

Deep learning-enabled accurate assessment of gait impairments in Parkinson's disease using smartphone videos

Received: 17 July 2025

Accepted: 2 November 2025

Cite this article as: Han, J., Tian, Z., Wu, J. *et al.* Deep learning-enabled accurate assessment of gait impairments in Parkinson's disease using smartphone videos. *npj Digit. Med.* (2025). <https://doi.org/10.1038/s41746-025-02150-8>

Jianda Han, Zihua Tian, Jialing Wu, Kai Zhang, Shaohua Li, Fahd Baig, Peipei Liu, Ravi Vaidyanathan, Francesca Morgante & Weiguang Huo

We are providing an unedited version of this manuscript to give early access to its findings. Before final publication, the manuscript will undergo further editing. Please note there may be errors present which affect the content, and all legal disclaimers apply.

If this paper is publishing under a Transparent Peer Review model then Peer Review reports will publish with the final article.

Deep learning-enabled accurate assessment of gait impairments in Parkinson's disease using smartphone videos

Jianda Han^{1,2}, Zhihua Tian¹, Jialing Wu³, Kai Zhang¹, Shaohua Li¹, Fahd Baig⁴,
Peipei Liu³, Ravi Vaidyanathan^{5,6}, Francesca Morgante⁴, Weiguang Huo^{1,2,*}

¹College of Artificial Intelligence, Nankai University, Tianjin, China

²Institute of Intelligence Technology and Robotic Systems, Shenzhen Research Institute,
Nankai University, Shenzhen, China

³Department of Neurology, Tianjin Huanhu Hospital, Tianjin, China

⁴Neurosciences and Cell Biology Institute, Neuromodulation and Motor Control Section, City St George's
University of London, London, UK

⁵Department of Mechanical Engineering, Imperial College London, London, UK

⁶Dementia Research Institute Care Research & Technology Centre, Imperial College London, London, UK

* E-mail: weiguang.huo@nankai.edu.cn

Abstract

Gait impairments are among the most prevalent and disabling symptoms in Parkinson's Disease (PD), featuring complex and highly heterogeneous manifestations. Here, we propose a deep learning-based framework to assess gait impairments using smartphone-recorded videos. This framework demonstrated high proficiency in predicting PD severity, with a micro-average area under the receiver operating characteristic curve (AUC) of 0.87 and an F1 score of 0.806, comparable to the average performance of three clinical specialists. Additionally, it effectively discerned the comprehensive efficacy of medications on gait impairments with a precision of 73.68%. In particular, it demonstrated the ability to discriminate medication-induced fine-granular gait changes beyond the resolution of the Unified Parkinson's Disease Rating Scale (UPDRS). Furthermore, our interpretable framework enabled the extraction of traditional clinically used motion markers and the discovery of novel digital biomarkers sensitive to disease progression and medication response. The findings underscore its great potential for efficiently assessing disease progression in both clinical and home settings, as well as evaluating disease-modifying effects in clinical trials to promote personalized therapies.

1

INTRODUCTION

2 Parkinson's disease (PD) is the second most prevalent progressive neurodegenerative dis-
3 order, affecting more than 10 million people worldwide, and its incidence increases signif-
4 icantly as the global population ages [1]. Gait impairments are among the most common
5 and disabling symptoms in PD, characterized by intricate underlying mechanisms and
6 substantial individual variation in clinical manifestations [2]. Current interventions, in-
7 cluding pharmacological [3], non-pharmacological [4], and neuromodulatory therapies [5],
8 still yield inconsistent effects on gait impairments across patients with PD and disease
9 stages [2]. Precise and routine assessment of PD-induced gait impairments is crucial
10 for elucidating underlying mechanisms, understanding disease progression, developing
11 personalized intervention strategies, and ultimately improving patient outcomes.

12 Leveraging objective gait parameters to assess the PD progression and effectiveness of
13 different interventions has been an emerging trend and a considerable challenge in clinical
14 practice [6], [7]. Although clinical rating scales, like the Unified Parkinson's Disease Rating
15 Scale (UPDRS), are still commonly used for PD diagnosis, their low sensitivity, inherent
16 subjectivity, and dependency on clinical specialists limit their utility in routinely assessing
17 gait impairments for monitoring disease progression and evaluating treatment responses
18 [8]. Recent clinical studies have employed a few easily calculable gait parameters, such
19 as gait speed and step length, to assess PD severity [9]–[13]. However, relying on single
20 or limited gait features fails to comprehensively portray the severity of gait impairments
21 or evaluate the treatment efficacy. Because PD-induced gait impairments display complex
22 and diverse spatiotemporal motion characteristics, such as slow gait speed, shortened step
23 length, reduced amplitude of arm swing, reduced smoothness of locomotion, increased
24 interlimb asymmetry, increased gait variability, and impaired rhythmicity [2]. While mo-
25 tion capture systems can accurately measure various gait features [14], [15], their high
26 costs and professional operation requirements hinder their acceptance as tools for routine
27 assessment. Consequently, there is a substantial need for objective and precise routine
28 assessment methods that are able to thoroughly reflect the severity of gait impairments
29 and illuminate the complex relationships between various gait parameters and individual
30 disease progression, alongside their responsiveness to interventions, which remain poorly
31 understood [16].

32 Machine learning combined with sensing devices has been a promising modality
33 for PD assessment [17]–[24]. Most existing methods primarily focused on quantifying
34 PD symptoms with specific and distinct characteristics of movement disorders, such as
35 tremors and bradykinesia [17], [22], [24], particularly in the upper extremities. A small
36 body of studies has developed systems to assess gait impairments using a single wearable

37 device embodied with an inertia measurement sensor, such as a smartwatch [17], [23]
38 and smartphone [21], or fixed radio sensors at home [25] for longitudinal monitoring.
39 Despite these advances, such approaches are restricted to monitoring limited gait features
40 linked to specific body parts and cannot track detailed movements across all PD-affected
41 regions. Conversely, video-based approaches can overcome this limitation by leveraging
42 their inherent capability to capture diverse body movements comprehensively [22], [26]–
43 [29]. Nonetheless, current video-based methods for assessing gait impairment severity
44 struggle to achieve clinician-level accuracy, whether employing RGB or depth cameras
45 [30]–[32]. Furthermore, existing methods face considerable challenges in 1) high-resolution
46 evaluations of gait impairments beyond clinical rating scales for accurately assessing
47 treatment efficacy; 2) efficient identification of various motion markers for elucidating
48 the evolutions of different gait parameters with disease progression and their interactions
49 with treatments [2]. Moreover, existing video-based approaches often require multiple
50 fixed cameras to simultaneously capture gait features to avoid the influence of body part
51 occlusions [30], [31], impeding their regular application.

52 In this study, we propose a deep learning-based framework that can efficiently extract
53 precise spatiotemporal motion characteristics of the entire body joints from gait videos
54 recorded with a single smartphone to accurately assess PD-induced gait impairments
55 (MDS-UPDRS Part III-Gait item) (Fig 1). By developing a novel Siamese contrastive
56 architecture, our framework can imitate clinician’s assessment to fuse gait videos recorded
57 from both left and right lateral perspectives during shuttle walks, ensuring accurate identifi-
58 cation of lateral motion characteristics and comparative analysis of whole-body movements.
59 Unlike existing methods that directly extract traditional clinical gait parameters, our inter-
60 pretable framework analyzes personalized joint impacts on gait impairment severity over
61 walking time. This approach allows us to not only extract traditional gait parameters but
62 also discover novel digital biomarkers.

63 We applied this framework to a dataset with a well-balanced distribution of PD
64 severities to highlight the model’s proficiency in accurately assessing disease severity. We
65 effectively extracted digital biomarkers most sensitive to disease progression by correlating
66 them with disease severity. Additionally, we demonstrated the framework’s validity in dis-
67 criminating medication-induced changes in gait impairments, particularly subtle responses
68 undetectable by the UPDRS. Furthermore, we showed the model’s capability of identifying
69 digital biomarkers exhibiting high responsiveness to medication across patients with PD
70 and pinpointed the body joints with associated motion biomarkers showing the highest
71 medication responsiveness. These findings enable us to quantitatively assess the progression
72 of PD based on smartphone videos and use the model outputs as potential outcomes in

73 disease-modifying clinical trials, promoting personalized therapies.

74 The proposed smartphone video-based framework can efficiently enable home-based,
75 objective, routine assessments of gait impairments in PD. Our results also benefit in
76 developing objective and personalized therapies for other neurological disorders [2], such
77 as stroke [33] and Alzheimer's disease [34], with similar motor symptoms and complexities
78 as PD.

79 RESULTS

80 Model development and evaluation

81 We developed and evaluated the deep learning model for assessing gait impairments based
82 on smartphone-recorded videos from 118 participants, including 87 patients with PD and
83 31 healthy elderly controls (Fig. 1 and Table I). The evaluation of the model performance
84 consists of two phases: 1) predicting the severity of gait impairments and extracting
85 motion markers that are sensitive to disease progression; 2) discriminating comprehensive
86 changes in gait impairments in response to medication and identifying individualized
87 motion markers with high responsiveness to medication interventions. The severity of
88 gait impairments for each patient with PD was rated according to the consensus of three
89 clinical experts based on the MDS-UPDRS Part III-Gait scales (Supplementary Figure 3).
90 When there is no consensus among all three experts, the agreement of two experts (i.e., the
91 majority vote) was selected as the ground truth. In our study, there are no cases where a
92 consensus of at least two experts was not reached. We trained the model using 558 videos
93 from 93 participants and tested it with an independent dataset of 25 participants (Fig. 1a).
94 Compared to assessing the severity of PD, the evaluation of the effect of medication on gait
95 impairments is more challenging in clinical practice due to the complicated representations
96 of gait impairments and patient-specific responses to medication. To further evaluate the
97 validity of the model, we performed a medication response assessment with 19 patients
98 with PD in the test dataset (Fig. 1d). The gold standard was also established through expert
99 consensus, where three clinical specialists evaluated patients' changes in gait impairments
100 during off- and on-medication states, according to the UPDRS, alongside a refined three-
101 level sub-UPDRS scoring approach. These evaluations highlight the model's capability
102 to perform a clinician-level accurate assessment of PD severity while also serving as an
103 effective assessment tool for tailoring personalized treatments.

104 Model performance in predicting gait impairment severity

105 We evaluated the performance of the model on a test dataset comprising 150 video segments
106 from 25 participants (six video segments per participant) (Fig. 1a). The model demonstrated

107 a highly accurate prediction of the severity of gait impairments, with a precision of 0.804,
108 a recall of 0.811, a specificity of 0.898 and an F1 score of 0.806 (Table II). These values
109 are comparable to the average results achieved by three clinical specialists (Table II and
110 Supplementary Table I). The model correctly predicted 86% cases with a score of 0,
111 70% with a score of 1, and 88% with a score of 2 (Fig. 2a). The receiver operating
112 characteristic (ROC) curve demonstrated robust model performance in various UPDRS
113 categories (Fig. 2b). The model achieved an area under the ROC curve (AUC) of 0.93 for
114 the UPDRS score of 0, 0.78 for the score of 1, and 0.92 for the score of 2, with a micro-
115 average AUC of 0.87. These high AUC values further highlighted the model's effectiveness
116 in accurately predicting the severity of gait impairments. Furthermore, we assessed the
117 model's performance on the validation dataset through a 5-fold cross-validation procedure,
118 wherein the model achieved an enhanced F1 score of 0.82 and an elevated micro-average
119 AUC of 0.92 (Fig. 3).

120 **Model alignment with scores rated by clinical experts**

121 We analyzed the alignments among the UPDRS scores rated by the AI model and the
122 three experts and the ground truth scores for the participants in the test dataset (Fig. 2).
123 The performance of the model closely resembled those of the experts, with the model
124 having five mismatched scores compared to Expert 1 with three, Expert 2 with eight and
125 Expert 3 with three (Fig. 2c-f). Notably, the deviations of all model's mismatched scores
126 were limited to one point, matching the experts' performances. Although we trained the
127 model based on the consensus of three experts rather than their individualized scores,
128 interestingly, four of the five mismatched scores for the model were also not correctly
129 rated by at least one expert. These comparisons demonstrated that the model performance
130 is on par with the average performance of the experts (Error rates: 0.20 vs 0.19, Fig. 2j),
131 indicating the effectiveness and reliability of the model. In addition, the results reveal the
132 difficulties in assessing gait impairments, particularly in patients with inconsistent ratings
133 among experts. Moreover, we observed that the scoring approach of the model is different
134 from any of the three experts (Fig. 2g, h, i), although most of the rated scores of the model
135 are the same as those of each expert. These findings underscore the model's robust capacity
136 to align with expert assessments, illustrating its great potential as an independent tool in
137 accurately assessing disease severity in clinical evaluations.

