline<sup>13</sup>  
80 (Supplementary Table S1) and TRIPOD statement<sup>14</sup> (Supplementary Table S2) for prediction  
81 model development were used to report this study. Data collection started on 1st January 2022  
82 and cut-off for inclusion in the current analysis was 5th April 2024. Data collection occurred in  
83 a research centre within a university-hospital.

84  
85 Study participants were recruited through a telephone interview and excluded if they met any  
86 exclusion criteria, as well as some excluded at the testing site (Supplementary Table S3). A  
87 total of 1273 participants were included in the study at cut-off point. These participants were  
88 initially guided through self-administered questionnaires by a study coordinator. Then they  
89 continued to a clinical examination by a trained medical doctor, which included physical  
90 examinations, questions, as well as a back shape and function test. The examinations and  
91 questionnaires took a total of 90 minutes to complete. Additionally, participants were offered a

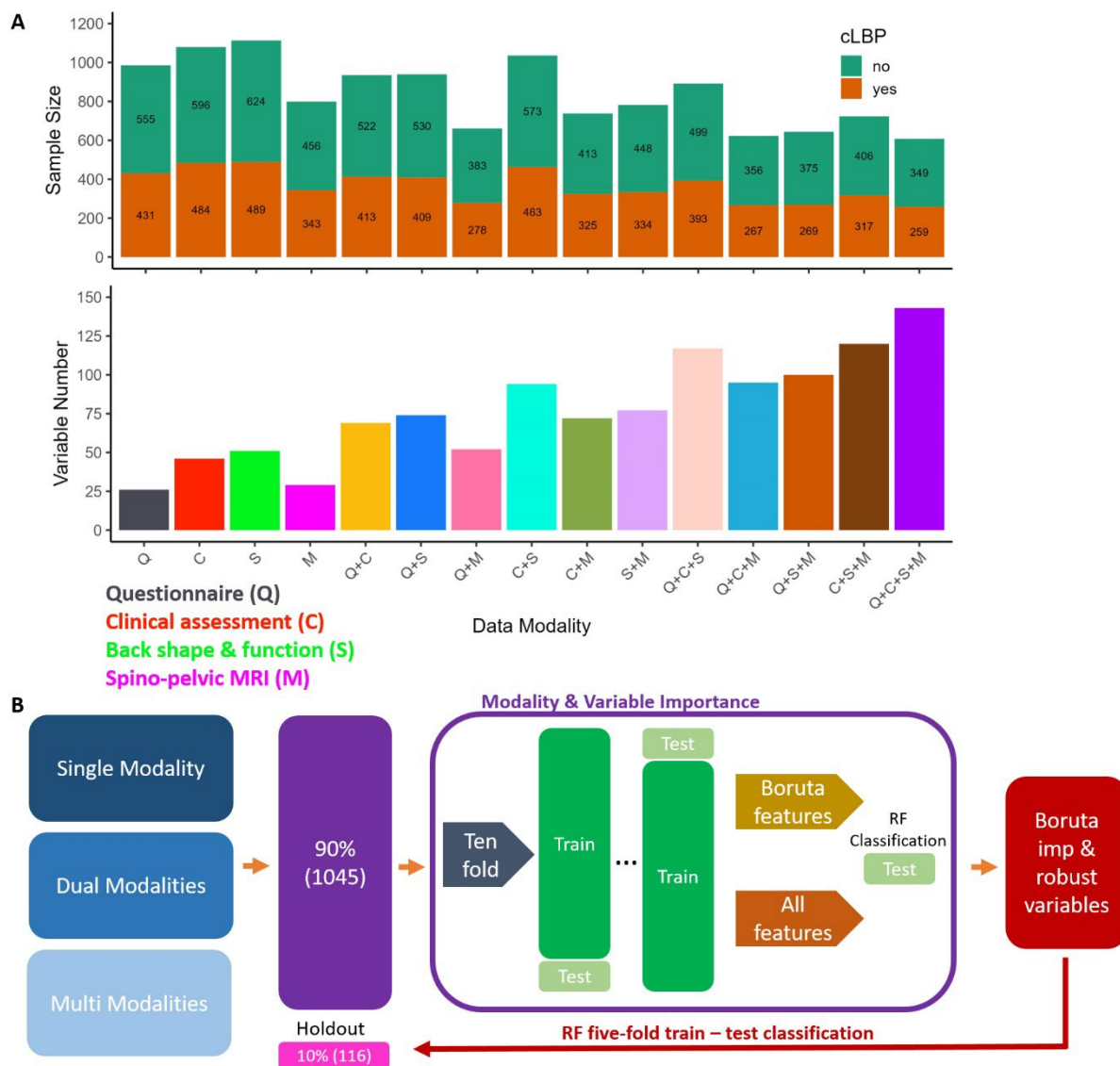
92 magnetic resonance imaging (MRI) within 14 days of the spino-pelvic region. During the  
93 clinical assessment the participants were classified by the clinician as asymptomatic (no back  
94 pain), symptomatic (cLBP), or previously suffering from cLBP. To ensure a more robust cLBP  
95 patient classification previous symptomatic subjects were removed from the sample. The  
96 criteria for designation as cLBP was daily back pain for more than three months. Furthermore,  
97 participants who revoked their inclusion in the study and those who were missing demographic  
98 data; age, sex, body mass index (BMI), and patient status were removed. This resulted in a  
99 study sample of 1161 subjects that included 649 asymptomatic (19 – 72 years old, mean age =  
100  $40.7 \pm 12.6$ , females = 353) and 512 cLBP (19 – 65 years old, mean age =  $43.5 \pm 11.7$ , females  
101 = 306) participants (Table 1). This sample was sub-divided into four modalities; questionnaires  
102 (Q), clinical physical assessment (C), back shape and function (S), and MRI (M). Each modality  
103 was combined with demographic data (age, sex, and BMI) and then joined with all  
104 combinations of the four modalities, resulting in 15 datasets (Figure 1A, and Supplementary  
105 Table S4 - S18). Both complete datasets (Fig. 2) and imputed datasets were used for cLBP  
106 classification (Fig. 1B). Imputation can lead to bias results, in particular when it comes to  
107 variable importance selection <sup>15</sup>. Therefore, presentation of variable importance was only  
108 considered for non-imputed datasets. Highly correlated ( $r > 0.9$ ) variables were removed from  
109 the datasets to reduce colinearity between dataset variables (Supplementary Figure S1 – S6).  
110 An overview of all Berlin Back Study variables (Supplementary Table S19), those removed  
111 during preprocessing (Supplementary Table S20), and a list of all variables (144) used for  
112 modelling is presented in Supplementary (Supplementary Table S21).

113 Table 1. Demographic data of entire Berlin Back dataset

	<b>Asymptomatic</b>	<b>Chronic low back pain</b>
Sample size	649	512
Females	353	306
Age (years)		
mean [sd]	40.65 [12.6]	43.52 [11.67]
Body mass index		
mean [sd]	23.55 [2.75]	23.55 [2.84]

Pain duration mean [sd]	549.68 [532.38]
Pain Intensity mean [sd]	2.87 [1.86]

