Interrater reliability of sensory testing in patients with diabetic neuropathy

Isabelle M.L.P. Kamm 5854318

Department of Plastic, Reconstructive and Hand Surgery University Medical Centre Utrecht

Supervisors: N. Boers, W.D. Rinkel, J.H. Coert Department of Plastic, Reconstructive and Hand Surgery – University Medical Centre Utrecht

July 18th 2022 – October 7th 2022

Abstract

Background: Peripheral neuropathy, occurring in 30-50% of diabetic patients, is the main risk factor for foot ulceration. In order to detect diabetic neuropathy, the continuous and dichotomous 39-item and the dichotomous 13-item Rotterdam Diabetic Foot Study test batteries (RDF-39-C, RDF-39-D, RDF-13-D) were developed. This study examined the interrater reliability of these test batteries in patients with diabetic neuropathy.

Methods: Interrater reliability of the 39-item test batteries was determined across one pair of raters, whereas reliability of the RDF-13-D was examined across two pairs of raters. Interrater agreement rates (IRA) and intraclass coefficients (ICC) were calculated to assess reliability.

Results: Sixty-five patients with diabetes and symptomatic neuropathy were included. Interrater agreement was found to be acceptable for the RDF-39-C test battery but low for the RDF-39-D (IRA = 84.6% and 73.3%, respectively). Regarding the RDF-13-D, agreement rates ranged from 64.3% for one pair of raters to 81.4% for the other pair of raters. ICCs were all above .80, indicating high correlation.

Conclusion: The results demonstrate that the continuous version of the RDF-39 test battery is reliable across different clinicians, but when measured dichotomously it is not. The level of interrater agreement of the 13-item RDF test battery is dependent on the pair of raters. Though raters do not always agree on the absolute total scores of both dichotomous test batteries, the results indicate that their ratings are highly consistent. Standardised training across clinicians may be important to improve reliability.

Key words: diabetes mellitus – peripheral neuropathy – sensibility test – interrater reliability

Abbreviations: DFU, diabetic foot ulcer; ICC, intraclass correlation coefficient; IRA, interrater agreement; κ, Cohen's kappa coefficient; M2PD, moving two-point discrimination; RDF-13-D, dichotomous 13-item Rotterdam Diabetic Foot Study test battery; RDF-39-C, continuous 39-item Rotterdam Diabetic Foot Study test battery; RDF-39-D, dichotomous 39-item Rotterdam Diabetic Foot Study test battery; S1PD, static one-point discrimination; S2PD, static two-point discrimination

Introduction

One in ten adults are affected by diabetes, accounting for approximately 537 million diabetic patients worldwide.¹ The development of peripheral neuropathy is one of the most common complications of diabetes, occurring in 30% to 50% of patients.^{2,3} As diabetic neuropathy causes sensory loss, it is the main risk factor for foot ulceration and eventually lower extremity amputations. The lifetime incidence of developing a diabetic foot ulcer (DFU) may be as high as 34%⁴, and up to 85% of all amputations in individuals with diabetes is preceded by foot ulcers.⁵ Early detection and management of diabetic neuropathy may delay or prevent foot complications. It has been shown that podiatric care, patient education and regular foot examinations on sensory loss reduce the risk of adverse outcomes.⁶ Hence, screening of diabetic patients to identify those at risk for a foot ulcer is of great importance.

Current international guidelines recommend annual examination with a 10 gram monofilament and a tuning fork to assess loss of sensation in diabetic patients.7 Previous research has, however, demonstrated that these tests only identify diabetic neuropathy in advanced stages.^{8,9} Recently, Rinkel et al. developed the so-called 39item Rotterdam Diabetic Foot (RDF-39) Study test battery to create an objective and reliable screening method for the detection of early diabetic neuropathy.^{10,11} The test battery allows a full evaluation of somatosensory function of the feet, as it includes tests on static and moving two-point discrimination (S2PD, M2PD), static one-point discrimination (S1PD), vibration sense, cold perception, Romberg's test and information about experienced numbness of the feet, prior foot ulceration and prior amputation. Data can be characterised as either continuous (RDF-39-C) or dichotomous (RDF-39-D). A shorter version consisting of 13 dichotomous items (RDF-13-D) was developed as well for clinical settings, as the RDF-39-C and -D are more time consuming. The RDF-13-D allows healthcare providers to perform a quick assessment of sensory loss and may serve as a complement to the current screening methods for diabetic neuropathy.¹⁰

