

# BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.

The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email [info.bmjopen@bmj.com](mailto:info.bmjopen@bmj.com)

# BMJ Open

## Test-retest reliability of the twenty-five-hole peg test in stroke patients

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2019-032560
Article Type:	Research
Date Submitted by the Author:	24-Jun-2019
Complete List of Authors:	Granström, Fredrik; Uppsala University, Centre for Clinical Research Sörmland Hedlund, Mattias; Umeå Universitet Medicinska fakulteten, Department of Community Medicine and rehabilitation, Physiotherapy Lindström, Britta; Umeå Universitet Medicinska fakulteten, Department of Community Medicine and Rehabilitation, Physiotherapy Eriksson, Staffan; Uppsala University, Centre for Clinical Research Sörmland; Uppsala University, Department of Neuroscience, Physiotherapy
Keywords:	Stroke < NEUROLOGY, measurement, Reliability, Clinical assessment, hand function, twenty-five-hole peg test

SCHOLARONE™  
Manuscripts

## Test-retest reliability of the twenty-five-hole peg test in stroke patients

Fredrik Granström<sup>1</sup>, Mattias Hedlund<sup>2</sup>, Britta Lindström<sup>2</sup>, Staffan Eriksson<sup>1,3</sup>

<sup>1</sup>Centre for Clinical Research Sörmland, Uppsala University, Sweden

<sup>2</sup>Department of Community Medicine and Rehabilitation, Physiotherapy, Umeå University, Sweden

<sup>3</sup>Department of Neuroscience, Physiotherapy, Uppsala University, Sweden

**CORRESPONDING AUTHOR:** Staffan Eriksson, Centre for Clinical Research Sörmland, Uppsala University, Kungsgatan 41, 631 88 Eskilstuna, Sweden. Telephone number: +46 16 105251, +46 73 949 99 33. Email: [staffan.eriksson@germed.umu.se](mailto:staffan.eriksson@germed.umu.se)

**KEY WORDS:** clinical assessment, measurement, reproducibility of results, reliability, hand function, stroke, twenty-five-hole peg test

**WORD COUNT** 3920 words, 4045 words including “strengths and limitations of this study”.

## ABSTRACT

**Objectives:** The twenty-five-hole peg test (TFHPT) is similar to the nine-hole peg test (NHPT), but the larger number of available pegs makes it straightforward to count the number of pegs inserted, during a stipulated time frame (50 seconds), as the result. The objective was to assess the test-retest reliability of the TFHPT when testing persons with stroke, a special focus was placed on the absolute reliability as quantified by the smallest real difference (SRD). Complementary aims were to investigate possible implications for the use of the TFHPT and for how the SRD of the TFHPT performance should be expressed.

**Design:** This study employed a test-retest design including 3 trials; the pause between trials was approximately 10-120 seconds.

**Participants, setting and outcome measure:** Thirty-one participants who had suffered a stroke were recruited from a group designated for constraint-induced movement therapy (CIMT) at outpatient clinics. The result of the TFHPT was expressed as the number of pegs inserted.

**Methods:** Absolute reliability was quantified by the SRD, including random and systematic error for a single trial,  $SRD_{2,1}$ , and for an average of three trials,  $SRD_{2,3}$ . For the SRD measures, the corresponding smallest real difference percentage (SRD%) measures were also reported.

**Results:** The differences in the number of pegs necessary to detect a change in the TFHPT for  $SRD_{2,1}$  and  $SRD_{2,3}$  were 4.0 and 2.3, respectively. The corresponding SRD% values for  $SRD_{2,1}$  and  $SRD_{2,3}$  were 36.5% and 21.3%, respectively.

**Conclusions:** The smallest change that can be detected in the TFHPT should be just above 2 pegs for a test procedure including an average of three trials when systematic error is also

1  
2  
3 considered (SRD<sub>2,3</sub>). The use of an average of three trials compared to a single trial (SRD<sub>2,1</sub>  
4  
5 vs SRD<sub>2,3</sub>) reduces the measurement error substantially.  
6  
7

8 **Trial registration:** ISRCTN registry, reference number ISRCTN24868616.  
9  
10

## 11 12 13 14 **ARTICLE SUMMARY**

### 15 16 17 **Strengths and limitations of this study**

- 18  
19  
20  
21 • There were some issues in this study regarding the generalizability of the results, as  
22  
23 the participants were selected because they should benefit from CIMT, and few scored  
24  
25 above 20 pegs during the 50-second trial duration.  
26
- 27  
28 • Among other measures of reliability, the SRD percentage was reported, which is a  
29  
30 good measure for comparisons between different tests, scales and populations.  
31
- 32  
33 • The results were presented with several different reliability measures, which helps to  
34  
35 gain some knowledge about the source of the measurement error.  
36
- 37  
38 • As the test-retest trials were performed within minutes, the possible day-to-day  
39  
40 variation was not captured.  
41
- 42  
43 • The intended practice trial was included as one of three trials in the analyses, which  
44  
45 appears to have contributed to the present learning effect.  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## INTRODUCTION

For a comprehensive upper limb assessment among persons with stroke it is important to combine a measure of proximal upper limb function with a measure of manual dexterity.[1] However, only approximately 25% of studies regarding upper limb interventions include a specific measure of manual dexterity.[1] The nine-hole peg test (NHPT) is a common test of fine manual dexterity and the most common in research.[1-3] The reliability of the NHPT has been investigated in two studies of stroke populations, revealing a large discrepancy in the reliability reported.[2, 3]

Reliability is a term that describes how the result of a measurement with an instrument is affected by measurement error.[4] The concept of absolute reliability refers to the consistency of measurements within individuals and can be quantified by, for example, the smallest real difference (SRD) and by SRD%.[5]

A weakness with the NHPT is that many persons with stroke cannot reach the lower limit; i.e., a floor effect arises. Furthermore, because there are only nine pegs, measures must be taken to avoid ceiling effects. Therefore, in the original test, the result is expressed as the time to complete the test, including inserting and removing all of the pegs.[2, 3, 6, 7] This, however, aggravates the floor effects because tests that are not completed during the stipulated time are excluded.[2, 3] The maximum time could be prolonged to include the great majority useful to test. However, this would be time consuming and possibly unethical due to the possibility of a non-completed test after a lengthy attempt. A modified NHPT is used to mitigate the floor effect while avoiding the ceiling effect, in which the result is expressed as the number of inserted pegs (not removed) per unit of time, i.e., the frequency;[8] however, this test has not been investigated for reliability.

1  
2  
3 In Sweden, a similar peg test, a twenty-five-hole peg test (TFHPT), has been used in clinical  
4  
5 practice. The larger number of available pegs makes it straightforward to count the number of  
6  
7 pegs inserted, during a stipulated time frame of 50 seconds, as the result. Thus, the TFHPT  
8  
9 measures the motor function on a numerical scale, with low floor effects and reasonable  
10  
11 ceiling effects. Concomitantly, in the two other studies in which the reliability of the NHPT  
12  
13 has been investigated, individuals with worse motor impairment, compared to what is possible  
14  
15 to test with the TFHPT, were excluded due to floor effects.[2, 3] The TFHPT is not previously  
16  
17 described in the literature, and its reliability has not been investigated. Due to the similarity of  
18  
19 the NHPT and the TFHPT, the underlying skill assessed with these tests is most likely the  
20  
21 same. However, since the tests have completely different stop criteria – a time limit for the  
22  
23 TFHP vs. all pegs inserted for the NHPT – equal reliability cannot be taken for granted.[6]  
24  
25  
26 Measurements with the TFHPT are quantified on a numerical scale and can be used on a large  
27  
28 portion of persons suffering from stroke. Thus, depending on the magnitude of measurement  
29  
30 error, this test may be useful, both in clinical practice and in research.  
31  
32  
33  
34  
35

36 The overall aim of this study was to assess the test-retest reliability of the TFHPT for persons  
37  
38 suffering from stroke; a special focus was placed on the absolute reliability, measured as the  
39  
40 SRD. Complementary aims were to investigate possible implications for the use of the  
41  
42 TFHPT and for how the SRD of the TFHPT performance should be expressed.  
43  
44  
45  
46  
47  
48  
49

## 50 **METHOD**

### 51 **Participants**

52  
53 The participants in this study were consecutively recruited in the process of screening patients  
54  
55 eligible for inclusion in a multicentre randomized controlled trial (RCT), reference number  
56  
57  
58  
59  
60

1  
2  
3 ISRCTN24868616 at the ISRCTN registry. The patients were considered for inclusion  
4  
5 because they were to undergo constraint-induced movement therapy (CIMT) at one of the  
6  
7 clinics participating in the RCT. The clinics were outpatient rehabilitation clinics in the public  
8  
9 health care system in Sweden. Data were collected at the clinics. The sample in this study  
10  
11 consisted of included and excluded participants in the multicentre RCT. The participants were  
12  
13 included if they had one stroke or more registered in the medical record and if TFHPT data  
14  
15 were available from three trials before and three trials after the CIMT. Moreover, related to  
16  
17 the outcome measure, a minimum of one peg and a maximum of 24 pegs inserted was  
18  
19 necessary for inclusion. This was to avoid an untrue low measurement error from participants  
20  
21 stable at 0 or 25 pegs inserted.  
22  
23  
24  
25

26  
27 A minimum of 30 participants were included to obtain a sufficient number for a reliability  
28  
29 study.[9]  
30  
31

### 32 **Procedure and measurements**

33  
34  
35 The TFHPT has twenty-five holes and pegs.[6] The test used in this study consisted of a  
36  
37 rectangular 21 cm × 45 cm board with a box containing pegs on one side and an elevated 18  
38  
39 cm × 18 cm area with holes on the other side. The holes were 9 mm wide, 18 mm deep and  
40  
41 spaced 20 mm apart. The box had a base of 13 cm × 18 cm and was 5 cm deep. The pegs were  
42  
43 40 mm long and 8 mm in diameter.  
44  
45

46  
47 The TFHPT was administered as the second test in a battery of different tests. The preceding  
48  
49 test required approximately 30-60 minutes to administer. The tests were administered in an  
50  
51 examination room, in which only the participant and the physiotherapist were present. For the  
52  
53 TFHPT, three trials were performed with each hand. The participants started with the less  
54  
55 affected hand, followed by the more affected hand, i.e., the hand of investigation in this test-  
56  
57 retest study. The pause between trials was approximately 10-120 seconds. The board was  
58  
59  
60



1  
2  
3 placed at a distance favoured by the participant with the centre row of holes centred towards  
4  
5 the navel and the box side oriented towards the tested hand. The starting position was with  
6  
7 both hands on the board, and the time keeping was begun upon first hand contact with the  
8  
9  
10 pegs. We gave the following instructions to the participants:

- 11  
12  
13 1. I want you to pick up one peg at a time and insert them in the holes of the board.
- 14  
15 2. Use only the right/left hand; you can only use the other hand to steady the board.
- 16  
17 3. You can fill the holes in any order you desire.
- 18  
19 4. We start with a practice trial.
- 20  
21 5. You have got 50 seconds to insert as many pegs as you can. After 50 seconds, the trial  
22  
23 is terminated.
- 24  
25 6. Are you ready? Ready, set, go!
- 26  
27 7. After the practice trial: This was practice, now come the two actual test trials where  
28  
29 the results are noted down. Repeat step 6.  
30  
31  
32  
33

34 The test-retest reliability of the TFHPT was assessed on two separate occasions, i.e., before  
35  
36 and after a two-week training period (the CIMT). The same procedure was used on both of  
37  
38 these occasions, and for each participant, all tests were administered by the same  
39  
40 physiotherapist. The assessment after the CIMT period was performed as an internal  
41  
42 validation.  
43  
44  
45

46 Two physiotherapists, SE and BL, administered the tests in this study. SE has general  
47  
48 experience with persons suffering from stroke and experience administering the original  
49  
50 NHPT. BL has extensive experience with persons suffering from stroke, including  
51  
52 administering the original NHPT.  
53  
54  
55

56 Background data were collected by the staff at the clinics, except for data on the dominant  
57  
58 hand before the stroke and the Fugl-Meyer test.[10]  
59  
60

## Statistics

All three trials were used in the analyses, although the first trial was introduced as a practice trial to the participants. The exception was the Bland-Altman plots, for which only trials two and three were used. Analyses of pre-intervention data and post-intervention data were performed separately.

Bland-Altman plots provided a graphic description of the variability of the data. The mean of trials two and three was plotted against the difference between trials three and two for each subject. The centre line displayed the mean difference for the group between trials three and two. The upper and lower confidence limits were calculated as the mean difference  $\pm$  standard deviation (SD) of the mean difference  $\times$  1.96.

Measurement error can be either random or systematic. In random error, there is no pattern to the variability, whereas in systematic error the measurement varies in a non-random way, i.e., the mean values between the trials differ.[5] To investigate whether there was a systematic error in test scores, one-way repeated measures ANOVA was used to detect potential between trial effects. Fisher's LSD post hoc tests between trials were performed when the main effect for trials was significant. The assumption of sphericity was met in all trials according to Mauchly's test, whereas the assumption of normal distribution was violated in trial two's pre-intervention according to the Kolmogorov-Smirnov test ( $p=0.043$ ).

Relative reliability refers to the consistency of the positions of measurements relative to those of others within the tested group and was quantified by using several intra-class correlation coefficients (ICCs).[5 11] Concomitantly, in ICCs, the within-subject variability is compared to the between-subject variability.[5] This makes ICCs sensitive to the degree of between-subject variability, and with all other things being equal, a more heterogeneous sample will

1  
2  
3 produce higher ICC values.[5, 11] In addition, it is difficult to draw statistical inference from  
4  
5 one sample to another.[12]  
6  
7

8 Three separate measures of relative reliability including 95% confidence interval (CIs),  
9  
10 namely,  $ICC_{2,1}$ ,  $ICC_{2,3}$ , and  $ICC_{3,3}$ , were calculated. This panel of measures was used to  
11  
12 compare the results representative of single and average measures and to obtain an estimate of  
13  
14 the influence of systematic error. The first figure in the ICC designation represents the type of  
15  
16 intraclass correlation coefficient (ICC model), while the second figure represents single or  
17  
18 average measures, where “1” represents single measures and “2” or higher represents the  
19  
20 number of trials from which the average is calculated.[5]  $ICC_{2,1}$  and  $ICC_{2,3}$  are calculated  
21  
22 from a two-way random effect model and incorporate both systematic and random error,  
23  
24 whereas  $ICC_{3,3}$  is calculated from a two-way fixed effect model and incorporates only random  
25  
26 error.[5, 13] Thus, the less the systematic error contributes to the total error, the closer  $ICC_{2,3}$   
27  
28 is to  $ICC_{3,3}$ . [5] Furthermore, ICCs were calculated for single measures, i.e., 2.1, and average  
29  
30 measures, i.e., 2.3 and 3.3. An ICC for single measures represents the reliability for a test  
31  
32 procedure in which the subject is tested with a single trial on a test occasion.[5] An ICC for  
33  
34 average measures represents the reliability for a test procedure in which the subject is tested  
35  
36 with two or more trials on a test occasion and the score is expressed as the average of these  
37  
38 trials.  
39  
40  
41  
42  
43  
44

45 To estimate the absolute reliability, the standard error of measurement (SEM), SRD and SRD  
46  
47 percentage (SRD%) were calculated. These three measures of absolute reliability were  
48  
49 calculated from each of the three different ICC-measures:  $SEM_{2,1}$ ,  $SRD_{2,1}$ ,  $SRD\%_{2,1}$ ,  $SEM_{2,3}$ ,  
50  
51  $SRD_{2,3}$ ,  $SRD\%_{2,3}$ ,  $SEM_{3,3}$ ,  $SRD_{3,3}$ , and  $SRD\%_{3,3}$ . SEM is the within-subject standard deviation  
52  
53 calculated from repeated tests.[11, 14] The variation of repeated tests can be thought of as the  
54  
55 error around a true value; concomitantly, the within-subject standard deviation is used as a  
56  
57 measure of measurement error.[14] The SEM was calculated according to  $SEM = SD$   
58  
59  
60

1  
2  
3  $\sqrt{(1 - ICC)}$ , where SD was calculated from the total sum of squares ( $SS_{TOTAL}$ ) in the  
4  
5 ANOVA table generated in the ICC analyses as  $\sqrt{SS_{total}/(n - 1)}$ .<sup>[5]</sup> SRD can be interpreted  
6  
7 as an extension of SEM, in which a 95% CI of the measurement error for both test occasions  
8  
9 in a test-retest situation has been incorporated in the measure.<sup>[5]</sup> The SRD can be seen as the  
10  
11 smallest difference that can be detected with 95% certainty in an individual using a test  
12  
13 instrument.<sup>[5]</sup> The SRD was calculated using the formula  $1.96 \times SEM \times \sqrt{2}$ , where 1.96 is  
14  
15 related to the 95% CI and  $\sqrt{2}$  refers to the error of two measurements.<sup>[5]</sup> The SRD% was  
16  
17 calculated by dividing the SRD value by the grand mean multiplied by 100.<sup>[2, 9]</sup> This value is  
18  
19 independent of measurement units and is indexed to the mean value of the observations from  
20  
21 which it was derived and is therefore a good measure for comparisons between different tests,  
22  
23 scales and populations.<sup>[9, 11, 12]</sup> An SRD% of 30% has been suggested as an acceptable  
24  
25 level of reliability.<sup>[15]</sup>

26  
27 Because estimates of absolute reliability vary with the type of ICC value, some caution is  
28  
29 warranted when comparing them with measures from other studies.<sup>[5]</sup> Therefore,  $SEM_{mean}$   
30  
31 square error term (MSE) =  $\sqrt{MSE}$  was also calculated, where MSE (this term is called residual error  
32  
33 by Hopkins and mean square residual in the SPSS-output) was taken from the ANOVA table  
34  
35 of the ICC calculation.<sup>[5, 11]</sup> This SEM measure represents the reliability of a test procedure  
36  
37 in which the subject is tested with a single trial on a test occasion and is a pure measure of  
38  
39 random error.<sup>[5]</sup>  $SRD_{MSE}$  and  $SRD\%_{MSE}$  were also derived from  $SEM_{MSE}$ .

40  
41 The analysis of test-retest reliability was pre-planned. SPSS version 21 was used to calculate  
42  
43 ICC and ANOVA. The alpha level was set to 0.05.

### 44 45 **Patient and public involvement**

46  
47 No patients or public were involved in the development or design of this study.

### 48 49 **RESULTS**

In this study, participants were recruited between January 2011 and September 2014. Of 60 eligible patients, 29 were excluded for any of the following reasons: not suffering a stroke, missing data, and yielding either below the minimum number of inserted pegs or above the maximum number of inserted pegs (Figure 1). This yielded 31 participants, 21 men and 10 women, for inclusion in the analysis, with a mean  $\pm$  SD age of  $66 \pm 9$  years (Table 1). The two eligible patients who were excluded because they exceeded the permitted maximum number of pegs inserted completed the 25 pegs in their best trial within 49.3 seconds and 39.4 seconds. Of the eight patients who were excluded because they fell below the minimum number of inserted pegs, five could insert at least one peg in one of the trials. Data were collected from 17 and 14 participants, respectively, by the two physiotherapists (first and last author) at seven clinics.

**Table 1.** Characteristics of participants at pre-intervention trials

Participants	N=31	
Age (years), mean $\pm$ SD <sup>a</sup>	66 $\pm$ 9	
Men/women, n <sup>b</sup>	21/10	
Time since stroke (months), median (IQR <sup>d</sup> ), (min–max)	17 (8–24), (2–70)	A graphic description of the variability of the data can be seen in the Bland-Altman plots (Figure 2a and b). A slight association between the random error and the magnitude of the measurements can
Previous dominant hand more affected by stroke, n	19	
TFHPT <sup>c</sup> , mean of three trials, mean $\pm$ SD, (min–max)	10.8 $\pm$ 6.8, (1–22.7)	
Fugl-Meyer test (score), median (IQR <sup>d</sup> ), (min–max)	46 (41–53), (29–62)	
More than one stroke, n	3	

<sup>a</sup>SD, <sup>b</sup>number of participants, <sup>c</sup>twenty-five-hole peg test, <sup>d</sup>interquartile range

be observed in the pre-intervention trials (Figure 2a).