114

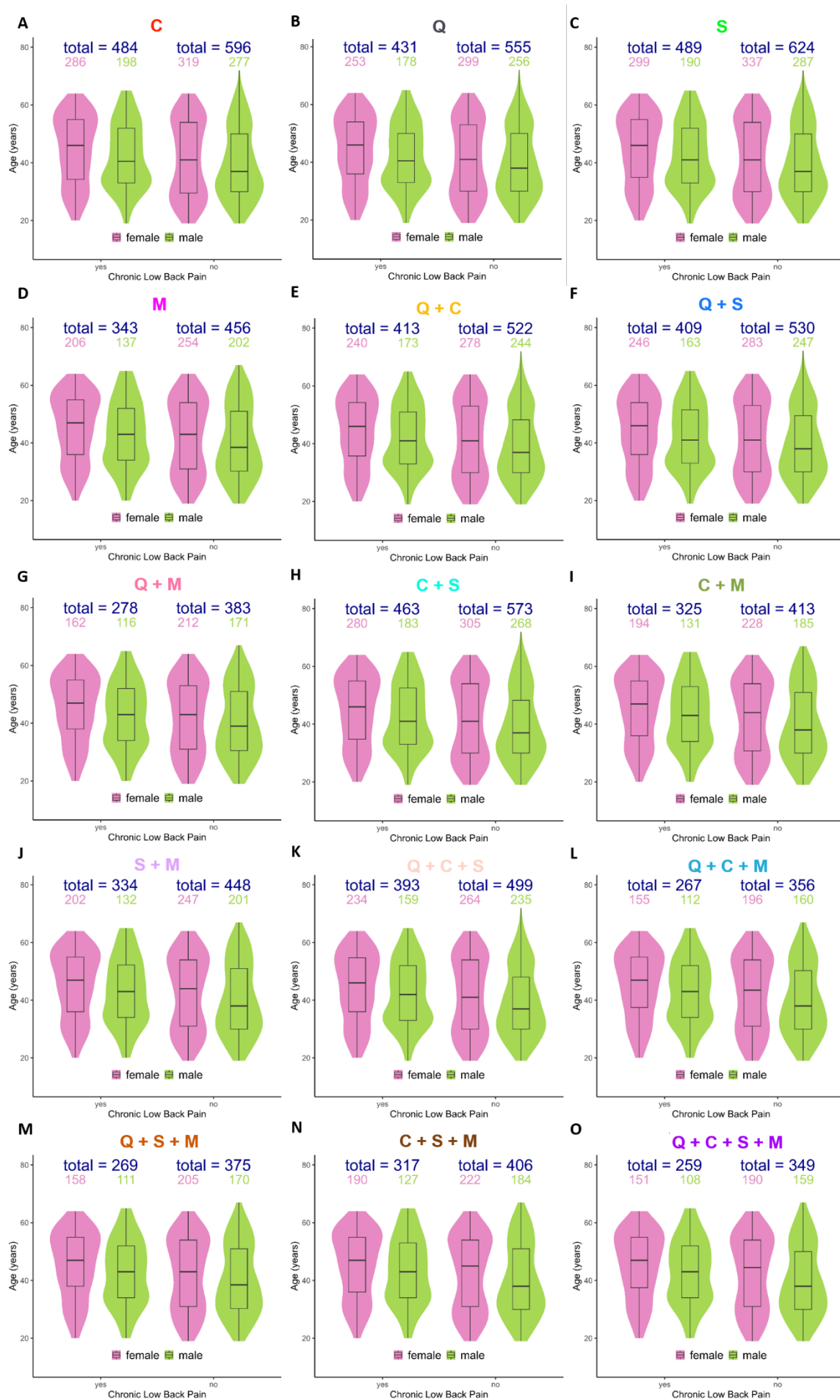


115

116 **Figure 1. Modality dataset distributions and machine learning workflow.** A – Top shows the chronic  
 117 low back pain (cLBP) sample size distribution across all 15 dataset modalities. Bottom presents the  
 118 number of variables used for cLBP classification and variable importance selection across the 15 dataset  
 119 modalities. B – Represents the machine learning workflow implemented to compare the different  
 120 modalities and determine the most important variables for cLBP patient delineation using a random  
 121 forest binary classification algorithm for training and testing.

122





124 **Figure 2. Modality datasets age, sex, and cLBP distributions.** Violin plots for non-imputed dataset  
125 A - questionnaire; B – clinical assessment: C – back shape and function; D – MRI; E – questionnaire +  
126 clinic; F – questionnaire + back shape and function; G – questionnaire + MRI; H – clinic + back shape  
127 and function; I – clinic + MRI; J – back shape and function + MRI; K – questionnaire + clinic + back  
128 shape and function; L – questionnaire + clinic + MRI; M – questionnaire + back shape and function +  
129 MRI; N – clinic + back shape and function + MRI; O – all datasets.  
130

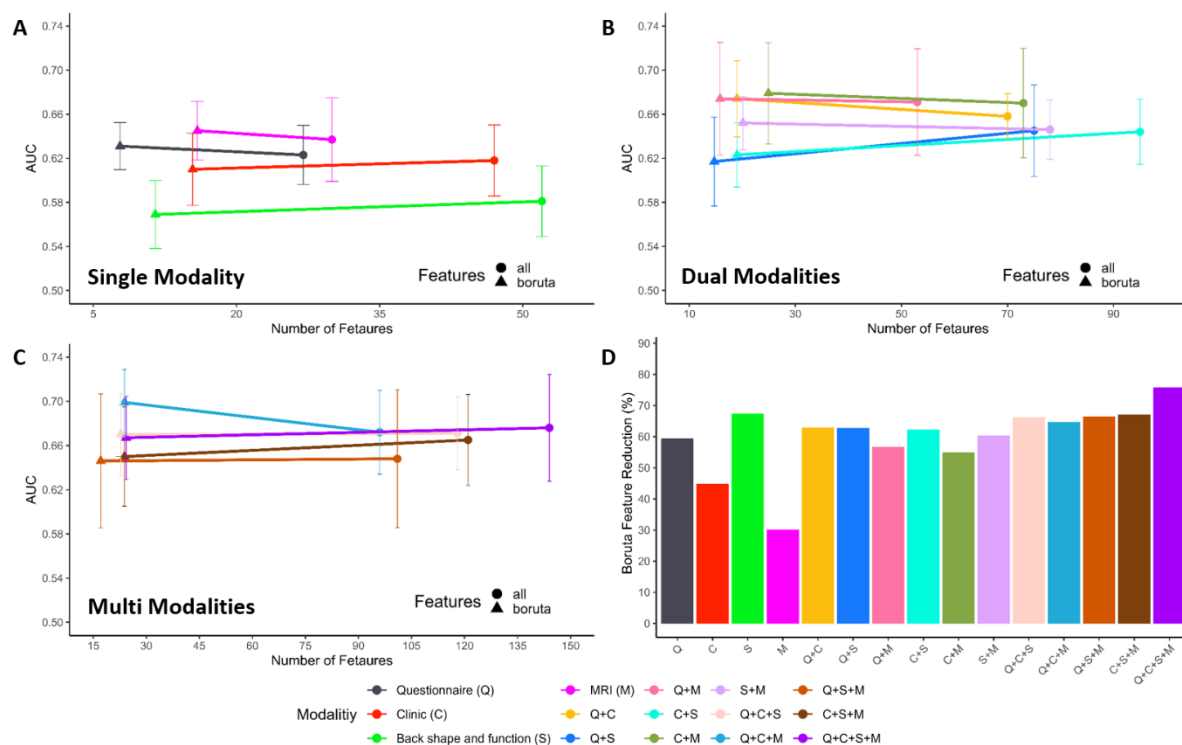
## 131 Chronic low back pain classification

132 Determining modality and variable importance was conducted on 90% (n = 1045, 19 – 72 y/o,  
133 mean age =  $41.71 \pm 12.28$ , cLBP = 469) of the whole dataset, with 10% (n = 116, 20 – 64 y/o,  
134 mean age =  $43.81 \pm 12.12$ , cLBP = 43) used as a hold-out sample to test cLBP patient  
135 classification performance of the most robust and important variables (Fig. 1B). Variable  
136 selection and importance was determined using the Boruta algorithm <sup>16</sup> on the training data  
137 within an iterative ten-fold train-test loop across all 15 modality datasets. Boruta is a wrapper  
138 function selects important variables and removes unimportant variables using random random  
139 forest (RF) <sup>17</sup>. A binary RF classification model was trained using the Boruta <sup>16</sup> selected  
140 variables as well as all variables to classify cLBP status in the test sample across the ten-fold  
141 train-test loop. Furthermore, imputation using missForest <sup>18</sup> was independently conducted on  
142 the training and test datasets within the ten-fold loop (Supplementary Table S23). Imputed  
143 datasets performed generally slightly worse than the full datasets. We used AUC (area under  
144 the receiver operating characteristic (ROC) curve) as the main model performance metric. AUC  
145 provides a good combination of sensitivity and specificity for comparing Boruta selected  
146 variables to all variables and the different modality datasets in their performance of cLBP  
147 patient delineation.

148

149 The best single modality for cLBP classification was Boruta selected variables from MRI (Fig.  
150 3A) with a mean AUC of 0.645 with a 95% confidence interval (CI) of 0.618 – 0.672 and  
151 accuracy of 0.657 (95% CI, 0.636 – 0.678). The Boruta selected questionnaire dataset modality

152 produced only minimally worse classification performance (mean AUC = 0.631, 95% CI =  
153 0.610 – 0.652) than Boruta reduced and all MRI variables (mean AUC = 0.637, 95% CI = 0.599  
154 – 0.675) using the least amount of variables (mean = 8) across all modality datasets (Fig. 3A-  
155 C). Dual and multi modalities generally performed better than single modality models in cLBP  
156 classification. Questionnaire, clinical physical assessment, and MRI (Q + C + M, Fig. 3C)  
157 modality with Boruta selected variables represents the best performing modality model with a  
158 mean AUC of 0.699 (95% CI, 0.669 – 0.729) and a mean accuracy of 0.709 (95% CI, 0.679 –  
159 0.739). Moreover, this model showed the highest sensitivity (mean = 0.622, 95% CI – 0.568 –  
160 0.676). The modality model showing the highest specificity was using all variables and all  
161 modalities (Q + C + S + M, mean = 0.840, 95% CI = 0.799 – 0.881). The three best dual  
162 modalities models, C + M (mean AUC = 0.679, 95% CI = 0.633 – 0.725), Q + M (mean AUC  
163 = 0.674, 95% CI = 0.623 – 0.725), and Q + C (mean AUC = 0.674, 95% CI = 0.639 – 0.709)  
164 all with Boruta selected variables (Fig. 3B), performed only slightly worse than the best model  
165 (Boruta – Q + C + M). Overall back shape and function dataset using Boruta selected variables  
166 produces the worst classification performance (Fig. 3A, AUC = 0.569, 95% CI = 0.538 – 0.60).  
167 Additionally, the dual modalities continuing back shape and function data always performed  
168 worse than those without (Fig. 3B), when using both all and Boruta selected variables.  
169 Classification performance metrics across all models is provided in Supplementary Table S7.  
170



171

172 **Figure 3. Boruta variable reduction performance.** A – C shows RF classification model performance  
 173 (AUC) following the reduction of variables using Boruta and all variables in the single, dual, and multi  
 174 data modalities respectively. This shows the change in performance follow variable reduction. Error  
 175 bars represent 95% CI in AUC over 10-fold train-test splits. D – Shows the amount of variable reduction  
 176 by using Boruta as a percentage of the total number of variables within each modality dataset.  
 177