138 **Extraction of motion markers sensitive to disease progression**

139 We first identified individualized joint contributions to the prediction of disease severity of
140 our model for participants in both training and test datasets using a dual maximum gradient-
141 weighted class activation mapping (DMGrad-CAM) method in conjunction with correlation

142 analysis (Fig.4b and Supplementary Figure 2). Across different levels of disease severity,
143 the foot, wrist, knee, and elbow were ranked as the primary body parts affected by the
144 disease (Fig.4b). This finding was consistent with typical gait characteristics observed in the
145 clinical assessment of gait impairments, such as the reduced amplitude of arm swing, gait
146 speed and step length as well as the diminished range of motion of the knee and ankle [2].
147 Based on the joint contributions, we first extracted the clinical commonly used biomarkers
148 [2] of gait impairments, including arm swing amplitude, gait speed, and step length (Fig.4c,
149 d, e). Arm swing amplitude demonstrated a significant correlation with disease severity,
150 yielding a correlation coefficient of -0.64 . Notably, there were significant differences
151 (Kruskal-Wallis test: $p < 0.05$) in arm swing amplitudes among participants with varying
152 UPDRS scores (Fig.4c). In addition, lower walking speeds and shorter step lengths were
153 significantly observed (Kruskal-Wallis test: $p < 0.05$) in patients with a UPDRS score of 2
154 compared to those with UPDRS scores of 0 and 1 (Fig.4d, e). These results align with the
155 PD progression of Parkinson's disease, where a reduced amplitude of arm swing appears
156 in the slight stage and worsens further in the mild state; slow speed and shortened step
157 length become common in the early stage [2]. Since increased cadence usually appears
158 in the moderate stage, we didn't observe this manifestation in the patients with early-
159 stage PD in our study (Supplementary Figure 4). In addition to these traditional motion
160 markers, we further discovered richer motion features consisting of the linear velocities
161 and accelerations of the skeletal joints and the joint angles (Fig.4a). We finally selected two
162 types of indicators that are sensitive to disease progression: the mean of each motion feature
163 and the variances in the means for six walking periods during three-time shuttle walk tests
164 (Supplementary Figure 5). The correlation analysis with the UPDRS score revealed that
165 the mean linear velocities of the foot and knee, the variances of the velocities of the
166 wrist, elbow and shoulder, the mean and the variances of the accelerations of the wrist and
167 elbow can especially serve as effective digital biomarkers to reflect disease progression. In
168 particular, the average linear velocity of the ankle, a newly identified digital marker, showed
169 a higher correlation with the gait impairment severity ($\rho = -0.66$) than all three traditional
170 motion markers. These body parts are also the ones with higher joint contributions across
171 different disease stages, compared to the other parts (Fig.4b). Moreover, the relatively
172 high correlation between the variances of average velocities of the upper limbs across six
173 walking periods and the severity of the disease aligns with clinical characteristics of gait
174 impairments, i.e., increased gait variability [2]. The model, combined with interpretable
175 joint contributions and extracted skeletal data, offers a promising approach for efficiently
176 identifying motion markers sensitive to individualized disease progression.

177 Model performance in discriminating medication effect on gait impairments

178 To evaluate the ability of the model to discriminate changes in gait impairments caused
179 by pharmacological interventions, we experimented with collecting gait videos from 19
180 patients with PD during off- and on-medication states (Table I). According to the consensus
181 of the three clinical specialists, the UPDRS scores (Part III-Gait scales) of seven patients
182 changed by one score on-medication (SwC cohort), with six patients' scores decreased and
183 one increased, while the scores of the others remained the same (SwoC cohort). For the
184 last cohort, the three experts further conducted a more granular sub-score comparison with
185 three outcomes, i.e. improvement, no change, and deterioration, to differentiate changes
186 in gait impairments of each patient between off- and on-medication states. This granular
187 rating scale reflected subtle gait alterations after medication that cannot be indicated by
188 a change in the UPDRS score. This comparison was repeated three times. First, we only
189 randomly provided the gait videos of each patient recorded during off- and on-medication
190 states to the experts without giving them the patient's medication states. Secondly, videos
191 aligned with corresponding medication states were provided, as well as the patient's
192 other medical information (Table I). Lastly, experts were also provided with the gait
193 characteristics extracted from the skeleton data, including stride length, gait speed, and
194 arm swing amplitude [2]. This staged evaluation benchmarks our model against clinicians
195 by comparing the model's video-only performance to that of experts receiving progressively
196 more data, demonstrating its ability to achieve a comparable or superior assessment with
197 more parsimonious data inputs. In addition, this graded approach can avoid preconceptions
198 caused by awareness of the state of the medication for clinicians. For the SwoC cohort,
199 the consensus rating, determined by the agreement of at least two experts in the final
200 evaluation, served as the gold standard for each patient. In instances where consensus was
201 not achieved (i.e., the three experts gave three different ratings), a "no-change" classification
202 was assigned, which occurred only once in our study. For the SwC cohort, the consensus
203 was still based on changes in UPDRS scores, which were further mapped to "improvement,"
204 "no change," or "deterioration." In addition, a non-expert clinician, certified in assessing
205 PD symptoms but possessing less clinical experience than the experts, conducted the same
206 evaluations for all 19 patients.

207 We introduced a comprehensive index based on the predicted scores along with their
208 confidence levels, as generated by the model, to identify changes in gait impairments
209 between off- and on-medication states for both the SwC and SwoC cohorts. The model
210 demonstrated a strong capability in discriminating the effects of medications on gait
211 impairments with an accuracy of 73.68%. This matches the best performance of the three
212 rounds of evaluations of two experts and the non-expert clinician and only falls short by

213 10.53% (two patients) compared to the highest expert performance (Fig.5a). Compared to
214 the rating results of the three experts in the first-round evaluation based on only gait videos,
215 the model's accuracy was even slightly higher than their average accuracy (70.17%). For the
216 SwC cohort, all experts exhibited a notable decrease in discrimination accuracy, indicating
217 greater challenges in distinguishing medication-induced changes in gait impairments using
218 UPDRS scores as opposed to merely assessing the severity of gait impairments (Fig.5c
219 and Table II). However, the comprehensive index derived from our model significantly
220 mitigated this issue, achieving a discrimination accuracy of 85.71%, thereby surpassing
221 all experts. For the SwoC cohort, all clinicians achieved their highest performance in the
222 third evaluation when receiving not only gait videos but also patient medication states and
223 quantitative gait characteristics. The discrimination accuracy of the model is less than the
224 best performance of the experts in the three-round evaluations; however, it matched the best
225 performance of the non-expert, with a discrimination accuracy of 66.67% (Fig.5c). In other
226 words, the model outperformed the experts in distinguishing relatively significant changes
227 in gait impairments (with a change in the UPDRS score) caused by medication based on
228 UPDRS scores and ensured the same level as the non-expert in distinguishing insignificant
229 gait changes (without a change in the UPDRS score). It is important to note that the clinical
230 commonly-used gait characteristics, including arm swing amplitude, step length, velocity,
231 and their combination, cannot effectively differentiate gait changes between off- and on-
232 medication based on their significance analysis in gait changes, giving a precision $\leq 42.11\%$
233 (Fig. 6a). Interestingly, the correlations between the outcomes of the medication in gait
234 impairments determined by the model, the clinicians, and the gait characteristics revealed
235 that the model's performance had more explicit relationships with the gait characteristics
236 compared to the clinicians' discrimination (Fig.5b). These findings demonstrate that our
237 model can efficiently extract representations of skeletal data that reflect medication-induced
238 gait changes; meanwhile, it can significantly outperform these traditional gait biomarkers to
239 perform expert-level discrimination of changes in gait impairments caused by medication.
240 This highlights its potential as a valuable tool for evaluating the effectiveness of medical
241 therapies.

242 **Identification of digital biomarkers with high responsiveness to medication**

243 In addition to comprehensively determining the outcomes of pharmacological interventions
244 using our model, we identified motion markers with high responsiveness to medication
245 across patients with PD. We used the percentage of patients ($N=19$) showing significant
246 changes ($p < 0.05$) in different motion markers between off- and on-medication states to
247 highlight the markers and the corresponding body joints sensitive to medication. We only
248 considered the cases in which significant changes in the motion markers aligned with the
249 consensus on medication efficacy rated by three clinical specialists.

250 For the newly extracted spatiotemporal digital biomarkers, the linear accelerations
251 of the neck and head, and the standard deviations of the elbow and neck joint angles,
252 demonstrated relatively high inter-patient response to medication outcomes (Fig. 6a). No-
253 tably, 42.11% of patients showed significant changes in these four digital biomarkers
254 consistent with medication outcomes rated by clinical experts. The inter-patient medication
255 responsiveness for these digital biomarkers was higher than that of traditional clinical
256 motion markers (31.58%), such as arm swing amplitude, walking speed, and step length
257 (Fig. 6a). For any single spatiotemporal marker or traditional motion marker, the percentage
258 of patients with significant changes was relatively low (any spatiotemporal marker: \leq
259 42.11%; any traditional motion marker: \leq 31.58%) In contrast, combining changes in sev-
260 eral key motion markers can more accurately reflect medication effects on gait impairments
261 across a broader patient group. It was shown that 63.16% of patients exhibited significant
262 changes in at least one of the four newly extracted spatiotemporal markers matching
263 their rated medication outcomes, compared to just 42.11% with substantial changes in the
264 traditional markers (Fig. 6a). Here, we used voting results to verify alignment with patient
265 medication outcomes if these markers displayed inconsistent significant changes. These
266 results highlight the challenge of identifying a universal digital biomarker for indicating
267 medication efficacy and underscore the importance of developing personalized treatments
268 based on comprehensive changes in multiple key motion markers, especially the extracted
269 spatiotemporal markers. Fig. 6b detailed the proportions of patients with PD exhibiting
270 significant changes ($p < 0.05$) in one or multiple markers among the Top-4 spatiotemporal
271 and Top-3 traditional markers. Only 15.79% and 21.05% of patients showed significant
272 changes in all Top-3 traditional and Top-4 spatiotemporal markers, respectively. However,
273 these percentages are still remarkably higher than those exhibiting significant changes in 1-
274 2 traditional or 1-3 spatiotemporal markers (\leq 10.53%), respectively. These results further
275 demonstrate the clinical heterogeneity of PD, indicating that medications have varying
276 effects on different motor symptoms across individuals.

277 Furthermore, we analyzed the effects of medication on motor abilities of different body
278 parts, taking into account both the Top-4 spatiotemporal and Top-3 traditional markers. Our
279 findings indicated that motion markers associated with the neck, head, elbow, knee, hip
280 showed higher inter patient medication responsiveness (with 57.86%, 52.63%, 43.37%,
281 42.11%, 42.11% patients, respectively), compared to other joints (Fig. 6c). Notably, the
282 average linear velocities and accelerations of the head and neck serve as key motion
283 markers with high response to medication across patients. The results also demonstrate
284 the capability of the proposed framework for identifying different medication effects on
285 various body parts. Detailed changes in these motion markers for all nineteen patients
286 between off- and on-medication states are presented in Supplementary Figure 7-10.

287

DISCUSSION

288 We proposed a deep learning model using smartphone videos to quantitatively assess PD-
289 induced gait impairments and discriminate the effect of pharmaceutical intervention on
290 gait impairments. Based on the proposed model, we further extracted motion markers
291 with explicit responsiveness to disease progression and medical treatment. The model was
292 trained on a dataset with 93 participants that are classified into three categories according
293 to the disease severity. To reduce the influence of subjectivity and inter-rater variability
294 (IRV), the consensus of three clinical experts was used to label the severity of the disease
295 for each participant according to the MDS-UPDRS Part III-Gait scales and to evaluate gait
296 impairment alterations between off- and on-medication states according to UPDRS scores
297 along with a specialized assessment score that offers greater resolution than UPDRS. The
298 ground truths (i.e., labeled UPDRS scores) of the dataset demonstrated an excellent inter-
299 rater reliability, with an intraclass correlation coefficient of 0.80 [19], [22]. The model
300 exhibited expert-level performance in predicting the UPDRS scores of an independent
301 test dataset with 25 participants and extracted various motion markers whose changes
302 are consistent with the clinical manifestations of PD disease at different severity stages.
303 The model also performed an effective identification of the patient's response on gait
304 impairments to medical treatment, with the same discrimination precision (73.68%) as
305 that of two of three experts (Fig. 5a). In particular, the model demonstrated the ability to
306 distinguish more granular changes in gait impairment that cannot be indicated using the
307 UPDRS score with the same precision as that of the non-expert clinician who used the
308 fine-grained rating method. For the SwC cohort, our model achieved 85.71% accuracy,
309 matching that of the non-expert clinician and outperforming the three experts (accuracies
310 of 57.14%–71.43%) who rated based on the UPDRS scale. For the SwoC cohort, although
311 the experts performed at a higher level of accuracy using the granular three-level criterion,
312 with scores ranging from 83.3% to 91.67% (the best performance), the model's accuracy
313 (66.67%) remained comparable to that of the non-expert (66.67%). These findings under-
314 score the great potential of the model in precisely tracking disease progression and response
315 to treatments over time, enhancing understanding of the fundamental mechanisms of gait
316 impairments, and paving the way for the development of personalized therapies. In addition,
317 we developed an online assessment system based on the model to allow home-based routine
318 assessment of gait impairments, efficiently addressing the challenges in routine quantitative
319 evaluation caused by the complicated representations of gait impairments [35].

320 Our model achieved an accurate and robust assessment of the severity of gait impair-
321 ment, which mainly lies in the design of the novel contrastive network architecture and a
322 weight-sharing mechanism, enabling our model to extract gait spatiotemporal features from

323 videos recorded from left and right perspectives simultaneously. Despite advancements
324 in sensing techniques for assessing gait impairment severity, existing methods [29]–[32],
325 [36] continue to face challenges in enhancing assessment accuracy. The best-reported
326 performances across studies are constrained with an F1 score of 0.77, a precision of
327 0.782, or an AUC of 0.83, while our model improved the F1 score, accuracy, and AUC
328 to 0.806, 0.811, and 0.93, respectively. For instance, a recent well-structured model [26],
329 [37] was developed to predict the 4-level severities (i.e., UPDRS scores of 0, 1, 2, 3)
330 of gait impairments using a 3D skeleton extracted from videos recorded from a frontal
331 perspective. That model, evaluated with a dataset consisting of 31 patients with PD and 23
332 controls without PD, reported an AUC of 0.8 and an F1 score of 0.76 during a leave-one-out
333 cross-validation. Given the inclusion of only four patients with a UPDRS score of 3 in that
334 study, we compared our model to that model. In contrast, our model demonstrated superior
335 performance, achieving a 0.12 boost in AUC and a 0.06 increase in F1 score during a 5-fold
336 cross-validation on a substantially larger dataset comprising 93 participants (Fig. 3). These
337 findings demonstrated the improved precision, efficiency and robustness of our model over
338 existing ones in assessing the severity of PD. To the best of our knowledge, we are the first
339 to demonstrate that a video-based model can achieve an assessment accuracy comparable to
340 that of clinical experts, highlighting its potential as an independent tool for gait impairment
341 diagnosis in clinical and home-based settings. The robustness of the model was further
342 emphasized by verifying it on a relatively large dataset with well-balanced categories
343 compared to those used in current methods [25], [30], [32], [38].

