

1
2
3
4 1 **Dual task turning in place: a reliable, valid and responsive outcome**
5
6 2 **measure of freezing of gait**
7
8

9 3 Nicholas D'Cruz, PhD,¹ Jana Seuthe, PhD,^{2,3} Clara de Somer, MSc,¹ Femke Hulzinga, MSc,¹
10 4 Pieter Ginis, PhD,¹ Christian Schlenstedt, PhD,^{2,3} Alice Nieuwboer, PhD¹
11
12

13
14 5 ¹ *KU Leuven, Department of Rehabilitation Sciences, Neurorehabilitation Research Group, B-3000 Leuven,*
15 6 *Belgium*
16
17

18
19 7 ² *Christian-Albrechts-University (CAU) Kiel, University Hospital Schleswig-Holstein, Department of Neurology,*
20 8 *Kiel, Germany*
21
22

23 9 ³ *MSH Medical School Hamburg, Section Exercise and Neurosciences, Hamburg, Germany*
24
25

26 10 **Corresponding Author:** Nicholas D'Cruz, Department of Rehabilitation Sciences, KU
27 11 Leuven, Tervuursevest 101 Bus 1501, B-3001 Leuven, Belgium. Ph: +3216376003 Email:
28 12 nicholas.dacruz@kuleuven.be
29
30
31

32
33 13 **Word count:** Abstract: 248 Manuscript: 3849
34

35
36 14 Figures: 2 Tables: 3 Supplementary tables: 5 Supplementary figures: 3
37
38

39 15 **Running Title:** Objective measurement of FOG
40
41

42 16 **Keywords:** Gait disorders, Dual tasking, Outcome measures, Split-belt treadmill, Interrater
43 17 reliability
44
45

46
47 18 **Competing interests:** The authors declare that there are no potential conflicts of interest in
48 19 regard to this work.
49
50

51
52 20 **Funding sources:** Research grant from the Jacques and Gloria Gossweiler Foundation
53
54
55

56 21
57
58
59
60

Abstract

Background: Freezing of gait (FOG) is a complex symptom in Parkinson's disease (PD) which is both elusive to elicit and varied in its presentation. These complexities present a challenge to measuring FOG in a sensitive and reliable way, precluding therapeutic advancement.

Objective: We investigated the reliability, validity and responsiveness of manual video-annotations of the turning in place task and compared it to the sensor-based FOG ratio.

Methods: Forty-five optimally medicated people with PD and FOG performed rapid alternating 360° turns without and with an auditory stroop dual task, thrice over two consecutive days. Tasks were recorded with video and inertial sensors placed on the lower back and shins. Interrater reliability between three raters, criterion validity with self-reported FOG, and responsiveness to single-session split-belt treadmill (SBT) training were investigated and contrasted with the sensor-based FOG ratio.

Results: Visual ratings showed excellent agreement between raters for the percent time frozen (%TF) (ICC = 0.99), the median duration of a FOG episode (ICC = 0.90) and the number of FOG episodes (ICC = 0.86). Dual tasking improved the sensitivity and validity of visual FOG ratings resulting in increased FOG detection, criterion validity with self-reported FOG ratings, and responsiveness to a short SBT intervention. The sensor-based FOG ratio on the other hand, showed complex FOG presentation-contingent relationships with visual and self-reported FOG ratings, and limited responsiveness to SBT training.

Conclusions: Manual video-annotations of FOG during dual task turning in place generate reliable, valid, and sensitive outcomes for investigating therapeutic effects on FOG.

46 **Introduction**

47 Freezing of gait (FOG) is a disabling symptom for persons with Parkinson's disease (PwPD).

48 FOG is defined as the "brief, episodic absence or marked reduction of forward progression of

49 the feet despite the intention to walk."¹ Due to these unpredictable interruptions to movement,

50 PwPD who experience FOG are at a higher risk of falling², and its unwanted consequences³.

51 In fact, a staggering 61% of falls in PD are directly attributable to FOG⁴. Rehabilitation in

52 conjunction with medical therapies^{5,6} aimed at reducing FOG severity or delaying its onset are

53 therefore urgently needed. A novel and relevant training for FOG involves split-belt treadmill

54 (SBT) whereby gait perturbations are imposed by driving the belts at different speeds.

55 Repeated exposure to such perturbations, even within one session, has been found to reduce

56 gait asymmetry, enhance adaptation to asymmetry and improve turning speed⁷⁻⁹, but it is

57 unclear whether it also reduces FOG.

58 In order to test intervention effects on FOG, reliable, valid and sensitive measures of FOG

59 severity are required. Quantification of FOG severity is frequently performed with subjective

60 rating scales such as the Freezing of Gait Questionnaire and the New Freezing of Gait

61 Questionnaire (NFOGQ)^{10,11}. Although these measures have shown adequate validity to detect

62 the presence of FOG, recent work suggests that retest reliability of subjective measures is

63 poor,¹² rendering them less suitable for measuring FOG severity. Given the sensitivity of

64 subjective measures to recall and expectation bias, objective measures of FOG severity are

65 increasingly being sought after^{13,14}.

66 Objective quantification of FOG uses FOG provoking tasks¹⁵, often under stress¹⁶, to elicit

67 freezing episodes which can be quantified using manual or automated methods. Manually

68 labelling FOG episodes from video recordings of Timed Up and Go tasks to quantify the

69 percentage of the task time with FOG (percent time frozen) is considered the gold-standard for

70 measuring FOG severity¹⁷. However, poor sensitivity to FOG and limited response to both

1
2
3 71 medication and training¹⁸ have raised questions about its utility as an outcome for clinical trials.
4
5 72 Subsequent work has shown that full and fast turning in place is more sensitive in eliciting
6
7 73 FOG^{19,20}, and this sensitivity can be further improved by dual tasking²¹. So far, determining
8
9 74 criteria for annotating FOG episodes while turning in place and their validation as an outcome
10
11 75 of FOG severity has not been undertaken. FOG annotation is challenging during alternating
12
13 76 fast 360° turns, as these do not show the same stepping patterns as straight-line walking tasks.
14
15 77 This is important, as reliability of labelling FOG episodes will have an impact on the
16
17 78 measurement error and subsequent usefulness of the FOG metric as an outcome for
18
19 79 intervention.

20
21
22
23
24 80 More recently, an automated sensor-based metric has been proposed to detect FOG during
25
26 81 turning in place, using automated algorithms based on temporal decomposition of the leg
27
28 82 movement signal²⁰. The FOG ratio has been used to quantify FOG severity in several
29
30 83 observational and interventional studies^{22–24} and has undergone initial validation²². However,
31
32 84 the FOG ratio may be affected by turning fragmentation and less affected by periods of
33
34 85 complete absence of movement²⁵, thus reducing its FOG sensitivity and thereby, its
35
36 86 responsiveness.

37
38
39
40
41 87 In this study, we developed standardized criteria for visually rating FOG episodes from 360°
42
43 88 turning in place videos and investigated the interrater reliability when these were applied to an
44
45 89 optimally medicated representative sample of PwPD and FOG. As a secondary objective, we
46
47 90 investigated criterion validity through associations with a validated subjective rating scale –
48
49 91 the NFOGQ, as well as automated sensor-based ratings of FOG severity. Finally, we compared
50
51 92 the responsiveness of the visual-rated and sensor-based metrics to the effects of a split-belt
52
53 93 treadmill intervention. We expected that visual-rated outcomes would show validity with both
54
55 94 subjective and sensor-based FOG measures, but owing to higher FOG-sensitivity, would be
56
57 95 more responsive than sensor-based outcomes of FOG severity.
58
59
60

96 **Methods**

97 2.1 Participants

98 Included participants were a subset of a larger multi-center study to investigate the short-term
99 effects of split-belt treadmill training on gait, turning and adaptation (ClinicalTrials.gov
100 NCT03725215). Forty-five people with Parkinson's disease and freezing of gait (PD+FOG)
101 were included in this study. Eligibility criteria included PD diagnosis based on the UK Brain
102 Bank Criteria, presence of FOG based on self-reported answer to the question "did you
103 experience freezing of gait in the past month", and the ability to walk unassisted for at least 5
104 minutes. Exclusion criteria included marked cognitive impairment ($MMSE \leq 24$),
105 cardiovascular risk for exercise, musculoskeletal disorders affecting gait, or recent changes in
106 Parkinson's medication or deep brain stimulation settings (< 1 month). Ethical approval was
107 obtained from the respective Institutional Review Boards and all participants provided written
108 informed consent prior to enrollment in the study. Measurements and training were performed
109 in the optimally medicated state and repeated measurements were standardized in relation to
110 medication intake.

111 2.2 Study design

112 The study utilized a randomized parallel design with one control and three active arms. FOG
113 provoking tests were performed at three moments over two days – pre and post intervention on
114 day one, and once the following day (retention). In addition, clinical questionnaires were
115 administered once to characterize participants' cognitive, balance and motor disease severity.
116 Subjective FOG severity was characterized using the NFOGQ, which enquires about the
117 severity and impact of FOG in the past month¹¹.