## 178 Boruta variable importance

179 Boruta variable importance selection resulted in a median of 62.7% reduction in the number of  
 180 variables across all 15 datasets (Fig. 3D). This large reduction in variables from Boruta  
 181 performed comparable (eight slightly worse and seven better) compared to using all variables,  
 182 with high overlap in confidence intervals. Furthermore, three of the top five performing  
 183 modality models where those employing Boruta selected variables (Supplementary Table S22).  
 184 The greatest performance improvement following Boruta variable selection was found in the Q  
 185 + C + M datasets with an AUC increase of 0.270 and an average reduction of 72.2 variables  
 186 (Fig. 3C). The smallest variable reduction was shown in MRI (30%) with an AUC increase of

187 0.008, while the largest variable reduction was found in the whole dataset (Q + C + S + M) with  
188 a reduction of 75.7% of the variables and an AUC decrease of 0.009.  
189  
190 Boruta selected variables represented the best patient delineation performance in single, dual,  
191 and multi-modality datasets (Fig. 3A-C). MRI was found to be the best single modality model  
192 using Boruta selected variables (Fig. 3A). The most important variables (Supplementary Table  
193 S27) were intervertebral disc (IVD) herniation L4 – L5, spinal canal width L2, IVD  
194 degeneration L3 – L4 and L4 – L5, and spinal canal width L1, showing a mean importance  
195 across the ten iterations of 19.49, 9.94, 9.46, 9.14, and 8.13 respectively. The best dual  
196 modalities cLBP patient stratification model, Boruta selected C + M (Fig. 3B), showed the  
197 second highest AUC (mean = 0.679). The MRI variable IVD herniation L4 – L5 was found to  
198 be the most important (mean = 13.72), with clinical mobility assessments of the hip, cervical  
199 spine, and the whole body showing high importance (Supplementary Table S31). The most  
200 important and robust variables of the best performing model (Boruta – Q + C + M, Fig. 3C)  
201 contained assessments from all three modalities (Supplementary Table S34). The Short-form  
202 36 Health Status Questionnaires (SF-36) psychological well-being, SF-36 social function, and  
203 hip pain presented a mean importance of 13.36, 13.07, and 9.71 respectively. The clinical  
204 assessment, cervical axial rotation (left) and MRI variable IVD herniation L4 – L5 each showed  
205 an importance of 9.47 and 9.12. These represent the top five most important variables and the  
206 rest are provided in Table 3. The questionnaire Boruta reduced model provides a sparse model  
207 with decent performance, meaning the modality performs comparably well with a small amount  
208 of data. On average, the questionnaire Boruta model used 8 variables with a mean AUC  
209 reduction of 0.068 compared to the best model (Q + C + M). The most robust and important  
210 variables (Supplementary Table S39) were SF-36 social function (mean = 21.42), SF-36  
211 psychological well-being (mean = 20.02), hip pain (mean = 15.0), smoking in pack years (mean  
212 = 11.8), and family history of back pain (mean = 7.27). All Boruta selected important variables

213 across the datasets are provided in Supplementary Table S24 – S38 and their correlation  
214 matrices in Supplementary Figure S7 – S21.

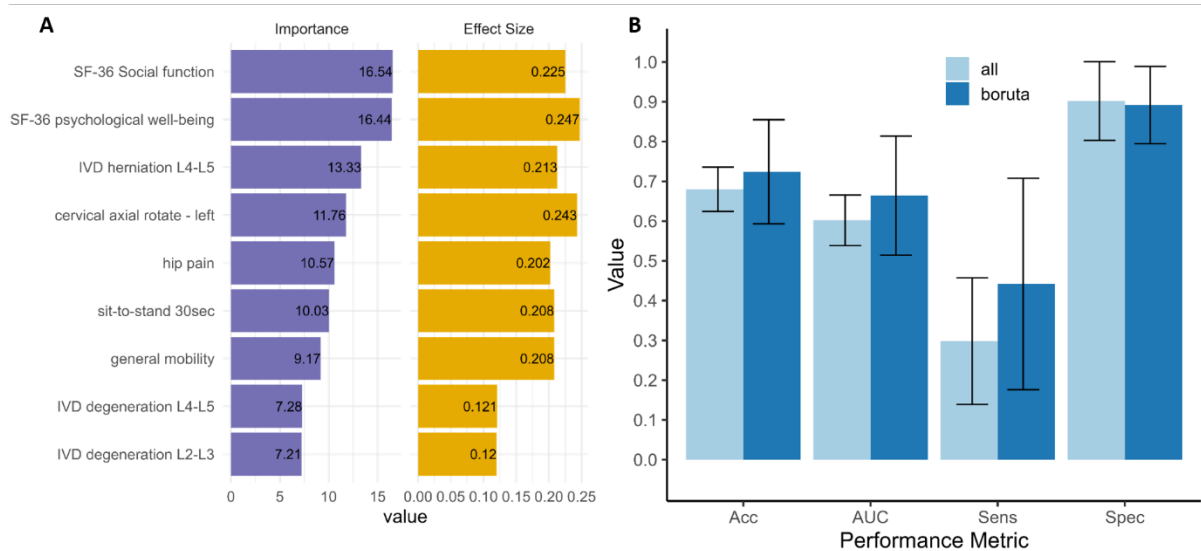
215

## 216 Robust and important variables

217 To select the most robust and important variables for cLBP patient classification, the percentage  
218 of selection and average importance score was calculated across the 15 datasets (Supplementary  
219 Table S39). The variables that were selected at every opportunity (100%) across the multiple  
220 iterations and datasets were defined as most important and solely implemented in cLBP  
221 classification compared to all variables. This resulted in nine robust and important variables  
222 (Fig. 4A). The variables importance score is the Z-score of the mean decrease accuracy measure  
223 and is computed by dividing the average accuracy loss by its standard deviation. Meaning this  
224 value shows how much the models accuracy would decrease without this variable included.  
225 These variables represented psychosocial factors, IVD herniation and degeneration of the lower  
226 lumbar spine, presence of hip pain, as well as mobility of the neck and general mobility of the  
227 whole body. The performance of these nine variables were compared to all variables utilising  
228 the hold-out dataset (Fig. 1B) in a five-fold train-test workflow to provide an unbiased  
229 comparison. The best nine variables showed better mean accuracy (Fig. 4B, Boruta = 0.724,  
230 95% CI = 0.593 – 0.855, All = 0.680, 95% CI = 0.625 – 0.736), AUC (Boruta = 0.664, 95% CI  
231 = 0.514 – 0.814, All = 0.602, 95% CI = 0.538 – 0.666), and sensitivity (Boruta = 0.442, 95%  
232 CI = 0.176 – 0.708, All = 0.298, 0.139 – 0.457), while all variables provided better specificity  
233 (Boruta = 0.892, 95% CI = 0.795 – 0.989, All = 0.902, 95% CI = 0.803 – 1.0). Moreover, all  
234 these nine variables showed significant univariate statistically significant differences between  
235 cLBP patients and asymptomatic controls in the questionnaire, clinical, and MRI datasets (Fig.  
236 4A). Utilising a Wilcoxon-Mann-Whitney test revealed reduced scores of the SF36 for social  
237 function ( $u = 183751$ ,  $z = -7.08$ ,  $p < 0.001$ , effect size  $r = -.225$ ) and psychological well-being

238 (u = 178506.5, z = -7.74, p < 0.001, r = -0.247) in people suffering from cLBP. In addition, the  
 239 occurrence of hip pain was also altered comparing people with and without cLBP ( $\chi^2$  (3, N=986)  
 240 = 40.17, p = 0.004,  $\omega$  = 0.202). Comparing clinical examinations further revealed reduced  
 241 cervical axial rotation to the left (u = 221595, z = -7.98, p < 0.001, r = -0.243), reduced sit to  
 242 stand 30 second repetition (u = 226818.5, z = -6.83, p < 0.001, r = -0.208), as well as altered  
 243 general mobility ( $\chi^2$  (2, N=1080) = 46.81, p = 0.006,  $\omega$  = 0.208) in cLBP patients. Regarding  
 244 MRI investigations, increased IVD degeneration at L2-L3 (u = 147168, z = 3.38, p = 0.011, r =  
 245 0.120) and L4-L5 (u = 147691, z = 3.43, p = 0.010, r = 0.121), as well as increased disc  
 246 herniation at L4-L5 (u = 153470.5, z = 6.02, p < 0.001, r = 0.213) were found in people suffering  
 247 from cLBP compared to asymptomatic controls. Univariate statistical results for all variable  
 248 can be found in Supplementary Table S40 – S47.

249



250

251 **Figure 4. Most robust and important variables for chronic low back pain classification**  
 252 **performance.** A – Presents a bar plot of the nine most robust variables in order of average Boruta  
 253 importance score (left). The right bar plot shows the absolute effect size (Cohen’s r or  $\omega$  depending on  
 254 data type) comparing controls and cLBP patients of the nine robust variables. B – Column plot showing  
 255 RF classification performance as mean of five-fold train-test iterations in hold-out set using Boruta  
 256 selected and all variables. Column plot error bars represent 95% CI. IVD – intervertebral disc, SF-36 –  
 257 short form 36 health status questionnaire, Acc – accuracy, AUC – Area under the receiver operating  
 258 characteristic curve, Sens – sensitivity, Spec – specificity.