For *pre-intervention* trials, the mean values  $\pm$  SDs for trials 1, 2 and 3 were  $10.0 \pm 6.5$ ,  $11.0 \pm 7.1$ , and  $11.5 \pm 6.9$ , respectively. The one-way repeated measures ANOVA revealed a main effect between trials with  $F(2, 29) = 10.9$  and  $p < 0.001$ . Post hoc tests revealed differences between trials 2 and 1 and between trials 3 and 1, with mean differences (95% CIs) of 1.0 (0.3–1.6) and 1.5 (0.9–2.2), respectively.

For *post-intervention* trials, the mean values  $\pm$  SD for trials 1, 2 and 3 were  $11.8 \pm 6.5$ ,  $12.4 \pm 6.7$ , and  $12.5 \pm 6.8$ , respectively. The one-way repeated measures ANOVA revealed a main effect between trials with  $F(2, 29) = 4.1$  and  $p = 0.027$ . Post hoc tests revealed a difference between trials 3 and 1, with a mean difference (95% CIs) of 0.6 (0.2–1.1).

For *pre-intervention* trials, the ICCs incorporating random and systematic error, the ICC<sub>2,1</sub> (95% CI) for single measures and the ICC<sub>2,3</sub> (95% CI) for average measures, were 0.96 (0.92–0.98) and 0.99 (0.97–0.99), respectively (Table 2). The SRDs incorporating random and systematic error, the SRD<sub>2,1</sub> for single measures and the SRD<sub>2,3</sub> for average measures, were 4.0 and 2.3 pegs, respectively. The corresponding SRD% values for SRD<sub>2,1</sub> and SRD<sub>2,3</sub> were 36.5% and 21.3%, respectively. The SRD only incorporating random error, the SRD<sub>3,3</sub> for average measures, was 2.0 pegs.

**Table 2.** Results of reliability measures for pre-intervention trials

	ICC <sup>a</sup> (95% CI)	SEM <sup>b</sup> , n <sup>c</sup>	SRD <sup>d</sup> , n <sup>c</sup>	SRD% <sup>e</sup>
ICC <sub>2,1</sub>	0.96 (0.90–0.98)	1.4	4.0	36.5
ICC <sub>2,3</sub>	0.99 (0.97–0.99)	0.8	2.3	21.3
ICC <sub>3,3</sub>	0.99 (0.98–0.99)	0.7	2.0	18.3

Derived from MSE<sup>f</sup> 1.3 3.5 32.1

<sup>a</sup>Intra-class correlation coefficient.

<sup>b</sup>Standard error of measurement derived from ICC<sub>2.1</sub>, 2.3, 3.3 and MSE.

<sup>c</sup>Number of pegs.

<sup>d</sup>Smallest real difference derived from ICC<sub>2.1</sub>, 2.3, 3.3 and MSE.

<sup>e</sup>SRD percentage derived from ICC<sub>2.1</sub>, 2.3, 3.3 and MSE.

<sup>f</sup>Mean square error term.

For *post-intervention* trials, the ICCs incorporating random and systematic error, the ICC<sub>2.1</sub> (95% CI) for single measures and the ICC<sub>2.3</sub> (95% CI) for average measures, were 0.97 (0.95–0.98) and 0.99 (0.98–1.0), respectively (Table 3). The SRDs incorporating random and systematic error, the SRD<sub>2.1</sub> for single measures and the SRD<sub>2.3</sub> for average measures, were 3.2 and 1.8 pegs, respectively. The corresponding SRD% values for SRD<sub>2.1</sub> and SRD<sub>2.3</sub> were 25.9% and 15.0%, respectively. The SRD only incorporating random error, the SRD<sub>3.3</sub> for average measures, was 1.8 pegs.

**Table 3.** Results of reliability measures for post-intervention trials

	ICC <sup>a</sup> (95% CI)	SEM <sup>b</sup> , n <sup>c</sup>	SRD <sup>d</sup> , n <sup>c</sup>	SRD% <sup>e</sup>
ICC <sub>2.1</sub>	0.97 (0.95–0.98)	1.1	3.2	25.9
ICC <sub>2.3</sub>	0.99 (0.98–1.0)	0.7	1.8	15.0
ICC <sub>3.3</sub>	0.99 (0.98–1.0)	0.7	1.8	15.0
Derived from MSE <sup>f</sup>		1.1	3.1	25.5

<sup>a</sup>Intra-class correlation coefficient.

<sup>b</sup>Standard error of measurement derived from ICC<sub>2.1</sub>, 2.3, 3.3 and MSE.

<sup>c</sup>Number of pegs.

<sup>d</sup>Smallest real difference derived from ICC<sub>2.1</sub>, 2.3, 3.3 and MSE.

<sup>e</sup>SRD percentage derived from ICC2.1, 2.3, 3.3 and MSE.

<sup>f</sup>Mean square error term.

## DISCUSSION

This study indicated that in a selected group of persons suffering from stroke, the use of an average of three trials reduced the measurement error substantially compared to a single trial (SRD<sub>2,3</sub> vs SRD<sub>2,1</sub>). Moreover, the absolute test-retest reliability of the TFHPT was at a level that can be considered acceptable for measures representing an average of three trials and incorporating systematic error, i.e., SRD<sub>2,3</sub> and SRD%<sub>2,3</sub>.

Comparing SRD<sub>2,1</sub> to SRD<sub>2,3</sub>, revealed that the use of an average of three trials reduced the measurement error by approximately 1.5 pegs compared to the use of a single trial.[5] The result of the ANOVA indicated the presence of systematic error. Comparing SRD<sub>2,3</sub> to SRD<sub>3,3</sub>, where SRD<sub>3,3</sub> incorporates only random error, revealed that the contribution of the systematic error was approximately 0.3 pegs of the total 2.3 pegs when the average of three trials was used.[5] Although the systematic error was small compared to the random error it was not small enough to be overlooked in the assessment of reliability. A measure of absolute reliability, expressed as an absolute number of pegs for the mean of a study population, can over- or underestimate the number of pegs necessary to demonstrate an improvement for an individual.[11 12]. This limitation arises because the random error of measurements often increases with the magnitude of the measurements (i.e. heteroscedasticity),[11 12] which also, to some extent, was evident in this study (Figure 2a). To remedy this, the use of a relative measure of absolute reliability, such as SRD%, has been proposed.[11 12] However, because the heteroscedasticity was modest (Figure 2a), for the TFHP, the plain SRD appears to be a better choice for use in individuals [12]. Thus, to capture the systematic error and express the absolute reliability as an absolute



1  
2  
3 number of pegs representing an average of three trials, the most accurate measure investigated  
4  
5 in this study for assessing the absolute reliability of the TFHPT is  $SRD_{2,3}$ .  
6  
7

8 The results for  $SRD_{2,3}$  and  $SRD\%_{2,3}$  were 2.3 pegs and 21.3%, respectively. The value of  
9  
10  $SRD\%_{2,3}$  fell within the 30% level that has been suggested as acceptable.[15] The 30% level  
11  
12 seems high in this context, with persons affected by stroke in a chronic stage; from a clinical  
13  
14 viewpoint, our opinion is that the results of 21.3% and 2.3 pegs in this study indicate a barely  
15  
16 acceptable level of absolute reliability. For a favourable level, we believe that a mean number  
17  
18 consisting of approximately 1.5 pegs is desired. The relative test-retest reliability, as measured  
19  
20 by  $ICC_{2,3}$ , was 0.99, which seems excellent. The discrepancy between the level of the relative  
21  
22 and the absolute reliability is most likely caused by the heterogeneity in this study population  
23  
24 (Figure 2a, Table 1) which inflates the relative reliability.[5]  
25  
26  
27  
28  
29

30 The level of the relative absolute test-retest reliability ( $SRD\%$ ), the most comparable measure,  
31  
32 observed for the TFHPT in this study (21,3%) is better than what Chen et al.[2] reported (54%),  
33  
34 and, is at approximately the same level as Ekstrand et al.[3] reported (24%) for the NHPT. Even  
35  
36 though the  $SRD\%$  measures reported in the studies by Chen et al. and by Ekstrand et al. were  
37  
38 calculated in different ways compared to the  $SRD\%_{2,3}$  reported in this study, the measures used  
39  
40 in these three studies are fairly equivalent.[5] Several methodological differences between these  
41  
42 3 studies could have affected the results.[2, 3] *First*, the results of the TFHPT and NHPT were  
43  
44 measured using different scales, where the use of time for completion of the test in the NHPT  
45  
46 should accommodate more variability compared to the peg count in the TFHPT. However, the  
47  
48  $SRD\%$  results should still be comparable between the TFHPT and the NHPT because this  
49  
50 relative measure of absolute reliability adjusts for different scales and study populations.[11,  
51  
52 12] *Second*, in this study of the TFHPT, the test and retest trials were performed within minutes  
53  
54 compared to within days in the studies of the NHPT, which may have resulted in seemingly  
55  
56 worse reliability for the NHPT because of possible random error from day-to-day variation in  
57  
58  
59  
60

1  
2  
3 performance.[11, 16] *Third*, the 3-5 days between test and retest trials in the study by Chen et  
4 al.[2] may also have resulted in seemingly worse reliability in that study because of systematic  
5 error. A systematic error may have originated in possible recovery from stroke because the time  
6 since stroke was 3 months or less for a quarter of the study sample [17].  
7  
8  
9

10  
11  
12 One advantage with the TFHPT, compared to the NHPT, is that persons with worse motor  
13 function can be tested.[2, 3] In the study by Ekstrand et al.[3], those who did not complete the  
14 NHPT in 180 seconds were excluded. This would correspond to inserting and removing a  
15 minimum of 2.5 pegs in 50 seconds, whereas 0 pegs inserted in 50 seconds is a valid result  
16 with the TFHPT.  
17  
18  
19

20  
21  
22 There was a tendency towards improved reliability after the CIMT period which, was due to  
23 decreased systematic error and decreased random error. The decreased systematic error can be  
24 observed in the elimination of the difference between the  $SRD_{2,3}$  that incorporates systematic  
25 error and the  $SRD_{3,3}$  that does not in the post-intervention trials and in the main effects of the  
26 trial in the ANOVA results.[5] The decreased systematic error is most likely due to a decreased  
27 learning effect, when the participants had previous experience in the test. This is indicated by  
28 the increases in the mean values over the trials, especially over trials 1-2, and by the less  
29 pronounced increase in the post intervention trials.[11, 12] The lower random error can be  
30 observed from the lower  $SRD_{3,3}$  results in the post-intervention trials.[5] The cause of the  
31 decreased random error is less clear, but it could also be attributed to the decreased systematic  
32 error.[11] Furthermore, it is likely that the  $SRD_{2,3}$  result of 2.3 pegs for TFHPT could, in reality,  
33 be adjusted downwards. A peg test is often used to evaluate a rehabilitation period; because the  
34 error is smaller in the post-intervention trials, the “true” SRD may be somewhere between those  
35 of the pre- and post-intervention trials (2.3 vs 1.8).  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57

58 There were some weaknesses in this study, including the relatively low number of participants,  
59 few observations above 20 pegs and a study population in which the participants were selected  
60

1  
2  
3 because they should benefit from CIMT.[9, 11] Both of these latter objections may, to some  
4 degree, hinder the generalization of the results to other groups of people suffering from stroke.  
5  
6 The SEM and SRD are not as population-independent as the SRD% but are still considered  
7 rather robust.[5] In addition, the intended practice trial was included as one of three trials in the  
8 analyses which appears to have contributed to systematic error through an increased learning  
9 effect, indicated by a large increase in the mean values between trials 1 and 2.[5, 11, 12] Thus,  
10 to mitigate the learning effect, a practice trial preceding regular trials is recommended.  
11  
12 Moreover, the possible day-to-day variation was not captured in the present study design. The  
13 advantage of this approach is that it yields a pure result for measurement error for the instrument  
14 in this population; the disadvantage is that the result is less clinically applicable.[11, 16]  
15  
16

17 In conclusion, our results suggest that the smallest difference that can be detected using a test  
18 procedure with an average of three trials ( $SRD_{2,3}$ ) conducted by a single tester should be just  
19 above 2 pegs with the TFHPT. Furthermore, to reach an acceptable level of measurement error,  
20 the use of the average of multiple trials is crucial. Future research should focus on optimizing  
21 the number of trials.  
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 **ACKNOWLEDGEMENTS**

49 The authors thank the staff at the clinics for the extra work associated with participating in the  
50 RCT. The authors thank Håkan Littbrand for input on the conception of the study and Monika  
51 Edström for administrative work with the study.  
52  
53  
54  
55  
56  
57  
58  
59

## 60 **AUTHOR STATEMENT**

1  
2  
3 Design of the study and data interpretation SE, FG, BL and MH. Acquisition of the data SE,  
4  
5 BL and MH. Statistical analysis SE and FG. Drafting and finalization of the manuscript SE.  
6  
7 Critical revision of the manuscript FG, BL and MH. Final approval of the submitted  
8  
9 manuscript FG, BL and MH.  
10  
11  
12  
13

#### 14 **COMPETING INTERESTS**

15  
16  
17 None declared.  
18  
19

#### 20 **FUNDING**

21  
22  
23 The Norrbacka-Eugenia Foundation; The Swedish STROKE-Association; The Stroke  
24  
25 Foundation of Northern Sweden; and the Centre for Clinical Research Sörmland, Uppsala  
26  
27 University.  
28  
29

#### 30 **DATA SHARING**

31  
32  
33 No additional data are available.  
34  
35

#### 36 **PATIENT CONSENT**

37  
38  
39 Written informed consent was obtained from the participants.  
40  
41

#### 42 **ETHICS APPROVAL**

43  
44 The Regional Ethical Review Board in Umeå, reference 09-104M, with additional approval  
45  
46 Dnr 2010/314-32M, Dnr 2011-244-32M, and 2012-235-32M.  
47

#### 48 **REFERENCES**

- 49  
50 1 Santisteban L, Teremetz M, Bleton JP, et al. Upper limb outcome measures used in  
51  
52 stroke rehabilitation studies: a systematic literature review. *PLoS One*  
53  
54 2016;11:e0154792.  
55  
56 2 Chen HM, Chen CC, Hsueh IP, et al. Test-retest reproducibility and smallest real  
57  
58 difference of 5 hand function tests in patients with stroke. *Neurorehabil Neural Repair*  
59  
60 2009;23:435-40.

- 1  
2  
3 Ekstrand E, Lexell J, Brogardh C. Test-retest reliability and convergent validity of  
4 three manual dexterity measures in persons with chronic stroke. *PM R* 2016;8:935-43.  
5  
6 Carter RE, Lubinsky J, Domholdt E. Rehabilitation Research: Principles and  
7 Applications. St. Louis: Elsevier-Saunders 2011: 237-239.  
8  
9 Weir JP. Quantifying test-retest reliability using the intraclass correlation coefficient  
10 and the SEM. *J Strength Cond Res* 2005;19:231-40.  
11  
12 Mathiowetz V, Weber K, Kashman N, et al. Adult norms for the nine hole peg test of  
13 finger dexterity. *Occup Ther J Res* 1985;5:25-38.  
14  
15 Grice KO, Vogel KA, Le V, et al. Adult norms for a commercially available nine hole  
16 peg test for finger dexterity. *Am J Occup Ther* 2003;57:570-73.  
17  
18 Heller A, Wade DT, Wood VA, et al. Arm function after stroke: measurement and  
19 recovery over the first three months. *J Neurol Neurosurg Psychiatry* 1987;50:714-19.  
20  
21 Lexell JE, Downham DY. How to assess the reliability of measurements in  
22 rehabilitation. *Am J Phys Med Rehabil* 2005;84:719-23.  
23  
24 Fugl-Meyer AR, Jaasko L, Leyman I, et al. The post-stroke hemiplegic patient. 1. A  
25 method for evaluation of physical performance. *Scand J Rehabil Med* 1975;7:13-31.  
26  
27 Hopkins WG. Measures of reliability in sports medicine and science. *Sports Med*  
28 2000;30:1-15.  
29  
30 Atkinson G, Nevill AM. Statistical methods for assessing measurement error  
31 (reliability) in variables relevant to sports medicine. *Sports Med* 1998;26:217-38  
32  
33 Shrout PE, Fleiss JL. Intraclass correlations: uses in assessing rater reliability. *Psychol*  
34 *Bull* 1979;86:420-8.  
35  
36 Bland JM, Altman DG. Measurement error. *BMJ* 1996;313:744.  
37  
38 Huang SL, Hsieh CL, Lin JH, et al. Optimal scoring methods of hand-strength tests in  
39 patients with stroke. *Int J Rehabil Res* 2011;34:178-80.  
40  
41 Sim J, Wright C. Research in Health Care: Concepts, Designs and Methods.  
42 Cheltenham: Nelson Thornes Ltd 2002: 133-134.  
43  
44 Kwakkel G, Kollen B, Twisk J. Impact of time on improvement of outcome after  
45 stroke. *Stroke* 2006;37:2348-53.  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
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 only

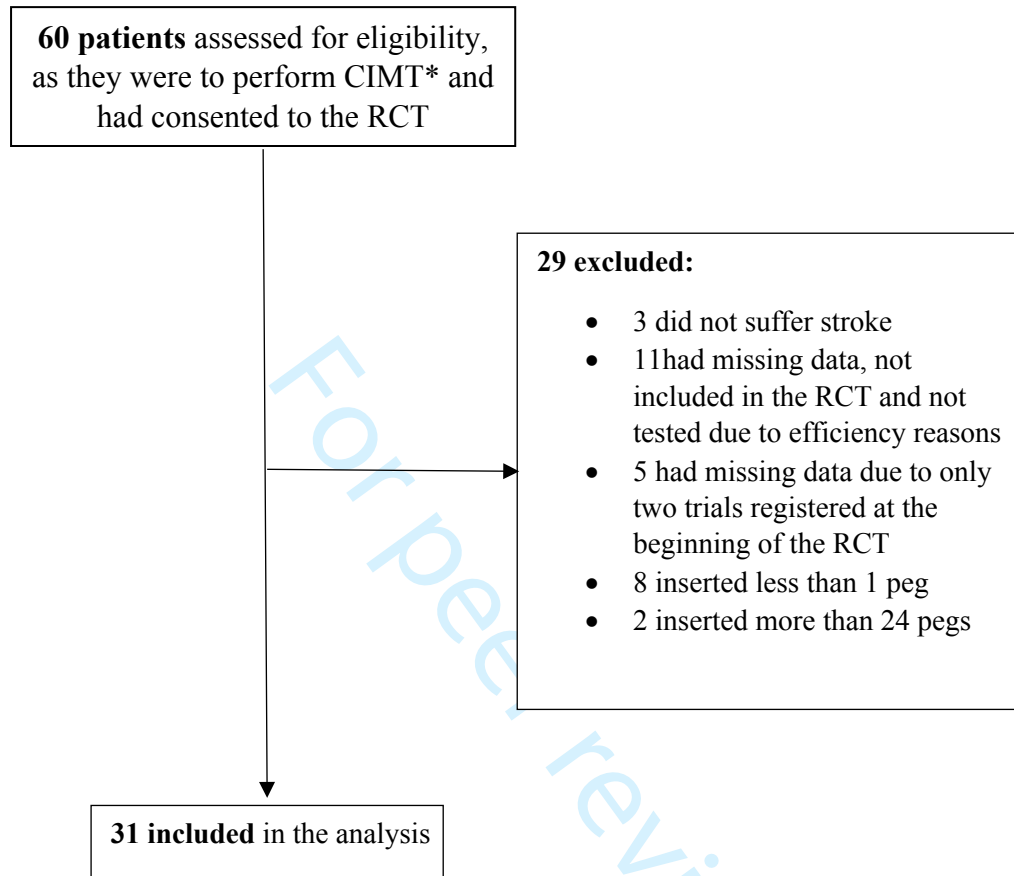
## FIGURE LEGENDS

**Figure 1.** Flowchart of the recruitment process in the study. \*Constraint-induced movement therapy.

1  
2  
3 **Figure 2.** Bland-Altman plots of numbers of pegs from pre-intervention (a) and post-  
4 intervention trials (b). The mean of trials 2 and 3 was plotted against the difference of trials 3  
5 and 2 for each subject. The centre line displays the mean difference for the group between  
6 trials 3 and 2. The upper and lower confidence limits were calculated as the mean difference  $\pm$   
7  
8  
9  
10  
11  
12 SD of the mean difference  $\times$  1.96.  
13  
14  
15  
16  
17  
18  
19  
20  
21  
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 only