118 2.3 FOG provoking task

1
2
3 119 Based on previous work²², the turning in place task for 60 seconds was performed to provoke
4
5 120 FOG. Participants were instructed to turn in place as quickly and as safely possible, alternating
6
7 121 direction after each full turn, and to take steps rather than pivot on one leg. The task was
8
9 122 performed without and with a cognitive dual task (in that order), namely the auditory stroop
10
11 123 task, delivered through a wireless headset. Inertial measurement units (Opals, APDM, Portland,
12
13 124 USA) were placed on the shins and lower back of the participants to capture objective turning
14
15 125 and FOG metrics. One video camera providing a single-angle (diagonal to starting position)
16
17 126 neck-to-foot view of the participant captured the trial for subsequent rating.
18
19
20
21

22 127 2.4 Visual FOG rating

23
24
25 128 Three raters (CDS, JS and ND) annotated the video recordings with the ELAN toolbox (version
26
27 129 5.8, Max Planck Institute for Psycholinguistics, The Netherlands), based on recent
28
29 130 recommendations²⁶. Raters were blinded to the time point, and whether it was single or dual
30
31 131 task (by muting sound). Criteria and definitions for all labels were developed in two iterations.
32
33 132 In the first iteration, using established criteria, 40 turning in place trials were rated between
34
35 133 two pairs of raters (20 trials for rater A and B, and 20 trials for rater B and C) and interrater
36
37 134 reliability was evaluated. Despite fair to excellent interrater reliability (ICC > 0.93, 95% CI
38
39 135 between 0.78 and 0.98 for both rater pairs), limits of agreement were large (18.1% and 37.7%
40
41 136 for the two rater pairs), therefore the decision was made to revise the criteria and re-annotate
42
43 137 the videos. Chief sources of variation in ratings were attributed to difficulty in labelling the
44
45 138 start and end of the task when it began or ended with akinetic FOG, as well as varied
46
47 139 interpretation of the “ineffective step” to label the start and end of the FOG episodes. These
48
49 140 criteria were therefore changed between iterations (Table 1). For the second iteration, 20 trials
50
51 141 with the largest interrater differences in the first iteration on the percent time frozen were rated
52
53 142 by all raters, and interrater reliability was reassessed. The remaining trials were randomly
54
55 143 distributed among raters and outcomes were calculated. Percent time frozen (%TF) was the
56
57
58
59
60

1
2
3 144 primary outcome, with number of FOG episodes and median duration of a FOG episode as
4
5 145 secondary outcomes. To minimize potential overfitting of criteria development to the 40 videos
6
7
8 146 earlier labelled, a random sample of 10 trials was re-evaluated after 15 months.
9

10 147 2.5 Sensor-based FOG rating

11
12
13 148 The FOG ratio as described by Mancini et al. (2017) was calculated from the power spectral
14
15
16 149 density (PSD) of the anterior-posterior acceleration signal in the shin sensors (sampled at
17
18 150 128Hz). Using a four-second Hanning window, the PSD was calculated with the Welch
19
20
21 151 method. The ratio of the square of the power within the freezing band (3 – 8 Hz) to the square
22
23 152 of the power within the movement band (0.5 to 3 Hz), averaged over the trial and over the two
24
25 153 legs, gave the FOG ratio. Further, the FOG ratio calculated from the mediolateral (task
26
27 154 movement direction) acceleration signal was also obtained. Objective measures of turning
28
29
30 155 performance were calculated from the lower back sensor, including mean and peak turning
31
32 156 speed (yaw angular velocity) and mediolateral jerkiness (measure of turning fluidity)²¹.
33
34

35 157 2.6 Training intervention

36
37
38 158 The training comprised of one session of 30 minutes of walking on either a split-belt treadmill
39
40 159 with both belts going at the same speed (tied-belt mode – TBT) or with each belt going at a
41
42 160 different speed (split-belt mode – SBT). Three SBT conditions were used to compare
43
44
45 161 effectiveness of the various modes, and participants were randomly assigned to receive TBT
46
47 162 or any one of the SBT conditions. Two SBT conditions differed in the extent of speed reduction
48
49 163 on the slow belt (25% or 50% slower) and one condition switched between the two speeds.
50
51

52 164 **Statistical analysis**

53 54 55 165 3.1 Interrater reliability analysis

56
57
58
59
60

1
2
3 166 Interrater reliability was assessed using the intra-class correlation coefficient (ICC) with a two-
4
5 167 way random effects analysis (random trials, random raters) for the absolute agreement between
6
7
8 168 raters for a single measurement (ICC(2,1)). Wilcoxon signed-rank tests, Bland-Altman plots
9
10 169 and limits of agreement were investigated for systematic bias between pairs of raters.
11
12

13 170 3.2 Criterion validity of visual FOG rating

14
15
16 171 To investigate criterion validity of the visual ratings, we performed correlation analyses with
17
18 172 the NFOG-Q and FOG ratio at pre-training. Distribution of the FOG outcomes were assessed
19
20 173 with histograms and found to be highly skewed to the right, so spearman rank correlation was
21
22 174 performed. Non-parametric bootstrapping (1000 resamples, unrestricted random sampling)
23
24 175 was used to estimate confidence intervals.
25
26
27

28 176 3.3 Responsiveness to treatment

29
30
31 177 To investigate if the FOG metrics were responsive to treatment, we tested whether the visual-
32
33 178 rated or sensor-based metrics demonstrated treatment effects. Constrained longitudinal data
34
35 179 analysis implemented in a linear mixed model framework was applied to investigate changes
36
37 180 within and between TBT and any SBT condition from pre-training to retention⁹. Data were
38
39 181 transformed with an inverse hyperbolic sine function to reduce skewness while allowing
40
41 182 inclusion of zero scores. Normality of model residuals were visually assessed with histograms
42
43 183 and QQ plots.
44
45
46
47

48 184 3.4 Exploratory analysis – Sensor-based turning metrics related to FOG over time

49
50
51 185 To investigate sensor-based turning metrics related to visual-rated FOG, we performed non-
52
53 186 parametric repeated measures correlation between visual-rated FOG and various turning
54
55 187 metrics. 1000 random permutations were used to obtain null distribution of the resulting Z
56
57 188 score and calculation of p-values²⁷.
58
59
60

189 **Results**

190 4.1 Participant characteristics and missing data

191 Forty-five PD+FOG, Hoehn & Yahr stages I - IV were included in this study. Participants
192 varied in their ages (mean: 68.6 years, range: 48 – 86), disease duration (mean: 12.8 years,
193 range: 1 – 38) and self-reported freezing severity (NFOGQ mean: 16.3, range: 6 - 29)
194 (Supplementary Table 1). No significant differences were found for any of the cognitive,
195 balance, disease-related or training-intensity measures between the four training groups
196 (Supplementary Table 2). Due to technical difficulties with video capture, seven participants
197 did not have recordings available, hence visual ratings are reported on the 38 remaining
198 participants.

199 4.2 Interrater reliability of visual FOG rating

200 Twenty videos from with varying amounts and presentations of FOG were rated to assess the
201 inter-rater reliability (mean and ranges of rater means – number of episodes: 4.55 (2 – 13);
202 median duration of an episode in seconds: 6.01 (0.4 - 27); %TF: 40.5 (1.35 - 90.1)). Intraclass
203 correlation coefficient for *absolute agreement* of %TF for a single rater (ICC (2,1)) using a 2-
204 way mixed effects model was 0.993 (95% CI: 0.986 – 0.997). Reliability was slightly lower
205 for other FOG outcomes, with the number of FOG episodes showing the lowest reliability.
206 Importantly, the standard error of measurement was less than 3% for the %TF (Table 2).

207 Bland-Altman plots revealed no systematic error across FOG severity, although scores tended
208 to be more spread out around medium severity (35 – 60 %TF) trials. Only number of episodes
209 showed a statistically significant difference between raters, mainly due to multiple short FOG
210 episodes being pooled into longer episodes (Rater A vs Rater C - Wilcoxon signed rank test Z
211 = -2.167, $p = 0.030$) (Figure 1). Re-evaluation of the annotations resulted in similarly excellent

212 reliability with smaller limits of agreement and lower measurement error (Supplementary Table
213 3 and Supplementary Figure 1).

214 4.3 FOG episode characteristics in all the annotated videos

215 Three hundred and sixty-four freezing episodes were annotated (ST: 159 episodes from 25
216 participants, DT: 205 episodes from 27 participants) (for detailed freezing characteristics see
217 Supplementary Table 4). Freezing most often occurred while initiating turns and within the
218 first 120° arc (ST: 58.5%, DT: 72% of the time) and the outer foot was more frequently unable
219 to initiate a step (ST: 68%, DT: 65% of the time). Notably, dual tasking led to a higher number
220 of freezing episodes (Wilcoxon signed rank test $Z = -2.295$, $p = 0.022$), longer episode duration
221 ($Z = -2.639$, $p = 0.008$) and higher %TF ($Z = -2.476$, $p = 0.013$) when pooled across time points
222 within participants showing freezing at any one measurement ($N = 29$).

223 4.4 Criterion validity of visual-rated FOG

224 Significant associations were found between the NFOGQ total score and the pre-training DT
225 FOG duration ($\rho = 0.44$, $p = 0.007$), DT %TF ($\rho = 0.47$, $p = 0.004$) and a trend for DT
226 number of episodes ($\rho = 0.30$, $p = 0.074$). NFGOQ sub-scores for turning FOG frequency
227 were associated with ST and DT number of episodes, and sub-scores for turning FOG duration
228 were associated with the ST and DT %TF, and DT number of episodes and duration. Sensor-
229 based AP and ML FOG ratio were only associated with DT number of episodes (Table 3).

230 4.5 Responsiveness to intervention

231 Only DT visual-rated FOG outcomes showed any indication of differences within ($ps: 0.047 -$
232 0.088) or between intervention groups ($ps: 0.026 - 0.067$) from pre-training to retention, with
233 two SBT groups showing significant reductions in number and duration of FOG episodes and
234 %TF compared to TBT (effect sizes d from t : SB75 = $-0.99 - -1.12$; SBCR = $-0.79 - -0.88$)

235 (Figure 2). Figure 2 also illustrates differing responsiveness between AP and ML FOG ratio,
236 where the ML ratio showed a more similar response to the %TF.

237 4.6 Exploratory analysis – visual FOG ratings and sensor-based turning in place metrics

238 Repeated measures associations showed that turning jerkiness was significantly associated with
239 number of episodes in both single and dual task, but not with duration of episodes. FOG ratio
240 in both AP and ML directions was significantly associated with number and duration of FOG
241 episodes and the %TF during single task, but not dual task (Supplementary Table 5).