## 259 Discussion

260

261 This study employs a large multi-modal dataset and a machine learning workflow to  
262 demonstrate the importance of using data from different domains in cLBP patient delineation.  
263 Increasing the number of modalities generally lead to a model performance improvement  
264 although it seems the inclusion of back shape and motion data resulted in little to no  
265 performance improvement. Utilising Boruta in our iterative selection workflow resulted in  
266 considerable variable reduction across all datasets (median = 62.7%), while model performance  
267 remained comparable. This may reflect many variables showing little difference between  
268 patients and controls, or the underlying mechanisms are more robustly captured by a small sub-  
269 set of variables. Both the best performing modality model (Q + C + M) and the most robust  
270 variables (Fig. 4), show the importance of measuring psychosocial factors, cervical axial  
271 rotation, general mobility, hip flexion, and lower lumbar spine disc degeneration and herniation  
272 ratings in cLBP patients.

273

274 The questionnaires probing the psychosocial factors, social function and psychological well-  
275 being, showed highest importance among all variables (Fig. 4A) and were the most important  
276 variables in the best sparse model (Boruta reduced questionnaire, Fig. 3A). Moreover, the  
277 questionnaires represent the most cost effective modality to the examiner, highlighting the  
278 clinical importance of psychosocial factors in cLBP diagnosis and treatment.

279

280 Social functioning describes the ability of a person to engage in social activities, which we have  
281 shown to be an important marker for delineating cLBP patients from asymptomatic controls.  
282 Using cross-sectional data from 180 chronic low back pain patients, Ge and colleagues<sup>19</sup>  
283 showed that these patients reported more limitations in performing (major life tasks and) social  
284 activities as compared to subjects without cLBP, even after adjusting for influencing factors,



285 such as socio-demographics, lifestyle and number of diseases. Furthermore, Tagliaferri et al. <sup>20</sup>  
286 were able to separate 4156 chronic back pain patients from the UK Biobank dataset into five  
287 sub-groups based on their scores of social isolation and depressive symptoms. Interestingly,  
288 increased social isolation was only a variable of three sub-groups, encompassing 26% of all  
289 back pain patients (n = 1085), while the remaining subgroups showed either no changes (4.1%;  
290 n = 776) or a reduction in social isolation scores (12%; n = 2296). This prevalence in patients  
291 with LBP may indicate that reduced social functioning was identified by some studies, while  
292 others did not find similar changes as compared to asymptomatic controls <sup>21</sup>. However, as levels  
293 of social function (here: social participation) were found to be correlated with self-perceived  
294 physical health status <sup>22</sup>, a direct impact of social functioning on personal functional  
295 impairments remains feasible.

296  
297 In addition to social function, the psychological health or well-being was shown to be an  
298 important variable in cLBP patient delineation. Using longitudinal data from the SwePain  
299 cohort, including 9361 participants with and without chronic pain, psychological well-being  
300 scores at baseline were able to predict pain intensity after 2 years <sup>23</sup>. Within this study, positive  
301 well-being was predictive of lower pain severity in participants without and with chronic pain.  
302 Similar conclusions were drawn from the comparison of back pain patients with different levels  
303 of mental distress, in which patients with higher mental distress showed, among other things,  
304 reduced psychological well-being and social function, and higher severe pain than patients with  
305 lower mental distress <sup>24</sup>. Furthermore, patients suffering from chronic pain exhibit significantly  
306 lower quality of life scores across all sub-domains, including psychological well-being <sup>25</sup>.  
307 Alterations in quality of life are stronger associated with changes in social functioning and  
308 psychological well-being (via pain catastrophizing) than pain intensity itself <sup>26</sup>, indicating the  
309 high importance of psycho-social aspects for daily living with painful conditions such as back  
310 pain.

311 Our findings highlight the potential importance of psychosocial factors in cLBP and suggest a  
312 need for change in clinical practice. In Berlin, and more widely, the integration of psychosocial  
313 assessment into routine cLBP care could improve treatment effectiveness. This is consistent  
314 with previous meta-analytic evidence <sup>27</sup> supporting a combined biopsychosocial therapeutic  
315 approach to optimise treatment efficacy. Early identification of psychosocial barriers such as  
316 anxiety, depression or stress should be prioritised, with interventions such as cognitive  
317 behavioural therapy, stress management and mental health support used alongside traditional  
318 physical therapies.

319  
320 We demonstrated that spinal herniation and degeneration observed on MRI may contribute to  
321 pain mechanisms in the bio-psychosocial model of cLBP, consistent with previous meta-  
322 analyses <sup>28</sup>. However, it is clear that MRI findings alone do not fully explain pain presence in  
323 cLBP, highlighting the need for caution when interpreting MRI results at the individual patient  
324 level. Interestingly, cervical spine rotation but not lumbar back motion assessments were shown  
325 to be robust important examinations for cLBP delineation. This contradictory finding is likely  
326 a result of two factors. First, the poor cervical rotation may be the result of neck pain that has  
327 high comorbidity with cLBP <sup>29</sup> and can lead to decreased axial rotation <sup>30</sup>. Second, poor  
328 psychological health has been associated with neck pain <sup>31</sup>, which we found to also present high  
329 importance in patient stratification and relates to the bio-physical-psychosocial interplay  
330 present in cLBP patients <sup>5</sup>. A systematic review on hip mobility in LBP patients <sup>32</sup> showed small  
331 to no changes in hip flexion compared to controls. As all studies had less than 110 subjects,  
332 they were likely under powered to uncover the decreased mobility we show here and may  
333 represent a diagnostic test for LBP. Previous studies have shown that clinical kinematic data  
334 can effectively stratify cLBP patients into high, low and intermediate risk groups <sup>33</sup>, suggesting  
335 that pain correlates with reduced physical function. Persistent nociceptive input from  
336 aggravated spinal joints/muscles may lead to reduced motor output and spinal cord excitability

337 <sup>34</sup>, potentially resulting in a reduced ability to recruit specific muscles and necessitating  
338 compensatory movement strategies. Our findings underscore that detailed movement analysis  
339 could serve as a diagnostic biomarker for LBP, potentially rivalling medical imaging in  
340 diagnostic accuracy and improving patient care by identifying sub-populations likely to respond  
341 well to specific therapies or at risk of adverse outcomes.

342  
343 Classification of cLBP patients has often been conducted on relatively small sample sizes (<  
344 200) as well as utilising a single data domain <sup>9</sup>. Performance of such models are subject to  
345 overfitting due to their small samples and would likely perform poorly at out-of-sample  
346 classification in unseen external datasets <sup>35</sup>. Furthermore, several studies have established  
347 classification models with high accuracy (> 0.8) at determining particular LBP symptoms <sup>36-43</sup>,  
348 although these lack clinical applicability in understanding the most appropriate variables and  
349 modalities in classification of cLBP as well as the classification of the disorder in general.  
350 Classification models created using large datasets (n > 1000) either contained psychosocial and  
351 demographic variables, without imaging and physical variables <sup>44,45</sup>. On the other hand, Jin-  
352 Heekun and colleagues <sup>46</sup> employed only physical variables without considering important  
353 psychosocial factors, which have shown to be important in previous research <sup>5,20,47,48</sup> as well as  
354 in our current study (Fig. 4). Our best model (Boruta – Q + C + M) performed slightly worse  
355 than Parsaeian et al. <sup>45</sup> (AUC 0.693 – 0.75) and compared to Shim et al. <sup>44</sup> (AUC, 0.693 – 0.716),  
356 we utilised more plentiful data points per subject in a significantly smaller sample size  
357 (approximately 34x and 6x smaller respectively) to address the clinically relevant question of  
358 what modalities and variables a best suited for cLBP classification.

359  
360 The multimodal nature of the data utilised and the amount of subjects that have participated in  
361 physical, imaging, and questionnaire measurements are major strengths in our current study.  
362 Utilising the “Berlin Back Study” dataset that contains more than 500 subjects in the different

363 domains highlighted as lacking in a recent review by our group <sup>9</sup>, enabling us to robustly  
364 investigate the importance of different modalities as well as specific variables in cLBP patient  
365 stratification. The large sample size enabled us to minimise model overfitting through cross-  
366 validation and hold-out testing of good sample size and distributions. A more accurate  
367 representation of a models performance is provided by out-of-sample testing that uses a new  
368 sample population containing comparable variables. This provides a test set with minimised  
369 sampling and dataset bias greatly improving the generalisability and applicability of the  
370 findings.

371  
372 The results of our study offer valuable insights into the potential of using a multimodal machine  
373 learning approach for the classification of individuals with chronic low back pain (cLBP).  
374 However, despite the promising nature of our dataset, the classification performance was  
375 moderate. One possible explanation is the size and diversity of the dataset. Although the sample  
376 size was substantial, it may not fully encompass the complex and multifactorial nature of cLBP,  
377 especially when considering the broad spectrum of psychosocial, environmental, socio-  
378 economic, and biological factors that contribute to the condition. Incorporating detailed  
379 information regarding socio-economic status, education, and other social factors could  
380 potentially improve model performance by capturing the broader context in which cLBP  
381 manifests.