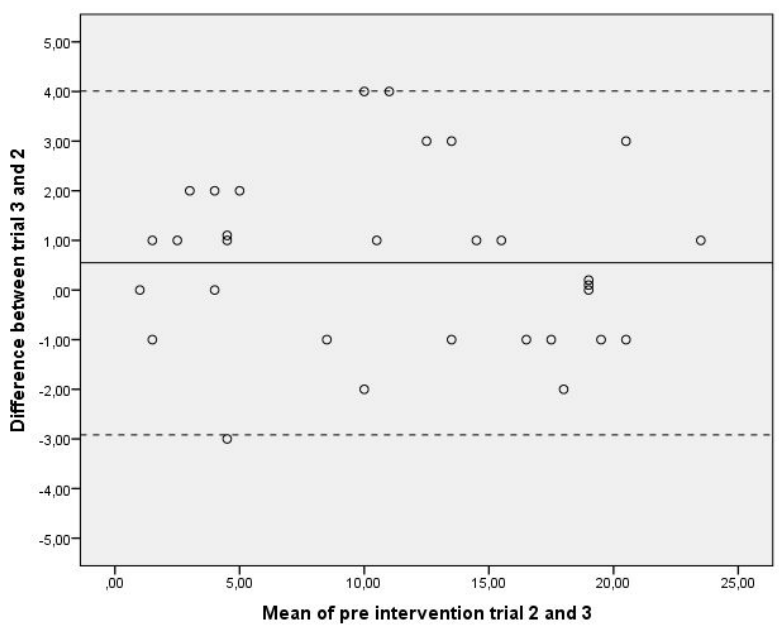
1  
2  
3 **Figure 1.**  
4  
5  
6





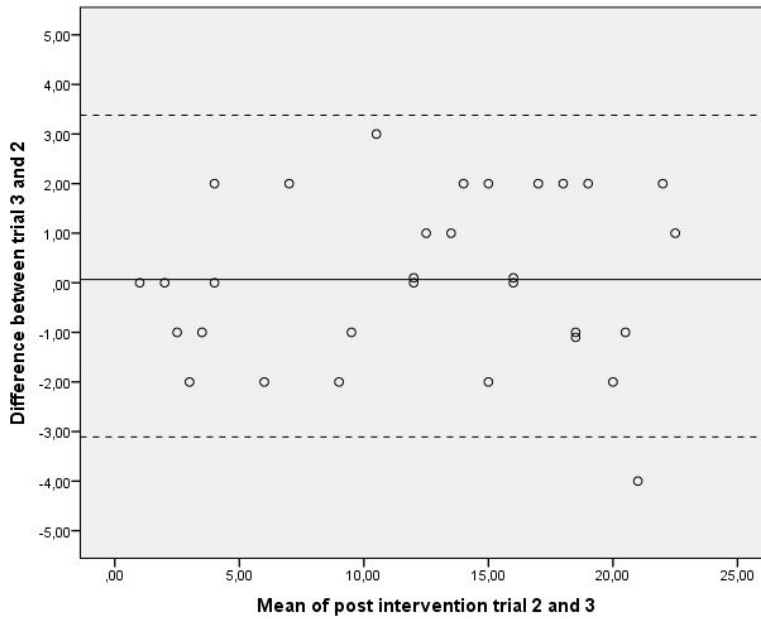
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Figure 2a.



er review only

Figure 2b.



Peer review only



# EDITORIAL CERTIFICATE

This document certifies that the manuscript listed below was edited for proper English language, grammar, punctuation, spelling, and overall style by one or more of the highly qualified native English speaking editors at American Journal Experts.

## Manuscript title:

Test-retest reliability of the twenty-five-hole peg test in stroke patients

## Authors:

Staffan Eriksson, Fredrik Granström, Mattias Hedlund, Britta Lindström

## Date Issued:

June 13, 2019

## Certificate Verification Key:

3D8C-EEE0-E006-60A1-591E



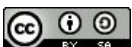
This certificate may be verified at [www.aje.com/certificate](http://www.aje.com/certificate). This document certifies that the manuscript listed above was edited for proper English language, grammar, punctuation, spelling, and overall style by one or more of the highly qualified native English speaking editors at American Journal Experts. Neither the research content nor the authors' intentions were altered in any way during the editing process. Documents receiving this certification should be English-ready for publication; however, the author has the ability to accept or reject our suggestions and changes. To verify the final AJE edited version, please visit our verification page. If you have any questions or concerns about this edited document, please contact American Journal Experts at [support@aje.com](mailto:support@aje.com).

Section & Topic	No	Item	Reported on page #
<b>TITLE OR ABSTRACT</b>			
	<b>1</b>	Identification as a study of diagnostic accuracy using at least one measure of accuracy (such as sensitivity, specificity, predictive values, or AUC)	(Reliability), 1
<b>ABSTRACT</b>			
	<b>2</b>	Structured summary of study design, methods, results, and conclusions (for specific guidance, see STARD for Abstracts)	2
<b>INTRODUCTION</b>			
	<b>3</b>	Scientific and clinical background, including the intended use and clinical role of the index test	4-5
	<b>4</b>	Study objectives and hypotheses	5
<b>METHODS</b>			
<i>Study design</i>	<b>5</b>	Whether data collection was planned before the index test and reference standard were performed (prospective study) or after (retrospective study)	10, registered, refnr: ISRCTN24868616
<i>Participants</i>	<b>6</b>	Eligibility criteria	5-6
	<b>7</b>	On what basis potentially eligible participants were identified (such as symptoms, results from previous tests, inclusion in registry)	5-6
	<b>8</b>	Where and when potentially eligible participants were identified (setting, location and dates)	Setting: 6 Dates: 11 Exact locations of the clinics not included.
	<b>9</b>	Whether participants formed a consecutive, random or convenience series	5
<i>Test methods</i>	<b>10a</b>	Index test, in sufficient detail to allow replication	Not applicable
	<b>10b</b>	Reference standard, in sufficient detail to allow replication	n.a.
	<b>11</b>	Rationale for choosing the reference standard (if alternatives exist)	n.a.
	<b>12a</b>	Definition of and rationale for test positivity cut-offs or result categories of the index test, distinguishing pre-specified from exploratory	n.a.
	<b>12b</b>	Definition of and rationale for test positivity cut-offs or result categories of the reference standard, distinguishing pre-specified from exploratory	n.a.
	<b>13a</b>	Whether clinical information and reference standard results were available to the performers/readers of the index test	n.a.
	<b>13b</b>	Whether clinical information and index test results were available to the assessors of the reference standard	7
<i>Analysis</i>	<b>14</b>	Methods for estimating or comparing measures of diagnostic accuracy	Reliability measures 8-10
	<b>15</b>	How indeterminate index test or reference standard results were handled	n.a.
	<b>16</b>	How missing data on the index test and reference standard were handled	n.a.
	<b>17</b>	Any analyses of variability in diagnostic accuracy, distinguishing pre-specified from exploratory	n.a.
	<b>18</b>	Intended sample size and how it was determined	6
<b>RESULTS</b>			
<i>Participants</i>	<b>19</b>	Flow of participants, using a diagram	11
	<b>20</b>	Baseline demographic and clinical characteristics of participants	10
	<b>21a</b>	Distribution of severity of disease in those with the target condition	Table 1, figure 2.
	<b>21b</b>	Distribution of alternative diagnoses in those without the target condition	n.a.
	<b>22</b>	Time interval and any clinical interventions between index test and reference standard	6
<i>Test results</i>	<b>23</b>	Cross tabulation of the index test results (or their distribution) by the results of the reference standard	Plots instead, figure 2.
	<b>24</b>	Estimates of diagnostic accuracy and their precision (such as 95% confidence intervals)	Page 12, table 2 and 3
	<b>25</b>	Any adverse events from performing the index test or the reference standard	n.a.
<b>DISCUSSION</b>			
	<b>26</b>	Study limitations, including sources of potential bias, statistical uncertainty, and generalisability	16-17

1		<b>27</b>	Implications for practice, including the intended use and clinical role of the index test	16
2	<b>OTHER INFORMATION</b>			
3				
4		<b>28</b>	Registration number and name of registry	ISRCTN registry; referene number: ISRCTN24868616
5				
6		<b>29</b>	Where the full study protocol can be accessed	At the registry (above), but not detailed.
7				
8		<b>30</b>	Sources of funding and other support; role of funders	18
9				

For peer review only

10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
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



# STARD 2015

---

## AIM

STARD stands for “Standards for Reporting Diagnostic accuracy studies”. This list of items was developed to contribute to the completeness and transparency of reporting of diagnostic accuracy studies. Authors can use the list to write informative study reports. Editors and peer-reviewers can use it to evaluate whether the information has been included in manuscripts submitted for publication.

---

## EXPLANATION

A **diagnostic accuracy study** evaluates the ability of one or more medical tests to correctly classify study participants as having a **target condition**. This can be a disease, a disease stage, response or benefit from therapy, or an event or condition in the future. A medical test can be an imaging procedure, a laboratory test, elements from history and physical examination, a combination of these, or any other method for collecting information about the current health status of a patient.

The test whose accuracy is evaluated is called **index test**. A study can evaluate the accuracy of one or more index tests. Evaluating the ability of a medical test to correctly classify patients is typically done by comparing the distribution of the index test results with those of the **reference standard**. The reference standard is the best available method for establishing the presence or absence of the target condition. An accuracy study can rely on one or more reference standards.

If test results are categorized as either positive or negative, the cross tabulation of the index test results against those of the reference standard can be used to estimate the **sensitivity** of the index test (the proportion of participants *with* the target condition who have a positive index test), and its **specificity** (the proportion *without* the target condition who have a negative index test). From this cross tabulation (sometimes referred to as the contingency or “2x2” table), several other accuracy statistics can be estimated, such as the positive and negative **predictive values** of the test. Confidence intervals around estimates of accuracy can then be calculated to quantify the statistical **precision** of the measurements.

If the index test results can take more than two values, categorization of test results as positive or negative requires a **test positivity cut-off**. When multiple such cut-offs can be defined, authors can report a receiver operating characteristic (ROC) curve which graphically represents the combination of sensitivity and specificity for each possible test positivity cut-off. The **area under the ROC curve** informs in a single numerical value about the overall diagnostic accuracy of the index test.

The **intended use** of a medical test can be diagnosis, screening, staging, monitoring, surveillance, prediction or prognosis. The **clinical role** of a test explains its position relative to existing tests in the clinical pathway. A replacement test, for example, replaces an existing test. A triage test is used before an existing test; an add-on test is used after an existing test.

Besides diagnostic accuracy, several other outcomes and statistics may be relevant in the evaluation of medical tests. Medical tests can also be used to classify patients for purposes other than diagnosis, such as staging or prognosis. The STARD list was not explicitly developed for these other outcomes, statistics, and study types, although most STARD items would still apply.

---

## DEVELOPMENT

This STARD list was released in 2015. The 30 items were identified by an international expert group of methodologists, researchers, and editors. The guiding principle in the development of STARD was to select items that, when reported, would help readers to judge the potential for bias in the study, to appraise the applicability of the study findings and the validity of conclusions and recommendations. The list represents an update of the first version, which was published in 2003.

More information can be found on <http://www.equator-network.org/reporting-guidelines/stard>.



# BMJ Open

## Test-retest reliability of the twenty-five-hole peg test in stroke patients

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2019-032560.R1
Article Type:	Original research
Date Submitted by the Author:	24-Oct-2019
Complete List of Authors:	Granström, Fredrik; Uppsala University, Centre for Clinical Research Sörmland Hedlund, Mattias; Umeå Universitet Medicinska fakulteten, Department of Community Medicine and rehabilitation, Physiotherapy Lindström, Britta; Umeå Universitet Medicinska fakulteten, Department of Community Medicine and Rehabilitation, Physiotherapy Eriksson, Staffan; Uppsala University, Centre for Clinical Research Sörmland; Uppsala University, Department of Neuroscience, Physiotherapy
<b>Primary Subject Heading</b>:	Rehabilitation medicine
Secondary Subject Heading:	Neurology
Keywords:	Stroke < NEUROLOGY, measurement, Reliability, Clinical assessment, hand function, twenty-five-hole peg test

SCHOLARONE™  
Manuscripts

## Test-retest reliability of the twenty-five-hole peg test in stroke patients

Fredrik Granström<sup>1</sup>, Mattias Hedlund<sup>2</sup>, Britta Lindström<sup>2</sup>, Staffan Eriksson<sup>1,3</sup>

<sup>1</sup>Centre for Clinical Research Sörmland, Uppsala University, Sweden

<sup>2</sup>Department of Community Medicine and Rehabilitation, Physiotherapy, Umeå University, Sweden

<sup>3</sup>Department of Neuroscience, Physiotherapy, Uppsala University, Sweden

### CORRESPONDING AUTHOR

Staffan Eriksson, Centre for Clinical Research Sörmland, Uppsala University, Kungsgatan 41, 631 88 Eskilstuna, Sweden. Telephone number: +46 16 105251, +46 73 949 99 33. Email: [staffan.eriksson@germed.umu.se](mailto:staffan.eriksson@germed.umu.se)

**KEY WORDS:** clinical assessment, measurement, reproducibility of results, reliability, hand function, stroke, twenty-five-hole peg test

**WORD COUNT** 3991 words, 4105 words including “strengths and limitations of this study”.



## ABSTRACT

**Objectives:** In the nine-hole peg test (NHPT), tests that are not completed within the stipulated time are excluded, resulting in floor effects. A modified NHPT exists in which the result is expressed as the number of inserted pegs per unit of time, which might be difficult to comprehend. In the twenty-five-hole peg test (TFHPT), the larger number of available pegs makes it straightforward to count the number of pegs inserted as the result. It thus provides a comprehensible result and low floor effects, as zero pegs is a valid result. The objective was to assess the test-retest reliability of the TFHPT when testing persons with stroke. A particular focus was placed on the absolute reliability, as quantified by the smallest real difference (SRD). Complementary aims were to investigate possible implications for how the TFHPT should be used and for how the SRD of the TFHPT performance should be expressed.

**Design:** This study employed a test-retest design including three trials.

**Participants, setting and outcome measure:** Thirty-one participants who had suffered a stroke were recruited from a group designated for constraint-induced movement therapy at outpatient clinics. The TFHPT result was expressed as the number of pegs inserted.

**Methods:** Absolute reliability was quantified by the SRD, including random and systematic error for a single trial,  $SRD_{2,1}$ , and for an average of three trials,  $SRD_{2,3}$ .

**Results:** The differences in the number of pegs necessary to detect a change in the TFHPT for  $SRD_{2,1}$  and  $SRD_{2,3}$  were 4.0 and 2.3, respectively.

**Conclusions:** The smallest change that can be detected in the TFHPT should be just above two pegs for a test procedure including an average of three trials. The use of an average of three trials compared to a single trial substantially reduces the measurement error.

**Trial registration:** ISRCTN registry, reference number ISRCTN24868616.

## ARTICLE SUMMARY

### Strengths and limitations of this study

- The generalizability of the results may be limited: the participants were selected because they should benefit from CIMT, and few scored above 20 pegs during the 50-second trial duration.
- Among other measures of reliability, the SRD percentage was reported; this is a good measure for comparisons between different tests, scales and populations.
- The results were presented with several different reliability measures to offer knowledge about the source of the measurement error.
- As the test-retest trials were performed within minutes, possible day-to-day variation was not captured.
- The intended practice trial was included as one of three trials in the analyses, which appears to have contributed to the learning effect.

## INTRODUCTION

For a comprehensive upper limb assessment among persons with stroke, it is important to combine a measure of proximal upper limb function with a measure of hand function.[1] However, only approximately 25% of studies regarding upper limb interventions include a specific measure of hand function.[1] The nine-hole peg test (NHPT) is a common test of hand function focusing on fine manual dexterity, and it is the most common such test used in research.[1-3] Two studies of stroke populations investigated the reliability of the NHPT, and there was a large discrepancy in the reliability reported: SRD% 24% vs. 52%.[2, 3] The NHPT has mostly shown moderate to excellent correlations (0.55-0.97) with other tests and self-reports focusing on hand function, including the Action Research Arm Test, the Jebsen-Taylor Hand Function Test, and the Stroke Impact Scale (hand function domain).[4, 5] The exception is the Motor Activity Log, for which a low correlation has been reported.[5]

A weakness of the NHPT is that many persons with stroke cannot reach the lower limit; i.e., a floor effect arises. Furthermore, if the number of completed pegs is used as an outcome measure, a test with only nine pegs can measure only a narrow range of hand function, resulting in profound ceiling effects.[6] Therefore, to widen the scale and avoid ceiling effects, the original NHPT expresses the result as the time needed to complete the test (including inserting and removing all the pegs).[2, 3, 7, 8] However, this approach aggravates the floor effects because tests that are not completed during the stipulated time (limits of 60 and 180 seconds has been used) are excluded.[2, 3] The maximum time could be prolonged; however, this would be time consuming, mentally strenuous and therefore possibly unethical due to the possibility of a non-completed test after a lengthy attempt. A modified NHPT is used to mitigate the floor effect while avoiding the ceiling effect; in this modified version, the result is expressed as the number of inserted pegs per unit of time, i.e., the frequency.[9] This modified test includes only peg insertion and not peg removal. It is thus possible also to

1  
2  
3 include tests that were not completed within the stipulated time limit and still measure  
4 performance on the same task across the entire range of hand function. However, it may be  
5 difficult both to interpret the frequency and to communicate it to other staff members and  
6 patients, especially to those suffering from a brain injury. The reliability of this modified test  
7 has not been investigated.  
8  
9

10  
11  
12 In Sweden, a similar peg test, a twenty-five-hole peg test (TFHPT), has been used in clinical  
13 practice. The larger number of available pegs makes it straightforward to count the number of  
14 pegs inserted during a stipulated time frame of 50 seconds, as the test result. Thus, the TFHPT  
15 measures fine manual dexterity on a numerical scale that is easy to comprehend, with low  
16 floor effects and presumably reasonable ceiling effects (based on pre-study data). Moreover,  
17 compared to the individuals whom the original NHPT can test, individuals with worse hand  
18 function can be tested with the TFHPT.[2, 3] Of the two studies investigating the reliability of  
19 the NHPT, the one with the most generous time limit excluded all tests that were not  
20 completed in 180 seconds.[3] This limit corresponds to inserting and removing a minimum of  
21 2.5 pegs in 50 seconds, whereas 0 pegs inserted in 50 seconds is a valid result with the  
22 TFHPT. The TFHPT has not been previously described in the literature, and its reliability has  
23 not been investigated. Due to the similarity of the NHPT and the TFHPT, the underlying skill  
24 assessed with these tests is most likely the same. However, since the tests have completely  
25 different stop criteria – a time limit for the TFHPT vs. the insertion of all pegs for the NHPT –  
26 equal reliability cannot be taken for granted.[7] Thus, if the size of the measurement error  
27 related to the TFHPT is shown to be acceptable, this test may be useful in both clinical  
28 practice and research.  
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 The overall aim of this study was to assess the test-retest reliability of the TFHPT for persons  
56 suffering from stroke. A particular focus was placed on the absolute reliability, as quantified  
57 by the smallest real difference (SRD). Complementary aims were to investigate possible  
58  
59  
60

1  
2  
3 implications for how the TFHPT should be used and for how the SRD of the TFHPT  
4  
5 performance should be expressed.  
6  
7

## 8 **METHOD**

9

### 10 **Participants**

11

12  
13  
14 The participants in this study were consecutively recruited in the process of screening patients  
15  
16 eligible for inclusion in a multicentre randomized controlled trial (RCT), reference number  
17  
18 ISRCTN24868616 at the ISRCTN registry. The patients were considered for inclusion  
19  
20 because they were to undergo constraint-induced movement therapy (CIMT) at one of the  
21  
22 clinics participating in the RCT. The clinics were outpatient rehabilitation clinics in the public  
23  
24 health care system in Sweden. Data were collected at the clinics. The sample in this study  
25  
26 consisted of included and excluded participants in the multicentre RCT. The participants were  
27  
28 included if they had one stroke or more registered in the medical record and if TFHPT data  
29  
30 were available from three trials before and three trials after the CIMT. Moreover, with regard  
31  
32 to the outcome measure, a minimum of one peg and a maximum of 24 pegs inserted was  
33  
34 necessary for inclusion. This was to avoid an untrue low measurement error from participants  
35  
36 stable at 0 or 25 pegs inserted. Because these two intervals are wider, measurements at these  
37  
38 intervals should be more stable.  
39  
40  
41  
42  
43