242 **Discussion**

243 This study investigated the reliability and validity of visual ratings of making 360° turns in
244 comparison to validated subjective and previously proposed objective automated FOG rating
245 methods. 360° turns were sensitive in provoking FOG when ON-medication, replicating earlier
246 work in the context of walking 180° turns^{21,28} in OFF. We found that application of the
247 proposed criteria for visual rating FOG resulted in higher interrater reliability compared to
248 previous studies (ICCs were 0.73¹⁷, 0.86¹⁸ and 0.9²⁹ respectively), and low measurement error,
249 making the percent time frozen during turning in place a very promising outcome measure.
250 Importantly, the dual task visual-rated metrics showed criterion validity through small but
251 significant associations with the validated NFOGQ and were responsive to the immediate
252 effects of intervention. This was not verified in relation to the FOG ratio particularly in the AP
253 direction, as this metric did not show consistent associations with visual-rated FOG, nor
254 responsiveness to intervention.

255 *Dual task turning in place as a sensitive and valid measure of FOG*

256 Like previous work¹⁹, 65% of self-professed freezers displayed clinically observable FOG
257 while performing alternating turns in place for one minute. In this respect, turning has shown
258 remarkable consistency^{29,30} and remains the most reliable trigger of FOG, possibly due to its

1
2
3 259 demands on coupling postural control and movement³¹ in the absence of external visual
4
5 260 strategies. Dual tasking not only improved this sensitivity (71%) as previously seen with
6
7
8 261 walking turns (180°)²⁸, but also revealed patterns of freezing that were more consistent with
9
10 262 the self-reported FOG severity, both in the number but particularly in the duration of FOG
11
12 263 episodes. Although large angle turns are less commonly encountered in daily life³², the 360°
13
14 264 turning task presents a greater motor challenge, revealing the extent of motor automaticity
15
16 265 deficits^{21,28}. In addition, the auditory stroop task presents an ongoing attentional demand with
17
18
19 266 response inhibition and set-switching components³³, loading strongly on prefrontal control
20
21 267 circuits, thereby limiting the ability to compensate for deficits in motor automaticity with
22
23 268 cognitive strategies. Functional near-infrared spectroscopy while turning in place showed
24
25
26 269 increased prefrontal activity without a dual task and decreased prefrontal activity with a dual
27
28 270 task in freezers³⁴, lending support to these observations. Dual task 360° turns therefore likely
29
30 271 reveal the “true” degree of FOG severity, making it not only a sensitive outcome measure, but
31
32 272 an ecologically valid one.

36 273 *Sensor-based FOG metrics show mixed relationship to FOG*

37
38
39 274 Consistent with earlier work²⁰, we found that the FOG ratios are sensitive to the occurrence of
40
41 275 freezing episodes, however results were inconsistent across tasks. Contrasting relationships
42
43 276 between the FOG ratio and visual ratings in single and dual task turning may be explained by
44
45 277 relative contributions of the number of FOG episodes and their duration to the percentage time
46
47 278 frozen. Interestingly, out of the two directions, the ML FOG ratio showed greater
48
49 279 responsiveness to the split-belt intervention (recall Figure 2), which may reflect a greater FOG
50
51 280 specificity of the ML FOG ratio during the 360° turning task, also supported by stronger
52
53 281 repeated measures correlations (see Supplementary Table 5). In contrast to previous work^{22,23},
54
55 282 relationships between both FOG ratios and the NFOGQ were absent, with scatter plots
56
57 283 suggesting an inverted – U relationship, apart from a few outlying points (see Supplementary

1
2
3 284 Figure 2). A similar distribution can be observed in a larger sample of freezers²³ (Peterson et
4
5 285 al., 2020 Supplementary Figure 1), corroborating our findings and suggesting that the FOG
6
7 286 ratio is only useful in the early disease stages when episodes are more frequent and shorter in
8
9 287 duration. However, once episodes become fewer and longer in duration, the FOG ratio
10
11 288 inadequately represents FOG severity. Hence, we advocate caution regarding the use of the
12
13 289 FOG ratio in its present form as a measure of FOG severity.
14
15
16
17

18 290 *Split-belt treadmill as a tool to reduce FOG severity*

19
20
21 291 Although not a self-evident choice for testing the responsiveness of the FOG outcomes, split-
22
23 292 belt treadmill training was found to improve turning performance in the same cohort⁹ and
24
25 293 therefore we expected to see similar effects on the freezing while turning. In line with this
26
27 294 hypothesis, we showed that the same SBT arms that showed the largest improvements on
28
29 295 turning speed (SB75 and SBCR), also improved FOG severity during the DT 360° turning task.
30
31 296 Mechanisms for the reduction of FOG are likely improved amplitude generation in the leg that
32
33 297 walked on the fast belt. Previously, circular treadmill walking^{35,36} as well as cued treadmill
34
35 298 walking to improve amplitude³⁷ have reported robust effects on freezing, likely via similar
36
37 299 mechanisms. Interestingly, improving spatiotemporal control of one limb may be sufficient to
38
39 300 reduce FOG³⁸.
40
41
42
43

44 301 *Challenges and future directions*

45
46
47 302 Rating FOG severity is particularly challenging due to the broad definition of “ineffective
48
49 303 stepping” that presents heterogeneously across participants³⁹. Previously, comparisons to
50
51 304 “relatively normal” steps^{29,30} were used to define effectiveness. However, during turning in
52
53 305 place, the stepping pattern is altered, making differentiation of effective small stepping patterns
54
55 306 from FOG-related shuffling particularly challenging. Furthermore, conflicting viewpoints exist
56
57 307 as to whether hesitations and shuffling or festination without a complete motor arrest should
58
59
60

1
2
3 308 be included within the calculation of the percentage time frozen⁴⁰. Treating these as equivalent
4
5 309 to complete and prolonged motor blocks, would likely lead to an overestimation of FOG
6
7
8 310 severity and burden. Thus, while we do not directly compare these approaches, here we provide
9
10 311 evidence that using a high-specificity approach provides reliable and valid FOG outcomes.
11
12 312 Future work may apply a weighting approach to integrate the various FOG-spectrum
13
14 313 presentations, similar to those used in clinical scoring tools^{22,41,42}. Additionally, testing these
15
16 314 methods and criteria in a larger sample of raters from multiple institutions may provide more
17
18
19 315 robust reliability estimates and lead to further refinement and validation of this approach.
20
21

22 316 Another limitation of this work is that we compared our clinical ratings to a summary measure
23
24 317 of the FOG ratio, rather than segment FOG episodes using automated algorithms⁴³. An attempt
25
26 318 to segment FOG episodes based on thresholding of the FOG ratio (>2.5) resulted in FOG
27
28 319 severity values that were unrelated to the visual-rated values (Supplementary Figure 3). Recent
29
30 320 work has shown moderate levels of reliability between automated FOG annotation models and
31
32 321 clinician ratings based on a similar two or three-sensor setup⁴⁴⁻⁴⁶. Critically, however, none of
33
34 322 these FOG segmentation algorithms have been validated for the turning in place task, and most
35
36 323 studies only included very severe freezers, off medication. We have highlighted some of the
37
38 324 challenges in manual annotation of FOG episodes – which are the basis for training the models,
39
40 325 as well as the inconsistent relationships between the FOG ratio and FOG severity in the context
41
42 326 of both FOG progression over the disease course (NFOGQ representing trait-FOG) and actual
43
44 327 task performance (manual annotations representing state-FOG), therefore these issues do not
45
46 328 appear to be trivial.
47
48
49
50

51
52
53 329 Finally, a third of the self-professed freezers did not freeze during the 360° turning task,
54
55 330 highlighting the scope for improving the sensitivity of FOG provoking tasks. Specifically,
56
57 331 studies validating these tasks should include milder freezers while ON medication to evaluate
58
59 332 the true sensitivity. Addition of anxiety-provoking components to these tasks through virtual
60

1
2
3 333 reality environments⁴⁷, may further overload the compensatory resources in early freezers.
4
5 334 Alternatively, longitudinal studies using non-episodic markers of FOG (such as turning
6
7 335 jerkiness³² or the ML FOG ratio) may serve to establish them as valid surrogates of FOG
8
9 336 severity and eliminate the need for eliciting FOG episodes completely.
10
11
12

13 337 *Manual rating of FOG severity – barriers and possibilities*

14
15
16 338 A recent systematic review of exercise effects on FOG⁶ found that only one¹⁸ out of the fifty
17
18 339 included studies annotated FOG episodes and used the percentage time frozen as an outcome
19
20 340 measure. Besides the challenges of eliciting FOG in the lab, the time and expertise to annotate
21
22 341 FOG is likely the biggest barrier to widespread adoption of these methods. Here, we used an
23
24 342 open-source software (<https://archive.mpi.nl/tla/elan>) having published methods²⁶ to perform
25
26 343 the annotations. Further, we developed standardized criteria to rate FOG trials which would
27
28 344 reduce the need for prior expertise. We believe that in doing so, the barriers to manually
29
30 345 annotating FOG would be lowered, resulting in more valid estimates of therapeutic effects on
31
32 346 FOG.
33
34
35
36

37 347 In summary, we showed that FOG severity during turning in place can be reliably rated using
38
39 348 our developed criteria. The resulting visual FOG ratings while dual tasking were both valid and
40
41 349 responsive to a split-belt treadmill intervention over a short timescale. Sensor-based ratings
42
43 350 were less favorable, showing a complex relationship to FOG that is contingent on specific
44
45 351 presentations of FOG. This work provides a robust clinical outcome to test potential therapeutic
46
47 352 interventions aimed at reducing the burden of FOG in PD.
48
49
50

51 353 **Acknowledgements**

52
53
54 354 We would like to thank all participants for their motivated and generous engagement in this
55
56 355 study. We are grateful to Demi Zoetewei and Prof. Colleen Canning for input on the annotation
57
58 356 criteria.
59
60

1
2
3 357 **Author's roles**
4
5

6 358 1) Research project: A. Conception, B. Organization, C. Execution; 2) Statistical Analysis: A.
7
8 359 Design, B. Execution, C. Review and Critique; 3) Manuscript: A. Writing of the first draft, B.
9
10 360 Review and Critique.