382  
383 While our study highlights the importance of psychosocial factors in the classification of cLBP,  
384 we acknowledge that the analysis of MRI data was limited. Several key MRI phenotypes,  
385 including disc bulging, spondylolisthesis, osteophytes, Modic lesions, endplate abnormalities,  
386 and high-intensity zones, were not included in the final analysis. It is also important to note that  
387 while psychosocial factors showed strong predictive power, the MRI analysis was not  
388 exhaustive, and a more complete MRI dataset would be necessary to make a fair comparison.

389 Future studies should aim to expand both the MRI and psychosocial data to provide a more  
390 comprehensive understanding of the relative contributions of anatomical and psychosocial  
391 factors in cLBP.

392  
393 A notable limitation of this study is the absence of key social determinants of health, such as  
394 income, education, and health insurance status, which are known to significantly influence both  
395 the risk of developing cLBP and its persistence. While we included job-related factors, such as  
396 posture during work, these social determinants were not available in our dataset. Incorporating  
397 these factors into future models could enhance the accuracy of our predictions and offer a more  
398 nuanced understanding of how socioeconomic and environmental factors intersect with clinical  
399 and psychosocial elements in the context of cLBP. The inclusion of such variables should be a  
400 priority in subsequent research to improve both model performance and clinical applicability.

401  
402 We acknowledge the challenge of translating these findings into clinical practice. The  
403 complexity and costs associated with collecting such comprehensive data may not be feasible  
404 in all healthcare settings. However, the value of this research lies in its potential to inform the  
405 development of more personalized, evidence-based treatment strategies. In future work, we aim  
406 to explore how these multimodal data could contribute to treatment recommendations or patient  
407 phenotyping, providing clinicians with more precise tools to tailor interventions to individual  
408 patient needs. By advancing the field toward personalized care, we believe the clinical utility  
409 of this approach will become more apparent. Furthermore, our study provides a cross-sectional  
410 investigation of cLBP, whereas a prospective design would be better suited to examine the  
411 causality of the disorder. Two different types of interviews were conducted: face-to-face  
412 interviews (clinical examination) and electronic interviews (questionnaires). The main  
413 difference is that body language, facial expressions and other non-verbal social cues are obvious  
414 to the interviewer in face-to-face interviews, whereas these aspects are absent in electronic

415 surveys. As both surveys have advantages and disadvantages, the answers of the study  
416 participants were weighted equally in this study. Additionally, the high number of questions  
417 could lead to a reduction in the participants' attention and concentration.

418

419 In conclusion, while our current model shows promise, there is room for improvement.  
420 Expanding the dataset, incorporating more detailed and diverse data sources, and exploring  
421 alternative machine learning models are potential next steps that could enhance the predictive  
422 power of the system. A key future application will be the use of baseline data to predict  
423 treatment response and the transition from acute to chronic pain. This predictive capability  
424 could enable early identification of high-risk patients, improving treatment outcomes and  
425 advancing personalised medicine in the management of cLBP. We also believe that as more  
426 data becomes available and our understanding of the complex interplay of factors contributing  
427 to cLBP deepens, the clinical utility of such models will become more apparent, leading to  
428 better patient stratification and more personalized treatment approaches.

429

## 430 Methods

### 431 Quantitative variables and data collection

432 Patients had to meet all of the following inclusion criteria: written informed consent to  
433 participate in the study, asymptomatic (no back pain) or symptomatic (cLBP) caucasian women  
434 and men aged 18–67 years, pain duration  $\geq 12$  weeks daily (cLBP only), pain localization in the  
435 lumbopelvic region (cLBP only). A telephone interview was conducted during recruiting and  
436 subjects were excluded if they met any exclusion criteria. However, some subjects came to  
437 testing that should have been excluded during the telephone interview, and were then excluded  
438 at the testing site (Supplementary Table S3). No minimal threshold for LBP intensity was  
439 defined. A list of all variables including number of missing values is shown in Supplementary  
440 Table S19.

### 441 Questionnaires

443 The localization, type, course, possible radiation, intensity, quality, duration, and any factors  
444 that may relieve or exacerbate the pain, as well as possible triggers or the patient's own  
445 explanations regarding the cause of the pain has been asked. The patient's medical history has  
446 been recorded, which includes any previous diseases and surgeries, and a detailed pain and  
447 general medication history as well as allergies, intolerances, and vaccination status among  
448 others. A family and social history (anamnesis) was taken (professional activity, family  
449 situation, diseases in the family, stressful situations, etc.). In addition, any past or present use  
450 of addictive substances has been asked (alcohol, nicotine, etc.).

451 The following questionnaires were completed within 30 minutes:

- 452 ➤ Pain intensity, pain duration, and pain-related disability: von Korff et al. <sup>49</sup>
- 453 ➤ Disability Questionnaire: Roland and Morris (RMDQ) <sup>50</sup>

454 ➤ Short-form 36 Health Status Questionnaire: SF-36<sup>51</sup>. The following four domains were  
455 considered: general mental health (psychological distress and well-being), limitations in  
456 usual role activities because of emotional problems, vitality (energy and fatigue), and  
457 general health perceptions.

458 ➤ International Physical Activity Questionnaire (IPAQ)<sup>52</sup>

459 ➤ Self-Report Behavioural Automaticity Index (SRBAI)<sup>53</sup>

460 ➤ Behavioural Regulation in Sport Questionnaire (BRSQ)<sup>54</sup>

461 ➤ Tampa Scale for Kinesiophobia (TSK-GV)<sup>55</sup>.

462 ➤ Fear-Avoidance Belief Questionnaire (FABQ)<sup>56</sup>.

463 The participants primarily answered the questionnaires in digital form using a survey program  
464 specially developed for the study. The data were collected under similar conditions (e.g., same  
465 room, same computer) for all subjects at the study centre.

466

## 467 Demographic data

468 During the clinical assessment, age, sex, body height, body weight, hip diameter, and waist  
469 diameter of the subjects were recorded. BMI was chosen instead of waist hip ratio to measure  
470 physical body size and health, as the BMI variable contained less missing values compared to  
471 waist hip ratio (Supplementary Table S19).

472

## 473 Clinical examination

474 The clinical examination included the evaluation of organ functions (inspection, palpation,  
475 percussion, and auscultation), the general impression, and the vital parameters of the patient  
476 (temperature, heart rate, blood pressure, etc.). Examination was performed by an experienced  
477 orthopaedic consultant. The neurological status was assessed by the examination of the  
478 coordination, reflexes, sensitivity, and motor function. The evaluation of the functional



479 parameters, that is, the assessment of posture, shape, orientation, and movement of the lumbar  
480 spine and pelvis, was based on current clinical standards (e.g., Ott and Schober test, 3-step  
481 hyperextension test, passive lumbar extension test, etc.) and self-assessment by the persons  
482 investigated. Data were documented according to their dimension using distances in cm,  
483 degrees of angle, and number of repetitions per defined time interval or bivalent whether pain  
484 provocation occurred. Self-assessment of functional restrictions of the back was recorded  
485 according to a scale from 1 (best) – 10 (worst).

486

## 487 Back shape and function

488 All study participants received measurements of the back shape in the sagittal and frontal planes  
489 during upright standing and sitting using the Idiag M360 (Idiag AG, Fehraltorf, Switzerland).  
490 The device measures segmental angles of the thoracic and lumbar spine. In both postures, study  
491 participants were measured upright, in flexed, extended, and in left and right lateral bending (3  
492 repetitions, ~10 sec each). Maximum upper body flexion, extension, as well as left and right  
493 lateral bending were performed with extended knees. During extension the arms were crossed  
494 in front of the body. The order of performed tasks was randomised. The measurements were  
495 performed by trained medical students. The validity and reliability were demonstrated in  
496 previous studies<sup>57-60</sup>.

497

## 498 Spino-pelvic MRI

499 MRIs were conducted using a 1.5 MRI scanner. Following sequences were evaluated: 1) Sag  
500 T1 (4 mm slices), 2) Sag T2 (4 mm slices), 3) Cor STIR-T2 (4 mm slices) and 4) Axial T2  
501 (3 mm slices). MRIs were evaluated for intervertebral disc degeneration (Pfirrmann  
502 classification<sup>61</sup>), disc herniation (Kramer classification<sup>62</sup>), facet joint arthrosis (Fujiwara  
503 classification<sup>63</sup>), osteochondrosis intervertebralis<sup>64</sup>, spondylolisthesis (Meyerding  
504 classification<sup>65</sup> and spinal canal stenosis (Schizas classification<sup>66</sup>) at each level of the lumbar

505 spine. The spino-pelvic MRI evaluation was performed blinded by two spine surgeons and a  
506 radiologist, all of whom have many years of experience in the evaluation of spinal pathologies.  
507 The inter-rater reliability was good-to-excellent for all measurement parameters.

508

## 509 Data storage

510 All data files electronically recorded during the study period were stored on a database server  
511 folder (SharePoint folder) hosted by Charite-Universitaetsklinikum. A data back-up for the  
512 database is run daily. Local study team members signed a non-disclosure agreement. They have  
513 access to the database using a personal password and are authorized only for entries depending  
514 on their function based on a role concept (investigator, statistician, monitor, administrator etc.).  
515 A multilevel data validation plan was developed to guarantee the correctness and consistency  
516 of the data. Data were entered only after a check for completeness and plausibility. Furthermore,  
517 data were cross-checked for plausibility with previously entered data for each participant.  
518 Questionnaires filled out on paper are stored in a lockable cabinet at the university.