344 The proposed quantification approach utilizing smartphone-recorded videos efficiently
345 balances two critical requirements for routine assessment of gait impairment severity:
346 capturing more motion features to enhance accuracy while minimizing sensors used to
347 improve acceptance and usability. Our approach can be effectively used for longitudinal
348 monitoring of gait impairments and analyzing the impact of significant daily events, such
349 as on-off medication. In addition, our model, based on the analysis of joint contributions
350 to disease severity, is capable of extracting not only highly correlated traditional motion
351 markers but also a broad spectrum of digital biomarkers to reflect disease stages, such as
352 linear accelerations of wrist and elbow, linear velocities of foot and knee, variances of
353 linear velocities of elbow, wrist and shoulder. The extracted traditional motion markers
354 of disease progression are consistent with the findings of recent studies using wearable
355 devices [17], [23] or radio waves [25]. For example, our analysis revealed significant
356 negative correlations between disease severity and arm swing amplitude (Spearman: $\rho =$
357 $-0.64, p < 0.001$), gait speed (Spearman: $\rho = -0.62, p < 0.001$), and step length
358 (Spearman: $\rho = -0.59, p < 0.001$). These findings align with previous work, which also re-
359 ported significant correlations between disease severity and arm swing amplitude (Pearson:

360 $r = -0.31, p = 0.008$) [39], step length (Pearson: $r = -0.519, p < 0.001$), and gait speed
361 (Pearson: $r = -0.433, p = 0.001$) [40]. The new digital motion markers can provide deeper
362 insights into personalized gait impairments associated with disease progression and clinical
363 treatments [2]. For instance, the wrist acceleration biomarker identified in our study aligns
364 with recent findings from a study on a large dataset (UK Biobank), which demonstrated
365 that wrist acceleration measured using a wrist-worn triaxial accelerometer, when analyzed
366 with a machine learning model, can better distinguish clinically diagnosed PD (with a
367 Area Under Precision Recall Curve (AUPRC) = 0.14 ± 0.04) and prodromal PD up to
368 seven years pre-diagnosis (AUPRC = 0.07 ± 0.03) from the general population, compared
369 to traditional biomarkers, such as genetics (AUPRC = 0.01 ± 0.00) and blood biochemistry
370 (AUPRC = 0.01 ± 0.00) [23]. Similarly, the correlation between wrist acceleration and
371 gait impairment severity was also observed in another study, which found a correlation
372 of $\rho = -0.46$ between arm swing acceleration and clinical gait scores using a wrist-worn
373 wearable device [17]. Compared to these methods, our model, leveraging smartphone-
374 recorded videos, successfully addresses the limitations of wearable-sensor-based methods
375 by enabling the comprehensive capture of motion features across diverse body regions
376 while ensuring patient acceptability.

377 To bridge the gap between the high demand for accessible and comprehensive analyses
378 of the effectiveness of pharmaceutical interventions on gait impairments and the critical
379 paucity of relevant assessment methods, we developed an approach to discriminate changes
380 in gait impairments caused by medication, leveraging the model and confidence values of
381 its outputs, without retraining it. As the low sensitivity of the UPDRS score limits its ability
382 to capture all the nuances of gait metrics and their abnormalities [8], it cannot be used
383 to adequately reflect changes in gait impairments [41], particularly medication-induced
384 changes in gait [8]. On the other hand, an expert clinician can often perceive fine-grained
385 improvements in movement quality that are not significant enough to change an UPDRS
386 score [42], [43]. Hence, we introduced a refined three-level sub-UPDRS scoring approach.
387 In our study, the UPDRS scores of 12 out of 19 patients did not change after medication,
388 while only one patient's score was still rated 'no change' when clinicians used a higher
389 resolution assessment method. So far, few motion markers (normally the speed) have been
390 used to evaluate the effect of interventions [25]. A major constraint behind this is the lack
391 of easily accessible means to perform a quantitative and comprehensive assessment. Our
392 framework demonstrated a capability to accurately identify changes in gait impairment in
393 patients with PD after medication with a precision of 73.68%, which slightly surpasses
394 the average proficiency of the clinical experts based on only gait videos and matches the
395 precision of two of them when supplemented with information on medication states and
396 gait characteristics. The comprehensive index, derived from our model, demonstrated an

397 accuracy of 85.71% in identifying medication-induced outcomes on gait impairments as
398 indicated by the changes of UPDRS scores, remarkably exceeding the average accuracy
399 of clinical experts, which was 61.9%. Furthermore, it accurately identified gait changes
400 undetectable by UPDRS scores, reaching the accuracy level of a non-expert clinician with
401 a fine-granular assessment method. Since UPDRS is limited by its subjective, coarse and
402 semi-quantitative nature [44], [45], it is insensitive to subtle but potentially significant motor
403 fluctuations [2], particularly in response to treatment [46], [47]. Our fine-grained assessment
404 method has the potential to equip clinicians with the ability to monitor fluctuations in the
405 clinical status of a patient more objectively [21], which in turn supports timely adjustments
406 to patient management and individualized treatment strategies [44], [46], [47].

407 Utilizing motion markers derived from our model, we discerned personalized med-
408 ication responses in patients with PD. Notably, four digital biomarkers, including linear
409 accelerations of the neck and head and standard deviations of the elbow and neck joint
410 angles, exhibited relatively high inter-patient responsiveness to medication, with 63.16%
411 of patients showing significant changes in at least one marker that aligned with the
412 patient's medication outcomes during the medication response test. These newly extracted
413 spatiotemporal biomarkers demonstrated higher inter-patient medication responsiveness
414 than traditional motion markers used in clinical studies, such as arm swing amplitude,
415 gait speed, and step length. Furthermore, motion markers associated with the head, neck,
416 elbow, knee, and hip showed higher inter-patient responsiveness than those linked to other
417 joints. So far, only a few studies have used quantitative gait measures to analyze the
418 medication response of patients with PD, mainly demonstrating that the gait speed [11],
419 [25] and step length [11] are sensitive to medication. Our model demonstrates the ability
420 to identify more digital biomarkers of medication response, which can further benefit
421 disease-modifying clinical studies and customized treatments.

422 There are some limitations to this study. The quality of the skeleton data, which
423 is the direct input to our model, significantly impacts the accuracy of the assessment.
424 Since the skeleton data come from deep learning extraction models, their accuracy directly
425 affects the final assessment. In our study, all gait videos were captured under indoor
426 lighting conditions. During indoor experiments, we found that subjects' clothing affects
427 keypoint extraction [48]. Loose clothing or pure black wear that obscures limb movements
428 can cause assessment failures. Patients are suggested not to wear this kind of clothing
429 during the assessment. We recorded the videos using smartphones with a resolution of
430 at least 720p. These conditions are suitable for the underlying DWPose pose estimation
431 model, indicating that factors such as poor lighting, phone type, or low resolution were
432 not main sources of error in our study. To ensure the generalizability of the model, we

433 implemented an assessment on a completely independent test set with 25 participants as
434 well as a five-fold cross-validation scheme. In addition, we evaluated the performance of
435 the model in discriminating the response of medication to gait impairments in 19 patients
436 with PD selected from the independent test set, which means that their data were never
437 used for training. Although the sample for this specific evaluation is not as large as that
438 for disease severity assessment, promising preliminary results demonstrate the model's
439 ability to discern medication effects and highlight its significant potential for future clinical
440 applications.

441 The confusion matrices for both the 5-fold cross-validation and the independent test
442 demonstrate that the model achieved a higher accuracy for classifying UPDRS scores
443 of 0 and 2 than classifying UPDRS scores of 1 (see Fig. 2a and Fig. 3a), suggesting
444 that the model encounters greater challenges in distinguishing gait impairment severity
445 with a UPDRS score of 1. This difficulty reflects an inherent challenge in the clinical
446 assessment itself, as the distinction between these mild severity levels (score of 1) is clini-
447 cally subtle, creating a blurry boundary that is challenging even for clinicians to delineate
448 consistently. There are no significant differences in sex (Fisher's Exact Test, test/training
449 set: $p=0.3217/0.4219$), age (Mann-Whitney U test, test/training set: $p=0.3755/0.9960$) and
450 height (Mann-Whitney U test, test/training set: $p=0.4131/0.2914$) between the correctly
451 predicted group and the mismatched group. With the test dataset, the model accurately
452 classified all cases with identical rated UPDRS scores among the three experts and nine of
453 the fourteen instances without a unanimous rated score (see Fig. 2), demonstrating a robust
454 capability to handle the IRV. On the other hand, the results showed that four of the five
455 mismatched scores (one-point difference) for the model on the test set were also incorrectly
456 rated in the same way by at least one expert and the model error was significantly linked to
457 the expert disagreement (Fisher's Exact Test, test/training set: $p=0.0464/0.0022$), indicating
458 that the model faces similar challenges in assessing gait impairments in complex cases with
459 inconsistent ratings among experts.

460 In our study, all data were from the same hospital and consisted exclusively of patients
461 of one ethnicity (Asian Chinese) with a mean age of 62.9 years (± 7.4 years), which may
462 not fully reflect the diverse gait characteristics of patients with PD. In the future, we will
463 expand the dataset by recruiting participants from different medical centers with a wider
464 range of races, ages and disease severity. We recorded gait videos under normal indoor
465 lighting conditions, and the use of smartphones allowed for a minimum video resolution of
466 720p. We will also include gait videos recorded with various types of smartphones under
467 different indoor and outdoor environments to further improve the model's generalization
468 before its widespread clinical use. In addition, the model mainly focuses on assessing

469 walking, allowing convenient application in the home setting. In the future, the model will
470 be further expanded to assess additional tasks, such as turning, Timed Up and Go, balance,
471 and freezing of gait, to obtain a more comprehensive diagnosis of the disease severity.
472 Moreover, the proposed video-based model can be further integrated with other modal data
473 from wearable sensors, such as kinematic data from inertial measurement units (IMUs)
474 [49], [50] and muscle activities from electromyography (EMG) sensors [51], to provide
475 more comprehensive assessment of PD symptoms related to features of specific body
476 parts (e.g., tremor and speech) and the whole body (e.g., gait, freezing of gait, balance).
477 Such an integration can also enable a deep insight into the relationships among different
478 symptoms and their complex responses to the same interventions. In addition, integrating
479 the gait videos with neuroimaging data, such as magnetic resonance imaging (MRI) and
480 functional MRI (fMRI) [26], can help further understand the mechanisms underlying gait
481 impairments. This multidimensional assessment will provide more comprehensive guidance
482 on Parkinson's disease progression and the effectiveness of therapeutic interventions.

483 In conclusion, we developed a smartphone video-based deep learning model that
484 accurately assessed the severity of PD-induced gait impairments (MDS-UPDRS Part III-
485 Gait item) and discriminated the effectiveness of pharmaceutical interventions on gait
486 impairments, including those that UPDRS scores failed to detect due to its low resolution.
487 In addition, the interpretability of the model enabled the extraction of valuable digital
488 biomarkers, which provide insights into disease progression and medication effects on gait
489 impairments.

490 METHODS

491 Participants and dataset

492 The study initially recruited 130 participants, including 99 patients with PD and 31 healthy
493 age-matched adults as controls. Each participant was instructed to perform the shuttle walk
494 test three times, covering a minimum distance of 10 meters for each test. Meanwhile,
495 we filmed the participants walking from the lateral perspective using a single smartphone
496 placed at a fixed point, capturing their full body motions. This perspective is more effective
497 for extracting key gait features, such as step length, arm swing amplitude, and gait speed,
498 compared to frontal views [26]. Twelve participants were excluded from the study due
499 to needing walking assistance ($N=3$), non-compliance with the experimental protocol by
500 raising their arms during walking ($N=2$), and failure in identifying parts of the skeletons
501 from the videos caused by the influence of loose black clothing ($N=7$). Consequently, the
502 final cohort of participants consisted of 87 patients with PD and 31 healthy controls (118
503 in total).

504 Gait videos from 118 participants, recorded during off-medication states, were inde-
505 pendently assessed and rated by three clinical experts according to the MDS-UPDRS Part
506 III-Gait exams. The consensus derived from the agreement of at least two experts was
507 established as the ground truth for each participant. In instances where consensus was not
508 reached, the average UPDRS scores of the three experts were used as the ground truth;
509 however, no such cases occurred in our study. We constructed a balanced dataset with a
510 distribution of UPDRS categories as follows: 38 participants with a score of 0, 39 with a
511 score of 1, and 41 with a score of 2 (Table I). Note that patients with a UPDRS of 4 cannot
512 walk independently and those with a score of 3 require assistance devices. The cohort also
513 ensured a rough gender balance with 63 males and 55 females ($p > 0.05$; Chi-Squared
514 Test). Detailed demographic information of the participants is presented in Table I.

515 We divided the dataset randomly into two groups: a training dataset consisting of
516 558 video segments from 93 participants and a test dataset comprising gait videos from 25
517 participants (Fig. 1a). All data splitting was performed at the participant level. This process
518 ensures that all video segments from a single participant were contained entirely within
519 one partition (i.e., the training set or the test set), thereby preventing data leakage. The
520 model was trained using the training dataset, and its performance was evaluated using the
521 independent test dataset. Initially, we assessed the model's effectiveness in predicting the
522 severity of gait impairments and compared it to the assessments made by individual clinical
523 experts. Subsequently, we calculated the personalized joint contributions to the prediction
524 of disease severity for each participant in the training and test datasets. By analyzing
525 the average joint contributions across all participants, we evaluated the model's ability to
526 extract both conventionally used clinical motion markers and novel digital biomarkers that
527 are sensitive to disease progression, based on their correlations with the UPDRS scores.