11
12
13 361 Nicholas D'Cruz: 1B, 1C, 2A, 2B, 3A

14
15
16 362 Jana Seuthe: 1B, 1C, 2C, 3B

17
18
19 363 Clara de Somer: 1B, 1C, 2C, 3B

20
21
22 364 Femke Hulzinga: 1B, 1C, 2C, 3B

23
24
25 365 Pieter Ginis: 1B, 1C, 2C, 3B

26
27
28 366 Christian Schlenstedt: 1A, 1B, 1C, 2C, 3B

29
30
31 367 Alice Nieuwboer: 1A, 1B, 2C, 3B

32
33
34 368 **Funding**
35

36
37 369 Support for this study was provided through a grant from the Jacques and Gloria Gossweiler
38
39 370 Foundation. The funding agency had no influence on the study design and implementation, or
40
41 371 in the interpretation of its results.
42
43
44

45 372 **Competing Interests**
46

47
48 373 The authors declare no relevant financial disclosures or conflicts of interest in regard to this
49
50 374 work.
51

52
53 375 **Financial Disclosures**
54

55
56 376 N.D., C.D.S., F.H., P.G., and A.N. are employed by the KU Leuven. J.S. and C.S. are employed
57
58 377 by the CAU Kiel. F.H. receives a doctoral fellowship from the Flanders Research Funds
59
60

1
2
3 378 (FWO). P.G. received research funding from the Internal Funds of the KU Leuven. C.S.
4
5 379 received research grants from the European Commission and the Jacques and Gloria
6
7 380 Gossweiler Foundation. A.N. received research grants from Flanders Research Funds (FWO);
8
9 381 Jacques and Gloria Gossweiler Foundation; King Baudouin Foundation; European
10
11 382 Commission; KU Leuven Internal Funds and MJ Fox Foundation.
12
13
14

15 383 **References**

- 16
17
18 384 1. Nutt JG, Bloem BR, Giladi N, Hallett M, Horak FB, Nieuwboer A. Freezing of gait: Moving
19 385 forward on a mysterious clinical phenomenon. *The Lancet Neurology*. 2011;10(8):734-744.
20 386 doi:10.1016/S1474-4422(11)70143-0
21
22 387 2. Paul SS, Canning CG, Sherrington C, Lord SR, Close JCT, Fung VSC. Three simple clinical tests to
23 388 accurately predict falls in people with Parkinson's disease. *Movement Disorders*.
24 389 2013;28(5):655-662. doi:10.1002/mds.25404
25
26 390 3. Bloem BR, Grimbergen YAM, Cramer M, Willemsen M, Zwinderman AH. Prospective
27 391 assessment of falls in Parkinson's disease. *Journal of Neurology*. 2001;248(11):950-958.
28 392 doi:10.1007/s004150170047
29
30 393 4. Pelicioni PHS, Menant JC, Latt MD, Lord SR. Falls in Parkinson's Disease Subtypes: Risk
31 394 Factors, Locations and Circumstances. *International Journal of Environmental Research and*
32 395 *Public Health*. 2019;16(12):2216. doi:10.3390/ijerph16122216
33
34 396 5. Cosentino C, Baccini M, Putzolu M, Ristori D, Avanzino L, Pelosin E. Effectiveness of
35 397 Physiotherapy on Freezing of Gait in Parkinson's Disease: A Systematic Review and
36 398 Meta-Analyses. *Movement Disorders*. 2020;35(4):523-536. doi:10.1002/mds.27936
37
38 399 6. Gilat M, Ginis P, Zoetewei D, et al. Effectiveness of exercise and training-based interventions
39 400 on Freezing of Gait in Parkinson's disease: A systematic review with meta-analysis. *npj*
40 401 *Parkinson's Disease*. Published online 2021.
41
42 402 7. Seuthe J, D'Cruz N, Ginis P, et al. The Effect of One Session Split-Belt Treadmill Training on
43 403 Gait Adaptation in People With Parkinson's Disease and Freezing of Gait. *Neurorehabilitation*
44 404 *and Neural Repair*. 2020;34(10):954-963. doi:10.1177/1545968320953144
45
46 405 8. Fasano A, Schlenstedt C, Herzog J, et al. Split-belt locomotion in Parkinson's disease links
47 406 asymmetry, dyscoordination and sequence effect. 2016;48:6-12. Accessed December 13,
48 407 2018. <https://www.sciencedirect.com/science/article/pii/S0966636216300376?via%3Dihub>
49
50 408 9. D'Cruz N, Seuthe J, Ginis P, Hulzinga F, Schlenstedt C, Nieuwboer A. Short-Term Effects of
51 409 Single-Session Split-Belt Treadmill Training on Dual-Task Performance in Parkinson's Disease
52 410 and Healthy Elderly. *Frontiers in Neurology*. 2020;11. doi:10.3389/fneur.2020.560084
53
54 411 10. Giladi N, Tal J, Azulay T, et al. Validation of the Freezing of Gait Questionnaire in Patients with
55 412 Parkinson's Disease. *Movement Disorders*. 2009;24(5):655-661. doi:10.1002/mds.21745
56
57
58
59
60

- 1
2
3 413 11. Nieuwboer A, Rochester L, Herman T, et al. Reliability of the new freezing of gait
4 414 questionnaire: Agreement between patients with Parkinson's disease and their carers. *Gait*
5 415 *and Posture*. 2009;30(4):459-463. doi:10.1016/j.gaitpost.2009.07.108
6
7 416 12. Hulzinga F, Nieuwboer A, Dijkstra BW, et al. The New Freezing of Gait Questionnaire:
8 417 Unsuitable as an Outcome in Clinical Trials? *Movement Disorders Clinical Practice*.
9 418 2020;7(2):199-205. doi:10.1002/mdc3.12893
10
11 419 13. Mancini M, Bloem BR, Horak FB, Lewis SJG, Nieuwboer A, Nonnekes J. Clinical and
12 420 methodological challenges for assessing freezing of gait: Future perspectives. *Movement*
13 421 *Disorders*. 2019;34(6):783-790. doi:10.1002/mds.27709
14
15 422 14. Delval A, Tard C, Rambour M, Defebvre L, Moreau C. Characterization and quantification of
16 423 freezing of gait in Parkinson's disease: Can detection algorithms replace clinical expert
17 424 opinion? *Neurophysiologie Clinique*. 2015;45(4-5):305-313. doi:10.1016/j.neucli.2015.09.009
18
19 425 15. Ziegler K, Schroeteler F, Ceballos-Baumann AO, Fietzek UM. A New Rating Instrument to
20 426 Assess Festination and Freezing Gait in Parkinsonian Patients. *Movement Disorders*.
21 427 2010;25(8):1012-1018. doi:10.1002/mds.22993
22
23 428 16. Nieuwboer A, Giladi N. Characterizing freezing of gait in Parkinson's disease: Models of an
24 429 episodic phenomenon. *Movement Disorders*. 2013;28(11):1509-1519.
25 430 doi:10.1002/mds.25683
26
27 431 17. Morris TR, Cho C, Dilda V, et al. A comparison of clinical and objective measures of freezing of
28 432 gait in Parkinson's disease. *Parkinsonism & Related Disorders*. 2012;18(5):572-577.
29 433 doi:10.1016/j.parkreldis.2012.03.001
30
31 434 18. Walton CC, Mowszowski L, Gilat M, et al. Cognitive training for freezing of gait in Parkinson's
32 435 disease: a randomized controlled trial. *npj Parkinson's Disease*. 2018;4(1):15.
33 436 doi:10.1038/s41531-018-0052-6
34
35 437 19. Snijders AH, Haaxma CA, Hagen YJ, Munneke M, Bloem BR. Freezer or non-freezer: Clinical
36 438 assessment of freezing of gait. *Parkinsonism and Related Disorders*. 2012;18(2):149-154.
37 439 doi:10.1016/j.parkreldis.2011.09.006
38
39 440 20. Zach H, Janssen AM, Snijders AH, et al. Identifying freezing of gait in Parkinson's disease
40 441 during freezing provoking tasks using waist-mounted accelerometry. Published online 2015.
41 442 doi:10.1016/j.parkreldis.2015.09.051
42
43 443 21. Bertoli M, Croce U della, Cereatti A, Mancini M. Objective measures to investigate turning
44 444 impairments and freezing of gait in people with Parkinson's disease. *Gait and Posture*.
45 445 Published online 2019. doi:10.1016/j.gaitpost.2019.09.001
46
47 446 22. Mancini M, Smulders K, Cohen RG, Horak FB, Giladi N, Nutt JG. The clinical significance of
48 447 freezing while turning in Parkinson's disease. *Neuroscience*. 2017;343(December):222-228.
49 448 doi:10.1016/j.neuroscience.2016.11.045
50
51 449 23. Peterson DS, Van Liew C, Stuart S, Carlson-Kuhta P, Horak FB, Mancini M. Relating Parkinson
52 450 freezing and balance domains: A structural equation modeling approach. *Parkinsonism and*
53 451 *Related Disorders*. 2020;79:73-78. doi:10.1016/j.parkreldis.2020.08.027
54
55
56
57
58
59
60