519

## 520 Potential sources of bias and minimisation

521 To generally reduce the possible location and assessor bias, measurements (clinical physical  
522 and questionnaire assessment and back shape and function) were administered by a few trained  
523 clinicians in the same room with the same lighting. The self-administered questionnaires (von  
524 Korff, RMDQ, SF-36, IPAQ, SRBAI, BRSQ, TSK-GV, and FABQ) were completed by the  
525 subjects under supervision by our trained study coordinator who provided explanations for  
526 unclear questions and mitigated possible lack of motivation to complete the questions by  
527 assuring the subjects of the importance to complete the questionnaires. Furthermore, generic  
528 questionnaires (SF-36 or IPAQ) were placed before specific ones (SRBAI, BRSQ) to minimize  
529 bias from order effects.

530

531 To minimise the bias in our classification and variable selection modelling, as well as the  
532 modality comparison, variables directly assessing back pain and questions heavily biased to  
533 pain patients were removed. Such variables assessed back pain during particular movements or  
534 upon physical manipulation. Furthermore, questions only back pain patients were asked for  
535 example, pain intensity, duration, and disability, as well as pain and health biased self-  
536 questionnaires (von Korff, RMDQ, TSK-GV, FABQ, therapies, and the SF-36 sub-categories  
537 regarding physical function, physical role function, physical pain, health perception, and  
538 vitality), and clinician administered questions regarding pain medication intake, previous spinal  
539 disorder diagnosis, and participants' subjective physical health assessment were removed as  
540 well to reduce model bias.

541

## 542 Outcome

543 The outcome target for our classification model is cLBP patients. All participants were assessed  
544 by a clinician and diagnosed as cLBP patient, asymptomatic control, or suffered from LBP in  
545 the past but not at present. We used the clinician diagnosis of either current cLBP patient or  
546 asymptomatic control as our two-class target outcome. The participants with LBP in the past  
547 were removed to enable better distinguishable groups for binary classification and variable  
548 importance selection.

549

## 550 Data handling, preprocessing, cleaning and missing data

551 Total and sub-scores for the self-administered questionnaires were used after removal of pain  
552 biased questionnaires (see *Potential sources of bias and minimisation*). This includes the SF-  
553 36<sup>51</sup>, the IPAQ<sup>52</sup>, the SRBAI<sup>53</sup>, and the BRSQ<sup>54</sup>. The SF-36 was used to collect statements  
554 related to the health domains 'emotional role limitation, (three items) 'social functioning' (two

555 items), and ‘mental health’ (five items). A scoring algorithm was used to convert the raw scores  
556 into these three domains. The scores were transformed to range from zero (worst possible  
557 health) and 100 (best possible health). The SRBAI and BRSQ collected ratings for multiple  
558 statements on a numerical scale from 1 (strongly disagree) to 6 (strongly agree). For the SRBAI,  
559 the total score was calculated by summing the numerical values across all 4 statements. In  
560 contrast, sub-scores were created for the BRSQ that related to ‘intrinsic motivation’, ‘integrated  
561 regulation’ and ‘external regulation’. These scores were calculated by summing the numerical  
562 values across two statements per sub-score. The IPAQ, recorded the average time spent per day  
563 over the past 7 days while ‘sitting’, ‘walking’, doing ‘moderate activities’ (e.g., heavy lifting,  
564 digging, aerobics, or fast bicycling), and ‘vigorous activities’ (e.g. carrying light loads,  
565 bicycling at a regular pace, or doubles tennis). To estimate the energy requirements for each  
566 activity type, the average time spent per day in minutes was multiplied by MET-score  
567 (metabolic equivalents) of 1.5, 3.3, 4, or 8 for sitting, walking, moderate and vigorous activities,  
568 respectively. This resulted in MET-minutes scores, describing the amount of energy in  
569 kilocalories required for a 60 kilogram person. Finally, the MET-minutes scores were summed  
570 across ‘sitting’, ‘walking’, ‘moderate activities’ and ‘vigorous activities’, resulting in a total  
571 MET-minutes score per subject.

572  
573 Modelling preprocessing was conducted by checking variables within each modality for very  
574 low variance and collinearity. Variables with close to zero (variance  $< 1$ ) variance were  
575 removed. Furthermore, variables presenting a Spearman correlation greater than 0.9 were also  
576 removed to reduce collinearity between variables. The decision on which of the correlated  
577 variables to removed, was the variable showing less correlation to the target, cLBP patient  
578 status. Following these cleaning steps, single data modalities were joined with demographic  
579 data and subsequently joined with other datasets for the dual- and multi-modalities datasets.  
580 Subjects having any missing values were removed from our main analyses and MissForest <sup>18</sup>

581 imputation was additionally conducted. Imputation was computed for the training and testing  
582 samples within each fold independently to avoid data leakage. Following cleaning and  
583 preprocessing modality dataset presented different number of subjects and variables with a  
584 similar age and sex distribution across cLBP patients and asymptomatic controls (Fig. 1A).

585

## 586 Univariate statistics

587 The univariate statistics were carried out separately for continuous, ordinal and nominal data to  
588 compare patients suffering from cLBP against asymptomatic controls using R (version 4.3.1;  
589 [www.r-project.org](http://www.r-project.org)). As most continuous variables did not follow a normal distribution  
590 according to the Anderson-Darling test <sup>67</sup>, we implemented the non-parametric Wilcoxon-  
591 Mann-Whitney test to determine significant difference between cLBP patients and  
592 asymptomatic controls for ordinal and continuous data. Hence, u-values, z-values, r-value  
593 (effect sizes), as well as the p-value are reported from. Nominal data were compared using the  
594 Chi-Square test and reported with  $\chi^2$  values, Cohen's  $\omega$ -values (effect size), and p-values.  
595 Statistical significance was determined at  $p \leq 0.05$  following family wise error (FWE) <sup>68</sup>  
596 correction for multiple comparisons within each modality (demographics, questionnaires,  
597 clinical examinations, superficial spine morphology and motion, and spino-pelvic MRI).

598

## 599 Machine learning

### 600 Boruta variable selection

601 We used the Boruta method <sup>16</sup> for importance variable selection. Boruta utilises random forest  
602 (RF) classification algorithm <sup>17</sup> with both the real variables and set of 'dummy' or 'shadow'  
603 variables, that are created by shuffling the variable values. This creates random variables

604 (dummy features) that have the same distribution as the original variables, although represent  
605 the classification accuracy of this variable randomly sampled. As these dummy variables  
606 represent random noise, they had their possible correlation to the target (cLBP) removed. All  
607 real and dummy variables are used to classify cLBP and variable importance is calculated.  
608 Variable importance is calculated as the Z-score of the mean decrease in classification accuracy  
609 following the removal of this variable from the model. The importance of the dummy variables  
610 can be used as a reference to test the variable importance of the real variables. Through an  
611 iterative process, real variables that have significantly greater importance than the maximum  
612 dummy variable importance are marked as important. The variables that have a significantly  
613 lower importance than the maximum dummy variables are deemed unimportant and removed  
614 for the next iteration of selection. Iterations are repeated until the importance is assigned to all  
615 variables or a user defined iteration number is reached. We used a max of 2000 iterations with  
616 a random forest containing 1000 trees. As there are occasionally a few variables not definitively  
617 identified as important or removed after 2000 iterations by Boruta, we only selected the  
618 variables that have been confirmed as important in our subsequent modelling and analyses.

619

## 620 Random forest classification algorithm

621 We utilised RF <sup>17</sup> implemented using the ranger package <sup>69</sup> in R (version 4.3.1, [www.r-](http://www.r-project.org)  
622 [project.org](http://www.r-project.org)) to classify cLBP patients and pain-free controls. As RF provides variable  
623 importance measures and can handle categorical, ordinal, and continuous data it represents the  
624 ideal choice to deal with the different data types present in the Berlin Back dataset. Ten-fold  
625 cross validation was conducted during model training and hyper-parameter tuning utilising  
626 1000 tree RF. We conducted hyper-parameter tuning using a tuning parameter search grid,  
627 which contained number of variables to be sampled at each split (mtry) of 1 – square root of

628 the number of variables, and a minimum node size of 5 and 10. Therefore, using a gini split rule  
629 a grid search was conducted with all combination of hyper-parameters to determine the best. A  
630 ten-fold train-test loop (Fig. 1B) was conducted to determine the cLBP classification  
631 performance within each of the 15 datasets utilising all and Boruta selected variables  
632 independently. Model performance was calculated as follows; accuracy =  $(TP + TN) / (TP +$   
633  $FP + TN + FN)$ , sensitivity =  $TP / (TP + FN)$ , specificity =  $TN / (TN + FP)$ , and AUC (area  
634 under the receiver operating characteristic (ROC) curve). Whereby, TP – true positive, TN –  
635 true negative, FP – false positive, and FN – false negative. The ROC curve represents the  
636 classification performance measured by sensitivity and specificity over a range [0, 1] of  
637 classification thresholds.