44  
45 A minimum of 30 participants were included to obtain a sufficient number for a reliability  
46  
47 study.[10]  
48  
49

### 50 **Procedure and measurements**

51

52  
53 The TFHPT has twenty-five holes and pegs (Figure 1). The test used in this study consisted of  
54  
55 a rectangular 21 cm × 45 cm board with a box containing pegs on one side and an elevated 18  
56  
57 cm × 18 cm area with holes on the other side. The holes were 9 mm wide and 18 mm deep,  
58  
59  
60

1  
2  
3 and they were spaced 20 mm apart. The box had a base of 13 cm × 18 cm and was 5 cm deep.  
4  
5 The pegs were 40 mm long and 8 mm in diameter.  
6  
7

8 The TFHPT was administered as the second test in a battery of different tests. The preceding  
9  
10 test, BL motor assessment, required approximately 30-60 minutes to administer.[11, 12] The  
11  
12 tests were administered in an examination room in which only the participant and the  
13  
14 physiotherapist were present. For the TFHPT, three trials were performed with each hand.  
15  
16 The participants started with the less affected hand, followed by the more affected hand, i.e.,  
17  
18 the hand of investigation in this test-retest study. The pause between trials was approximately  
19  
20 10-120 seconds. The board was placed at a distance favoured by the participant with the  
21  
22 centre row of holes centred towards the navel and the box side oriented towards the tested  
23  
24 hand. The starting position was with both hands on the board, and time keeping began upon  
25  
26 first hand contact with the pegs. We gave participants the following instructions:  
27  
28  
29  
30  
31

- 32 1. I want you to pick up one peg at a time and insert them in the holes of the board.
  - 33 2. Use only the right/left hand; you can only use the other hand to steady the board.
  - 34 3. You can fill the holes in any order you desire.
  - 35 4. We start with a practice trial.
  - 36 5. You have 50 seconds to insert as many pegs as you can. After 50 seconds, the trial is  
37 terminated.
  - 38 6. Are you ready? Ready, set, go!
  - 39 7. After the practice trial: This was practice; now come the two actual test trials, where  
40 the results are recorded. Repeat step 6.
- 41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52

53 The test-retest reliability of the TFHPT was assessed on two separate occasions, i.e., before  
54 and after a two-week training period (the CIMT). The same procedure was used on both of  
55 these occasions, and for each participant, all tests were administered by the same  
56  
57  
58  
59  
60

1  
2  
3 physiotherapist. The assessment after the CIMT period was performed as an internal  
4  
5 validation.  
6  
7

8 Two physiotherapists, SE and BL, administered the tests in this study, including the Fugl-  
9  
10 Meyer test.[13] SE has general experience with persons suffering from stroke and experience  
11  
12 administering the original NHPT. BL has extensive experience with persons suffering from  
13  
14 stroke, including administering the original NHPT. Background data from medical records  
15  
16 were collected by staff at the clinics.  
17  
18

### 19 20 **Statistics**

21  
22  
23 All three trials were used in the analyses, although the first trial was introduced to the  
24  
25 participants as a practice trial. Analyses of pre-intervention data and post-intervention data  
26  
27 were performed separately.  
28  
29

30  
31 Bland-Altman plots of trials one and two provided a graphic description of the data  
32  
33 variability. The mean of trials one and two was plotted against the difference between trials  
34  
35 two and one for each subject. Heteroscedasticity – i.e., an association between the random  
36  
37 error and the magnitude of measurements [14] – was investigated with pairwise comparisons  
38  
39 of trials using Koenker's [15] studentized test, which is useful for small samples and skewed  
40  
41 data. Heteroscedasticity is indicated by a significant result.  
42  
43

44  
45 Measurement error can be either random or systematic. In random error, there is no pattern of  
46  
47 variability between trials, whereas in systematic error, the measurements varies in a non-  
48  
49 random way; i.e., the mean values between the trials differ.[16] To investigate whether there  
50  
51 was a systematic error in test scores, one-way repeated measures ANOVA was used to detect  
52  
53 potential between trial effects.  
54  
55

56  
57 Reliability is a term that describes how the measurement result of an instrument is affected by  
58  
59 measurement error.[6, 14] Reliability can be quantified as either relative or absolute.[6]  
60

1  
2  
3 *Relative reliability* refers to the consistency of the positions of measurements relative to those  
4 of others within the tested group, and it is quantified using several intra-class correlation  
5 coefficients (ICCs).[16, 17] In ICCs, between-subject variability is related to the within-  
6 subject variability by a ratio.[16] Thus, ICCs are sensitive to the degree of between-subject  
7 variability, and with all other things being equal, a more heterogeneous sample (i.e., a larger  
8 between-subject variability) produces higher ICC values.[16, 17] The concept of *absolute*  
9 *reliability* refers to the consistency of measurements within individuals.[6, 16] Measurement  
10 error, quantified as within-subject standard deviations in repeated tests, is a common measure  
11 of absolute reliability [6, 14, 16-18] and is called the standard error of measurement (SEM).  
12 SRD is an extension of the SEM, and it can be seen as the smallest detectable difference, with  
13 95% certainty, using a test instrument on an individual.[16]

14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29 Three separate measures of *relative reliability*, i.e.,  $ICC_{2,1}$ ,  $ICC_{2,3}$ , and  $ICC_{3,3}$ , including 95%  
30 confidence intervals (CIs), were calculated. This panel of measures was used to compare the  
31 results representative of single and average measures and to obtain an estimate of the  
32 influence of systematic error. *The first* figure in the ICC designation represents the type of  
33 ICC model.[16]  $ICC_{2,1}$  and  $ICC_{2,3}$  are calculated from a two-way random effect model and  
34 incorporate both systematic and random error, whereas  $ICC_{3,3}$  is calculated from a two-way  
35 fixed effect model and incorporates only random error.[16, 19] Thus, the less systematic error  
36 contributes to the total error, the closer  $ICC_{2,3}$  is to  $ICC_{3,3}$ . [16] *The second* figure in the ICC  
37 designation represents single or average measures, where “1” represents single measures and  
38 “2” or higher represents the number of trials from which the average is calculated.[16]  $ICC_{2,1}$   
39 represents the reliability of a test procedure in which the subject is tested with a single trial on  
40 a test occasion.[16]  $ICC_{2,3}$  and  $ICC_{3,3}$  represents the reliability of a test procedure in which the  
41 subject is tested with three trials on a test occasion and the score is expressed as the average  
42 of these trials.



To estimate *absolute reliability*, the SEM, SRD and SRD percentage (SRD%) were calculated for each of the three different ICC measures (ICC<sub>2,1</sub>, ICC<sub>2,3</sub>, and ICC<sub>3,3</sub>), resulting in the corresponding properties SEM<sub>2,1</sub>, SRD<sub>2,1</sub>, SRD%<sub>2,1</sub>, SEM<sub>2,3</sub>, SRD<sub>2,3</sub>, SRD%<sub>2,3</sub>, SEM<sub>3,3</sub>, SRD<sub>3,3</sub>, and SRD%<sub>3,3</sub>. The SEM was calculated according to  $SEM = SD\sqrt{(1 - ICC)}$ , where SD was calculated from the total sum of squares (SS<sub>TOTAL</sub>) in the ANOVA table generated in the ICC analyses as  $\sqrt{SS_{total}/(n - 1)}$ .<sup>[16]</sup> The SRD was calculated using the formula  $1.96 \times SEM \times \sqrt{2}$ , where 1.96 is related to the 95% CI and  $\sqrt{2}$  refers to the error of two measurements.<sup>[16]</sup> The SRD% was calculated by dividing the SRD value by the grand mean multiplied by 100.<sup>[2, 10]</sup> This value is independent of measurement units and is indexed to the mean value of the observations from which it was derived. It is therefore a good measure for comparisons between different tests, scales and populations.<sup>[10, 14, 17]</sup> An SRD% of 30% has been suggested as an acceptable level of reliability.<sup>[20]</sup>

Because estimates of absolute reliability vary with the type of ICC value, some caution is warranted when comparing them with measures from other studies.<sup>[16]</sup> Therefore, SEM<sub>mean square error term (MSE)</sub> =  $\sqrt{MSE}$  was also calculated, where MSE (this term is called residual error by Hopkins and the mean square residual in the SPSS output) was taken from the ANOVA table of the ICC calculation.<sup>[16, 17]</sup> This SEM measure represents the reliability of a test procedure in which the subject is tested with a single trial on a test occasion, and it is a pure measure of random error.<sup>[16]</sup> SRD<sub>MSE</sub> and SRD%<sub>MSE</sub> were also derived from SEM<sub>MSE</sub>.

The analysis of test-retest reliability was pre-planned. SPSS version 21 was used to calculate ICC and ANOVA. The alpha level was set to 0.05.

### **Patient and public involvement**

No patients or members of the public were involved in the development or design of this study.

## RESULTS

In this study, participants were recruited between January 2011 and September 2014. Of 60 eligible patients, 29 were excluded for any of the following reasons: not suffering a stroke, missing data, and yielding either below the minimum or above the maximum number of inserted pegs (Figure 2). This yielded 31 participants (21 men and 10 women) for inclusion in the analysis, with a mean  $\pm$  SD age of  $66 \pm 9$  years (Table 1). The two eligible patients who were excluded because they exceeded the permitted maximum number of pegs inserted completed the 25 pegs in their best trial within 49.3 seconds and 39.4 seconds. Of the ten patients who were excluded because they fell below the minimum number of inserted pegs, five inserted at least one peg in one of the trials. Data were collected from 17 and 14 participants by the two physiotherapists (third and last author, respectively) at seven clinics.

**Table 1.** Characteristics of participants at pre-intervention trials

Participants	N=31
Age (years), mean $\pm$ SD <sup>a</sup>	66 $\pm$ 9
Men/women, n <sup>b</sup>	21/10
Time since stroke (months), median (IQR <sup>d</sup> ), (min-max)	17 (8–24), (2–70)
Previous dominant hand more affected by stroke, n	19
TFHPT <sup>c</sup> , mean of three trials (number of pegs), mean $\pm$ SD, (min-max)	10.8 $\pm$ 6.8, (1–22.7)
Fugl-Meyer test (score), median (IQR <sup>d</sup> ), (min-max)	46 (41–53), (29–62)
More than one stroke, n	3

<sup>a</sup>Standard deviation, <sup>b</sup>number of participants, <sup>c</sup>twenty-five-hole

peg test, <sup>d</sup>interquartile range

A graphic description of the data variability can be seen in the Bland-Altman plots (Figures 3 and 4). According to Koenker's studentized test, the measurement error was not affected by heteroscedasticity (Table 2).

**Table 2.** Results of the Koenker's studentized test, n=31

Pairwise test of trials	Pre-intervention trials		Post-intervention trials	
	Chi-square	P-value	Chi-square	P-value
1-2	1.33	0.25	0.41	0.52
2-3	0.05	0.83	1.38	0.24
1-3	0.28	0.60	0.20	0.66

For *pre-intervention* trials, the mean values  $\pm$  SDs for trials 1, 2 and 3 were  $10.0 \pm 6.5$ ,  $11.0 \pm 7.1$ , and  $11.5 \pm 6.9$ , respectively. The one-way repeated measures ANOVA revealed a main effect between trials with  $F(2, 29) = 10.9$  and  $p < 0.001$ . Post hoc tests revealed differences between trials 2 and 1 and between trials 3 and 1, with mean differences (95% CIs) of 1.0 (0.3–1.6) and 1.5 (0.9–2.2), respectively.

For *post-intervention* trials, the mean values  $\pm$  SD for trials 1, 2 and 3 were  $11.8 \pm 6.5$ ,  $12.4 \pm 6.7$ , and  $12.5 \pm 6.8$ , respectively. The one-way repeated measures ANOVA revealed a main effect between trials with  $F(2, 29) = 4.1$  and  $p = 0.027$ . Post hoc tests revealed a difference between trials 3 and 1, with a mean difference (95% CIs) of 0.6 (0.2–1.1).

For *pre-intervention* trials,  $ICC_{2,3}$  (95% CI) was 0.99 (0.97–0.99) (Table 3). The SRDs incorporating random and systematic error,  $SRD_{2,1}$  and  $SRD_{2,3}$ , were 4.0 and 2.3 pegs,

respectively. The corresponding SRD% values for SRD<sub>2,1</sub> and SRD<sub>2,3</sub> were 36.5% and 21.3%, respectively. The SRD incorporating only random error, SRD<sub>3,3</sub>, was 2.0 pegs.

For *post-intervention* trials, SRD<sub>2,1</sub> and SRD<sub>2,3</sub> were 3.2 and 1.8 pegs, respectively (Table 4). SRD<sub>3,3</sub>, was 1.8 pegs.

**Table 3.** Results of reliability measures for pre-intervention trials

	ICC <sup>a</sup> (95% CI)	SEM <sup>b</sup> , n <sup>c</sup>	SRD <sup>d</sup> , n <sup>c</sup>	SRD% <sup>e</sup>
ICC <sub>2,1</sub>	0.96 (0.90–0.98)	1.4	4.0	36.5
ICC <sub>2,3</sub>	0.99 (0.97–0.99)	0.8	2.3	21.3
ICC <sub>3,3</sub>	0.99 (0.98–0.99)	0.7	2.0	18.3
Derived from MSE <sup>f</sup>		1.3	3.5	32.1

<sup>a</sup>Intra-class correlation coefficient.

<sup>b</sup>Standard error of measurement derived from ICC<sub>2,1</sub>, <sub>2,3</sub>, <sub>3,3</sub> and MSE.

<sup>c</sup>Number of pegs.

<sup>d</sup>Smallest real difference derived from ICC<sub>2,1</sub>, <sub>2,3</sub>, <sub>3,3</sub> and MSE.

<sup>e</sup>SRD percentage derived from ICC<sub>2,1</sub>, <sub>2,3</sub>, <sub>3,3</sub> and MSE.

<sup>f</sup>Mean square error term.

**Table 4.** Results of reliability measures for post-intervention trials

	ICC <sup>a</sup> (95% CI)	SEM <sup>b</sup> , n <sup>c</sup>	SRD <sup>d</sup> , n <sup>c</sup>	SRD% <sup>e</sup>
ICC <sub>2,1</sub>	0.97 (0.95–0.98)	1.1	3.2	25.9
ICC <sub>2,3</sub>	0.99 (0.98–1.0)	0.7	1.8	15.0
ICC <sub>3,3</sub>	0.99 (0.98–1.0)	0.7	1.8	15.0
Derived from MSE <sup>f</sup>		1.1	3.1	25.5

<sup>a</sup>Intra-class correlation coefficient.

<sup>b</sup>Standard error of measurement derived from ICC2.1, 2.3, 3.3 and MSE.

<sup>c</sup>Number of pegs.

<sup>d</sup>Smallest real difference derived from ICC2.1, 2.3, 3.3 and MSE.

<sup>e</sup>SRD percentage derived from ICC2.1, 2.3, 3.3 and MSE.

<sup>f</sup>Mean square error term.

## DISCUSSION

This study indicated that in a selected group of persons suffering from stroke, the absolute test-retest reliability of the TFHPT was at a level that can be considered acceptable for measures representing an average of three trials and incorporating systematic error.

To assess implications for the use of the TFHPT and to determine which SRD measure best captures the absolute reliability, three issues were considered: 1) whether to use single or average measures, 2) whether to include systematic error in the assessments, and 3) whether to take heteroscedasticity into account.

Comparing SRD<sub>2,1</sub> to SRD<sub>2,3</sub> revealed that the use of an average of three trials reduced the measurement error by approximately 1.5 pegs compared to the use of a single trial. This finding suggests that the reliability of the TFHPT is substantially improved when an average of three trials is used.

Comparing SRD<sub>2,3</sub> to SRD<sub>3,3</sub>, where SRD<sub>3,3</sub> incorporates only random error, revealed that the contribution of the systematic error was approximately 0.3 pegs of the total 2.3 pegs when the average of three trials was used.[16] Although the systematic error was small compared to the random error it was not small enough to be overlooked in the assessment of reliability. Therefore, SRD<sub>2,3</sub> is preferable to SRD<sub>3,3</sub> for measuring the reliability of the TFHPT.

1  
2  
3 The choice of SRD% instead of SRD is dependent on whether the measurement error is affected  
4 by heteroscedasticity. A measure of absolute reliability, expressed as an absolute number of  
5 pegs, can over- or underestimate the number of pegs necessary to demonstrate an improvement  
6 for an individual.[14, 17] The reason is that the random error of measurements often increases  
7 with the magnitude of the measurements (i.e., heteroscedasticity).[14, 17] As a remedy, the use  
8 of a relative measure of absolute reliability, such as SRD%, has been proposed.[14, 17]  
9  
10 However, the lack of heteroscedasticity detection suggests that both  $SRD_{2,3}$  and  $SRD\%_{2,3}$  are  
11 appropriate measures of reliability for the TFHPT.[14]

12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22 The results for  $SRD_{2,3}$  and  $SRD\%_{2,3}$  were 2.3 pegs and 21.3%, respectively. The value of  
23  $SRD\%_{2,3}$  fell within the 30% level, which has been suggested as acceptable.[20] The 30% level  
24 seems high in this context, with persons affected by stroke in a chronic stage; from a clinical  
25 viewpoint, our opinion is that the results of 21.3% and 2.3 pegs in this study indicate a barely  
26 acceptable level of absolute reliability. For a favourable level, we believe that a mean number  
27 consisting of approximately 1.5 pegs is desirable.