- 1
2
3 452 24. Silva-Batista C, de Lima-Pardini AC, Nucci MP, et al. A Randomized, Controlled Trial of Exercise
4 453 for Parkinsonian Individuals With Freezing of Gait. *Movement disorders : official journal of the*
5 454 *Movement Disorder Society*. 2020;(5):1-12. doi:10.1002/mds.28128
- 7 455 25. Cockx H, Nonnekes J, Bastiaan R, Radboudumc B, van Wezel R, Wang Y. Dealing with the
8 456 Heterogeneous Presentations of Freezing of Gait: How Reliable are the Freezing Index and
9 457 Heart Rate for Freezing Detection? Published online July 30, 2021. doi:10.21203/rs.3.rs-
11 458 735366/v1
- 13 459 26. Gilat M. How to Annotate Freezing of Gait from Video: A Standardized Method Using Open-
14 460 Source Software. *Journal of Parkinson's Disease*. 2019;9(4):821-824. doi:10.3233/JPD-191700
- 16 461 27. Mohr DL, Marcon RA. Testing for a 'within-subjects' association in repeated measures data.
17 462 *Journal of Nonparametric Statistics*. 2005;17(3):347-363. doi:10.1080/10485250500038694
- 19 463 28. Spildooren J, Vercruyse S, Desloovere K, Vandenberghe W, Kerckhofs E, Nieuwboer A.
20 464 Freezing of gait in Parkinson's disease: The impact of dual-tasking and turning. *Movement*
21 465 *Disorders*. 2010;25(15):2563-2570. doi:10.1002/mds.23327
- 23 466 29. Shine JM, Moore ST, Bolitho SJ, et al. Assessing the utility of Freezing of Gait Questionnaires
24 467 in Parkinson's Disease. *Parkinsonism and Related Disorders*. 2012;18(1):25-29.
26 468 doi:10.1016/j.parkreldis.2011.08.002
- 28 469 30. Schaafsma JD, Balash Y, Gurevich T, Bartels AL, Hausdorff JM, Giladi N. Characterization of
29 470 freezing of gait subtypes and the response of each to levodopa in Parkinson's disease.
30 471 *European Journal of Neurology*. 2003;10(4):391-398. doi:10.1046/j.1468-1331.2003.00611.x
- 32 472 31. Earhart GM. Dynamic control of posture across locomotor tasks. *Movement Disorders*.
33 473 Published online 2013. doi:10.1002/mds.25592
- 35 474 32. Mancini M, Weiss A, Herman T, Hausdorff JM. Turn around freezing: Community-living
36 475 turning behavior in people with Parkinson's disease. *Frontiers in Neurology*. 2018;9(JAN).
37 476 doi:10.3389/fneur.2018.00018
- 39 477 33. McFadyen BJ, Hegeman J, Duysens J. Dual task effects for asymmetric stepping on a split-belt
40 478 treadmill. *Gait and Posture*. 2009;30(3):340-344. doi:10.1016/j.gaitpost.2009.06.004
- 42 479 34. Belluscio V, Stuart S, Bergamini E, Vannozzi G, Mancini M. The Association between
43 480 Prefrontal Cortex Activity and Turning Behavior in People with and without Freezing of Gait.
44 481 *Neuroscience*. 2019;416:168-176. doi:10.1016/j.neuroscience.2019.07.024
- 46 482 35. Hong M, Earhart GM. Rotating treadmill training reduces freezing in Parkinson disease:
47 483 Preliminary observations. *Parkinsonism & Related Disorders*. 2008;14(4):359-363.
49 484 doi:10.1016/j.parkreldis.2007.07.003
- 51 485 36. Cheng FY, Yang YR, Wu YR, Cheng SJ, Wang RY. Effects of curved-walking training on curved-
52 486 walking performance and freezing of gait in individuals with Parkinson's disease: A
53 487 randomized controlled trial. *Parkinsonism and Related Disorders*. Published online 2017.
54 488 doi:10.1016/j.parkreldis.2017.06.021
- 56 489 37. Frazzitta G, Pezzoli G, Bertotti G, Maestri R. Asymmetry and freezing of gait in parkinsonian
57 490 patients. *Journal of Neurology*. 2013;260(1):71-76. doi:10.1007/s00415-012-6585-4
- 59
60

- 1
2
3 491 38. Spildooren J, Vercruyse S, Meyns P, et al. Turning and unilateral cueing in Parkinson's
4 492 disease patients with and without freezing of gait. *Neuroscience*. 2012;207:298-306.
5 493 doi:10.1016/j.neuroscience.2012.01.024
6
7 494 39. Giladi N, Nieuwboer A. Understanding and treating freezing of gait in Parkinsonism, proposed
8 495 working definition, and setting the stage. *Movement Disorders*. 2008;23(SUPPL. 2):423-425.
9 496 doi:10.1002/mds.21927
10
11 497 40. Fietzek UM, Zwosta J, Schroeteler FE, Ziegler K, Ceballos-Baumann AO. Levodopa changes the
12 498 severity of freezing in Parkinson's disease. *Parkinsonism and Related Disorders*.
13 499 2013;19(10):894-896. doi:10.1016/j.parkreldis.2013.04.004
14
15 500 41. Fling BW, Cohen RG, Mancini M, Nutt JG, Fair DA, Horak FB. Asymmetric pedunculo-pontine
16 501 network connectivity in parkinsonian patients with freezing of gait. *Brain*. 2013;136(8):2405-
17 502 2418. doi:10.1093/brain/awt172
18
19 503 42. Zoetewei D, Herman T, Brozgol M, et al. *Protocol for the DeFOG Trial: A Randomized*
20 504 *Controlled Trial on the Effects of Smartphone-Based, on-Demand Cueing for Freezing of Gait*
21 505 *in Parkinson's Disease.*; 2021. doi:10.1016/j.conctc.2021.100817
22
23 506 43. Moore ST, Yungher DA, Morris TR, et al. Autonomous identification of freezing of gait in
24 507 Parkinson's disease from lower-body segmental accelerometry. *Journal of NeuroEngineering*
25 508 *and Rehabilitation*. 2013;10(1):19. doi:10.1186/1743-0003-10-19
26
27 509 44. Mancini M, Shah V V., Stuart S, et al. Measuring freezing of gait during daily-life: an open-
28 510 source, wearable sensors approach. *Journal of NeuroEngineering and Rehabilitation*.
29 511 2021;18(1):1. doi:10.1186/s12984-020-00774-3
30
31 512 45. Reches T, Dagan M, Herman T, et al. Using Wearable Sensors and Machine Learning to
32 513 Automatically Detect Freezing of Gait during a FOG-Provoking Test. Published online 2020.
33 514 doi:10.3390/s20164474
34
35 515 46. O'Day J, Syrkin-Nikolau J, Anidi C, Kidzinski L, Delp S, Bronte-Stewart H. The turning and
36 516 barrier course reveals gait parameters for detecting freezing of gait and measuring the
37 517 efficacy of deep brain stimulation. Barbieri FA, ed. *PLOS ONE*. 2020;15(4):e0231984.
38 518 doi:10.1371/journal.pone.0231984
39
40 519 47. Economou K, Quek D, MacDougall H, Lewis SJG, Ehgoetz Martens KA. Heart rate changes
41 520 prior to freezing of gait episodes are related to anxiety. *Journal of Parkinson's Disease*.
42 521 2021;11(1):271-282. doi:10.3233/JPD-202146
43
44 522
45
46
47
48

49 523 Figure Legends:

50
51
52 524 Figure 1. Bland Altman plots for the three visual-rated outcomes. Dots represent the difference
53
54 525 in scores for each rater pair on the Y axis, plotted against the mean score from the three raters
55
56 526 on the X axis. No systematic error across severity was seen, and no significant differences
57
58
59
60

1
2
3 527 between raters was seen, apart from for number of FOG episodes between raters A and C
4
5 528 (Wilcoxon signed ranks $p = 0.03$). LOA - limits of agreement
6
7

8 529 Figure 2. LOESS curves fit to model predicted values for DT %TF and AP and ML FOG ratios
9
10 530 in the four intervention groups. Only %TF showed significant differences from pre-training to
11
12 retention, although ML FOG ratio also partially captured this pattern. Shaded regions depict
13 531 95% confidence intervals. LOESS – Locally weighted scatterplot smoothing, DT – dual task,
14
15 532 AP – anterior-posterior, ML – mediolateral, %TF – percent time frozen
16
17 533
18
19

20
21 534
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

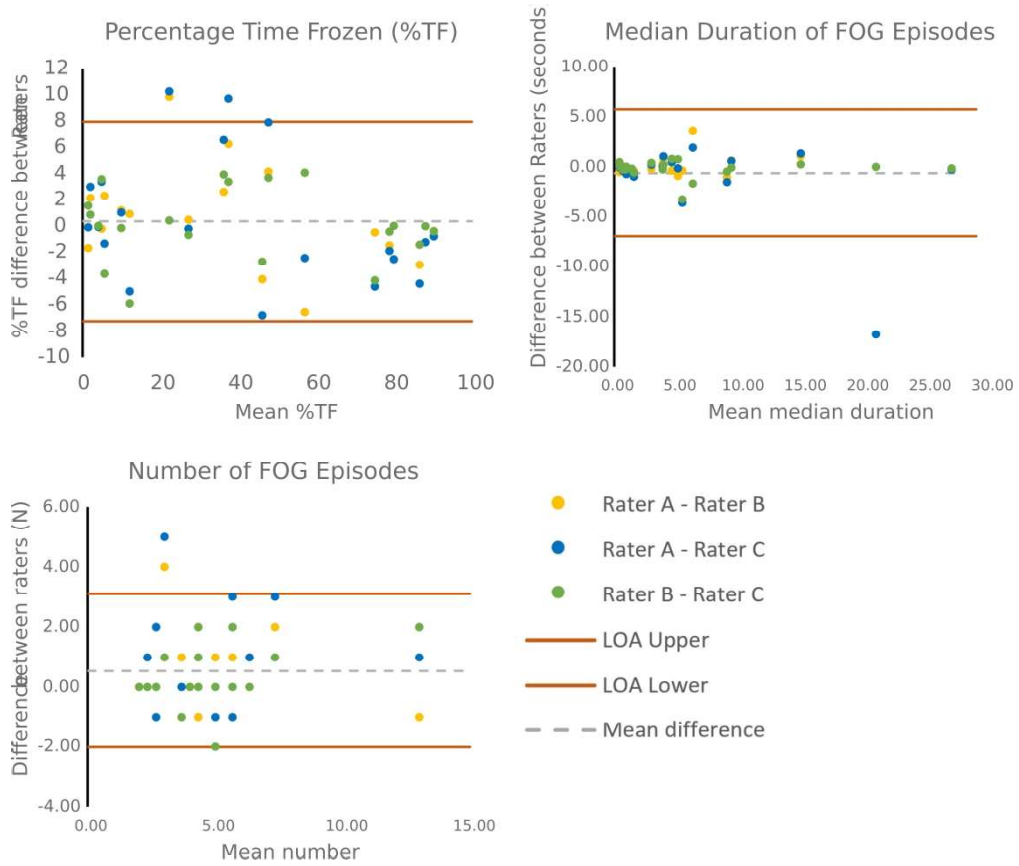


Figure 1. Bland Altman plots for the three visual-rated outcomes. Dots represent the difference in scores for each rater pair on the Y axis, plotted against the mean score from the three raters on the X axis. No systematic error across severity was seen, and no significant differences between raters was seen, apart from for number of FOG episodes between raters A and C (Wilcoxon signed ranks $p = 0.03$). LOA - limits of agreement