528 We evaluated the model's performance in distinguishing comprehensive gait impair-
529 ment outcomes in response to medication based on gait videos from the 19 patients with
530 PD during both off- and on-medication states in the test dataset. This evaluation was
531 accomplished by integrating the UPDRS scores predicted by the model, along with their
532 confidence levels. To accurately evaluate the effectiveness of pharmacological interventions
533 (Table I) on gait impairments, we developed a fine-granular assessment approach with a
534 higher resolution than UPDRS scores. Three clinical experts assessed the gait videos of 19
535 patients with PD in the test dataset (Fig. 1), recorded in both off- and on-medication states,
536 according to the UPDRS and the fine-granular approach. Regarding medication regimens
537 during the medication response test, sixteen patients were treated with Madopar alone, one
538 patient received a combination of Madopar and Sifrol, and two patients were administered
539 Madopar and Benzhexol. First, the UPDRS scores of the 19 patients were independently

540 rated by three experts. According to the consensus of experts, changes in UPDRS scores
541 after the medication were used to indicate alterations in gait impairments caused by the
542 medication. For those whose UPDRS scores remained unchanged, three experts further
543 identified changes in gait impairments using a more granular three-level sub-score criterion
544 (i.e., improvement, no change, and deterioration). To capture the nuances in patients' gaits
545 with higher resolution, we asked the three specialists to independently compare all gait
546 metrics of the same patient based on the videos recorded during off and on medication
547 states to identify the differences between them as much as possible, then leading to one
548 of three categories: improvement, no change, and deterioration. To avoid preconceptions
549 caused by medication state awareness, clinicians first evaluated randomly sorted gait videos
550 without associated medication states. Afterwards, they assessed the videos with the patient's
551 medication states and medical history, with a new sorted sequence. Finally, we provided
552 them with the gait characteristics derived from the modified skeleton data, including stride
553 length, gait speed, and arm swing amplitude [2] in addition to the videos and previous
554 information with a new order. The consensus rating, determined by the agreement of at
555 least two experts in the final evaluation, served as the gold standard for each patient. In
556 instances where consensus was not achieved, a "no-change" classification was assigned,
557 which occurred only once in our study. Moreover, a non-expert clinician conducted the
558 three rounds of fine-granular evaluations for all 19 patients. Apart from discriminating the
559 comprehensive changes in gait impairments, we further identified individualized motion
560 markers with high responsiveness to medication interventions by analyzing their significant
561 changes between off and on-medication states.

562 **Ethics declaration**

563 Participants were recruited from Tianjin Huanhu Hospital in Tianjin, China. They were
564 screened for Parkinson's disease, and their clinical tests were assessed by expert neurol-
565 ogists specializing in movement disorders. The inclusion criterion for the patient group
566 was a confirmed diagnosis of Parkinson's disease and the capability to walk without assis-
567 tance. Healthy controls, matched by age, were also included in the study. All participants
568 were fully informed of the experimental procedures and provided their written consent
569 to participate in this study. All ethical and experimental procedures and protocols were
570 approved by the Institutional Review Board of Tianjin Huanhu Hospital, Tianjin, China,
571 under Approval No. ChiCTR1900025372.

572 **Data preprocessing**

573 We developed an automated video segmenting approach to segment long shuttle walking
574 videos into short clips that only capture directional walking to enable simultaneous analysis

575 of left- and right-side motion characteristics. This approach works by identifying the
576 segmentation locations via detecting peaks and valleys in the horizontal profiles of the
577 neck keypoint and analyzing their magnitudes and the horizontal distances between them.
578 Each gait video was automatically split into six segments of unidirectional movement
579 (left-to-right or right-to-left), ensuring that only walking segments were included. This
580 segmentation was consistently applied to model training, inference, and feature extraction.

581 We extracted the body skeleton from smartphone videos using DWPose [52], consid-
582 ering its superior effectiveness and precision compared to the commonly used OpenPose
583 [53]. Although the accuracy and robustness the DWPose have been verified on a large-scale
584 public dataset (COCO-WholeBody [54]), we further verify its accuracy in extracting human
585 skeletons in the setting of our study (see Fig.1). We conducted a comparison experiment
586 using a motion capture system (Vicon, UK) as the ground truth. Two healthy young subjects
587 were asked to perform a five-meter shuttle walk test five times. We compared the key points
588 extracted using DWPose and Vicon without considering the depth information of Vicon
589 along the frontal axis since DWPose provides only 2D skeleton data. To align the pixel
590 coordinate system of DWPose with the world coordinate system of Vicon, we performed a
591 perspective transformation using a nonlinear fitting with a quadratic function. In addition,
592 time scale alignment [55] was performed to ensure synchronization between the key points
593 extracted from both systems. We calculated the Spearman correlation coefficients for each
594 joint's 2D trajectories estimated by DWPose and Vicon across 16 video segments for
595 the two subjects. The results showed a high correlation between their estimated joint's
596 trajectories, with a $\rho \geq 0.85 \pm 0.06$ in the sagittal plane and a ρ in between 0.67 ± 0.07
597 and 0.99 ± 0.01 in the transverse plane (Supplementary Figure 11).

598 Skeleton extraction models such as DWPose and OpenPose face intrinsic challenges
599 in precisely identifying keypoints for each leg in videos recorded in lateral perspective
600 when the legs alternate between front and back positions, stemming from their current
601 technological limitations. Misidentification of leg keypoints can result in significant errors
602 in feature extraction and consequently cause incorrect UPDRS classification. To address
603 this, we developed a modification algorithm to correct misdetections of leg keypoints
604 for each video segment. This algorithm analyzes the trajectory of the horizontal distance
605 between the left and right ankles, as estimated by the skeleton extraction model, to correct
606 keypoints. To eliminate potential errors arising from confusion between the left and right
607 ankle keypoints during foot alternation, our algorithm implemented iterative adjustments
608 to the curve trend on a per-frame basis. For each video frame, it evaluates two scenarios:
609 one in which the model directly identifies the keypoints and another in which they are

610 swapped. By selecting the scenario that maintains a smoother and more consistent motion
611 trajectory based on the predicted curve and the second-order and third-order differences of
612 the predicted curves, the algorithm can efficiently identify and rectify incorrect swaps. This
613 algorithm iteratively updated the distance between the left and right ankles across all video
614 frames to correct the confusion between them. This approach largely reduced the number
615 of incorrect identifications of left and right ankles and refined motion trajectories in video
616 frames. Consequently, it provided a more accurate representation of the ankle joint motion
617 curve. Supplementary Figure 1 demonstrates the algorithm's effectiveness in improving the
618 reliability of tracking key points. Furthermore, the detected keypoints for the occluded arms
619 are unreliable when using DWPose or OpenPose, as they rely on statistical estimates when
620 the limbs are occluded. These estimates are often similar to those of healthy individuals
621 rather than patients with PD, leading to classification errors. To mitigate this, we set the
622 unreliable keypoints for the occluded arms to zero in our model, focusing primarily on the
623 motion of the arm nearest to the smartphone side.

624 **Deep learning-based model**

625 To accurately extract motion characteristics from gait videos recorded from both left- and
626 right-lateral perspectives and efficiently fuse them to perform a comprehensive assessment
627 of gait impairments, we developed a novel Siamese contrastive network architecture (Fig.
628 7), inspired by the model in [56]. This architecture can concurrently extract features from
629 videos recorded in left- and right-lateral perspectives. By using a weight-sharing mechanism
630 in feature extraction and a feature fusion strategy, the proposed Siamese network can
631 not only efficiently combine all joints' features, but also reflect the asymmetry of joint
632 movements in spatiotemporal domains, ensuring an accurate prediction of disease severity.
633 Note that interlimb asymmetry is typically specific to PD and often appears in its early stage
634 [2]. Instead of using the heavy Transformer-based MotionBERT network [57], [58], we
635 utilized a lightweight spatial-temporal graph convolutional network (ST-GCN) [59] as the
636 backbone of the architecture to guarantee high efficiency. Unlike our previous model [27],
637 which required both gait energy images and a complex convolutional neural network, our
638 present model uses only skeleton data and eliminates the need for complex convolutional
639 neural networks, leading to a significant reduction in computational resources. This enables
640 convenient home-based online prediction of disease severity.

641 The model was trained using 558 video segments from 93 participants (Fig. 1). In
642 addition, we enhanced the robustness of the model by further performing data augmentation
643 to expand the effective size of the dataset [60], [61]. Spatial data augmentation was
644 carried out by randomly scaling and rotating coordinate values of all joints to introduce

645 data variability. We did not perform temporal data augmentation due to the constraint of
646 the patients' zero gait velocity at the beginning of walking. To prevent overfitting, a 5-
647 fold cross-validation strategy was implemented. We selected the five-fold cross-validation
648 procedure to ensure both an efficient model training and an effective evaluation of model
649 performance. Note that we split the folds at the patient level, ensuring that all data from
650 a single patient was contained entirely within one fold and a given patient's data never
651 appeared in both the training and validation sets simultaneously in any iteration. The data
652 augmentation was restricted to the training phase to avoid data leakage, ensuring that no
653 augmentation occurred during the validation or independent test phases, thus maintaining
654 the integrity of the validation/test data with the original data.

655 To evaluate the model's performance in predicting disease severity, we conducted tests
656 using an independent dataset comprising 150 video segments from 25 participants (Fig. 1).
657 For each participant, six gait video segments were randomly paired with counterparts
658 recorded from the opposite side, generating a pair of dual-perspective video inputs for the
659 model. We utilized five models derived from the five-fold training during the evaluation.
660 For each video segment, the majority vote among the outputs of five models was selected
661 as the predicted result. Similarly, the majority vote of the predicted results for six video
662 segments was chosen as the final predicted result for each participant. Evaluation metrics
663 were calculated based on these final results of all participants in the test dataset (Table II and
664 Fig. 2). Additionally, to evaluate the performance in discriminating the effect of medication
665 on gait impairments, we performed the same evaluation for videos recorded during off-
666 and on-medication states. It is noteworthy that the same model trained with the training
667 dataset was used for discriminating medication response in gaits rather than retraining a
668 new model.

669 **Digital biomarker extraction**

670 Before extracting motion makers, we identified the participant's body parts with spa-
671 tiotemporal features strongly correlating with disease severity by examining each joint's
672 contribution to the prediction of UPDRS scores in our model. Since two side-view video
673 segments were used to predict PD severity, we calculated the joint contributions based
674 on both videos. We proposed a dual maximum gradient-weighted class activation mapping
675 (DMGrad-CAM) method based on the traditional Grad-CAM [62], tailored for our Siamese
676 contrastive network architecture. This method allows us to observe and assess the gait
677 features in PD from two perspectives. We used the final spatiotemporal graph convolutional
678 network (ST-GCN) layer as the target for the DMGrad-CAM analysis. To generate the
679 DMGrad-CAM heatmap, we applied Grad-CAM to the model's final layer, producing a

680 spatiotemporal heatmap that assigns an importance score to each joint for every frame.
 681 After generating the Grad-CAM heatmaps for videos recorded from both left and right
 682 perspectives, we selected the joints with the maximum Grad-CAM heatmap for the utility
 683 of videos on both sides. By comparing the results from both perspectives, we calculated
 684 the normalized contribution ratio for each joint in both perspectives. We further used a
 685 sliding average filter with a window size of w ($w = 5$) to process the heatmap matrix and
 686 selected the maximum value of the filtered heatmap for each joint. For each participant,
 687 the joint contributions were estimated based on all three pairs of video segments by taking
 688 the average across them. These contributions were then normalized and aggregated across
 689 categories to assess their correlations with different levels of PD severity.

690 We extracted two types of motion markers associated with joints that contribute most
 691 significantly: traditional clinical markers such as arm swing amplitude, gait speed, and step
 692 length, and novel digital biomarkers indicative of the joint's spatiotemporal movements,
 693 including linear velocity, linear acceleration, joint angles, and the standard deviation of
 694 joint angles. We calculated the means of these gait parameters and their variances across
 695 six video segments, revealing not only the average gait performance but also the gait
 696 variability associated with PD, such as slower walking speed or shorter step length over
 697 walking time. To transform pixel coordinates into real-world measurements, the forearm
 698 length was assumed to constitute a proportion $p = 0.1608$ of the total body height, h_{real} ,
 699 according to the anthropometric analysis [63]. Utilizing the median pixel length, L_i^{pix} , in
 700 the i -th video frame, the pixel-to-real-world scaling factor, s , was determined as

$$701 \quad L_i^{\text{pix}} = \text{median}_i (\| \mathbf{l}_e(i) - \mathbf{l}_w(i) \|) \quad (1)$$

$$702 \quad s = \frac{L_{\text{pix}}}{ph_{\text{real}}} \quad (2)$$

704 where $\mathbf{l}_e(i)$ and $\mathbf{l}_w(i)$ are the 2D pixel coordinates of the elbow and wrist joints in the i -th
 705 frame, respectively.

706 Then, the pixel data $\mathbf{l}(i)$ was scaled into real-world coordinates as

$$707 \quad \mathbf{j}(i) = s\mathbf{l}(i) \quad (3)$$

708 where $\mathbf{j}(i)$ is the real-world 2D coordinates in the i -th frame.