638

## 639 Robust variable selection workflow

640 Selecting important variables and comparing the 15 different dataset modalities was conducted  
641 on 90% (n = 1045, 19 – 72 y/o, mean age =  $41.71 \pm 12.28$ , cLBP = 469) of the preprocessed  
642 Berlin Back study dataset, with a 10% (n = 116, 20 – 64 y/o, mean age =  $43.81 \pm 12.12$ , cLBP  
643 = 43) hold-out set used to evaluate the most robust and important variables (Fig. 1B). The hold-  
644 out sample comprised of participants that had data for all variables. Each 15 modality datasets  
645 went through the iterative ten-fold train-test loop, where Boruta variable selection was  
646 conducted on the training set. Followed by RF training using all and Boruta selected variables  
647 implemented on the training dataset. Model performance was calculated using the test set.  
648 Across the ten loops, the percentage a variable was selected and average importance value was  
649 calculated across all 15 datasets. Finally, the variables that were selection within every train-  
650 test loop and within every possible dataset were provided as the most robust and important  
651 variables for cLBP classification. The demographic variables can be selected a maximum of 15

652 times, while the modality-specific variables a maximum of 8 times across the 15 dataset  
653 modalities. These variables were then used in a five-fold train test loop on the hold-out dataset  
654 and compared to a model trained using all variables.



## 655 **Disclosures**

656 **Funding:** This study is part of the Research Unit FOR 5177 funded by the German Research  
657 Foundation (DFG), Hendrik Schmidt: SCHM 2572/11-1, SCHM 2572/12-1, SCHM 2572/13-  
658 1; Sandra Reitmaeier: RE 4292/3-1, Matthias Pumberger: PU762/1-1. The analyses and  
659 contribution from the Hochschule für Gesundheit were funded, in part, by grant number  
660 50WK2273A (to DLB) from the German AeroSpace Center (DLR).

661 **Conflicts of interest:** All authors declare no conflict of interests.

662

## 663 **Acknowledgments**

664 We would like to thank all patients and healthy participants for their selfless participation in  
665 this study and the participating companies for informing their employees about this study.

666

## 667 **Data and code availability**

668 All results in this study are provided in the (Supplementary) tables. The Berlin Back study is  
669 currently ongoing (end date 31/12/2025) and therefore the raw data used in this manuscript  
670 cannot be provided. The raw data will be openly released from the Berlin Back Study as per  
671 agreement with the funding agency following the completion of the data acquisition. A link to  
672 the raw data will be provided on the Github repository where the analysis code is located  
673 ([https://github.com/viko18/BerlinBack\\_FeatImp/](https://github.com/viko18/BerlinBack_FeatImp/)) when it is made available.

674

## 675 **Author Contributions**

676 **CRedit Contributions:**

677 Conceptualization – SV, HS, DLB

678 Methodology – SV, HS, DLB, MA, CK

- 679 Software – SV, FJ
- 680 Validation – SV
- 681 Formal Analysis – SV, FJ
- 682 Investigation – SV, LAB, NT, MP, SR
- 683 Resources – HS, DLB
- 684 Data Curation – SV FJ
- 685 Writing original draft – SV, RD, FJ, DLB, HS
- 686 Writing review & editing – All co-authors
- 687 Visualizations – SV, FJ
- 688 Supervision – DLB, HS, MA, CK
- 689 Project Administration – HS, SV
- 690 Funding Acquisition – HS, DLB
- 691

## 692 **References**

- 693
- 694 1. Andersson, G. B. Epidemiologic aspects on low-back pain in industry. *Spine (Phila Pa*
- 695 *1976)* **6**, 53–60 (1981).
- 696 2. Meucci, R. D., Fassa, A. G. & Faria, N. M. X. Prevalence of chronic low back pain:
- 697 systematic review. *Rev Saude Publica* **49**, 1 (2015).
- 698 3. Pastorino, R. *et al.* Benefits and challenges of Big Data in healthcare: an overview of the
- 699 European initiatives. *Eur J Public Health* **29**, 23–27 (2019).
- 700 4. Shilo, S., Rossman, H. & Segal, E. Axes of a revolution: challenges and promises of big
- 701 data in healthcare. *Nat Med* **26**, 29–38 (2020).
- 702 5. Tagliaferri, S. D. *et al.* Relative contributions of the nervous system, spinal tissue and
- 703 psychosocial health to non-specific low back pain: Multivariate meta-analysis. *Eur J Pain*
- 704 (2021) doi:10.1002/ejp.1883.
- 705 6. Grotle, M. *et al.* Lumbar spine surgery across 15 years: trends, complications and
- 706 reoperations in a longitudinal observational study from Norway. *BMJ Open* **9**, e028743
- 707 (2019).
- 708 7. Lee, W. *et al.* Identifying and Assessing Interesting Subgroups in a Heterogeneous
- 709 Population. *Biomed Res Int* **2015**, 462549 (2015).
- 710 8. Lötsch, J. & Ultsch, A. Machine learning in pain research. *Pain* **159**, 623–630 (2018).
- 711 9. Tagliaferri, S. D. *et al.* Artificial intelligence to improve back pain outcomes and lessons
- 712 learnt from clinical classification approaches: three systematic reviews. *npj Digital*
- 713 *Medicine* **3**, 93 (2020).
- 714 10. Noroozi, Z., Orooji, A. & Erfannia, L. Analyzing the impact of feature selection methods
- 715 on machine learning algorithms for heart disease prediction. *Sci Rep* **13**, 22588 (2023).
- 716 11. Mwangi, B., Tian, T. S. & Soares, J. C. A review of feature reduction techniques in
- 717 neuroimaging. *Neuroinformatics* **12**, 229–244 (2014).

- 718 12. World Medical Association. World Medical Association Declaration of Helsinki: ethical  
719 principles for medical research involving human subjects. *JAMA* **310**, 2191–2194 (2013).
- 720 13. von Elm, E. *et al.* The Strengthening the Reporting of Observational Studies in  
721 Epidemiology (STROBE) Statement: Guidelines for reporting observational studies.  
722 *Preventive Medicine* **45**, 247–251 (2007).
- 723 14. Collins, G. S., Reitsma, J. B., Altman, D. G. & Moons, K. G. M. Transparent Reporting of  
724 a multivariable prediction model for Individual Prognosis or Diagnosis (TRIPOD): the  
725 TRIPOD statement. *Ann Intern Med* **162**, 55–63 (2015).
- 726 15. Shadbahr, T. *et al.* The impact of imputation quality on machine learning classifiers for  
727 datasets with missing values. *Commun Med* **3**, 1–15 (2023).
- 728 16. Kursa, M. B. & Rudnicki, W. R. Feature Selection with the Boruta Package. *Journal of*  
729 *Statistical Software* **36**, 1–13 (2010).
- 730 17. Breiman, L. Random Forests. *Machine Learning* **45**, 5–32 (2001).
- 731 18. Stekhoven, D. J. & Bühlmann, P. MissForest—non-parametric missing value imputation  
732 for mixed-type data. *Bioinformatics* **28**, 112–118 (2012).
- 733 19. Ge, L., Pereira, M. J., Yap, C. W. & Heng, B. H. Chronic low back pain and its impact on  
734 physical function, mental health, and health-related quality of life: a cross-sectional study  
735 in Singapore. *Sci Rep* **12**, 20040 (2022).
- 736 20. Tagliaferri, S. D. *et al.* Chronic back pain sub-grouped via psychosocial, brain and  
737 physical factors using machine learning. *Sci Rep* **12**, 15194 (2022).
- 738 21. Igti, A. M., Guimarães, M. & Barros, M. B. A. Health-related quality of life (SF-36) in  
739 back pain: a population-based study, Campinas, São Paulo State, Brazil. *Cad Saude*  
740 *Publica* **37**, e00206019 (2021).
- 741 22. Takeyachi, Y. *et al.* Correlation of low back pain with functional status, general health  
742 perception, social participation, subjective happiness, and patient satisfaction. *Spine*  
743 (*Phila Pa 1976*) **28**, 1461–1466; discussion 1467 (2003).

- 744 23. Larsson, B., Dragioti, E., Gerdle, B. & Björk, J. Positive psychological well-being  
745 predicts lower severe pain in the general population: a 2-year follow-up study of the  
746 SwePain cohort. *Ann Gen Psychiatry* **18**, 8 (2019).
- 747 24. Hnatešen, D. *et al.* Quality of Life and Mental Distress in Patients with Chronic Low  
748 Back Pain: A Cross-Sectional Study. *Int J Environ Res Public Health* **19**, 10657 (2022).
- 749 25. Hadi, M. A., McHugh, G. A. & Closs, S. J. Impact of Chronic Pain on Patients' Quality of  
750 Life: A Comparative Mixed-Methods Study. *J Patient Exp* **6**, 133–141 (2019).
- 751 26. Lamé, I. E., Peters, M. L., Vlaeyen, J. W. S., Kleef, M. v & Patijn, J. Quality of life in  
752 chronic pain is more associated with beliefs about pain, than with pain intensity. *Eur J*  
753 *Pain* **9**, 15–24 (2005).
- 754 27. Kamper, S. J. *et al.* Multidisciplinary biopsychosocial rehabilitation for chronic low back  
755 pain: Cochrane systematic review and meta-analysis. *BMJ* **350**, h444 (2015).
- 756 28. Brinjikji, W. *et al.* MRI Findings of Disc Degeneration are More Prevalent in Adults with  
757 Low Back Pain than in Asymptomatic Controls: A Systematic Review and Meta-  
758 Analysis. *AJNR Am J Neuroradiol* **36**, 2394–2399 (2015).
- 759 29. von der Lippe, E. *et al.* Prevalence of back and neck pain in Germany. Results from the  
760 BURDEN 2020 Burden of Disease Study. *J Health Monit* **6**, 2–14 (2021).
- 761 30. Rampazo, É. P. *et al.* Sensory, Motor, and Psychosocial Characteristics of Individuals  
762 With Chronic Neck Pain: A Case Control Study. *Physical Therapy* **101**, pzab104 (2021).
- 763 31. Mansfield, M. *et al.* The association between psychosocial factors and mental health  
764 symptoms in cervical spine pain with or without radiculopathy on health outcomes: a  
765 systematic review. *BMC Musculoskeletal Disorders* **24**, 235 (2023).
- 766 32. Avman, M. A., Osmotherly, P. G., Snodgrass, S. & Rivett, D. A. Is there an association  
767 between hip range of motion and nonspecific low back pain? A systematic review.  
768 *Musculoskeletal Science and Practice* **42**, 38–51 (2019).