728x623mm (600 x 600 DPI)

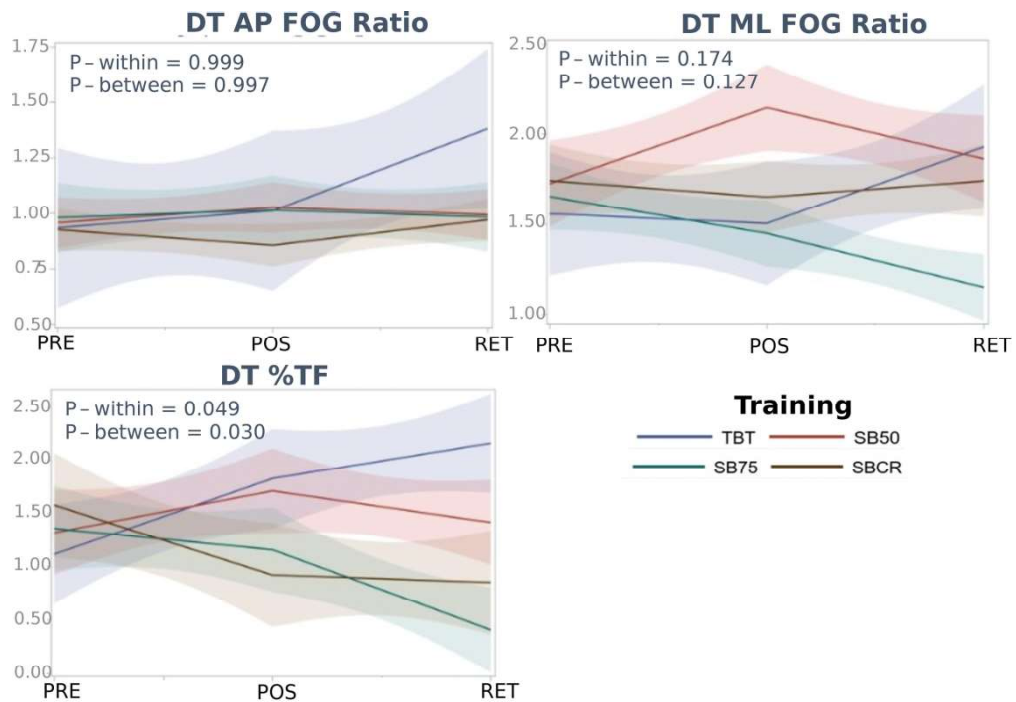


Figure 2. LOESS curves fit to model predicted values for DT %TF and AP and ML FOG ratios in the four intervention groups. Only %TF showed significant differences from pre-training to retention, although ML FOG ratio also partially captured this pattern. Shaded regions depict 95% confidence intervals. LOESS – Locally weighted scatterplot smoothing, DT – dual task, AP – anterior-posterior, ML – mediolateral, %TF – percent time frozen

213x148mm (300 x 300 DPI)

Table 1*Original and Final Criteria for Visual-Rating of FOG*

Label	Tiers	Criteria - original	Criteria - final
Start and end of the trial	Trial	From the first intention to move to the last completed full turn	From the first intention to move until 60 seconds later or until the end of the video, whichever is earlier
Start and end of a turn & direction	Turn Right Turn Left	From the first intention to move in one direction to the end of the last step (heel strike) in the same direction; OR if two whole turns are completed in the same direction – the end of the first turn is when the start position is crossed	From the first intention to move in one direction to the end of the last step (heel strike) in the same direction; OR if two whole turns are completed in the same direction – the end of the first turn is when the start position is crossed. For turns that undershoot or overshoot within 90° of the starting position, these are still considered as full turns
Start and end of a FOG episode		From intention to move or the heel/toe-off of the first ineffective step to the toe-off of the first of two effective steps	From the attempted initiation of movement of the unsuccessful step to the initiation of the first of two successful voluntary steps (one on each side)
Type of FOG	Trembling Akinetic	Trembling when accompanied by high frequency movements preceding or during the FOG episode; akinetic when these movements are not visible	Trembling when accompanied by high frequency movements preceding or during the FOG episode; akinetic when these movements are not visible
Position of FOG	Transition First 120 Mid 120 End 120	Transition when FOG occurs at the change of direction (last effective step of previous turn to first effective step of next turn). Use clock position (6 – 10 – 2 – 6 when front facing camera or 12 – 4 – 8 – 12 when rear facing camera) as reference for the 120-degree arcs	Transition when FOG occurs at the change of direction (last effective step of previous turn to first effective step of next turn). Use clock position (6 – 10 – 2 – 6 when front facing camera or 12 – 4 – 8 – 12 when rear facing camera) as reference for the 120-degree arcs
FOG leg	Outer Inner	Leg that first shows sequential reductions in excursion	Leg that fails to initiate movement at the beginning of the FOG episode

Note: Annotated labels included trial duration, direction and duration of each turn, FOG episode duration, type (trembling or akinetic), position (transition or any of three 120° arcs), and leg (inner or outer). Criteria and definitions for all labels were developed in two iterations. In the first iteration, two sets of 20 turning in place trials were rated by two pairs of raters (set 1 – rater A & B, set 2 – rater B & C) using the original criteria. Of these, 20 trials with the largest interrater variability on the percent time frozen were selected (by ND) for the second and final rating. Three of these trials were first discussed between the raters to refine criteria and align operationalization (final criteria). Interrater reliability between the three raters was calculated on these 20 trials. Freezing is defined as lack of voluntary stepping despite the *intention* to move, which is determined by the movements of the upper body, arms, and opposite leg. In case the participant moves the freezing leg due to turning momentum (usually trailing behind the hip), freezing is labelled so long as the freezing leg does not initiate a *voluntary* step. This might be seen as dragging of the leg without foot clearance, or an unexpected lack of progression that is followed by a corrective balance or stepping response. Festination or shuffling *without movement arrests* is not classified as freezing as per this definition and is not included in the rating. Video playback speed was set at 80% to capture shorter episodes and volume was muted to avoid bias for dual tasking. Trunk rotation was used as reference for labelling turn completion, and intention to move was determined from rotation of the shoulders and trunk, from arm movements, or from foot, knee or hip lifting.

Table 2*Interrater Reliability for Visual-Rated FOG Outcomes*

Outcome	ICC	LCL	UCL	P	α	Mean Δ	LOA	SEM
% Time Frozen	0.993	0.986	0.997	<0.001	0.998	0.35	± 7.59	2.68
Number of Episodes	0.859	0.723	0.937	<0.001	0.955	0.53	± 2.55	0.97
Duration of Episodes(s)	0.908	0.820	0.959	<0.001	0.967	-0.64	± 6.34	2.27

Note: ICC – Intraclass Correlation Coefficient, LCL – 95% lower confidence limit, UCL – 95% upper confidence limit, α – Cronbach's alpha, Δ – mean difference between the three pairs of raters, LOA – limits of agreement, SEM – standard error of measurement

For Peer Review

Table 3*Criterion validity of pre-training visual FOG metrics*

Outcome	%TF			FOG number			FOG duration		
	rho	p	CI	rho	p	CI	rho	p	CI
ST									
NFOGQ total	0.22	0.209	(-0.15 - 0.53)	0.23	0.177	(-0.14 - 0.54)	0.16	0.346	(-0.18 - 0.50)
Q3 turn FOG frequency	0.29	0.086	(-0.02 - 0.56)	0.34	0.044	(0.02 - 0.58)	0.24	0.153	(-0.07 - 0.52)
Q4 turn FOG duration	0.14	0.469	(-0.23 - 0.50)	0.12	0.510	(-0.25 - 0.48)	0.10	0.602	(-0.28 - 0.46)
FOG Ratio AP	0.12	0.505	(-0.26 - 0.47)	0.12	0.487	(-0.24 - 0.49)	0.10	0.560	(-0.25 - 0.46)
FOG Ratio ML	0.31	0.074	(-0.05 - 0.62)	0.37	0.029	(-0.01 - 0.62)	0.25	0.141	(-0.15 - 0.58)
DT									
NFOGQ	0.47	0.004	(0.14 - 0.73)	0.30	0.074	(-0.07 - 0.61)	0.44	0.007	(0.11 - 0.71)
Q3 turn FOG frequency	0.34	0.043	(0.01 - 0.63)	0.32	0.059	(-0.04 - 0.6)	0.27	0.107	(-0.08 - 0.57)
Q4 turn FOG duration	0.52	0.002	(0.23 - 0.76)	0.39	0.032	(0.03 - 0.68)	0.57	0.001	(0.26 - 0.77)
FOG Ratio AP	0.25	0.145	(-0.11 - 0.59)	0.41	0.014	(0.07 - 0.66)	0.23	0.179	(-0.16 - 0.58)
FOG Ratio ML	0.34	0.045	(-0.02 - 0.64)	0.57	<0.001	(0.23 - 0.77)	0.25	0.149	(-0.10 - 0.59)

Note: Spearman rank correlation coefficients, p-values and bootstrap 95% percentile confidence intervals (CI) are presented. ST – single task, DT – dual task, AP – anterior-posterior, ML – mediolateral, %TF – percent time frozen