709 Based on the calibrated coordinates, we can easily calculate the spatiotemporal motion
 710 markers. For the traditional clinical gait features, step length was defined as the maximum
 711 horizontal distance between the ankles, and arm swing amplitude as the maximum hori-
 712 zontal displacement between the wrist and hip. An inter-quartile range (IQR [64]) filter

713 was applied to the traditional features to remove outliers before determining the maximum
 714 value. Finally, the mean and variance of all features were calculated across the video
 715 segments to quantify overall gait performance and variability.

716 **Discrimination of medication effectiveness on gait impairment**

717 To effectively discriminate the effect of medical interventions in gait impairments, we
 718 proposed a novel fine-granular assessment score (FGAS) by integrating the UPDRS scores
 719 predicted by the model along with their associated confidence levels. Medication outcomes
 720 (MO) in gait impairments can be predicted as

$$721 \quad \text{MO} = \begin{cases} \text{Improvement} & \text{if } S_c \geq \delta, \\ \text{Deterioration} & \text{if } S_c \leq -\delta, \\ \text{No change} & \text{otherwise.} \end{cases} \quad (4)$$

722 with

$$723 \quad S_c = \alpha_1(S_1^{post} - S_1^{pre}) + \alpha_2(S_2^{post} - S_2^{pre}) + \alpha_3(S_3^{post} - S_3^{pre}) \quad (5)$$

724 where S_c represents the FGAS. S_1, S_2, S_3 denote the confidence levels for the predicted
 725 UPDRS scores 0, 1, and 2, respectively. The superscripts *post* and *pre* indicate post- and
 726 pre-medication states, respectively. The weights assigned to each confidence level, $\alpha_1, \alpha_2,$
 727 and α_3 , were set at 1, 2, and 3, respectively. δ is the threshold to determine the significance
 728 of changes in the gait impairments, which was established at 0.02 to efficiently identify
 729 the medication outcomes.

730 Benchmarking the consensus of the three experts on the medication outcomes, we
 731 calculated the agreement rates of the predicted medication outcomes using our model and
 732 the assessment results of each expert, along with a non-expert clinician. For all experts,
 733 the changed UPDRS scores, rated before and after medication, were first selected as the
 734 medication outcome evaluation for participants who experienced a change in their UPDRS
 735 scores. Subsequently, for participants without a UPDRS score change, the FGAS sub-score
 736 evaluations were used as the medication outcomes in gait impairments.

737 We evaluated the ability of extracted motion markers to characterize personalized
 738 medication responses. We statistically analyzed all motion markers extracted from modified
 739 skeleton data from 19 patients with PD in the test dataset (Table I). We assessed the
 740 normality of the distributions of the values of each marker across 19 patients using
 741 the Shapiro-Wilk test. The independent samples *t*-test and the Kruskal-Wallis test were
 742 employed to analyze the differences in values of motion markers obtained between the
 743 off- and on-medication states for those with and without normal distribution, respectively.

744 We compared the significant changes ($p < 0.05$) in motion markers to the consensus of
745 experts on the medication outcomes in gait impairments for each participant, selecting the
746 top three joints with markers showing the highest agreement rates. Finally, we performed
747 a statistical analysis on the differences in all associated motion markers of the top three
748 joints between the off- and on-medication states for all 19 patients.

749 **Statistical Analysis**

750 Statistical analyses were performed using Python version 3.8 (Python Software Founda-
751 tion). Boxplots are shown with a central mark at the median, bottom, and top edges
752 of the boxes at the 25th and 75th percentiles, respectively, and whiskers extending to
753 the most extreme points within 1.5 times the interquartile range. To assess the inter-
754 group in participant characteristics such as gender, age, height and time since diagnosis
755 across three groups with different UPDRS scores, two-sided χ^2 tests were performed
756 for genders and one-way analysis of variance (ANOVA) for age, height and time since
757 diagnosis. The Kruskal-Wallis test was performed to calculate the statistical significance
758 of extracted traditional clinical motion markers among participants with varying disease
759 stages. A two-sided t-test was conducted to analyze the statistical significance in the
760 extracted motion markers with a normal distribution between off- and on-medication states,
761 while the Kruskal-Wallis test was employed for markers without a normal distribution.
762 Shapiro-Wilk test was used to validate the distribution of each motion marker. Spearman
763 correlation coefficients were calculated to evaluate the correlations between extracted
764 motion markers and UPDRS scores, as well as to compare the medication outcomes rated
765 by the clinicians, the model, and the clinical motion markers. To investigate the sources
766 of model error reflected in the confusion matrices, Fisher's Exact Test was utilized to
767 examine sex differences between accurately predicted and mismatched groups, and the
768 Mann-Whitney U test was employed to assess age and height variations. Additionally,
769 Fisher's Exact Test was also applied to determine whether model errors were significantly
770 associated with discrepancies among expert ratings. The mean and standard deviation of
771 the Spearman correlation coefficients between the coordinates of body joints measured
772 using the Vicon and the DWPose are shown in Supplementary Figure 11.

773 **Report summary**

774 Further information on research design is available in the Nature Research Reporting
775 Summary linked to this article.

776

DATA AVAILABILITY

777 All data associated with this study are presented in the paper or the supplementary materials
778 (Supplementary Information and Data 1). We are unable to provide the original dataset

779 (i.e., gait videos) due to the patient data privacy policy. We provide a toy dataset, including
780 the gait features extracted from the skeleton data and the features extracted using DMGrad-
781 CAM, which are available at the Code Ocean capsule associated with the article (capsule
782 number 7700825), enabling the reproduction of results based on these features.

783 **CODE AVAILABILITY**

784 The code for data preprocessing, model training, validation, testing, DMGrad-CAM calcu-
785 lation, and joint feature extraction is available at Code Ocean (capsule number 7700825).
786 The developed online assessment system (Fast Assessment of Gait Impairments in Parkin-
787 son's Disease, FAGI-PD) is available at gaitanalysis.simplaj.fun. A demonstration of this
788 online system is provided in Supplementary Movie 1.

789 **ACKNOWLEDGEMENTS**

790 This study was supported in part by the National Key Research and Development Pro-
791 gram of China (2022YFB4700200), the National Natural Science Foundation of China
792 (62373202), the Science and Technology Program of Tianjin (23JCYBJC01200 and 24JCZX
793 JC00340), and the Fundamental Research Funds for the Central Universities of China.

794 **AUTHOR CONTRIBUTIONS STATEMENT**

795 J.H., Z.T., and J.W. contributed equally to this work. W.H. conceptualized this study.
796 J.H., Z.T., J.W., and W.H. designed the model framework and experimental protocol. Z.T.,
797 J.W., and P.L. were responsible for recruiting patients and collecting gait data. Z.T., F.B.,
798 R.V., F.M., and W.H. analyzed the data. Z.T., J.H., F.B., R.V., F.M., and W.H. interpreted
799 and discussed the experimental results. Z.T., J.H., K.Z., S.L., and W.H. were responsible
800 for drafting the initial manuscript of the article. Z.T., F.B., F.M., and W.H. revised the
801 manuscript. All authors read the manuscript and provided valuable suggestions for revision.
802 J.H., J.W., and W.H. provided equipment, venues, and software and hardware environments.
803 W.H. provided financial support and led this study, ensuring that the research project
804 proceeded as planned. All authors accept the responsibility to submit it for publication.

805 **COMPETING INTERESTS**

806 Francesca Morgante declares the following competing interests: Research support from
807 NIHR, Innovate UK; Speaking honoraria from Abbvie, Medtronic, Boston Scientific, Bial,
808 Merz; Travel grants from the International Parkinson's disease and Movement Disorder
809 Society; Advisory board fees from Merz and Boston Scientific; Consultancies fees from
810 Boston Scientific, Merz and Bial; Research support from Boston Scientific, Merz; Royalties
811 for the book "Disorders of Movement" from Springer; member of the editorial board of
812 Movement Disorders, Movement Disorders Clinical Practice, European Journal of Neurol-
813 ogy, European Journal of Neurology. The remaining authors declare no competing interests.

- 815 [1] Aarsland, D. *et al.* Parkinson disease-associated cognitive impairment. *Nat. Rev. Dis. Primers.* **7**, 47 (2021).
- 816 [2] Mirelman, A. *et al.* Gait impairments in Parkinson's disease. *Lancet Neurol.* **18**, 697–708 (2019).
- 817 [3] Armstrong, M. J. & Okun, M. S. Diagnosis and treatment of Parkinson disease: a review. *JAMA* **323**, 548–560
- 818 (2020).
- 819 [4] Forte, R., Tocci, N. & De Vito, G. The impact of exercise intervention with rhythmic auditory stimulation to
- 820 improve gait and mobility in Parkinson Disease: an umbrella review. *Brain Sci.* **11**, 685 (2021).
- 821 [5] Mao, Z. *et al.* Comparison of efficacy of deep brain stimulation of different targets in Parkinson's disease: a
- 822 network meta-analysis. *Front. Aging Neurosci.* **11**, 23 (2019).
- 823 [6] Demrozi, F., Bacchin, R., Tamburin, S., Cristani, M. & Pravadelli, G. Toward a wearable system for predicting
- 824 freezing of gait in people affected by Parkinson's disease. *IEEE J. Biomed. Health. Inf.* **24**, 2444–2451 (2019).
- 825 [7] Sabo, A., Mehdizadeh, S., Iaboni, A. & Taati, B. Estimating parkinsonism severity in natural gait videos of older
- 826 adults with dementia. *IEEE J. Biomed. Health. Inf.* **26**, 2288–2298 (2022).
- 827 [8] Bloem, B. R. *et al.* Measurement instruments to assess posture, gait, and balance in Parkinson's disease: Critique
- 828 and recommendations. *Movement Disorders* **31**, 1342–1355 (2016).
- 829 [9] Baron, E. I., Miller Koop, M., Streicher, M. C., Rosenfeldt, A. B. & Alberts, J. L. Altered kinematics of arm
- 830 swing in Parkinson's disease patients indicates declines in gait under dual-task conditions. *Parkinsonism Relat.*
- 831 *Disord.* **48**, 61–67 (2018).
- 832 [10] Rochester, L., Galna, B., Lord, S. & Burn, D. The nature of dual-task interference during gait in incident Parkinson's
- 833 disease. *Neurosci.* **265**, 83–94 (2014).
- 834 [11] Curtze, C., Nutt, J. G., Carlson-Kuhta, P., Mancini, M. & Horak, F. B. Levodopa is a double-edged sword for
- 835 balance and gait in people with Parkinson's disease. *Mov. Disord.* **30**, 1361–1370 (2015).
- 836 [12] Samotus, O., Parrent, A. & Jog, M. Spinal cord stimulation therapy for gait dysfunction in advanced Parkinson's
- 837 disease patients. *Mov. Disord.* **33**, 783–792 (2018).
- 838 [13] Schlick, C. *et al.* Visual cues combined with treadmill training to improve gait performance in Parkinson's disease:
- 839 a pilot randomized controlled trial. *Clin. Rehabil.* **30**, 463–471 (2016).
- 840 [14] Roiz, R. d. M. *et al.* Gait analysis comparing Parkinson's disease with healthy elderly subjects. *Arq. Neuro-*
- 841 *Psiquiatr.* **68**, 81–86 (2010).
- 842 [15] Pachoulakis, I. & Kourmoulis, K. Building a gait analysis framework for Parkinson's disease patients: motion
- 843 capture and skeleton 3d representation. In *2014 International Conference on Telecommunications and Multimedia*
- 844 *(TEMU)*, 220–225 (IEEE, 2014).
- 845 [16] Pistacchi, M. *et al.* Gait analysis and clinical correlations in early Parkinson's disease. *Funct Neurol.* **32**, 28–34
- 846 (2017).
- 847 [17] Burq, M. *et al.* Virtual exam for Parkinson's disease enables frequent and reliable remote measurements of motor
- 848 function. *NPJ Digit. Med.* **5**, 65 (2022).
- 849 [18] Elm, J. J. *et al.* Feasibility and utility of a clinician dashboard from wearable and mobile application Parkinson's
- 850 disease data. *NPJ Digit. Med.* **2**, 95 (2019).
- 851 [19] Islam, M. S. *et al.* Using ai to measure Parkinson's disease severity at home. *NPJ Digit. Med.* **6**, 156 (2023).
- 852 [20] Yang, Y. *et al.* Artificial intelligence-enabled detection and assessment of Parkinson's disease using nocturnal
- 853 breathing signals. *Nat. Med.* **28**, 2207–2215 (2022).
- 854 [21] Zhan, A. *et al.* Using smartphones and machine learning to quantify Parkinson Disease severity: the mobile
- 855 Parkinson Disease score. *JAMA Neurol.* **75**, 876–880 (2018).
- 856 [22] Morinan, G. *et al.* Computer vision quantification of whole-body parkinsonian bradykinesia using a large multi-site
- 857 population. *NPJ Parkinsons Dis.* **9**, 10 (2023).
- 858 [23] Schalkamp, A.-K., Peall, K. J., Harrison, N. A. & Sandor, C. Wearable movement-tracking data identify Parkinson's
- 859 disease years before clinical diagnosis. *Nat. Med.* **29**, 2048–2056 (2023).
- 860 [24] Powers, R. *et al.* Smartwatch inertial sensors continuously monitor real-world motor fluctuations in Parkinson's
- 861 disease. *Sci. Transl. Med.* **13**, eabd7865 (2021).