28  
29  
30  
31  
32  
33  
34  
35  
36  
37 The relative test-retest reliability, as measured by  $ICC_{2,3}$ , was 0.99, which seems excellent. The  
38 discrepancy between the level of the relative and the absolute reliability is most likely caused  
39 by the heterogeneity in this study population (Figure 2a, Table 1) which inflates the relative  
40 reliability.[16]

41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
The level of the relative absolute test-retest reliability (SRD%), the most comparable measure,  
observed for the TFHPT in this study (21.3%) is better than what Chen et al.[2] reported (54%),  
and, is at approximately the same level as Ekstrand et al.[3] reported (24%) for the NHPT.  
Although the SRD% measures reported in the studies by Chen et al. and by Ekstrand et al. were  
calculated in different ways than the  $SRD\%_{2,3}$  reported in this study, the measures used in these  
three studies are fairly equivalent.[16] Several methodological differences between these three  
studies could have affected the results.[2, 3] *First*, the results of the TFHPT and NHPT were

1  
2  
3 measured using different scales, where the use of time for completion of the test in the NHPT  
4 should accommodate more variability than the peg count used in the TFHPT. However, the  
5 SRD% results should still be comparable between the TFHPT and the NHPT because this  
6 relative measure of absolute reliability adjusts for different scales and study populations.[14,  
7 17] *Second*, in this study of the TFHPT, the test and retest trials were performed within minutes  
8 compared to within days in the studies of the NHPT. Thus, the TFHPT may seem more reliable  
9 because of possible random error from day-to-day variation in performance which was not  
10 captured in this study.[17, 21] *Third*, the longer time since stroke in this study of the TFHPT  
11 compared to the study of NHPT by Chen et al.[2] may have resulted in seemingly better  
12 reliability for the TFHPT because of a more stable level of hand function. In the study by Chen  
13 et al., a systematic error may have originated in recovery from stroke in the 3-5 days between  
14 the test and retest trials because the time since stroke was 3 months or less for a quarter of the  
15 study sample.[22]

16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34 It seems that the ceiling effects in the TFHPT can be considered acceptable. Only two of the  
35 persons assessed for eligibility inserted 25 pegs, and only one of them actually did hit the  
36 ceiling because according to this individual's best times for completion of the 25 pegs, he/she  
37 would have been able to insert more pegs if available. This occurred in a sample where  
38 approximately a quarter of the included participants suffered from a mild impairment of arm  
39 and hand function, as judged by the Fugl-Meyer test.[23, 24]

40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
There was a tendency towards improved reliability after the CIMT period, which was due to  
decreased systematic error and decreased random error. The decreased systematic error can be  
observed in the main effects of trial in the ANOVA results.[16] The decreased systematic error  
is most likely due to a decreased learning effect when the participants had previous experience  
in the test. The learning effect is indicated by the increases in the mean values over the trials,  
especially over trials 1-2, and the decreased learning effect is indicated by the less pronounced

1  
2  
3 increase in the post-intervention trials.[14, 17] The lower random error can be observed from  
4  
5 the lower SRD<sub>3,3</sub> results in the post-intervention trials.[16] The cause of the decreased random  
6  
7 error is less clear, but it could also be attributed to the decreased systematic error.[17] This is  
8  
9 because, the magnitude of the learning effect probably differs between individuals, which will  
10  
11 show as random error. Furthermore, it is likely that the SRD<sub>2,3</sub> result of 2.3 pegs for TFHPT  
12  
13 could, in reality, be adjusted downwards. A peg test is often used to evaluate a rehabilitation  
14  
15 period; because the error is smaller in the post-intervention trials, the “true” SRD may be  
16  
17 somewhere between the SRDs of the pre- and post-intervention trials (2.3 vs 1.8).  
18  
19  
20

21  
22 Four weaknesses of this study should be considered. The sample included a relatively low  
23  
24 number of participants with few observations above 20 pegs, and participants that were selected  
25  
26 because they should benefit from CIMT.[10, 17] These sample qualities may thus, to some  
27  
28 degree, hinder the generalization of the results to other groups of people suffering from stroke.  
29  
30 In addition, the intended practice trial was included as one of three trials in the analyses which  
31  
32 appears to have contributed to systematic error through an increased learning effect, indicated  
33  
34 by a large increase in the mean values between trials 1 and 2.[14, 16, 17] Thus, to mitigate the  
35  
36 learning effect, a practice trial preceding regular trials is recommended. Moreover, the possible  
37  
38 day-to-day variation was not captured in the present study design. The advantage of this  
39  
40 approach is that it yields a pure result for measurement error for the instrument in this  
41  
42 population; the disadvantage is that the result is less clinically applicable.[17, 21] Finally, in  
43  
44 this study, sensitivity to change and validity were not examined. However, the criterion validity  
45  
46 for NHPT has mostly shown a moderate to excellent level [4, 5] and the underlying skill  
47  
48 assessed with the TFHPT is most likely the same. A high reliability level is a prerequisite for a  
49  
50 high validity, and because the reliability of the TFHPT was at the same level as that of the  
51  
52 NHPT, the criterion validity should also be similar.[21]  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 In conclusion, our results suggest that the smallest detectable difference between two  
4 assessments using a test procedure with an average of three trials conducted by a single tester  
5 should be just above two pegs with the TFHPT. Furthermore, to reach an acceptable level of  
6 measurement error, the use of the average of multiple trials is crucial. Future research should  
7 focus on optimizing the number of trials.  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19

## 20 **ACKNOWLEDGEMENTS**

21  
22 The authors thank the staff at the clinics for the extra work associated with participating in the  
23 RCT. The authors thank Håkan Littbrand for input on the conception of the study and Monika  
24 Edström for administrative work with the study.  
25  
26  
27  
28  
29  
30

## 31 **AUTHOR CONTRIBUTIONS**

32  
33 Study design and data interpretation: SE, FG, BL and MH. Data acquisition: SE, BL and MH.  
34 Statistical analysis: SE and FG. Drafting and finalization of the manuscript: SE. Critical  
35 revision of the manuscript: FG, BL and MH. Final approval of the submitted manuscript: FG,  
36 BL and MH.  
37  
38  
39  
40  
41  
42  
43  
44

## 45 **COMPETING INTERESTS**

46  
47 None declared.  
48  
49  
50

## 51 **FUNDING**

52  
53 The Norrbacka-Eugenia Foundation; The Swedish STROKE-Association; The Stroke  
54 Foundation of Northern Sweden; and the Centre for Clinical Research Sörmland, Uppsala  
55 University.  
56  
57  
58  
59  
60

1  
2  
3  
4  
5 **DATA SHARING**  
6

7 No additional data are available.  
8  
9

10 **PATIENT CONSENT**  
11

12 Written informed consent was obtained from the participants.  
13  
14

15 **ETHICS APPROVAL**  
16

17 The Regional Ethical Review Board in Umeå, reference 09-104M, with additional approval  
18 Dnr 2010/314-32M, Dnr 2011-244-32M, and 2012-235-32M.  
19  
20  
21  
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

**REFERENCES**

1. Santisteban L, Teremetz M, Bleton JP, et al. Upper limb outcome measures used in stroke rehabilitation studies: a systematic literature review. *PLoS One* 2016;11:e0154792.
2. Chen HM, Chen CC, Hsueh IP, et al. Test-retest reproducibility and smallest real difference of 5 hand function tests in patients with stroke. *Neurorehabil Neural Repair* 2009;23:435–40.
3. Ekstrand E, Lexell J, Brogardh C. Test-retest reliability and convergent validity of three manual dexterity measures in persons with chronic stroke. *PM R* 2016;8:935–43.
4. Beebe JA, Lang CE. Relationships and responsiveness of six upper extremity function tests during the first six months of recovery after stroke. *J Neurol Phys Ther* 2009;33:96–103.
5. Lin KC, Chuang LL, Wu CY, et al. Responsiveness and validity of three dexterous function measures in stroke rehabilitation. *J Rehabil Res Dev* 2010;47:563–71.
6. Carter RE, Lubinsky J, Domholdt E. *Rehabilitation Research: Principles and Applications*. St. Louis, MO: Elsevier-Saunders, 2016:239–244.
7. Mathiowetz V, Weber K, Kashman N, et al. Adult norms for the Nine-Hole Peg Test of Finger Dexterity. *Occup Ther J Res* 1985;5:25–38.
8. Grice KO, Vogel KA, Le V, et al. Adult norms for a commercially available Nine Hole Peg Test for finger dexterity. *Am J Occup Ther* 2003;57:570–3.
9. Heller A, Wade DT, Wood VA, et al. Arm function after stroke: measurement and recovery over the first three months. *J Neurol Neurosurg Psychiatry* 1987;50:714–9.
10. Lexell JE, Downham DY. How to assess the reliability of measurements in rehabilitation. *Am J Phys Med Rehabil* 2005;84:719–23.

11. Lindmark B, Hamrin E. Evaluation of functional capacity after stroke as a basis for active intervention. Validation of a modified chart for motor capacity assessment. *Scand J Rehabil Med* 1988;20:111–5.
12. Lindmark B, Hamrin E. Evaluation of functional capacity after stroke as a basis for active intervention. Presentation of a modified chart for motor capacity assessment and its reliability. *Scand J Rehabil Med* 1988;20:103–9.
13. Fugl-Meyer AR, Jaasko L, Leyman I, et al. The post-stroke hemiplegic patient. 1. A method for evaluating of physical performance. *Scand J Rehabil Med* 1975;7:13-31.
14. Atkinson G, Nevill AM. Statistical methods for assessing measurement error (reliability) in variables relevant to sports medicine. *Sports Med* 1998;26:217–38.
15. Koenker R. A note on studentizing a test for heteroscedasticity. *J Econom* 1981;17:107–12.
16. Weir JP. Quantifying test-retest reliability using the intraclass correlation coefficient and the SEM. *J Strength Cond Res* 2005;19:231–40.
17. Hopkins WG. Measures of reliability in sports medicine and science. *Sports Med* 2000;30:1–15.
18. Bland JM, Altman DG. Measurement error. *BMJ* 1996;313:744.
19. Shrout PE, Fleiss JL. Intraclass correlations: uses in assessing rater reliability. *Psychol Bull* 1979;86:420–8.
20. Huang SL, Hsieh CL, Lin JH, et al. Optimal scoring methods of hand-strength tests in patients with stroke. *Int J Rehabil Res* 2011;34:178–80.
21. Sim J, Wright C. *Research in Health Care: Concepts, Designs and Methods*. Cheltenham, England: Nelson Thornes Ltd 2002: 123–33.
22. Kwakkel G, Kollen B, Twisk J. Impact of time on improvement of outcome after stroke. *Stroke* 2006;37:2348–53.

- 1  
2  
3 23. Duncan PW, Goldstein LB, Horner RD, et al. Similar motor recovery of upper and  
4 lower extremities after stroke. *Stroke* 1994;25:1181–8.  
5  
6  
7 24. Woytowicz EJ, Rietschel JC, Goodman RN, et al. Determining levels of upper  
8 extremity movement impairment by applying a cluster analysis to the fugl-meyer  
9 assessment of the upper extremity in chronic stroke. *Arch Phys Med Rehabil*  
10 2017;98:456–62.  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
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 only

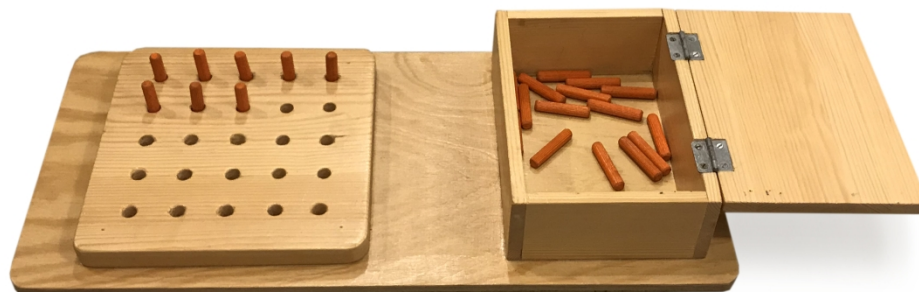
## FIGURE LEGENDS

**Figure 1.** The twenty-five-hole peg test.

**Figure 2.** Flowchart of the recruitment process in the study. \*Constraint-induced movement therapy.

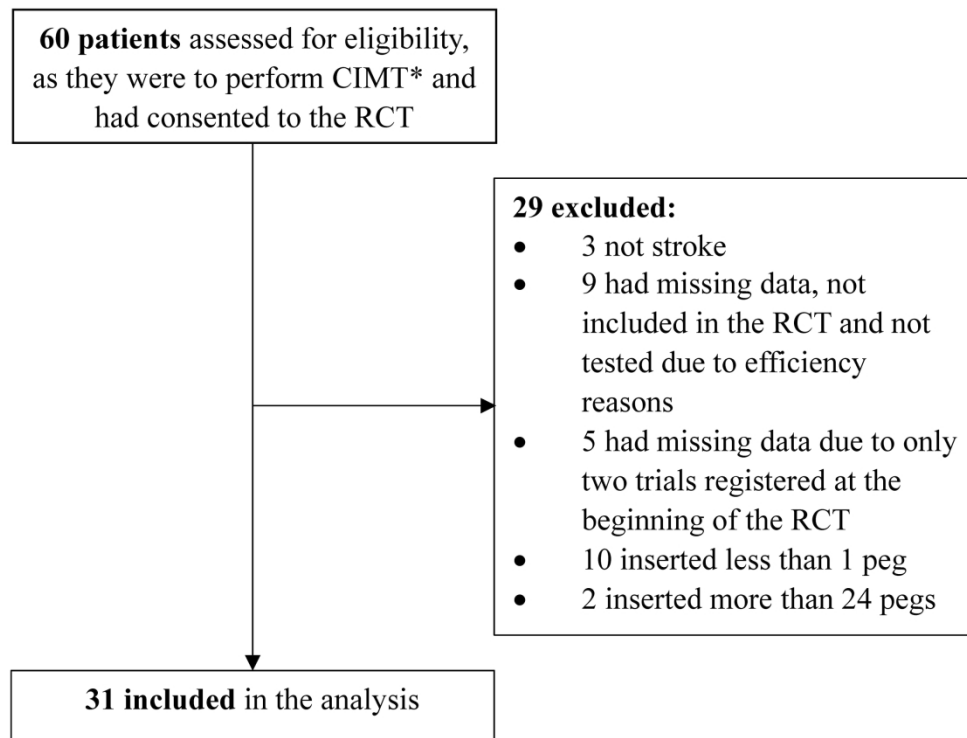
**Figure 3.** Bland-Altman plots of numbers of pegs from pre-intervention trials. The mean of trials 1 and 2 was plotted against the difference of trials 2 and 1 for each subject. The centre line displays the mean difference for the group between trials 2 and 1. The upper and lower confidence limits were calculated as the mean difference  $\pm$  SD of the mean difference  $\times$  1.96.

**Figure 4.** Bland-Altman plots of numbers of pegs from post-intervention trials. The mean of trials 1 and 2 was plotted against the difference of trials 2 and 1 for each subject. The centre line displays the mean difference for the group between trials 2 and 1. The upper and lower confidence limits were calculated as the mean difference  $\pm$  SD of the mean difference  $\times$  1.96.



The twenty-five-hole peg test.

208x140mm (300 x 300 DPI)

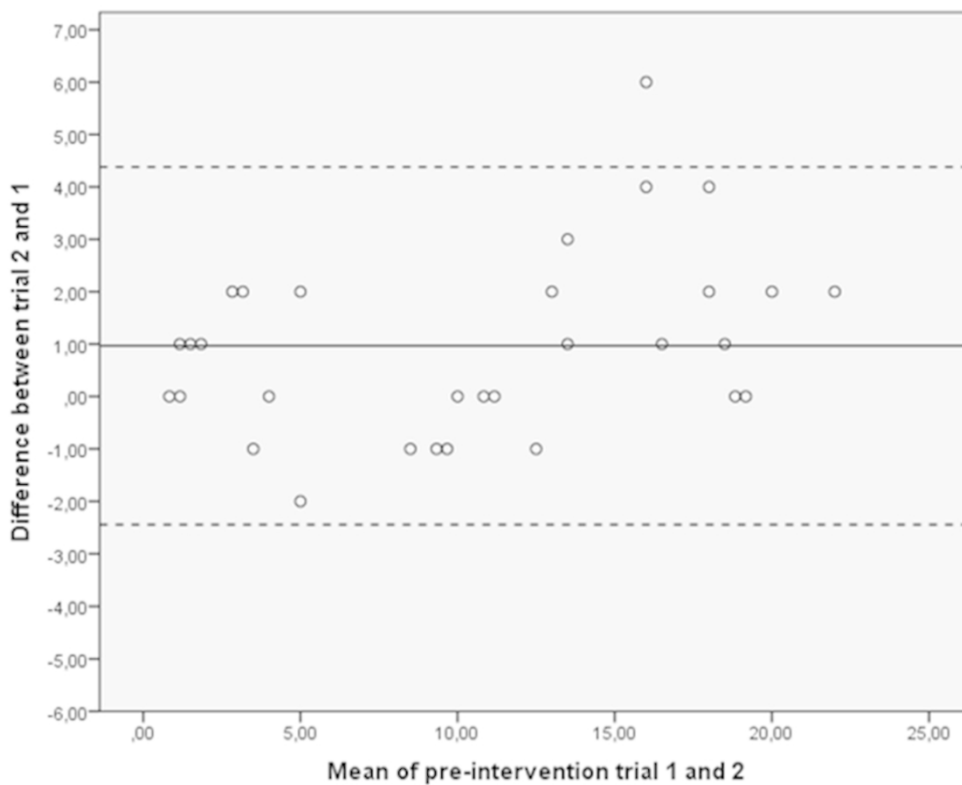


31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Flowchart of the recruitment process in the study. \*Constraint-induced movement therapy.

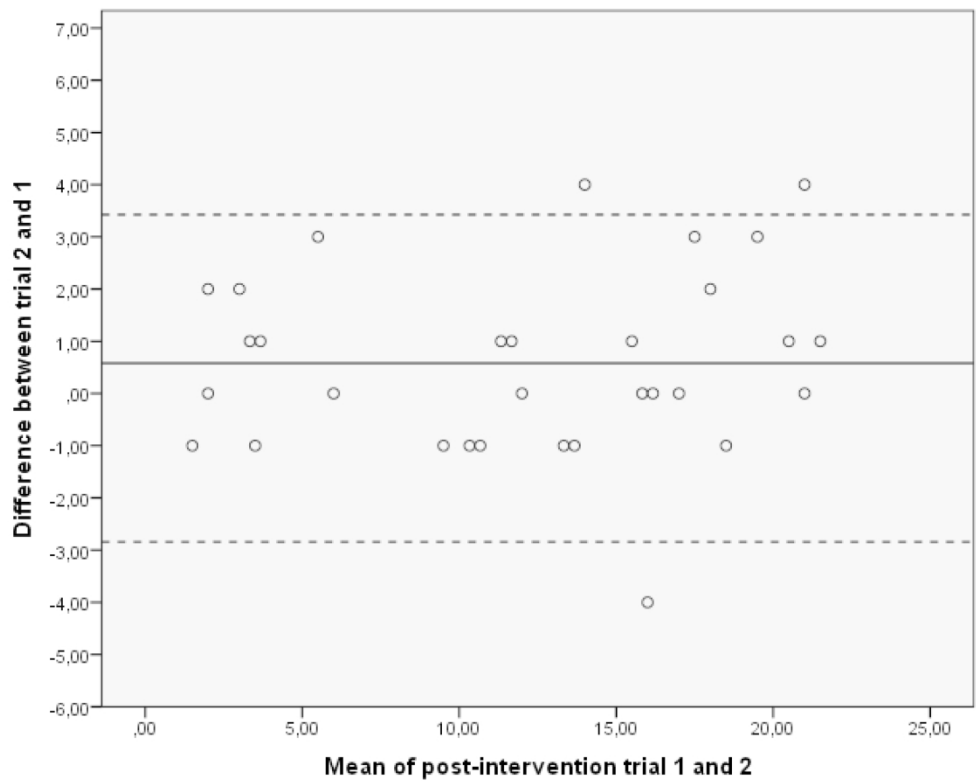
121x92mm (600 x 600 DPI)





Bland-Altman plots of numbers of pegs from pre-intervention trials. The mean of trials 1 and 2 was plotted against the difference of trials 2 and 1 for each subject. The centre line displays the mean difference for the group between trials 2 and 1. The upper and lower confidence limits were calculated as the mean difference  $\pm$  SD of the mean difference  $\times$  1.96.

140x112mm (600 x 600 DPI)



Bland-Altman plots of numbers of pegs from post-intervention trials. The mean of trials 1 and 2 was plotted against the difference of trials 2 and 1 for each subject. The centre line displays the mean difference for the group between trials 2 and 1. The upper and lower confidence limits were calculated as the mean difference  $\pm$  SD of the mean difference  $\times$  1.96.