- 769 33. Abdollahi, M. *et al.* Using a Motion Sensor to Categorize Nonspecific Low Back Pain  
770 Patients: A Machine Learning Approach. *Sensors* **20**, 3600 (2020).
- 771 34. Nijs, J. *et al.* Nociception affects motor output: a review on sensory-motor interaction  
772 with focus on clinical implications. *Clin J Pain* **28**, 175–181 (2012).
- 773 35. Rajput, D., Wang, W.-J. & Chen, C.-C. Evaluation of a decided sample size in machine  
774 learning applications. *BMC Bioinformatics* **24**, 48 (2023).
- 775 36. Al Imran, A., Rifat, M. R. I. & Mohammad, R. Enhancing the classification performance  
776 of lower back pain symptoms using genetic algorithm-based feature selection. in 455–469  
777 (Springer, 2020).
- 778 37. Abdullah, A. A., Yaakob, A. & Ibrahim, Z. Prediction of Spinal Abnormalities Using  
779 Machine Learning Techniques. in 1–6 (IEEE, 2018).
- 780 38. Riveros, N. A. M., Espitia, B. A. C. & Pico, L. E. A. Comparison between K-means and  
781 self-organizing maps algorithms used for diagnosis spinal column patients. *Informatics in*  
782 *Medicine Unlocked* **16**, 100206 (2019).
- 783 39. Sandag, G. A., Tedry, N. E. & Lolong, S. Classification of lower back pain using K-  
784 Nearest Neighbor algorithm. in 1–5 (IEEE, 2018).
- 785 40. Karabulut, E. M. & Ibrikci, T. Effective automated prediction of vertebral column  
786 pathologies based on logistic model tree with SMOTE preprocessing. *Journal of Medical*  
787 *Systems* **38**, 50 (2014).
- 788 41. Mathew, B., Norris, D., Hendry, D. & Waddell, G. Artificial intelligence in the diagnosis  
789 of low-back pain and sciatica. *Spine* **13**, 168–172 (1988).
- 790 42. Vaughn, M. L., Cavill, S. J., Taylor, S. J., Foy, M. A. & Fogg, A. J. Direct explanations  
791 for the development and use of a multi-layer perceptron network that classifies low-back-  
792 pain patients. *International Journal of Neural Systems* **11**, 335–347 (2001).
- 793 43. Zhang, W. *et al.* Deep learning-based detection and classification of lumbar disc  
794 herniation on magnetic resonance images. *JOR SPINE* **6**, e1276 (2023).

- 795 44. Shim, J.-G. *et al.* Machine Learning Approaches to Predict Chronic Lower Back Pain in  
796 People Aged over 50 Years. *Medicina* **57**, 1230 (2021).
- 797 45. Parsaeian, M., Mohammad, K., Mahmoudi, M. & Zeraati, H. Comparison of logistic  
798 regression and artificial neural network in low back pain prediction: second national  
799 health survey. *Iranian Journal of Public Health* **41**, 86 (2012).
- 800 46. Jin-Heeku. Analysis of sitting posture using wearable sensor data and support vector  
801 machine model. *Medico-Legal Update* **1**, 334–338 (2018).
- 802 47. Tagliaferri, S. D. *et al.* Brain structure, psychosocial, and physical health in acute and  
803 chronic back pain: a UKBioBank study. *Pain* **163**, 1277–1290 (2022).
- 804 48. Tagliaferri, S. D. *et al.* Towards data-driven biopsychosocial classification of non-specific  
805 chronic low back pain: a pilot study. *Sci Rep* **13**, 13112 (2023).
- 806 49. Von Korff, M., Ormel, J., Keefe, F. J. & Dworkin, S. F. Grading the severity of chronic  
807 pain. *PAIN* **50**, 133 (1992).
- 808 50. Roland, M. & Fairbank, J. The Roland–Morris Disability Questionnaire and the Oswestry  
809 Disability Questionnaire. *Spine* **25**, 3115 (2000).
- 810 51. Ware, J. E. & Sherbourne, C. D. The MOS 36-item short-form health survey (SF-36). I.  
811 Conceptual framework and item selection. *Med Care* **30**, 473–483 (1992).
- 812 52. Craig, C. L. *et al.* International Physical Activity Questionnaire: 12-Country Reliability  
813 and Validity. *Medicine & Science in Sports & Exercise* **35**, 1381 (2003).
- 814 53. Gardner, B., Abraham, C., Lally, P. & de Bruijn, G.-J. Towards parsimony in habit  
815 measurement: Testing the convergent and predictive validity of an automaticity subscale  
816 of the Self-Report Habit Index. *International Journal of Behavioral Nutrition and*  
817 *Physical Activity* **9**, 102 (2012).
- 818 54. Lonsdale, C., Hodge, K. & Rose, E. A. The behavioral regulation in sport questionnaire  
819 (BRSQ): Instrument development and initial validity evidence. *Journal of Sport &*  
820 *Exercise Psychology* **30**, 323–355 (2008).

- 821 55. Rusu, A. C., Kreddig, N., Hallner, D., Hülsebusch, J. & Hasenbring, M. I. Fear of  
822 movement/(Re)injury in low back pain: confirmatory validation of a German version of  
823 the Tampa Scale for Kinesiophobia. *BMC Musculoskelet Disord* **15**, 280 (2014).
- 824 56. Waddell, G., Newton, M., Henderson, I., Somerville, D. & Main, C. J. A Fear-Avoidance  
825 Beliefs Questionnaire (FABQ) and the role of fear-avoidance beliefs in chronic low back  
826 pain and disability. *PAIN* **52**, 157 (1993).
- 827 57. Dreischarf, B. *et al.* Comparison of three validated systems to analyse spinal shape and  
828 motion. *Sci Rep* **12**, 10222 (2022).
- 829 58. Guermazi, M. *et al.* [Validity and reliability of Spinal Mouse to assess lumbar flexion].  
830 *Ann Readapt Med Phys* **49**, 172–177 (2006).
- 831 59. Topalidou, A., Tzagarakis, G., Souvatzis, X., Kontakis, G. & Katonis, P. Evaluation of  
832 the reliability of a new non-invasive method for assessing the functionality and mobility  
833 of the spine. *Acta Bioeng Biomech* **16**, 117–124 (2014).
- 834 60. Barrett, E., McCreesh, K. & Lewis, J. Reliability and validity of non-radiographic  
835 methods of thoracic kyphosis measurement: a systematic review. *Man Ther* **19**, 10–17  
836 (2014).
- 837 61. Pfirrmann, C. W., Metzdorf, A., Zanetti, M., Hodler, J. & Boos, N. Magnetic resonance  
838 classification of lumbar intervertebral disc degeneration. *Spine* **26**, 1873–8 (2001).
- 839 62. Kraemer, J. Natural course and prognosis of intervertebral disc diseases. International  
840 Society for the Study of the Lumbar Spine Seattle, Washington, June 1994. *Spine (Phila*  
841 *Pa 1976)* **20**, 635–639 (1995).
- 842 63. Fujiwara, A. *et al.* The relationship between facet joint osteoarthritis and disc  
843 degeneration of the lumbar spine: an MRI study. *Eur Spine J* **8**, 396–401 (1999).
- 844 64. Modic, M. T., Steinberg, P. M., Ross, J. S., Masaryk, T. J. & Carter, J. R. Degenerative  
845 disk disease: assessment of changes in vertebral body marrow with MR imaging.  
846 *Radiology* **166**, 193–199 (1988).



- 847 65. Meyerding, H. W. Spondyloptosis. *Surgery, Gynecology & Obstetrics* 371–377 (1932).
- 848 66. Schizas, C. *et al.* Qualitative grading of severity of lumbar spinal stenosis based on the  
849 morphology of the dural sac on magnetic resonance images. *Spine (Phila Pa 1976)* **35**,  
850 1919–1924 (2010).
- 851 67. Anderson, T. W. & Darling, D. A. Asymptotic Theory of Certain ‘Goodness of Fit’  
852 Criteria Based on Stochastic Processes. *The Annals of Mathematical Statistics* **23**, 193–  
853 212 (1952).
- 854 68. Holm, S. A Simple Sequentially Rejective Multiple Test Procedure. *Scandinavian Journal*  
855 *of Statistics* **6**, 65–70 (1979).
- 856 69. Wright, M. N. & Ziegler, A. ranger: A Fast Implementation of Random Forests for High  
857 Dimensional Data in C++ and R. *Journal of Statistical Software* **77**, 1–17 (2017).
- 858