- 862 [25] Liu, Y. *et al.* Monitoring gait at home with radio waves in Parkinson's disease: a marker of severity, progression,
863 and medication response. *Sci. Transl. Med.* **14**, eadc9669 (2022).
- 864 [26] Endo, M. *et al.* Data-driven discovery of movement-linked heterogeneity in neurodegenerative diseases. *Nat. Mach.*
865 *Intell.* 1034–1045 (2024).
- 866 [27] Zeng, Q. *et al.* Video-based quantification of gait impairments in Parkinson's disease using skeleton-silhouette
867 fusion convolution network. *IEEE Trans. Neural Syst. Rehabilitation Eng.* **31**, 2912–2922 (2023).
- 868 [28] Guo, R., Shao, X., Zhang, C. & Qian, X. Multi-scale sparse graph convolutional network for the assessment of
869 parkinsonian gait. *IEEE Trans. Multimedia* **24**, 1583–1594 (2021).
- 870 [29] Lu, M. *et al.* Quantifying Parkinson's disease motor severity under uncertainty using MDS-UPDRS videos. *Med.*
871 *Image Anal.* **73**, 102179 (2021).
- 872 [30] Kim, K., Lyu, S., Mantri, S. & Dunn, T. W. Tulip: Multi-camera 3d precision assessment of Parkinson's disease.
873 In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 22551–22562 (IEEE, 2024).
- 874 [31] Tian, H. *et al.* Cross-spatiotemporal graph convolution networks for skeleton-based parkinsonian gait MDS-UPDRS
875 score estimation. *IEEE Trans. Neural Syst. Rehabilitation Eng.* **32**, 412–421 (2024).
- 876 [32] Krajushkina, A. *et al.* Gait analysis based approach for Parkinson's disease modeling with decision tree classifiers.
877 In *2018 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, 3720–3725 (IEEE, 2018).
- 878 [33] Peters, D. M. *et al.* Utilization of wearable technology to assess gait and mobility post-stroke: a systematic review.
879 *Journal of neuroengineering and rehabilitation* **18**, 1–18 (2021).
- 880 [34] Muir, S. W. *et al.* Gait assessment in mild cognitive impairment and alzheimer's disease: the effect of dual-task
881 challenges across the cognitive spectrum. *Gait & posture* **35**, 96–100 (2012).
- 882 [35] Guerra, A., D'Onofrio, V., Ferreri, F., Bologna, M. & Antonini, A. Objective measurement versus clinician-based
883 assessment for Parkinson's disease. *Expert Rev. Neurother.* **23**, 689–702 (2023).
- 884 [36] Vasquez-Correa, J. C. *et al.* Multi-view representation learning via gcca for multimodal analysis of Parkinson's
885 disease. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2966–2970
886 (IEEE, 2017).
- 887 [37] Endo, M. *et al.* Gaitforemer: Self-supervised pre-training of transformers via human motion forecasting for few-shot
888 gait impairment severity estimation. In *Medical Image Computing and Computer Assisted Intervention (MICCAI)*,
889 130–139 (Springer Nature Switzerland, Cham, 2022).
- 890 [38] Lu, M. *et al.* Vision-based estimation of MDS-UPDRS gait scores for assessing Parkinson's disease motor severity.
891 In *Medical Image Computing and Computer Assisted Intervention (MICCAI)*, 637–647 (Springer, 2020).
- 892 [39] Mirelman, A. *et al.* Arm swing as a potential new prodromal marker of parkinson's disease. *Movement Disorders*
893 **31**, 1527–1534 (2016).
- 894 [40] Mei, J. *et al.* The comparison of gait disorders among different motor subtypes in Parkinson's disease patients
895 during the early and middle stages. *Clinical Parkinsonism & Related Disorders* **12**, 100309 (2025).
- 896 [41] Goetz, C. G. *et al.* Movement disorder society-sponsored revision of the unified parkinson's disease rating scale
897 (mds-updrs): scale presentation and clinimetric testing results. *Movement disorders: official journal of the Movement*
898 *Disorder Society* **23**, 2129–2170 (2008).
- 899 [42] Goetz, C. G. *et al.* The unified parkinson's disease rating scale (UPDRS): status and recommendations. *Mov.*
900 *Disord.* **18**, 738–750 (2003).
- 901 [43] Kenny, L. *et al.* Inter-rater reliability of hand motor function assessment in parkinson's disease: Impact of clinician
902 training. *Clinical Parkinsonism & Related Disorders* **11**, 100278 (2024).
- 903 [44] Yin, W. *et al.* Gait analysis in the early stage of parkinson's disease with a machine learning approach. *Frontiers*
904 *in Neurology* **15**, 1472956 (2024).
- 905 [45] Kim, K., Lyu, S., Mantri, S. & Dunn, T. W. Tulip: Multi-camera 3d precision assessment of parkinson's disease.
906 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 22551–22562
907 (2024).
- 908 [46] Yang, K. *et al.* Objective and quantitative assessment of motor function in parkinson's disease—from the perspective
909 of practical applications. *Annals of translational medicine* **4**, 90 (2016).

- 910 [47] Narayanaswami, P. The spectrum of functional rating scales in neurology clinical trials. *Neurotherapeutics* **14**,
911 161–175 (2017).
- 912 [48] Johnson, S. & Everingham, M. Learning effective human pose estimation from inaccurate annotation. In *Proc.*
913 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 1465–1472 (IEEE, 2011).
- 914 [49] Celik, Y., Stuart, S., Woo, W. L., Sejdic, E. & Godfrey, A. Multi-modal gait: A wearable, algorithm and data
915 fusion approach for clinical and free-living assessment. *Information Fusion* **78**, 57–70 (2022).
- 916 [50] Formstone, L. *et al.* Quantification of motor function post-stroke using novel combination of wearable inertial and
917 mechanomyographic sensors. *IEEE Transactions on Neural Systems and Rehabilitation Engineering* **29**, 1158–1167
918 (2021).
- 919 [51] Guo, Z., Wang, Z., Wang, Y., Huo, W. & Han, J. Continuous estimation of swallowing motion with emg and mmg
920 signals. *IEEE Transactions on Neural Systems and Rehabilitation Engineering* **33**, 787–797 (2025).
- 921 [52] Yang, Z., Zeng, A., Yuan, C. & Li, Y. Effective whole-body pose estimation with two-stages distillation. In *Proc.*
922 *IEEE International Conference on Computer Vision (ICCV)*, 4210–4220 (IEEE, 2023).
- 923 [53] Cao, Z., Simon, T., Wei, S.-E. & Sheikh, Y. Realtime multi-person 2d pose estimation using part affinity fields.
924 In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 7291–7299 (IEEE, 2017).
- 925 [54] Xu, L. *et al.* Zoomnas: Searching for whole-body human pose estimation in the wild. *IEEE Transactions on*
926 *Pattern Analysis and Machine Intelligence* **45**, 5296–5313 (2023).
- 927 [55] Salvador, S. & Chan, P. Toward accurate dynamic time warping in linear time and space. *Intell. Data Anal.* **11**,
928 561–580 (2007).
- 929 [56] Koch, G., Zemel, R., Salakhutdinov, R. *et al.* Siamese neural networks for one-shot image recognition. In
930 *International Conference on Machine Learning (ICML) Workshops*, vol. 2, 1–30 (Lille, 2015).
- 931 [57] Vaswani, A. *et al.* Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*,
932 vol. 30 (Curran, 2017).
- 933 [58] Zhu, W. *et al.* Motionbert: A unified perspective on learning human motion representations. In *Proc. IEEE*
934 *International Conference on Computer Vision (ICCV)*, 15085–15099 (IEEE, 2023).
- 935 [59] Yan, S., Xiong, Y. & Lin, D. Spatial temporal graph convolutional networks for skeleton-based action recognition.
936 In *Proc. AAAI Conference on Artificial Intelligence (AAAI)*, vol. 32 (AAAI Pres, 2018).
- 937 [60] Rebuffi, S.-A. *et al.* Data augmentation can improve robustness. In *Advances in Neural Information Processing*
938 *Systems (NeurIPS)*, vol. 34, 29935–29948 (Curran, 2021).
- 939 [61] Li, B., Hou, Y. & Che, W. Data augmentation approaches in natural language processing: a survey. *Ai Open* **3**,
940 71–90 (2022).
- 941 [62] Selvaraju, R. R. *et al.* Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proc.*
942 *IEEE International Conference on Computer Vision (ICCV)*, 618–626 (IEEE, 2017).
- 943 [63] Hoang, K.-L. H. & Mombaur, K. Adjustments to de leva-anthropometric regression data for the changes in body
944 proportions in elderly humans. *Journal of biomechanics* **48**, 3732–3736 (2015).
- 945 [64] Laurikkala, J. *et al.* Informal identification of outliers in medical data. In *Fifth international workshop on intelligent*
946 *data analysis in medicine and pharmacology*, vol. 1, 20–24 (2000).

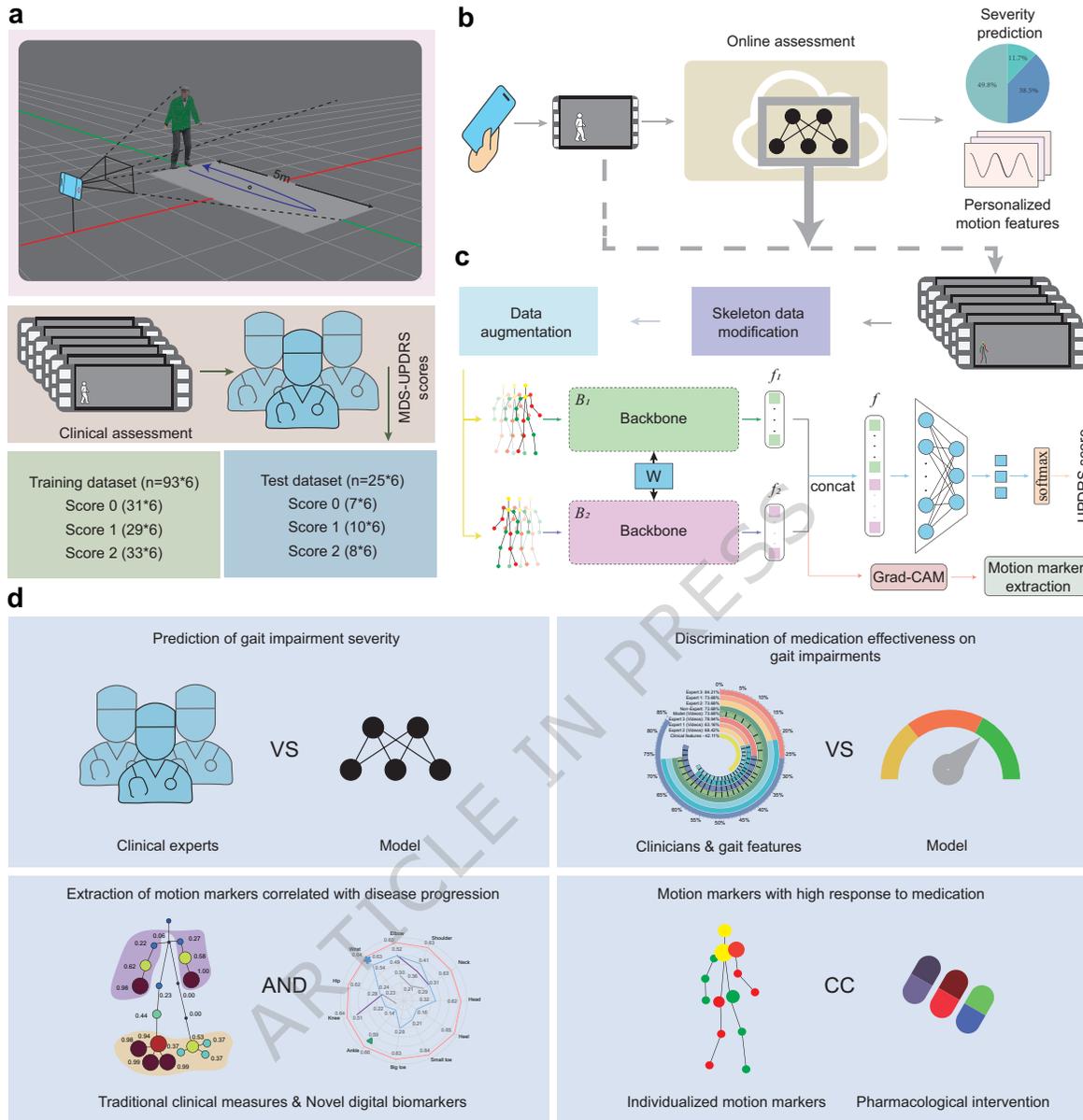


Fig. 1: Overview of the AI-based assessment of gait impairments using smartphone videos. **a**, The participants performed the shuttle walk three times over a 5-meter distance (i.e., six repetitions), and we used a smartphone to film their whole-body movements from the lateral perspective. Three clinical specialists independently rated the severity of gait impairments for each participant according to the MDS-UPDRS Part III-Gait exams. We randomly selected the gait videos from 93 participants (6 video segments per participant) to train the model and used the remaining ones from 25 participants for testing. **b**, We developed an online assessment system for clinical and home-based assessments of gait impairments based on smartphone-recorded gait videos. **c**, We designed a Siamese contrastive deep-learning network framework for predicting the UPDRS scores and extracting digital biomarkers. The recorded videos were automatically segmented into six parts, corresponding to six walking repetitions. The model was trained using the UPDRS scores rated by clinicians and augmented skeleton data extracted from the video segments with spatial augmentation. Skeleton data from videos recorded from both left and right perspectives were inputs for the two identical backbone networks. **d**, We evaluated the model's capabilities to 1) predict gait impairment severity, 2) discriminate medication effect on gait impairments, 3) extract motion markers correlated with disease progression, and 4) identify motion markers with high response (i.e., high correlation coefficients (CC)) to medication.