140x112mm (600 x 600 DPI)

Section & Topic	No	Item	Reported on page #
<b>TITLE OR ABSTRACT</b>			
	<b>1</b>	Identification as a study of diagnostic accuracy using at least one measure of accuracy (such as sensitivity, specificity, predictive values, or AUC)	(Reliability), 1
<b>ABSTRACT</b>			
	<b>2</b>	Structured summary of study design, methods, results, and conclusions (for specific guidance, see STARD for Abstracts)	2
<b>INTRODUCTION</b>			
	<b>3</b>	Scientific and clinical background, including the intended use and clinical role of the index test	4-5
	<b>4</b>	Study objectives and hypotheses	5-6
<b>METHODS</b>			
<i>Study design</i>	<b>5</b>	Whether data collection was planned before the index test and reference standard were performed (prospective study) or after (retrospective study)	10, registered, refnr: ISRCTN24868616
<i>Participants</i>	<b>6</b>	Eligibility criteria	6
	<b>7</b>	On what basis potentially eligible participants were identified (such as symptoms, results from previous tests, inclusion in registry)	6
	<b>8</b>	Where and when potentially eligible participants were identified (setting, location and dates)	Setting: 6 Dates: 11 Exact locations of the clinics not included.
	<b>9</b>	Whether participants formed a consecutive, random or convenience series	6
<i>Test methods</i>	<b>10a</b>	Index test, in sufficient detail to allow replication	Not applicable
	<b>10b</b>	Reference standard, in sufficient detail to allow replication	n.a.
	<b>11</b>	Rationale for choosing the reference standard (if alternatives exist)	n.a.
	<b>12a</b>	Definition of and rationale for test positivity cut-offs or result categories of the index test, distinguishing pre-specified from exploratory	n.a.
	<b>12b</b>	Definition of and rationale for test positivity cut-offs or result categories of the reference standard, distinguishing pre-specified from exploratory	n.a.
	<b>13a</b>	Whether clinical information and reference standard results were available to the performers/readers of the index test	n.a.
	<b>13b</b>	Whether clinical information and index test results were available to the assessors of the reference standard	7
<i>Analysis</i>	<b>14</b>	Methods for estimating or comparing measures of diagnostic accuracy	Reliability measures 8-10
	<b>15</b>	How indeterminate index test or reference standard results were handled	n.a.
	<b>16</b>	How missing data on the index test and reference standard were handled	n.a.
	<b>17</b>	Any analyses of variability in diagnostic accuracy, distinguishing pre-specified from exploratory	n.a.
	<b>18</b>	Intended sample size and how it was determined	6
<b>RESULTS</b>			
<i>Participants</i>	<b>19</b>	Flow of participants, using a diagram	11
	<b>20</b>	Baseline demographic and clinical characteristics of participants	11
	<b>21a</b>	Distribution of severity of disease in those with the target condition	Table 1, figure 3 and 4.
	<b>21b</b>	Distribution of alternative diagnoses in those without the target condition	n.a.
	<b>22</b>	Time interval and any clinical interventions between index test and reference standard	7
<i>Test results</i>	<b>23</b>	Cross tabulation of the index test results (or their distribution) by the results of the reference standard	Plots instead, figure 3 and 4.
	<b>24</b>	Estimates of diagnostic accuracy and their precision (such as 95% confidence intervals)	Page 12, table 3 and 4
	<b>25</b>	Any adverse events from performing the index test or the reference standard	n.a.
<b>DISCUSSION</b>			
	<b>26</b>	Study limitations, including sources of potential bias, statistical uncertainty, and generalisability	17

1		<b>27</b>	Implications for practice, including the intended use and clinical role of the index test	14-17
2	<b>OTHER INFORMATION</b>			
3				
4		<b>28</b>	Registration number and name of registry	ISRCTN registry; referene number: ISRCTN24868616
5				
6		<b>29</b>	Where the full study protocol can be accessed	At the registry (above), but not detailed.
7				
8		<b>30</b>	Sources of funding and other support; role of funders	18
9				

For peer review only

10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
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



# STARD 2015

---

## AIM

STARD stands for “Standards for Reporting Diagnostic accuracy studies”. This list of items was developed to contribute to the completeness and transparency of reporting of diagnostic accuracy studies. Authors can use the list to write informative study reports. Editors and peer-reviewers can use it to evaluate whether the information has been included in manuscripts submitted for publication.

---

## EXPLANATION

A **diagnostic accuracy study** evaluates the ability of one or more medical tests to correctly classify study participants as having a **target condition**. This can be a disease, a disease stage, response or benefit from therapy, or an event or condition in the future. A medical test can be an imaging procedure, a laboratory test, elements from history and physical examination, a combination of these, or any other method for collecting information about the current health status of a patient.

The test whose accuracy is evaluated is called **index test**. A study can evaluate the accuracy of one or more index tests. Evaluating the ability of a medical test to correctly classify patients is typically done by comparing the distribution of the index test results with those of the **reference standard**. The reference standard is the best available method for establishing the presence or absence of the target condition. An accuracy study can rely on one or more reference standards.

If test results are categorized as either positive or negative, the cross tabulation of the index test results against those of the reference standard can be used to estimate the **sensitivity** of the index test (the proportion of participants *with* the target condition who have a positive index test), and its **specificity** (the proportion *without* the target condition who have a negative index test). From this cross tabulation (sometimes referred to as the contingency or “2x2” table), several other accuracy statistics can be estimated, such as the positive and negative **predictive values** of the test. Confidence intervals around estimates of accuracy can then be calculated to quantify the statistical **precision** of the measurements.

If the index test results can take more than two values, categorization of test results as positive or negative requires a **test positivity cut-off**. When multiple such cut-offs can be defined, authors can report a receiver operating characteristic (ROC) curve which graphically represents the combination of sensitivity and specificity for each possible test positivity cut-off. The **area under the ROC curve** informs in a single numerical value about the overall diagnostic accuracy of the index test.

The **intended use** of a medical test can be diagnosis, screening, staging, monitoring, surveillance, prediction or prognosis. The **clinical role** of a test explains its position relative to existing tests in the clinical pathway. A replacement test, for example, replaces an existing test. A triage test is used before an existing test; an add-on test is used after an existing test.

Besides diagnostic accuracy, several other outcomes and statistics may be relevant in the evaluation of medical tests. Medical tests can also be used to classify patients for purposes other than diagnosis, such as staging or prognosis. The STARD list was not explicitly developed for these other outcomes, statistics, and study types, although most STARD items would still apply.

---

## DEVELOPMENT

This STARD list was released in 2015. The 30 items were identified by an international expert group of methodologists, researchers, and editors. The guiding principle in the development of STARD was to select items that, when reported, would help readers to judge the potential for bias in the study, to appraise the applicability of the study findings and the validity of conclusions and recommendations. The list represents an update of the first version, which was published in 2003.

More information can be found on <http://www.equator-network.org/reporting-guidelines/stard>.



# BMJ Open

## Test-retest reliability of the twenty-five-hole peg test in stroke patients

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2019-032560.R2
Article Type:	Original research
Date Submitted by the Author:	13-Nov-2019
Complete List of Authors:	Granström, Fredrik; Uppsala University, Centre for Clinical Research Sörmland Hedlund, Mattias; Umeå Universitet Medicinska fakulteten, Department of Community Medicine and rehabilitation, Physiotherapy Lindström, Britta; Umeå Universitet Medicinska fakulteten, Department of Community Medicine and Rehabilitation, Physiotherapy Eriksson, Staffan; Uppsala University, Centre for Clinical Research Sörmland; Uppsala University, Department of Neuroscience, Physiotherapy
<b>Primary Subject Heading</b>:	Rehabilitation medicine
Secondary Subject Heading:	Neurology
Keywords:	Stroke < NEUROLOGY, measurement, Reliability, Clinical assessment, hand function, twenty-five-hole peg test

SCHOLARONE™  
Manuscripts

## Test-retest reliability of the twenty-five-hole peg test in stroke patients

Fredrik Granström<sup>1</sup>, Mattias Hedlund<sup>2</sup>, Britta Lindström<sup>2</sup>, Staffan Eriksson<sup>1,3</sup>

<sup>1</sup>Centre for Clinical Research Sörmland, Uppsala University, Sweden

<sup>2</sup>Department of Community Medicine and Rehabilitation, Physiotherapy, Umeå University, Sweden

<sup>3</sup>Department of Neuroscience, Physiotherapy, Uppsala University, Sweden

### CORRESPONDING AUTHOR

Staffan Eriksson, Centre for Clinical Research Sörmland, Uppsala University, Kungsgatan 41, 631 88 Eskilstuna, Sweden. Telephone number: +46 16 105251, +46 73 949 99 33. Email: [staffan.eriksson@germed.umu.se](mailto:staffan.eriksson@germed.umu.se)

**KEY WORDS:** clinical assessment, measurement, reproducibility of results, reliability, hand function, stroke, twenty-five-hole peg test

**WORD COUNT** 4060 words, 4176 words including “strengths and limitations of this study”.

## ABSTRACT

**Objectives:** Weaknesses of the nine-hole peg test include high floor effects and a result that might be difficult to interpret. In the twenty-five-hole peg test (TFHPT), the larger number of available pegs allows for the straightforward counting of the number of pegs inserted as the result. The TFHPT provides a comprehensible result and low floor effects. The objective was to assess the test-retest reliability of the TFHPT when testing persons with stroke. A particular focus was placed on the absolute reliability, as quantified by the smallest real difference (SRD). Complementary aims were to investigate possible implications for how the TFHPT should be used and for how the SRD of the TFHPT performance should be expressed.

**Design:** This study employed a test-retest design including three trials. The pause between trials was approximately 10-120 seconds.

**Participants, setting and outcome measure:** Thirty-one participants who had suffered a stroke were recruited from a group designated for constraint-induced movement therapy at outpatient clinics. The TFHPT result was expressed as the number of pegs inserted.

**Methods:** Absolute reliability was quantified by the SRD, including random and systematic error for a single trial,  $SRD_{2,1}$ , and for an average of three trials,  $SRD_{2,3}$ . For the SRD measures, the corresponding smallest real difference percentage (SRD%) measure was also reported.

**Results:** The differences in the number of pegs necessary to detect a change in the TFHPT for  $SRD_{2,1}$  and  $SRD_{2,3}$  were 4.0 and 2.3, respectively. The corresponding SRD% values for  $SRD_{2,1}$  and  $SRD_{2,3}$  were 36.5% and 21.3%, respectively.

**Conclusions:** The smallest change that can be detected in the TFHPT should be just above two pegs for a test procedure including an average of three trials. The use of an average of three trials compared to a single trial substantially reduces the measurement error.



1  
2  
3 **Trial registration:** ISRCTN registry, reference number ISRCTN24868616.  
4  
5  
6  
7  
8

9 **ARTICLE SUMMARY**  
10

11  
12 **Strengths and limitations of this study**  
13

- 14
- 15 • The generalizability of the results may be limited: the participants were selected  
16 because they should benefit from CIMT, and few scored above 20 pegs during the 50-  
17 second trial duration.  
18
  - 19 • Among other measures of reliability, the SRD percentage was reported; this is a good  
20 measure for comparisons between different tests, scales and populations.  
21
  - 22 • The results are presented with several different reliability measures to offer knowledge  
23 about the source of the measurement error.  
24
  - 25 • As the test-retest trials were performed within minutes, possible day-to-day variation  
26 was not captured.  
27
  - 28 • The intended practice trial was included as one of three trials in the analyses, which  
29 appears to have contributed to the learning effect.  
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

## INTRODUCTION

For a comprehensive upper limb assessment among persons with stroke, it is important to combine a measure of proximal upper limb function with a measure of hand function.[1]

However, only approximately 25% of studies regarding upper limb interventions include a specific measure of hand function.[1] The nine-hole peg test (NHPT) is a common test of hand function focusing on fine manual dexterity, and it is the most common such test used in research.[1-3] Two studies of stroke populations investigated the reliability of the NHPT, and there was a large discrepancy in the reliability reported: the smallest real difference percentage (SRD%) 24% vs. 52%.[2, 3] The NHPT has mostly shown moderate to excellent correlations (0.55-0.97) with other tests and self-reports focusing on hand function, including the Action Research Arm Test, the Jebsen-Taylor Hand Function Test, and the Stroke Impact Scale (hand function domain).[4, 5] The exception is the Motor Activity Log, for which low correlations have been reported (0.23-0.33).[5]

A weakness of the NHPT is that many persons with stroke cannot reach the lower limit; i.e., a floor effect arises. Furthermore, if the number of completed pegs is used as an outcome measure, a test with only nine pegs can measure only a narrow range of hand function, resulting in profound ceiling effects.[6] Therefore, to widen the scale and avoid ceiling effects, the original NHPT expresses the result as the time needed to complete the test (including inserting and removing all the pegs).[2, 3, 7, 8] However, this approach aggravates the floor effects because tests that are not completed during the stipulated time (limits of 60 and 180 seconds have been used) are excluded.[2, 3] The maximum time could be prolonged; however, this would be time consuming, mentally strenuous and therefore possibly unethical due to the possibility of a non-completed test after a lengthy attempt. A modified NHPT is used to mitigate the floor effect while avoiding the ceiling effect; in this modified version, the result is expressed as the number of inserted pegs per unit of time, i.e., the frequency.[9] This

1  
2  
3 modified test includes only peg insertion and not peg removal. It is thus possible also to  
4  
5 include tests that were not completed within the stipulated time limit and still measure  
6  
7 performance on the same task across the entire range of hand function. However, it may be  
8  
9 difficult both to interpret the frequency and to communicate it to other staff members and  
10  
11 patients, especially to those suffering from a brain injury. The reliability of this modified test  
12  
13 has not been investigated.  
14  
15

16  
17 In Sweden, a similar peg test, a twenty-five-hole peg test (TFHPT), has been used in clinical  
18  
19 practice. The larger number of available pegs makes it straightforward to count the number of  
20  
21 pegs inserted during a stipulated time frame of 50 seconds as the test result. Thus, the TFHPT  
22  
23 measures fine manual dexterity on a numerical scale that is easy to comprehend, with low  
24  
25 floor effects and presumably reasonable ceiling effects (based on pre-study data). Moreover,  
26  
27 compared to the individuals whom the original NHPT can test, individuals with worse hand  
28  
29 function can be tested with the TFHPT.[2, 3] Of the two studies investigating the reliability of  
30  
31 the NHPT, the one with the most generous time limit excluded all tests that were not  
32  
33 completed in 180 seconds.[3] This limit corresponds to inserting and removing a minimum of  
34  
35 2.5 pegs in 50 seconds, whereas 0 pegs inserted in 50 seconds is a valid result with the  
36  
37 TFHPT. The TFHPT has not been previously described in the literature, and its reliability has  
38  
39 not been investigated. Due to the similarity of the NHPT and the TFHPT, the underlying skill  
40  
41 assessed with these tests is most likely the same. However, since the tests have completely  
42  
43 different stop criteria – a time limit for the TFHPT vs. the insertion of all pegs for the NHPT –  
44  
45 equal reliability cannot be taken for granted.[7] Thus, if the size of the measurement error  
46  
47 related to the TFHPT is shown to be acceptable, this test may be useful in both clinical  
48  
49 practice and research.  
50  
51

52  
53 The overall aim of this study was to assess the test-retest reliability of the TFHPT for persons  
54  
55 suffering from stroke. A particular focus was placed on the absolute reliability, as quantified  
56  
57  
58  
59  
60

1  
2  
3 by the smallest real difference (SRD). Complementary aims were to investigate possible  
4  
5 implications for how the TFHPT should be used and for how the SRD of the TFHPT  
6  
7 performance should be expressed.  
8  
9

## 10 **METHOD**

### 11 **Participants**

12  
13  
14  
15  
16 The participants in this study were consecutively recruited in the process of screening patients  
17  
18 eligible for inclusion in a multicentre randomized controlled trial (RCT), reference number  
19  
20 ISRCTN24868616 at the ISRCTN registry. The patients were considered for inclusion  
21  
22 because they were to undergo constraint-induced movement therapy (CIMT) at one of the  
23  
24 clinics participating in the RCT. The clinics were outpatient rehabilitation clinics in the public  
25  
26 health care system in Sweden. Data were collected at the clinics. The sample in this study  
27  
28 consisted of included and excluded participants in the multicentre RCT. The participants were  
29  
30 included if they had one stroke or more registered in the medical record and if TFHPT data  
31  
32 were available from three trials before and three trials after the CIMT. Moreover, with regard  
33  
34 to the outcome measure, a minimum of one peg and a maximum of 24 pegs inserted was  
35  
36 necessary for inclusion. This was to avoid an untrue low measurement error from participants  
37  
38 stable at 0 or 25 pegs inserted. These two intervals are wider, a person can be far below the  
39  
40 floor or high over the ceiling, so measurements at these intervals should be more stable.  
41  
42  
43  
44  
45

46  
47 A minimum of 30 participants were included to obtain a sufficient number for a reliability  
48  
49 study.[10]  
50

### 51 **Procedure and measurements**

52  
53  
54  
55 The TFHPT has twenty-five holes and pegs (Figure 1). The test used in this study consisted of  
56  
57 a rectangular 21 cm × 45 cm board with a box containing pegs on one side and an elevated 18  
58  
59 cm × 18 cm area with holes on the other side. The holes were 9 mm wide and 18 mm deep,  
60

1  
2  
3 and they were spaced 20 mm apart. The box had a base of 13 cm × 18 cm and was 5 cm deep.  
4  
5 The pegs were 40 mm long and 8 mm in diameter.  
6  
7

8 A battery of different tests was administered in this study, including the Fugl-Meyer test [11]  
9  
10 and the Birgitta Lindmark motor assessment (BL motor assessment).[12, 13] The TFHPT  
11  
12 was administered as the second test. The preceding test, the BL motor assessment, required  
13  
14 approximately 30-60 minutes to administer. The tests were administered in an examination  
15  
16 room in which only the participant and the physiotherapist were present. For the TFHPT,  
17  
18 three trials were performed with each hand. The participants started with the less affected  
19  
20 hand, followed by the more affected hand, i.e., the hand of investigation in this test-retest  
21  
22 study. The pause between trials was approximately 10-120 seconds. The board was placed at  
23  
24 a distance favoured by the participant with the centre row of holes centred towards the navel  
25  
26 and the box side oriented towards the tested hand. The starting position was with both hands  
27  
28 on the board, and time keeping began upon first hand contact with the pegs. We gave  
29  
30 participants the following instructions:  
31  
32  
33  
34  
35

- 36 1. I want you to pick up one peg at a time and insert them in the holes of the board.
- 37 2. Use only the right/left hand; you can only use the other hand to steady the board.
- 38 3. You can fill the holes in any order you desire.
- 39 4. We start with a practice trial.
- 40 5. You have 50 seconds to insert as many pegs as you can. After 50 seconds, the trial is  
41 terminated.
- 42 6. Are you ready? Ready, set, go!
- 43 7. After the practice trial: This was practice; now come the two actual test trials, where  
44 the results are recorded. Repeat step 6.