1
2
3 **Supplementary Material**
4

5 **Supplementary Table 1**
6

7 *Demographic and clinical profile of the included sample of people with PD*
8

9

Measure	Mean (Range)
Demographics, cognition & balance	
Age (Years)	68.62 (48 - 86)
Gender (%F)	26.7
MMSE (/30)	28.38 (24 - 30)
MOCA (/30)	24.52 (17 - 30)
FAB (/18)	15.77 (11 - 18)
Mini-BEST (/28)	20.42 (4 - 28)
FES-I (/64)	28.5 (16 - 54)
Fallers (%)	52.3
Disease characteristics	
MDS-UPDRS III (/132)	35.88 (6 - 81)
Disease Duration (Years)	12.89 (1 - 38)
Hoehn & Yahr (% I/II/III/IV)	2.2/40/44.4/13.3
LEDD (mg)	818.6 (175 - 1698.5)
NFOGQ (/30)	16.33 (6 - 29)

10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

31 *Note: Means and ranges or percentages are reported. MMSE – Mini Mental Status Examination, MOCA*
32 *– Montreal Cognitive Assessment, FAB – Frontal Assessment Battery, FES-I – International version of*
33 *the Falls Efficacy Scale, MDS-UPDRS III – Motor subscale of the Movement Disorders Society sponsored*
34 *revision of the Unified Parkinson’s Disease Rating Scale, LEDD – Daily Levodopa Equivalent Dose,*
35 *NFOGQ – New Freezing of Gait Questionnaire*
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Supplementary Table 2

Demographic, clinical and training information for the four groups included in the study

<i>Demographics, Cognition and Balance</i>						
N	10	12	11	12		
Age	67 (53.5 - 77.5)	66.5 (57.75 - 75)	71 (65 - 78)	71.5 (63.25 - 78.5)	NS	0.680
Sex (% Female)	10	33.3	27.3	33.3	NS	0.523
MMSE	28.5 (26.5 - 30)	28 (28 - 29)	28 (28 - 29)	29 (28 - 30)	NS	0.551
MOCA	26.5 (20.5 - 27)	24 (22 - 29)	24 (20 - 26)	25.5 (23.25 - 28)	NS	0.611
FAB	16 (15 - 17.25)	16 (13 - 18)	16 (15 - 17)	16 (15.25 - 17.75)	NS	0.938
FES-I	22 (20.75 - 42.75)	23 (17.25 - 34.75)	28 (20 - 36)	26.5 (22.25 - 33)	NS	0.764
Retrospective falls (N)	1 (0 - 9.5)	2 (0 - 24)	1 (0 - 5)	0 (0 - 2)	NS	0.413
Mini-BEST	23 (17.75 - 25)	20.5 (17 - 24.75)	18 (17 - 22)	24 (17 - 26)	NS	0.619
<i>Disease-specific Scales</i>						
Disease Duration (years)	9.75 (5 - 17.75)	13.5 (7.75 - 16.75)	14 (8.5 - 15)	11.5 (6 - 15.5)	NS	0.882
H&Y % I/II/III/IV	10/40/40/10	0/41.7/41.7/16.7	0/27.3/54.5/18.2	0/50/41.7/8.3	NS	0.857
LEDD	773 (576.25 - 1189.25)	814 (516.5 - 933.75)	810 (704 - 958.75)	805 (543.75 - 994.37)	NS	0.994
MDS-UPDRS Part III	32.5 (22 - 54.75)	43 (23 - 53)	35 (25 - 41)	33.5 (26 - 43.75)	NS	0.827
NFOGQ	16 (9 - 21.5)	15 (13.5 - 23.5)	15 (10 - 20)	17.5 (13.25 - 18.75)	NS	0.833
<i>Training Engagement and Intensity</i>						
Training Velocity (m/s)	0.98 (0.77 - 1.30)	1.09 (1.04 - 1.23)	0.95 (0.80 - 1.21)	1.12 (0.86 - 1.27)	NS	0.870
Training Duration (min)	30 (24.75 - 30)	30 (29.25 - 30)	30 (25.1 - 30)	30 (27.75 - 30)	NS	0.979
Training Reduced (% Yes)	30	18.2	27.3	27.3	NS	0.888
Handrail use (% Yes)	20	16.7	36.4	16.7	NS	0.655
Borg During	13 (11.75 - 13.5)	12.5 (10.25 - 14.5)	13 (12.75 - 14.25)	13 (11 - 15)	NS	0.702
Borg Post	14 (13 - 15.5)	13 (12.25 - 14.75)	15 (13 - 15.25)	14 (13 - 16.75)	NS	0.643
VAS Mental Pre	2 (1.15 - 3.55)	2.5 (1 - 5.8)	1.55 (0.47 - 3.12)	1.9 (0.425 - 3)	NS	0.566
VAS Mental Post	4.5 (2 - 6.25)	3.85 (1.35 - 7.4)	4.85 (2.87 - 6.3)	3.2 (2 - 5.35)	NS	0.793
VAS Physical Pre	3.1 (0.75 - 5.9)	3.1 (1.5 - 7.5)	4 (1.27 - 5.01)	1.35 (0.8 - 4.75)	NS	0.574
VAS Physical Post	4.6 (1.95 - 7.3)	4 (3.42 - 6.1)	4.9 (4 - 6.12)	3.4 (1.27 - 5.95)	NS	0.620

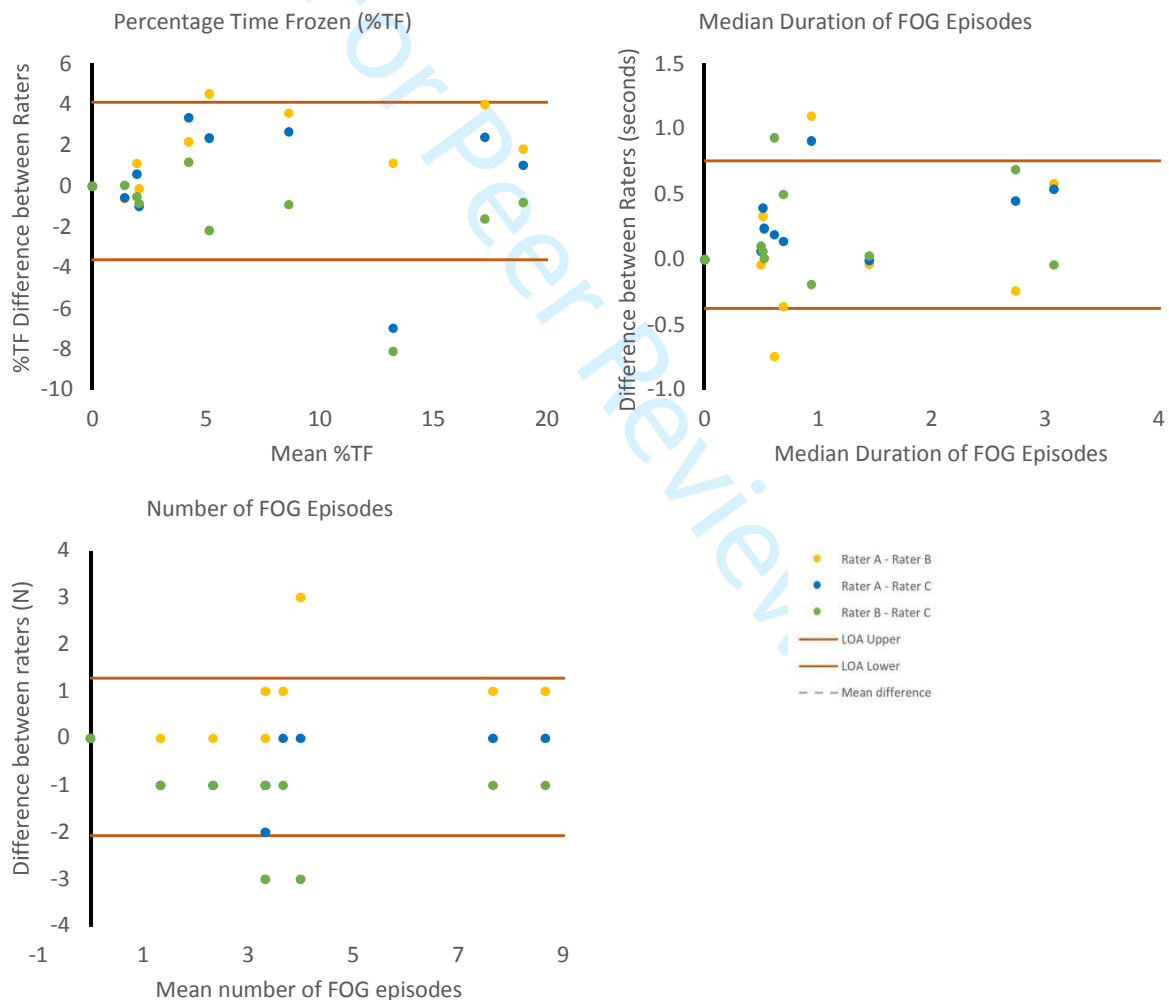
Note: Median values (Q1 – Q3) or percentages are reported. KW – Kruskal-Wallis test, MMSE – Mini Mental Status Examination, MOCA – Montreal Cognitive Assessment, FAB – Frontal Assessment Battery, FES-I – International version of the Falls Efficacy Scale, Mini-BEST – LEDD – Daily Levodopa Equivalent Dose, MDS-UPDRS part III – Motor subscale of the Movement Disorders Society sponsored revision of the Unified Parkinson's Disease Rating Scale, NFOGQ – New Freezing of Gait Questionnaire, Pre – Pre-training, During – Halfway through training (>3 blocks), Post – Post-training, VAS – Visual Acuity Scale

Supplementary Table 3

Re-evaluation of Interrater Reliability for Visual-Rated FOG Outcomes after 15 months

Outcome	ICC	LCL	UCL	P	α	Mean Δ	LOA	SEM
% Time Frozen	0.928	0.805	0.980	<0.001	0.979	0.25	± 3.87	1.90
Number of Episodes	0.901	0.613	0.975	<0.001	0.982	0.4	± 1.67	0.87
Duration of Episodes(s)	0.915	0.772	0.976	<0.001	0.974	0.19	± 0.57	0.30