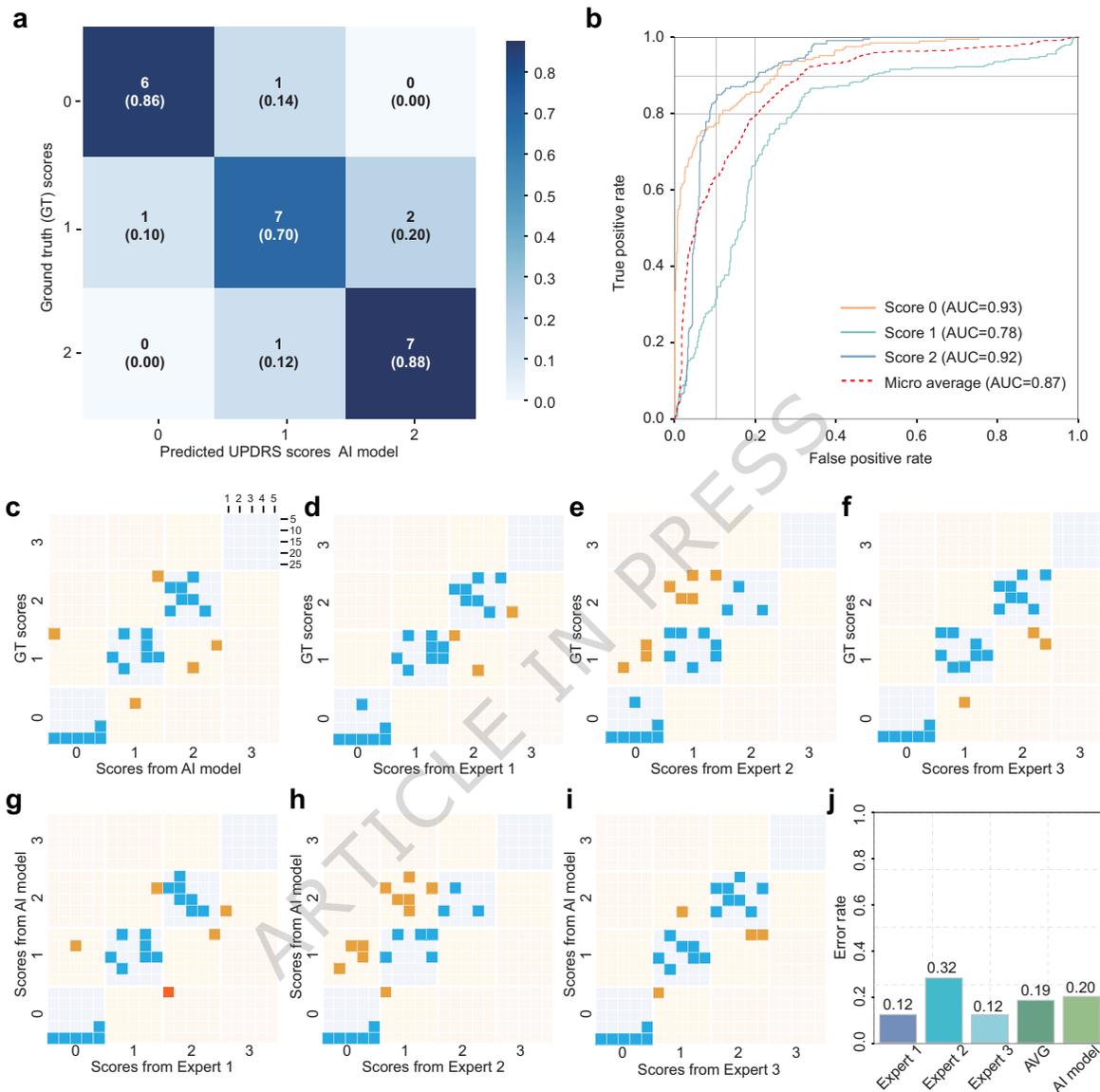


Fig. 2: Model alignments with ground truth and UPDRS scores rated by experts. **a**, Confusion matrix of the UPDRS score prediction on the test dataset. **b**, The receiver operating characteristic (ROC) curve for each severity category. The model achieved a high area under the ROC curve (AUC) with a micro-average AUC of 0.87. **c-f**, Agreements between the ground truth scores and those rated by the AI model and three experts for the test dataset. Within the large squares, each small square represents a participant as depicted in the top right corner of **c**. Blue, yellow and red squares indicate perfect agreements, one-score discrepancy, and two-score discrepancies, respectively. **g-i**, Agreements between the scores predicted by the AI model and the ones rated by the three experts. **j**, The error rates of predicted UPDRS scores of the experts and the AI model, compared to the ground truth. AVG represents the average error rate of three experts. The error rate of the AI model is close to the average one of three experts (0.2 vs 0.19).

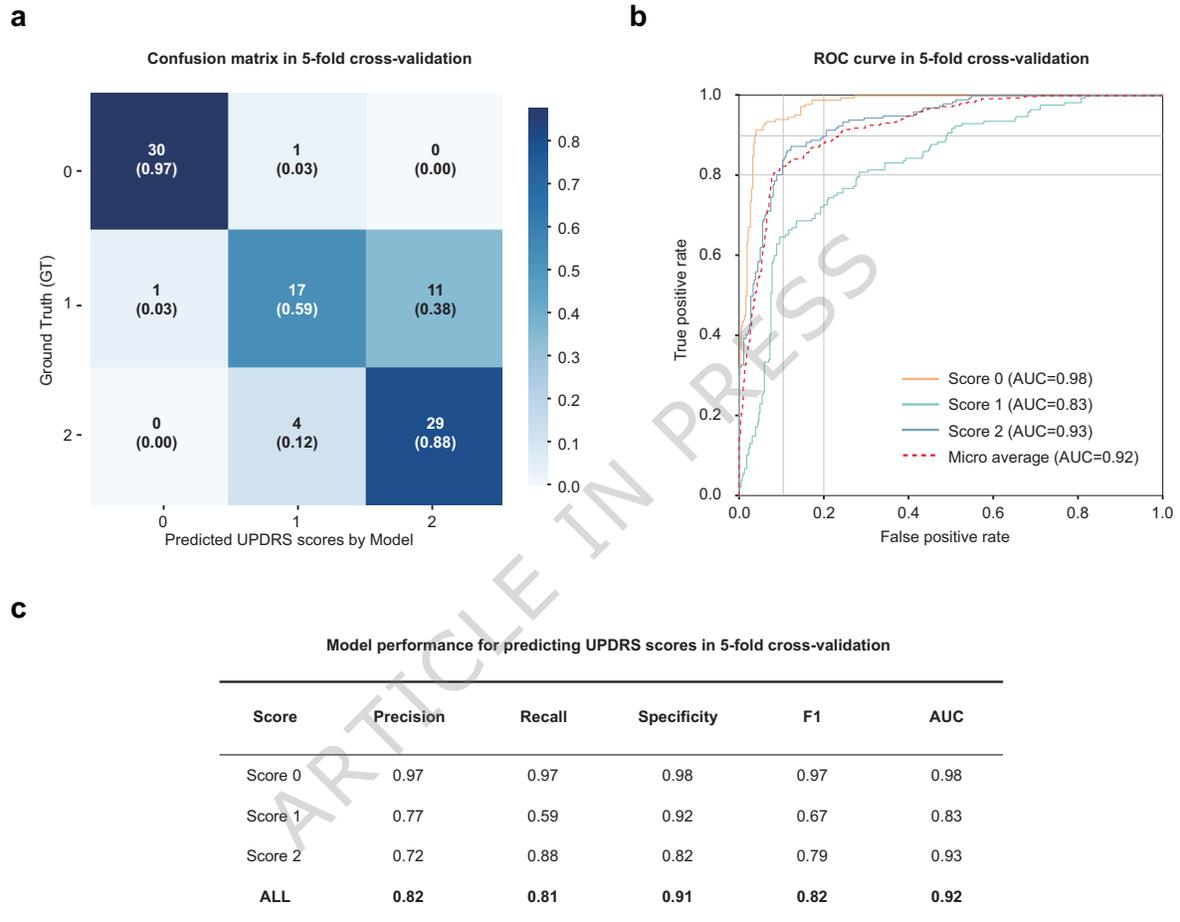


Fig. 3: Model performance in predicting gait impairment severity among 93 participants within the training dataset. **a**, The confusion matrix of the UPDRS score prediction. **b**, The receiver operating characteristic (ROC) curve for each severity category, as well as the area under the ROC curve (AUC) and the micro-average AUC. **c**, Performance metrics including macro precision, recall, specificity, F1 score and AUC. Except for recall, which remained constant, all other performance metrics exhibited slight improvements compared to those with the test dataset (Table II and Fig.2). Notably, the AUC increased by 0.05.

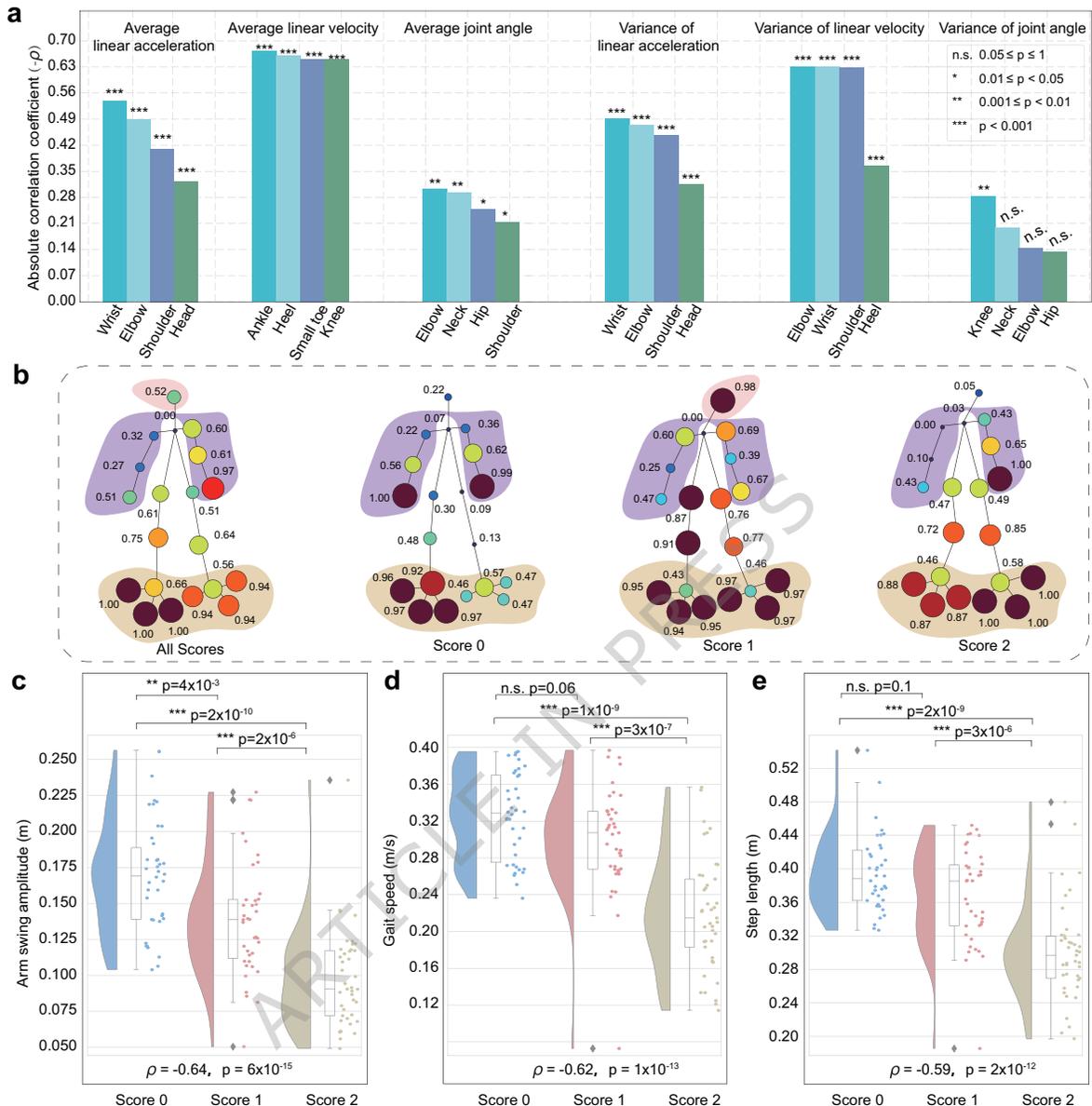


Fig. 4: Model performance in extracting motion markers of disease progression. **a**, The four joints with the highest Spearman correlations to UPDRS scores among all 20 joints for each type of extracted spatiotemporal biomarkers (additional information on the remaining joints are shown in Supplementary Figure 6). The ranges of the p -values of the Spearman correlations for each joint are displayed on the corresponding bars. We selected the largest correlation coefficients between the left and right joints for bilaterally symmetrical joints. **b**, Joint contributions to predict the severity of gait impairment cross all UPDRS scores and each score. The joint contributions were normalized to represent the relative contribution ratios among joints for severity prediction, indicated through points with different sizes and colors. Larger and red points indicate greater normalized contributions, whereas smaller and blue joints denote lesser normalized contributions. A contribution value of 0 means the lowest contribution ratio within the analyzed joints rather than the absence of contribution. **c-e**, The Spearman correlations between the UPDRS scores and the extracted motion markers, along with the significant differences in these markers across three disease severity categories. Spearman correlation coefficients and their p -values are displayed at the bottom of each graph. We used the Kruskal-Wallis test to analyze the statistical significance. p -values of the significance analyses between different groups are presented at the top of each graph.