45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58 The test-retest reliability of the TFHPT was assessed on two separate occasions, i.e., before  
59  
60 and after a two-week training period (the CIMT). The same procedure was used on both of

1  
2  
3 these occasions, and for each participant, all tests were administered by the same  
4  
5 physiotherapist. The assessment after the CIMT period was performed as an internal  
6  
7 validation.  
8  
9

10 Two physiotherapists, SE and BL, administered the tests in this study. SE has general  
11  
12 experience with persons suffering from stroke and experience administering the original  
13  
14 NHPT. BL has extensive experience with persons suffering from stroke, including  
15  
16 administering the original NHPT. Background data from medical records were collected by  
17  
18 staff at the clinics.  
19  
20  
21  
22

### 23 **Statistics**

24  
25 All three trials were used in the analyses, although the first trial was introduced to the  
26  
27 participants as a practice trial. Analyses of pre-intervention data and post-intervention data  
28  
29 were performed separately.  
30  
31  
32

33 Bland-Altman plots of trials one and two provided a graphic description of the data  
34  
35 variability. The mean of trials one and two was plotted against the difference between trials  
36  
37 two and one for each subject. Heteroscedasticity – i.e., an association between the random  
38  
39 error and the magnitude of measurements [14] – was investigated with pairwise comparisons  
40  
41 of trials using Koenker's [15] studentized test, which is useful for small samples and skewed  
42  
43 data. Heteroscedasticity is indicated by a significant result.  
44  
45  
46  
47

48 Measurement error can be either random or systematic. In random error, there is no pattern of  
49  
50 variability between trials, whereas in systematic error, the measurements vary in a non-  
51  
52 random way; i.e., the mean values between the trials differ.[16] To investigate whether there  
53  
54 was a systematic error in test scores, one-way repeated measures ANOVA was used to detect  
55  
56 potential between trial effects.  
57  
58  
59  
60

1  
2  
3 Reliability is a term that describes how the measurement result of an instrument is affected by  
4 measurement error.[6, 14] Reliability can be quantified as either relative or absolute.[6]

7 *Relative reliability* refers to the consistency of the positions of measurements relative to those  
8 of others within the tested group, and it is quantified using several intra-class correlation  
9 coefficients (ICCs).[16, 17] In ICCs, between-subject variability is related to the within-  
10 subject variability by a ratio.[16] Thus, ICCs are sensitive to the degree of between-subject  
11 variability, and with all other things being equal, a more heterogeneous sample (i.e., a larger  
12 between-subject variability) produces higher ICC values.[16, 17] The concept of *absolute*  
13 *reliability* refers to the consistency of measurements within individuals.[6, 16] Measurement  
14 error, quantified as within-subject standard deviations in repeated tests, is a common measure  
15 of absolute reliability [6, 14, 16-18] and is called the standard error of measurement (SEM).  
16 SRD is an extension of the SEM, and it can be seen as the smallest detectable difference, with  
17 95% certainty, using a test instrument on an individual.[16]

21 Three separate measures of *relative reliability*, i.e.,  $ICC_{2,1}$ ,  $ICC_{2,3}$ , and  $ICC_{3,3}$ , including 95%  
22 confidence intervals (CIs), were calculated. This panel of measures was used to compare the  
23 results representative of single and average measures and to obtain an estimate of the  
24 influence of systematic error. *The first* figure in the ICC designation represents the type of  
25 ICC model.[16]  $ICC_{2,1}$  and  $ICC_{2,3}$  are calculated from a two-way random effect model and  
26 incorporate both systematic and random error, whereas  $ICC_{3,3}$  is calculated from a two-way  
27 fixed effect model and incorporates only random error.[16, 19] Thus, the less systematic error  
28 contributes to the total error, the closer  $ICC_{2,3}$  is to  $ICC_{3,3}$ . [16] *The second* figure in the ICC  
29 designation represents single or average measures, where “1” represents single measures and  
30 “2” or higher represents the number of trials from which the average is calculated.[16]  $ICC_{2,1}$   
31 represents the reliability of a test procedure in which the subject is tested with a single trial on  
32 a test occasion.[16]  $ICC_{2,3}$  and  $ICC_{3,3}$  represent the reliability of a test procedure in which the

1  
2  
3 subject is tested with three trials on a test occasion and the score is expressed as the average  
4  
5 of these trials.  
6  
7

8 To estimate *absolute reliability*, the SEM, SRD and SRD percentage (SRD%) were calculated  
9  
10 for each of the three different ICC measures (ICC<sub>2,1</sub>, ICC<sub>2,3</sub>, and ICC<sub>3,3</sub>), resulting in the  
11  
12 corresponding properties SEM<sub>2,1</sub>, SRD<sub>2,1</sub>, SRD%<sub>2,1</sub>, SEM<sub>2,3</sub>, SRD<sub>2,3</sub>, SRD%<sub>2,3</sub>, SEM<sub>3,3</sub>,  
13  
14 SRD<sub>3,3</sub>, and SRD%<sub>3,3</sub>. The SEM was calculated according to  $SEM = SD\sqrt{(1 - ICC)}$ , where  
15  
16 SD was calculated from the total sum of squares (SS<sub>TOTAL</sub>) in the ANOVA table generated in  
17  
18 the ICC analyses as  $\sqrt{SS_{total}/(n - 1)}$ . [16] The SRD was calculated using the formula  $1.96 \times$   
19  
20  $SEM \times \sqrt{2}$ , where 1.96 is related to the 95% CI and  $\sqrt{2}$  refers to the error of two  
21  
22 measurements. [16] The SRD% was calculated by dividing the SRD value by the grand mean  
23  
24 multiplied by 100. [2, 10] This value is independent of measurement units and is indexed to  
25  
26 the mean value of the observations from which it was derived. It is therefore a good measure  
27  
28 for comparisons between different tests, scales and populations. [10, 14, 17] An SRD% of  
29  
30 30% has been suggested as an acceptable level of reliability. [20]  
31  
32  
33  
34  
35

36 Because estimates of absolute reliability vary with the type of ICC value, some caution is  
37  
38 warranted when comparing them with measures from other studies. [16] Therefore, SEM<sub>mean</sub>  
39  
40 square error term (MSE) =  $\sqrt{MSE}$  was also calculated, where MSE (this term is called residual error  
41  
42 by Hopkins and the mean square residual in the SPSS output) was taken from the ANOVA  
43  
44 table of the ICC calculation. [16, 17] This SEM measure represents the reliability of a test  
45  
46 procedure in which the subject is tested with a single trial on a test occasion, and it is a pure  
47  
48 measure of random error. [16] SRD<sub>MSE</sub> and SRD%<sub>MSE</sub> were also derived from SEM<sub>MSE</sub>.  
49  
50  
51  
52  
53

54 The analysis of test-retest reliability was pre-planned. SPSS version 21 was used to calculate  
55  
56 ICC and ANOVA. The alpha level was set to 0.05.  
57  
58

## 59 Patient and public involvement

60



No patients or members of the public were involved in the development or design of this study.

## RESULTS

In this study, participants were recruited between January 2011 and September 2014. Of 60 eligible patients, 29 were excluded for any of the following reasons: not suffering a stroke, missing data, and yielding either below the minimum or above the maximum number of inserted pegs (Figure 2). This yielded 31 participants (21 men and 10 women) for inclusion in the analysis, with a mean  $\pm$  SD age of  $66 \pm 9$  years (Table 1). The two eligible patients who were excluded because they exceeded the permitted maximum number of pegs inserted completed the 25 pegs in their best trial within 49.3 seconds and 39.4 seconds. Of the ten patients who were excluded because they fell below the minimum number of inserted pegs, six inserted at least one peg in one of the trials. Data were collected from 17 and 14 participants by the two physiotherapists (third and last author, respectively) at seven clinics.

**Table 1.** Characteristics of participants at pre-intervention trials

Participants	N=31
Age (years), mean $\pm$ SD <sup>a</sup>	66 $\pm$ 9
Men/women, n <sup>b</sup>	21/10
Time since stroke (months), median (IQR <sup>d</sup> ), (min-max)	17 (8–24), (2–70)
Previous dominant hand more affected by stroke, n	19
TFHPT <sup>c</sup> , mean of three trials (number of pegs), mean $\pm$ SD, (min-max)	10.8 $\pm$ 6.8, (1–22.7)
Fugl-Meyer test (score), median (IQR <sup>d</sup> ), (min-max)	46 (41–53), (29–62)

More than one stroke, n 3

<sup>a</sup>Standard deviation, <sup>b</sup>number of participants, <sup>c</sup>twenty-five-hole  
peg test, <sup>d</sup>interquartile range

A graphic description of the data variability can be seen in the Bland-Altman plots (Figures 3 and 4). According to Koenker's studentized test, the measurement error was not affected by heteroscedasticity (Table 2).

**Table 2.** Results of the Koenker's studentized test, n=31

Pairwise test of trials	Pre-intervention trials		Post-intervention trials	
	Chi-square	P-value	Chi-square	P-value
1-2	1.33	0.25	0.41	0.52
2-3	0.05	0.83	1.38	0.24
1-3	0.28	0.60	0.20	0.66

For *pre-intervention* trials, the mean values  $\pm$  SDs for trials 1, 2 and 3 were  $10.0 \pm 6.5$ ,  $11.0 \pm 7.1$ , and  $11.5 \pm 6.9$ , respectively. The one-way repeated measures ANOVA revealed a main effect between trials with  $F(2, 29) = 10.9$  and  $p < 0.001$ . Post hoc tests revealed differences between trials 2 and 1 and between trials 3 and 1, with mean differences (95% CIs) of 1.0 (0.3–1.6) and 1.5 (0.9–2.2), respectively.

For *post-intervention* trials, the mean values  $\pm$  SD for trials 1, 2 and 3 were  $11.8 \pm 6.5$ ,  $12.4 \pm 6.7$ , and  $12.5 \pm 6.8$ , respectively. The one-way repeated measures ANOVA revealed a main effect between trials with  $F(2, 29) = 4.1$  and  $p = 0.027$ . Post hoc tests revealed a difference between trials 3 and 1, with a mean difference (95% CIs) of 0.6 (0.2–1.1).

For *pre-intervention* trials, ICC<sub>2,3</sub> (95% CI) was 0.99 (0.97–0.99) (Table 3). The SRDs incorporating random and systematic error, SRD<sub>2,1</sub> and SRD<sub>2,3</sub>, were 4.0 and 2.3 pegs, respectively. The corresponding SRD% values for SRD<sub>2,1</sub> and SRD<sub>2,3</sub> were 36.5% and 21.3%, respectively. The SRD incorporating only random error, SRD<sub>3,3</sub>, was 2.0 pegs.

For *post-intervention* trials, SRD<sub>2,1</sub> and SRD<sub>2,3</sub> were 3.2 and 1.8 pegs, respectively (Table 4). SRD<sub>3,3</sub>, was 1.8 pegs.

**Table 3.** Results of reliability measures for pre-intervention trials

	ICC <sup>a</sup> (95% CI)	SEM <sup>b</sup> , n <sup>c</sup>	SRD <sup>d</sup> , n <sup>c</sup>	SRD% <sup>e</sup>
ICC <sub>2,1</sub>	0.96 (0.90–0.98)	1.4	4.0	36.5
ICC <sub>2,3</sub>	0.99 (0.97–0.99)	0.8	2.3	21.3
ICC <sub>3,3</sub>	0.99 (0.98–0.99)	0.7	2.0	18.3
Derived from MSE <sup>f</sup>		1.3	3.5	32.1

<sup>a</sup>Intra-class correlation coefficient.

<sup>b</sup>Standard error of measurement derived from ICC<sub>2,1</sub>, 2,3, 3,3 and MSE.

<sup>c</sup>Number of pegs.

<sup>d</sup>Smallest real difference derived from ICC<sub>2,1</sub>, 2,3, 3,3 and MSE.

<sup>e</sup>SRD percentage derived from ICC<sub>2,1</sub>, 2,3, 3,3 and MSE.

<sup>f</sup>Mean square error term.

**Table 4.** Results of reliability measures for post-intervention trials

	ICC <sup>a</sup> (95% CI)	SEM <sup>b</sup> , n <sup>c</sup>	SRD <sup>d</sup> , n <sup>c</sup>	SRD% <sup>e</sup>
ICC <sub>2,1</sub>	0.97 (0.95–0.98)	1.1	3.2	25.9
ICC <sub>2,3</sub>	0.99 (0.98–1.0)	0.7	1.8	15.0
ICC <sub>3,3</sub>	0.99 (0.98–1.0)	0.7	1.8	15.0
Derived from MSE <sup>f</sup>		1.1	3.1	25.5

<sup>a</sup>Intra-class correlation coefficient.

<sup>b</sup>Standard error of measurement derived from ICC2.1, 2.3, 3.3 and MSE.

<sup>c</sup>Number of pegs.

<sup>d</sup>Smallest real difference derived from ICC2.1, 2.3, 3.3 and MSE.

<sup>e</sup>SRD percentage derived from ICC2.1, 2.3, 3.3 and MSE.

<sup>f</sup>Mean square error term.

## DISCUSSION

This study indicated that in a selected group of persons suffering from stroke, the absolute test-retest reliability of the TFHPT was at a level that can be considered acceptable for measures representing an average of three trials and incorporating systematic error.

To assess implications for the use of the TFHPT and to determine which SRD measure best captures the absolute reliability, three issues were considered: 1) whether to use single or average measures, 2) whether to include systematic error in the assessments, and 3) whether to take heteroscedasticity into account.

Comparing  $SRD_{2,1}$  to  $SRD_{2,3}$  revealed that the use of an average of three trials reduced the measurement error by approximately 1.5 pegs compared to the use of a single trial. This finding suggests that the reliability of the TFHPT is substantially improved when an average of three trials is used.

Comparing  $SRD_{2,3}$  to  $SRD_{3,3}$ , where  $SRD_{3,3}$  incorporates only random error, revealed that the contribution of the systematic error was approximately 0.3 pegs of the total 2.3 pegs when the average of three trials was used.[16] Although the systematic error was small compared to the random error it was not small enough to be overlooked in the assessment of reliability. Therefore,  $SRD_{2,3}$  is preferable to  $SRD_{3,3}$  for measuring the reliability of the TFHPT.

1  
2  
3 The choice of SRD% instead of SRD is dependent on whether the measurement error is affected  
4 by heteroscedasticity. A measure of absolute reliability, expressed as an absolute number of  
5 pegs, can over- or underestimate the number of pegs necessary to demonstrate an improvement  
6 for an individual.[14, 17] The reason is that the random error of measurements often increases  
7 with the magnitude of the measurements (i.e., heteroscedasticity).[14, 17] As a remedy, the use  
8 of a relative measure of absolute reliability, such as SRD%, has been proposed.[14, 17]  
9  
10 However, the lack of heteroscedasticity detection suggests that both  $SRD_{2,3}$  and  $SRD\%_{2,3}$  are  
11 appropriate measures of reliability for the TFHPT.[14]

12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22 The results for  $SRD_{2,3}$  and  $SRD\%_{2,3}$  were 2.3 pegs and 21.3%, respectively. The value of  
23  $SRD\%_{2,3}$  fell within the 30% level, which has been suggested as acceptable.[20] The 30% level  
24 seems high in this context, with persons affected by stroke in a chronic stage; from a clinical  
25 viewpoint, our opinion is that the results of 21.3% and 2.3 pegs in this study indicate a barely  
26 acceptable level of absolute reliability. For a favourable level, we believe that a mean number  
27 consisting of approximately 1.5 pegs is desirable.

28  
29  
30  
31  
32  
33  
34  
35  
36  
37 The relative test-retest reliability, as measured by  $ICC_{2,3}$ , was 0.99, which seems excellent. The  
38 discrepancy between the level of the relative and the absolute reliability is most likely caused  
39 by the heterogeneity in this study population (Figure 3, Table 1) which inflates the relative  
40 reliability.[16]

41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
The level of the relative absolute test-retest reliability (SRD%), the most comparable measure,  
observed for the TFHPT in this study (21.3%) is better than what Chen et al.[2] reported (54%)  
and is at approximately the same level as Ekstrand et al.[3] reported (24%) for the NHPT.  
Although the SRD% measures reported in the studies by Chen et al. and by Ekstrand et al. were  
calculated in different ways than the  $SRD\%_{2,3}$  reported in this study, the measures used in these  
three studies are fairly equivalent.[16] Several methodological differences between these three  
studies could have affected the results.[2, 3] *First*, the results of the TFHPT and NHPT were

1  
2  
3 measured using different scales, and the use of time for completion of the test in the NHPT  
4 should accommodate more variability than the peg count used in the TFHPT. However, the  
5 SRD% results should still be comparable between the TFHPT and the NHPT because this  
6 relative measure of absolute reliability adjusts for different scales and study populations.[14,  
7 17] *Second*, in this study of the TFHPT, the test and retest trials were performed within minutes  
8 compared to within days in the studies of the NHPT. Thus, the TFHPT may seem more reliable  
9 because of possible random error from day-to-day variation in performance which was not  
10 captured in this study.[17, 21] *Third*, the longer time since stroke in this study of the TFHPT  
11 compared to the study of NHPT by Chen et al.[2] may have resulted in seemingly better  
12 reliability for the TFHPT because of a more stable level of hand function. In the study by Chen  
13 et al., a systematic error may have originated in recovery from stroke in the 3-5 days between  
14 the test and retest trials because the time since stroke was 3 months or less for a quarter of the  
15 study sample.[22]

16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34 The implications of the results of this study are that the TFHPT can be used in a clinical situation  
35 to detect changes in a patient's hand function. The test procedure should employ an average of  
36 three trials on each occasion, and a change of 2.3 pegs or more between two occasions should  
37 be considered real improvement/worsening. Furthermore, it seems that the ceiling effects in the  
38 TFHPT can be considered acceptable. Only two of the persons assessed for eligibility inserted  
39 25 pegs, and only one of them actually hit the ceiling because according to this individual's  
40 best times for completion of the 25 pegs, he/she would have been able to insert more pegs if  
41 available. This occurred in a sample where approximately a quarter of the included participants  
42 suffered from mild impairment of arm and hand function as judged by the Fugl-Meyer test.[23,  
43 24]

44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
There was a tendency towards improved reliability after the CIMT period, which was due to  
decreased systematic error and decreased random error. The decreased systematic error can be

1  
2  
3 observed in the main effects of trial in the ANOVA results.[16] The decreased systematic error  
4  
5 is most likely due to a decreased learning effect when the participants had previous experience  
6  
7 in the test. The learning effect is indicated by the increases in the mean values over the trials,  
8  
9 especially over trials 1-2, and the decreased learning effect is indicated by the less pronounced  
10  
11 increase in the post-intervention trials.[14, 17] The lower random error can be observed from  
12  
13 the lower SRD<sub>3,3</sub> results in the post-intervention trials.[16] The cause of the decreased random  
14  
15 error is less clear, but it could also be attributed to the decreased systematic error.[17] This is  
16  
17 because the magnitude of the learning effect probably differs between individuals, which will  
18  
19 show as random error. Furthermore, it is likely that the SRD<sub>2,3</sub> result of 2.3 pegs for TFHPT  
20  
21 could, in reality, be adjusted downwards. A peg test is often used to evaluate a rehabilitation  
22  
23 period; because the error is smaller in the post-intervention trials, the “true” SRD may be  
24  
25 somewhere between the SRDs of the pre- and post-intervention trials (2.3 vs 1.8).  
26  
27  
28  
29  
30

31 Four weaknesses of this study should be considered. The sample included a relatively low  
32  
33 number of participants with few observations above 20 pegs and participants who were selected  
34  
35 because they should benefit from CIMT.[10, 17] These sample qualities may thus, to some  
36  
37 degree, hinder the generalization of the results to other groups of people suffering from stroke.  
38  
39 In addition, the intended practice trial was included as one of three trials in the analyses which  
40  
41 appears to have contributed to systematic error through an increased learning effect, indicated  
42  
43 by a large increase in the mean values between trials 1 and 2.[14, 16, 17] Thus, to mitigate the  
44  
45 learning effect, a practice trial preceding regular trials is recommended. Moreover, the possible  
46  
47 day-to-day variation was not captured in the present study design. The advantage of this  
48  
49 approach is that it yields a pure result for measurement error for the instrument in this  
50  
51 population; the disadvantage is that the result is less clinically applicable.[17, 21] Finally, in  
52  
53 this study, sensitivity to change and validity were not examined. However, the criterion validity  
54  
55 for NHPT has mostly shown a moderate to excellent level [4, 5] and the underlying skill  
56  
57  
58  
59  
60

1  
2  
3 assessed with the TFHPT is most likely the same. A high reliability level is a prerequisite for  
4 high validity, and because the reliability of the TFHPT was at the same level as that of the  
5 NHPT, the criterion validity should also be similar.[21]  
6  
7  
8

9  
10 In conclusion, our results suggest that the smallest detectable difference between two  
11 assessments using a test procedure with an average of three trials conducted by a single tester  
12 should be just above two pegs with the TFHPT. Furthermore, to reach an acceptable level of  
13 measurement error, the use of the average of multiple trials is crucial. Future research should  
14 focus on optimizing the number of trials.  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

## 27 **ACKNOWLEDGEMENTS**

28  
29 The authors thank the staff at the clinics for the extra work associated with participating in the  
30 RCT. The authors thank Håkan Littbrand for input on the conception of the study and Monika  
31 Edström for administrative work with the study.  
32  
33  
34  
35  
36  
37  
38

## 39 **AUTHOR CONTRIBUTIONS**

40  
41 Study design and data interpretation: SE, FG, BL and MH. Data acquisition: SE, BL and MH.  
42  
43 Statistical analysis: SE and FG. Drafting and finalization of the manuscript: SE. Critical  
44 revision of the manuscript: FG, BL and MH. Final approval of the submitted manuscript: FG,  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 53 **COMPETING INTERESTS**

54  
55 None declared.