Note: ICC – Intraclass Correlation Coefficient, LCL – 95% lower confidence limit, UCL – 95% upper confidence limit, α – Cronbach’s alpha, Δ – mean difference between the three pairs of raters, LOA – limits of agreement, SEM – standard error of measurement



Supplementary Figure 1 – Re-evaluation of visual annotations. Bland Altman plots for the three visual-rated outcomes for 10 additional trials rated 15 months after criteria development. Dots represent the difference in scores for each rater pair on the Y axis, plotted against the mean score from the three raters on the X axis. Interrater reliability was similar to the second iteration with smaller LOA and measurement error. LOA - limits of agreement, FOG – Freezing of gait

Supplementary Table 4*Distribution of freezing episodes among participants during ST and DT turning in place*

ID	FOG in ST	FOG in DT	ST number of episodes	DT number of episodes
001	YES	YES	11	16
002	YES	YES	11	16
003	YES	YES	1	4
004	YES	YES	7	5
005	NO	YES	0	2
006	YES	YES	14	11
007	NO	YES	0	1
008	YES	YES	11	12
009	NO	NO	0	0
010	YES	YES	13	12
011	YES	NO	2	0
012	YES	YES	8	12
013	NO	NO	0	0
014	YES	YES	4	6
015	NO	NO	0	0
016	YES	YES	5	6
017	YES	YES	19	9
018	YES	YES	4	4
019	NO	NO	0	0
020	NO	NO	0	0
021	NO	YES	0	1
022	YES	YES	2	2
023	YES	YES	13	25
024	YES	YES	1	2
025	NO	NO	0	0
026	NO	NO	0	0
027	NO	YES	0	2
028	YES	YES	2	1
029	YES	NO	1	0
030	YES	YES	3	6
031	YES	YES	11	17
032	YES	YES	4	4
033	YES	YES	1	1
034	NO	NO	0	0
035	NO	NO	0	0
036	YES	YES	1	5
037	YES	YES	6	19
038	NO	NO	0	0
039	YES	YES	4	4
Total	25	27	159	205

Note: Video data was captured from 39 people with Parkinson's disease and FOG. FOG was provoked in 29 of the participants – 23 in both tasks, 25 in ST and 27 in DT. FOG – Freezing of gait, ST – single task, DT – dual task

Supplementary Table 5

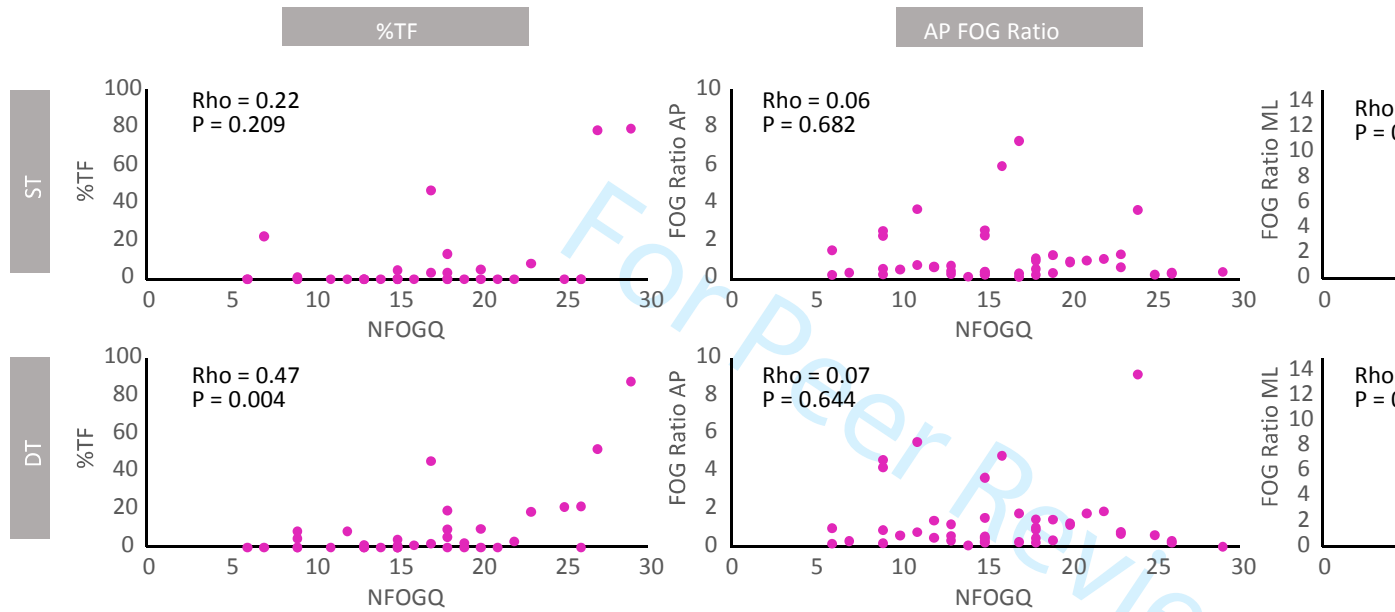
Associations between turning metrics and visual-rated FOG over time

Outcome	Task	Episode duration	Episode number	%TF
Turning				
Mean turning speed	ST	-1.51 (0.120)	-0.35 (0.723)	-1.05 (0.263)
	DT	-0.53 (0.595)	-1.69 (0.088)	-1.73 (0.067)
Peak turning speed	ST	-0.38 (0.702)	0.11 (0.913)	-0.54 (0.611)
	DT	-0.94 (0.331)	-1.95 (0.045)	-2.15 (0.03)
Mean Jerk	ST	1.75 (0.078)	2.25 (0.023)	2.38 (0.019)
	DT	-0.76 (0.425)	2.05 (0.036)	1.04 (0.315)
FOG Ratio AP	ST	2.14 (0.029)	2.72 (0.004)	2.61 (0.009)
	DT	0.90 (0.371)	1.61 (0.095)	1.35 (0.168)
FOG Ratio ML	ST	2.14 (0.029)	2.88 (0.002)	2.92 (0.004)
	DT	1.31 (0.210)	1.94 (0.036)	1.61 (0.099)

Note: Non-parametric within-subject repeated measures correlation Z score and associated p-values (in brackets) are shown for visual FOG rating with objective measures of turning. Measures with a significant p-value are shown in bold. AP - antero-posterior, ML - mediolateral

Supplementary Figure 2

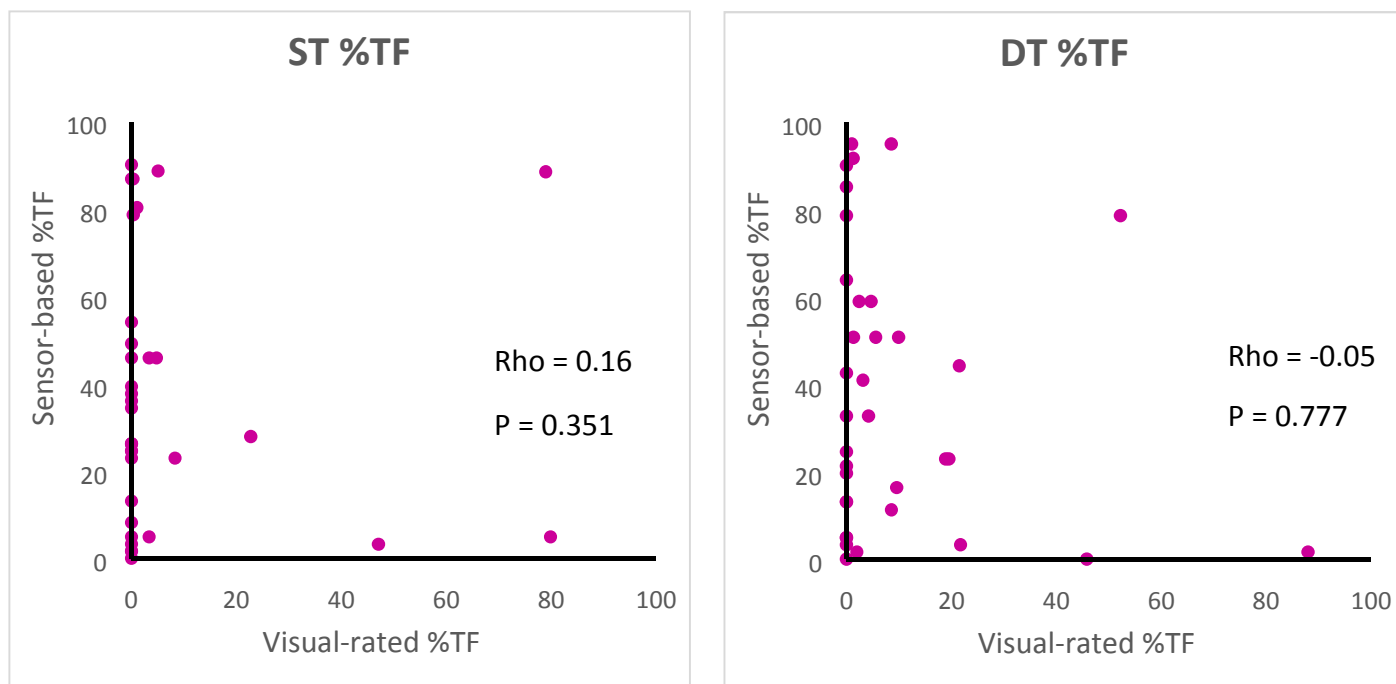
Scatter plots for the pre-training single and dual task objective FOG measures and the subjective NFOGQ



Note: Spearman rank correlations (rho) and associated p-values are presented in each panel. Dual task performance was significantly associated with the subjective rating. One outlying value (all FOG ratios > 28) was omitted from the analysis. ST – single task, DT – dual task, AP – anterior-posterior, ML – mediolateral, %TF – percent time

Supplementary Figure 3

Scatter plots for the sensor-based and visual-rated percent time frozen at baseline



Note: Sensor-based percent time frozen was calculated by thresholding the AP FOG ratio > 2.5 to segment periods of possible-freezing from normal turning. Scatter plots and spearman correlations revealed no relationship between these two metrics.