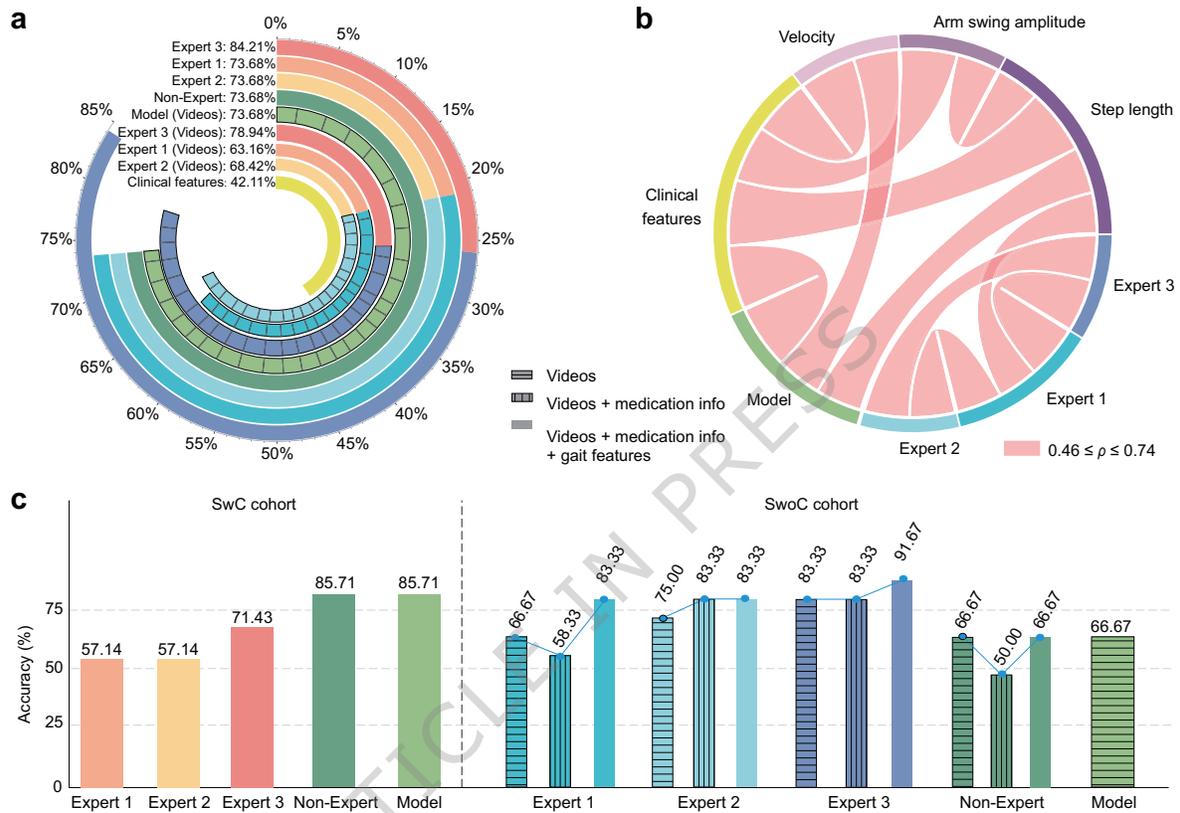


Fig. 5: Model performance in discriminating the medication effect on gait impairments. **a**, The accuracies in discriminating the patient outcomes by three expert clinicians, a non-expert clinician, the AI model, and the significant changes in the traditional clinical motion markers. For the sub-UPDRS score assessment, three experts performed three rounds of independent assessments of the patient outcomes with different sets of information: 1) solely gait videos, 2) gait videos accompanied by medication status and disease details (Table I), 3) gait videos with medication status and disease details, supplemented by values of traditional motion markers measured in both off and on medication states (Fig. 4). **b**, The Spearman correlation coefficients (ρ) between the patient outcomes discriminated by clinicians, AI model, and the clinical motion markers, with thicker chords indicating stronger correlations. **c**, Discrimination accuracies of the three experts, the non-expert clinician, and the AI model for different cohorts: 1) 7 patients with changed UPDRS scores after medication (SwC cohort), 2) 12 patients without changed UPDRS scores after medication (SwoC cohort). Note that for assessing the medication's effect on gait impairments in the SwC cohort, the non-expert clinician used a granular three-level sub-score rating approach. The three experts' ratings were still derived from changes in the standard UPDRS scale.

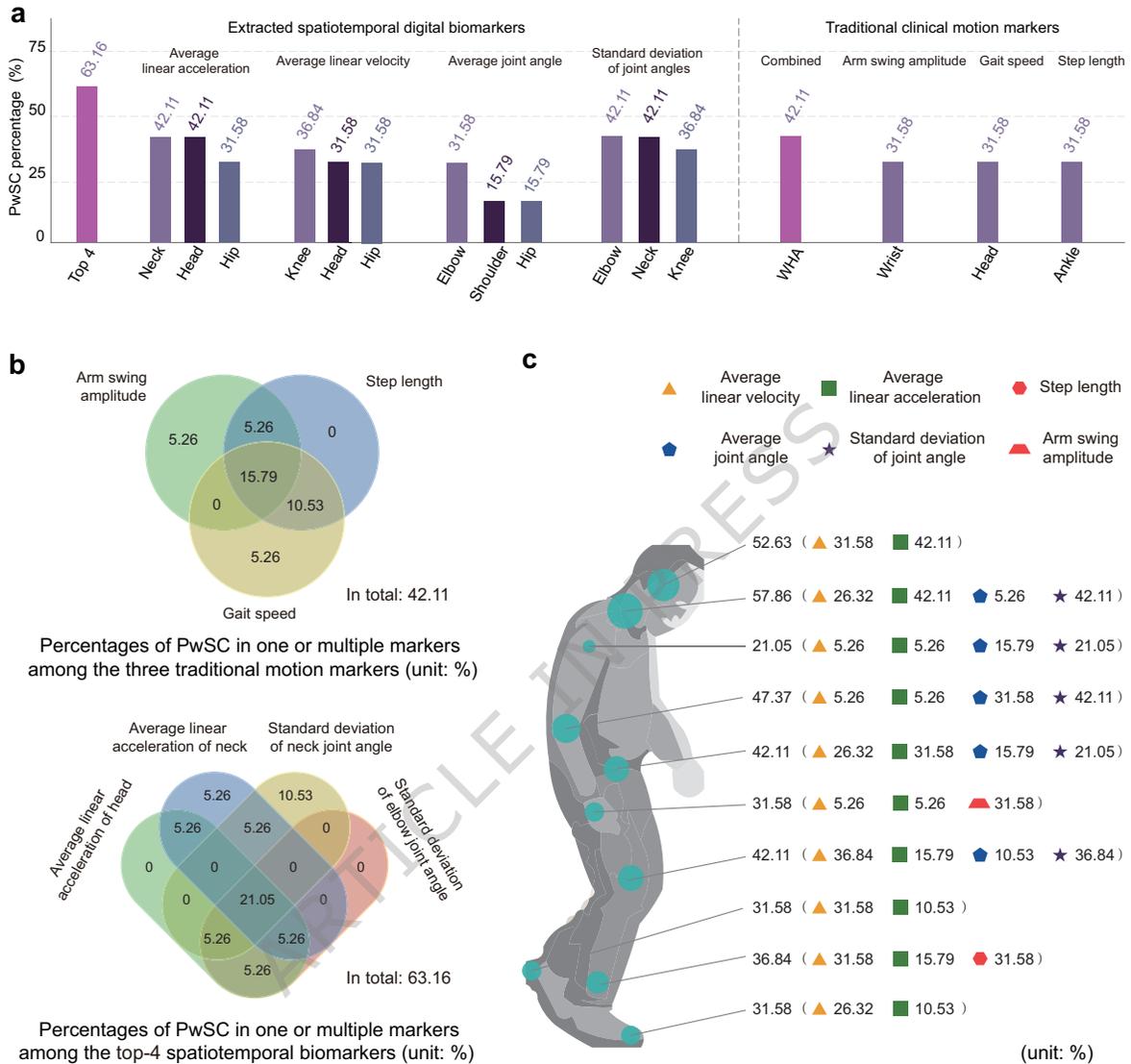


Fig. 6: Model performance in identifying digital biomarkers of medication response. **a**, Percentages of patients with PD ($N=19$) with significant changes (PwSC) ($p < 0.05$) in different motion markers between off- and on-medication states. Meanwhile, we only considered those with significant changes that agree with the medication efficacy rated by clinical specialists. For each type of newly extracted spatiotemporal digital biomarker, we present the top three joint points with the highest percentage of patients. The Top-4 indicates the total percentage of patients with significant changes in at least one of the top four spatiotemporal markers. In contrast, the WHA represents those for the combined traditional clinical motion markers. A voting was used when the markers displayed inconsistent significant changes. **b**, Proportions of patients with PD ($N=19$) exhibiting significant changes ($p < 0.05$) in one or multiple markers among the Top-4 spatiotemporal and WHA traditional markers. **c**, Percentages of patients ($N=19$) showing significant changes ($p < 0.05$) in motion markers associated with different joints (the motion markers' changes should be consistent with the rated medication outcomes).

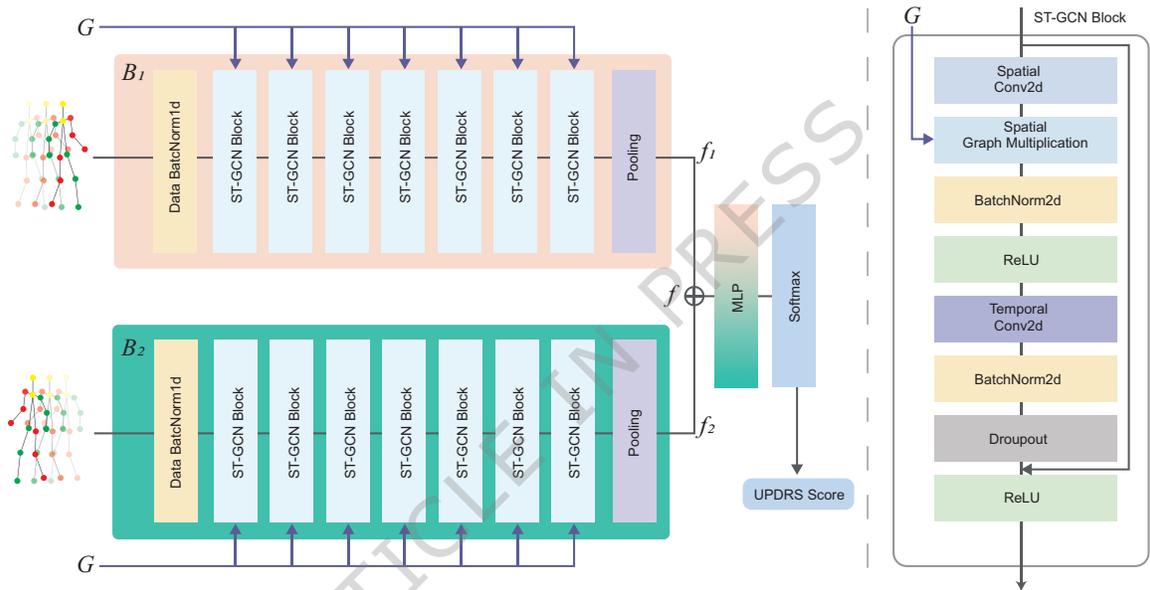


Fig. 7: **Siamese contrastive network architecture.** Skeleton data extracted from gait videos recorded from the left and right perspectives are fed into two identical backbone networks, B_1 and B_2 , which share the same weights. G represents the spatial topology of the data, capturing the spatial dependencies among nodes. These dependencies are propagated through the network via graph convolution operations. Each backbone network extracts feature vectors, denoted as f_1 and f_2 , which are subsequently concatenated into a unified vector, f . This composite vector is then processed through fully connected layers followed by a softmax layer, yielding a probabilistic distribution across three distinct classes. The right part shows the structure of the spatial-temporal graph convolutional (ST-GCN) network.

TABLE I: Characteristics of the participants

Gait impairment score	All	0	1	2	<i>P</i> value
<i>N</i>	118	38	39	41	
Sex, (<i>F/M</i>)	55/63	21/17	15/24	19/22	0.34
Age, years (mean±SD)	62.9 ± 7.4	61.1 ± 5.0	62.8 ± 7.9	64.7 ± 8.4	0.09
Height, cm (mean±SD)	166.2 ± 8.4	165.4 ± 7.2	168.0 ± 8.8	165.1 ± 9.0	0.26
Off/On Medication Test (<i>N</i>)	19	0	11	8	

We recruited 87 patients with PD and 31 healthy participants and classified them into three categories based on their gait impairment severity, assessed according to the MDS-UPDRS Part III-Gait scales. Here, we considered the UPDRS scores of healthy participants to be zero. In clinical practice, the severity of gait impairments of some patients with mild PD can be diagnosed as a UPDRS score of zero. To align with this, we incorporated six patients with a UPDRS score of 0 into the dataset. We used two-sided χ^2 tests to analyze the statistical significance of inter-group differences for genders and one-way analysis of variance (ANOVA) for ages and heights. The *p* values show that the dataset has a balanced distribution of these categorical items across three groups ($p > 0.05$). Nineteen patients with PD in the independent test dataset further participated in an Off/On-medication test, in which their gait videos were recorded during off and on medication states, respectively.

TABLE II: Comparison of gait impairment severity predictions between the model and three clinical experts

Examiner	Precision	Recall	Specificity	F1 score
Expert 1	0.926	0.892	0.961	0.904
Expert 2	0.761	0.692	0.833	0.668
Expert 3	0.897	0.886	0.939	0.885
Expert-Average	0.861	0.823	0.911	0.819
AI Model	0.804	0.811	0.898	0.806

Evaluation with an independent test dataset of 25 participants. Performance was evaluated using macro F1 score, precision, recall, and specificity. The consensus of three clinical specialists on the UPDRS scores of participants' gaits served as the gold standard (Table I). Expert-Average represents the average value of each performance metric across three experts. The F1 score shows a comprehensive performance of the predictions.

948

ADDITIONAL INFORMATION

949 **Correspondence and requests for materials** should be addressed to Weiguang Huo.

ARTICLE IN PRESS