## 57 **FUNDING**



1  
2  
3 The Norrbacka-Eugenia Foundation; The Swedish STROKE-Association; The Stroke  
4  
5 Foundation of Northern Sweden; and the Centre for Clinical Research Sörmland, Uppsala  
6  
7 University.  
8  
9

## 10 11 12 **DATA SHARING**

13  
14 No additional data are available.  
15

## 16 17 **PATIENT CONSENT**

18  
19 Written informed consent was obtained from the participants.  
20

## 21 22 **ETHICS APPROVAL**

23 The Regional Ethical Review Board in Umeå, reference 09-104M, with additional approval  
24 Dnr 2010/314-32M, Dnr 2011-244-32M, and 2012-235-32M.  
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

**REFERENCES**

1. Santisteban L, Teremetz M, Bleton JP, et al. Upper limb outcome measures used in stroke rehabilitation studies: a systematic literature review. *PLoS One* 2016;11:e0154792.
2. Chen HM, Chen CC, Hsueh IP, et al. Test-retest reproducibility and smallest real difference of 5 hand function tests in patients with stroke. *Neurorehabil Neural Repair* 2009;23:435–40.
3. Ekstrand E, Lexell J, Brogardh C. Test-retest reliability and convergent validity of three manual dexterity measures in persons with chronic stroke. *PM R* 2016;8:935–43.
4. Beebe JA, Lang CE. Relationships and responsiveness of six upper extremity function tests during the first six months of recovery after stroke. *J Neurol Phys Ther* 2009;33:96–103.
5. Lin KC, Chuang LL, Wu CY, et al. Responsiveness and validity of three dexterous function measures in stroke rehabilitation. *J Rehabil Res Dev* 2010;47:563–71.
6. Carter RE, Lubinsky J, Domholdt E. *Rehabilitation Research: Principles and Applications*. St. Louis, MO: Elsevier-Saunders, 2016:239–244.
7. Mathiowetz V, Weber K, Kashman N, et al. Adult norms for the Nine-Hole Peg Test of Finger Dexterity. *Occup Ther J Res* 1985;5:25–38.
8. Grice KO, Vogel KA, Le V, et al. Adult norms for a commercially available Nine Hole Peg Test for finger dexterity. *Am J Occup Ther* 2003;57:570–3.
9. Heller A, Wade DT, Wood VA, et al. Arm function after stroke: measurement and recovery over the first three months. *J Neurol Neurosurg Psychiatry* 1987;50:714–9.
10. Lexell JE, Downham DY. How to assess the reliability of measurements in rehabilitation. *Am J Phys Med Rehabil* 2005;84:719–23.
11. Fugl-Meyer AR, Jaasko L, Leyman I, et al. The post-stroke hemiplegic patient. 1. A method for evaluating of physical performance. *Scand J Rehabil Med* 1975;7:13-31.

12. Lindmark B, Hamrin E. Evaluation of functional capacity after stroke as a basis for active intervention. Validation of a modified chart for motor capacity assessment. *Scand J Rehabil Med* 1988;20:111–5.
13. Lindmark B, Hamrin E. Evaluation of functional capacity after stroke as a basis for active intervention. Presentation of a modified chart for motor capacity assessment and its reliability. *Scand J Rehabil Med* 1988;20:103–9.
14. Atkinson G, Nevill AM. Statistical methods for assessing measurement error (reliability) in variables relevant to sports medicine. *Sports Med* 1998;26:217–38.
15. Koenker R. A note on studentizing a test for heteroscedasticity. *J Econom* 1981;17:107–12.
16. Weir JP. Quantifying test-retest reliability using the intraclass correlation coefficient and the SEM. *J Strength Cond Res* 2005;19:231–40.
17. Hopkins WG. Measures of reliability in sports medicine and science. *Sports Med* 2000;30:1–15.
18. Bland JM, Altman DG. Measurement error. *BMJ* 1996;313:744.
19. Shrout PE, Fleiss JL. Intraclass correlations: uses in assessing rater reliability. *Psychol Bull* 1979;86:420–8.
20. Huang SL, Hsieh CL, Lin JH, et al. Optimal scoring methods of hand-strength tests in patients with stroke. *Int J Rehabil Res* 2011;34:178–80.
21. Sim J, Wright C. *Research in Health Care: Concepts, Designs and Methods*. Cheltenham, England: Nelson Thornes Ltd 2002: 123–33.
22. Kwakkel G, Kollen B, Twisk J. Impact of time on improvement of outcome after stroke. *Stroke* 2006;37:2348–53.
23. Duncan PW, Goldstein LB, Horner RD, et al. Similar motor recovery of upper and lower extremities after stroke. *Stroke* 1994;25:1181–8.

- 1  
2  
3 24. Woytowicz EJ, Rietschel JC, Goodman RN, et al. Determining levels of upper  
4 extremity movement impairment by applying a cluster analysis to the fugl-meyer  
5 assessment of the upper extremity in chronic stroke. *Arch Phys Med Rehabil*  
6  
7  
8 2017;98:456–62.  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
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 only

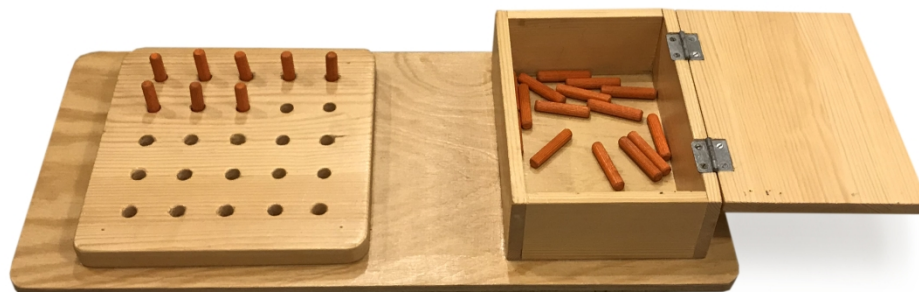
## FIGURE LEGENDS

**Figure 1.** The twenty-five-hole peg test.

**Figure 2.** Flowchart of the recruitment process in the study. \*Constraint-induced movement therapy.

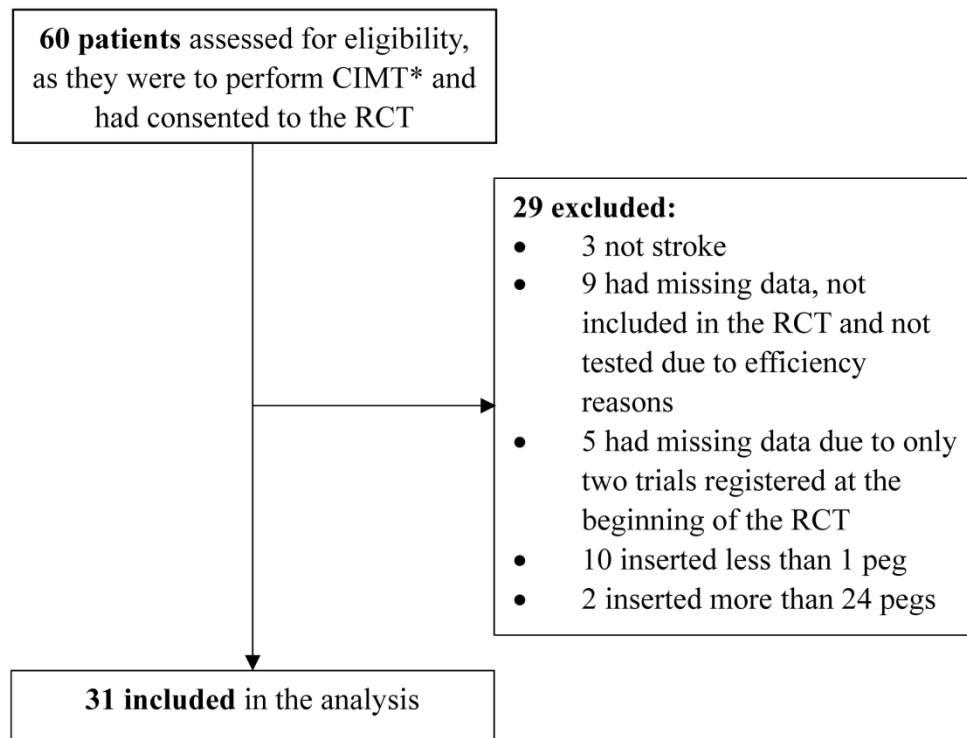
**Figure 3.** Bland-Altman plots of numbers of pegs from pre-intervention trials. The mean of trials 1 and 2 was plotted against the difference of trials 2 and 1 for each subject. The centre line displays the mean difference for the group between trials 2 and 1. The upper and lower confidence limits were calculated as the mean difference  $\pm$  SD of the mean difference  $\times$  1.96.

**Figure 4.** Bland-Altman plots of numbers of pegs from post-intervention trials. The mean of trials 1 and 2 was plotted against the difference of trials 2 and 1 for each subject. The centre line displays the mean difference for the group between trials 2 and 1. The upper and lower confidence limits were calculated as the mean difference  $\pm$  SD of the mean difference  $\times$  1.96.



The twenty-five-hole peg test.

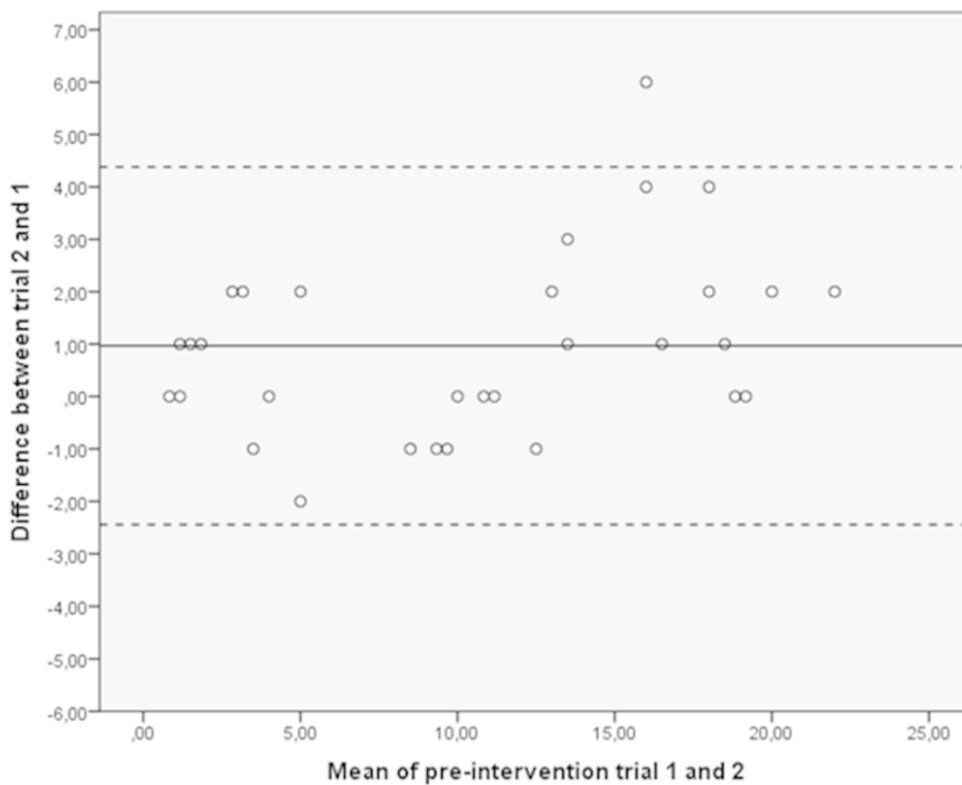
208x140mm (300 x 300 DPI)



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

Flowchart of the recruitment process in the study. \*Constraint-induced movement therapy.

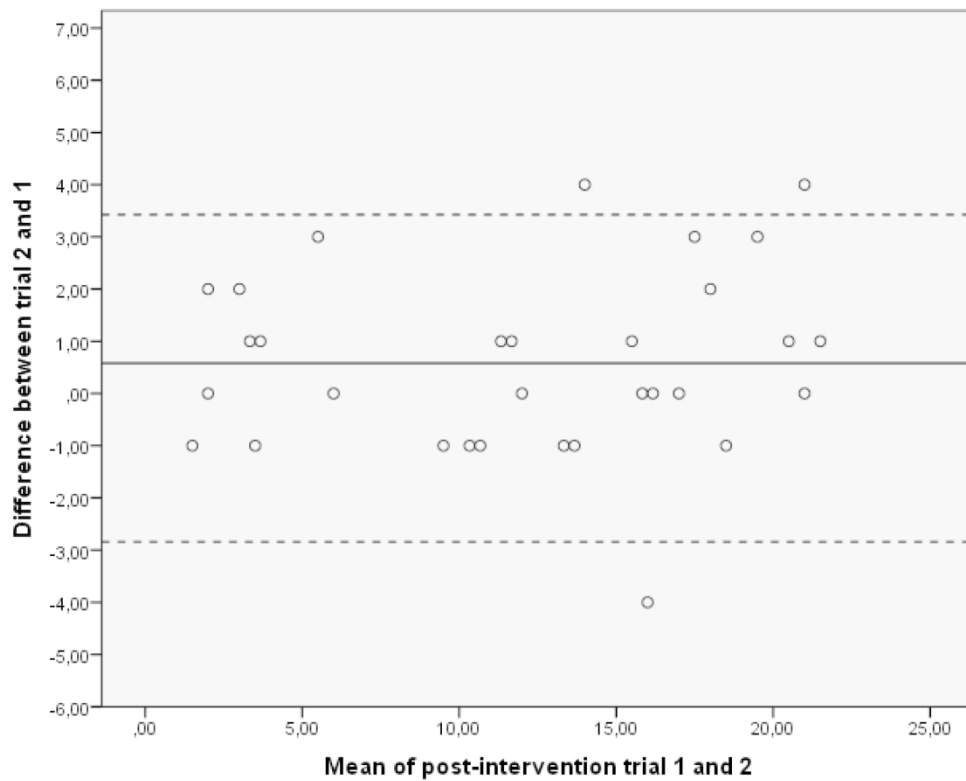
121x92mm (600 x 600 DPI)



Bland-Altman plots of numbers of pegs from pre-intervention trials. The mean of trials 1 and 2 was plotted against the difference of trials 2 and 1 for each subject. The centre line displays the mean difference for the group between trials 2 and 1. The upper and lower confidence limits were calculated as the mean difference  $\pm$  SD of the mean difference  $\times$  1.96.

140x112mm (600 x 600 DPI)





Bland-Altman plots of numbers of pegs from post-intervention trials. The mean of trials 1 and 2 was plotted against the difference of trials 2 and 1 for each subject. The centre line displays the mean difference for the group between trials 2 and 1. The upper and lower confidence limits were calculated as the mean difference  $\pm$  SD of the mean difference  $\times$  1.96.

140x112mm (600 x 600 DPI)

Section & Topic	No	Item	Reported on page #
<b>TITLE OR ABSTRACT</b>			
	<b>1</b>	Identification as a study of diagnostic accuracy using at least one measure of accuracy (such as sensitivity, specificity, predictive values, or AUC)	(Reliability), 1
<b>ABSTRACT</b>			
	<b>2</b>	Structured summary of study design, methods, results, and conclusions (for specific guidance, see STARD for Abstracts)	2
<b>INTRODUCTION</b>			
	<b>3</b>	Scientific and clinical background, including the intended use and clinical role of the index test	4-5
	<b>4</b>	Study objectives and hypotheses	5-6
<b>METHODS</b>			
<i>Study design</i>	<b>5</b>	Whether data collection was planned before the index test and reference standard were performed (prospective study) or after (retrospective study)	10, registered, refnr: ISRCTN24868616
<i>Participants</i>	<b>6</b>	Eligibility criteria	6
	<b>7</b>	On what basis potentially eligible participants were identified (such as symptoms, results from previous tests, inclusion in registry)	6
	<b>8</b>	Where and when potentially eligible participants were identified (setting, location and dates)	Setting: 6 Dates: 11 Exact locations of the clinics not included.
	<b>9</b>	Whether participants formed a consecutive, random or convenience series	6
<i>Test methods</i>	<b>10a</b>	Index test, in sufficient detail to allow replication	Not applicable
	<b>10b</b>	Reference standard, in sufficient detail to allow replication	n.a.
	<b>11</b>	Rationale for choosing the reference standard (if alternatives exist)	n.a.
	<b>12a</b>	Definition of and rationale for test positivity cut-offs or result categories of the index test, distinguishing pre-specified from exploratory	n.a.
	<b>12b</b>	Definition of and rationale for test positivity cut-offs or result categories of the reference standard, distinguishing pre-specified from exploratory	n.a.
	<b>13a</b>	Whether clinical information and reference standard results were available to the performers/readers of the index test	n.a.
	<b>13b</b>	Whether clinical information and index test results were available to the assessors of the reference standard	7
<i>Analysis</i>	<b>14</b>	Methods for estimating or comparing measures of diagnostic accuracy	Reliability measures 8-10
	<b>15</b>	How indeterminate index test or reference standard results were handled	n.a.
	<b>16</b>	How missing data on the index test and reference standard were handled	n.a.
	<b>17</b>	Any analyses of variability in diagnostic accuracy, distinguishing pre-specified from exploratory	n.a.
	<b>18</b>	Intended sample size and how it was determined	6
<b>RESULTS</b>			
<i>Participants</i>	<b>19</b>	Flow of participants, using a diagram	11
	<b>20</b>	Baseline demographic and clinical characteristics of participants	11
	<b>21a</b>	Distribution of severity of disease in those with the target condition	Table 1, figure 3 and 4.
	<b>21b</b>	Distribution of alternative diagnoses in those without the target condition	n.a.
	<b>22</b>	Time interval and any clinical interventions between index test and reference standard	7
<i>Test results</i>	<b>23</b>	Cross tabulation of the index test results (or their distribution) by the results of the reference standard	Plots instead, figure 3 and 4.
	<b>24</b>	Estimates of diagnostic accuracy and their precision (such as 95% confidence intervals)	Page 12, table 3 and 4
	<b>25</b>	Any adverse events from performing the index test or the reference standard	n.a.
<b>DISCUSSION</b>			
	<b>26</b>	Study limitations, including sources of potential bias, statistical uncertainty, and generalisability	17

1		<b>27</b>	Implications for practice, including the intended use and clinical role of the index test	16
2	<b>OTHER INFORMATION</b>			
3				
4		<b>28</b>	Registration number and name of registry	ISRCTN registry; referene number: ISRCTN24868616
5				
6		<b>29</b>	Where the full study protocol can be accessed	At the registry (above), but not detailed.
7				
8		<b>30</b>	Sources of funding and other support; role of funders	18-19
9				

For peer review only

# STARD 2015

---

## AIM

STARD stands for “Standards for Reporting Diagnostic accuracy studies”. This list of items was developed to contribute to the completeness and transparency of reporting of diagnostic accuracy studies. Authors can use the list to write informative study reports. Editors and peer-reviewers can use it to evaluate whether the information has been included in manuscripts submitted for publication.

---

## EXPLANATION

A **diagnostic accuracy study** evaluates the ability of one or more medical tests to correctly classify study participants as having a **target condition**. This can be a disease, a disease stage, response or benefit from therapy, or an event or condition in the future. A medical test can be an imaging procedure, a laboratory test, elements from history and physical examination, a combination of these, or any other method for collecting information about the current health status of a patient.

The test whose accuracy is evaluated is called **index test**. A study can evaluate the accuracy of one or more index tests. Evaluating the ability of a medical test to correctly classify patients is typically done by comparing the distribution of the index test results with those of the **reference standard**. The reference standard is the best available method for establishing the presence or absence of the target condition. An accuracy study can rely on one or more reference standards.

If test results are categorized as either positive or negative, the cross tabulation of the index test results against those of the reference standard can be used to estimate the **sensitivity** of the index test (the proportion of participants *with* the target condition who have a positive index test), and its **specificity** (the proportion *without* the target condition who have a negative index test). From this cross tabulation (sometimes referred to as the contingency or “2x2” table), several other accuracy statistics can be estimated, such as the positive and negative **predictive values** of the test. Confidence intervals around estimates of accuracy can then be calculated to quantify the statistical **precision** of the measurements.

If the index test results can take more than two values, categorization of test results as positive or negative requires a **test positivity cut-off**. When multiple such cut-offs can be defined, authors can report a receiver operating characteristic (ROC) curve which graphically represents the combination of sensitivity and specificity for each possible test positivity cut-off. The **area under the ROC curve** informs in a single numerical value about the overall diagnostic accuracy of the index test.

The **intended use** of a medical test can be diagnosis, screening, staging, monitoring, surveillance, prediction or prognosis. The **clinical role** of a test explains its position relative to existing tests in the clinical pathway. A replacement test, for example, replaces an existing test. A triage test is used before an existing test; an add-on test is used after an existing test.

Besides diagnostic accuracy, several other outcomes and statistics may be relevant in the evaluation of medical tests. Medical tests can also be used to classify patients for purposes other than diagnosis, such as staging or prognosis. The STARD list was not explicitly developed for these other outcomes, statistics, and study types, although most STARD items would still apply.

---

## DEVELOPMENT

This STARD list was released in 2015. The 30 items were identified by an international expert group of methodologists, researchers, and editors. The guiding principle in the development of STARD was to select items that, when reported, would help readers to judge the potential for bias in the study, to appraise the applicability of the study findings and the validity of conclusions and recommendations. The list represents an update of the first version, which was published in 2003.

More information can be found on <http://www.equator-network.org/reporting-guidelines/stard>.