However, the application of these test batteries by different healthcare providers could potentially cause discrepancies in test results, which may lead to conflicting interpretations of the sensory status of patients' feet. Establishing the reliability of the three test batteries has crucial implications for their generalizability and suitability in both clinical settings and research, especially when used as a screening tool. Previous studies have investigated the interrater reliability of some test instruments across examiners in patients with diabetic neuropathy^{12–15}, but the reliability of the three test batteries described by Rinkel and colleagues has not yet been studied. Therefore, the aim of this study was to determine the interrater reliability of the RDF-39-C, RDF-39-D and RDF-13-D Study test batteries in patients with diabetic neuropathy.

Methods

Study design and patient population

This study was part of the DeCompression (DECO) study, an ongoing multicentre, randomised controlled trial which investigates the (cost-)effectiveness of surgical decompression of lower limb compression neuropathy compared to the standard, nonsurgical care in diabetic patients. The DECO study is carried out in 11 hospitals in the Netherlands (University Medical Centre, Utrecht; Diakonessenhuis, Utrecht; Franciscus Gasthuis & Vlietland, Rotterdam; Maasstad Hospital, Rotterdam; Jeroen Bosch Hospital, Den Bosch; Isala Hospital, Zwolle; St. Antonius Hospital, Nieuwegein; Meander Medical Centre, Amersfoort; University Medical Centre Amsterdam, VUmc, Amsterdam; University Medical Centre Amsterdam, AMC, Amsterdam; OLVG Hospital, Amsterdam). In the current study, participants were recruited during their initial or follow-up visit to the outpatient clinic of one of the participating hospitals of the DECO trial between July 14th and October 4th 2022. All patients with diabetes type 1 or 2 and symptomatic neuropathy were eligible for inclusion. Further inclusion and exclusion criteria of the DECO trial were not applicable to the current study.

Procedure

The interrater reliability of the RDF-39 and RDF-13-D test batteries was determined across a researcher with one year experience in performing the tests (rater 1) and a final-year medical student (rater 2). The medical student had no previous experience in conducting the test batteries, but the protocol was thoroughly explained and all tests were demonstrated beforehand by rater 1. In addition, the interrater reliability of the RDF-13-D was tested across the medical student (rater 2) and a research nurse with six months experience (rater 3). All raters were blinded to the test results from each other. The order of raters conducting the testing sessions was randomised in a non-predetermined manner. Prior to the first testing session, the test battery was explained to the patient. Either the RDF-39-C, RDF-39-D or RDF-13-D test battery was used, depending on patients' availability. Taking into account patient travel time and costs, participants were re-tested the same day, with a few minutes interval between the first and the second testing session.

Data-collection and sensory testing

Demographic data and information on patients' height, weight, type and duration of diabetes, glycated hemoglobin value and blood pressure were retrieved from the electronic files of the participants. For current analysis, baseline data from all patients were used.

The 39-item Rotterdam Diabetic Foot Study test battery (continuous and dichotomous)

The RDF-39 consists of 39 items to assess the sensory status of the feet. It includes static and moving two-point discrimination (S2PD, M2PD), static one-point discrimination (S1PD), vibration sense, perception of cold, Romberg's balance test, and information on complaints of numbness and history of foot ulceration or amputation. The tests are

performed on both feet. The continuous version of the RDF-39 test battery (RDF-39-C) takes about 30 minutes to complete. The time needed to complete the dichotomised version (RDF-39-D) is approximately ten minutes.

In this study, tests of S2PD and M2PD were performed using a Disk-Criminator. To evaluate S2PD, the skin of the test site was randomly touched with one or two metal spikes, whereas for M2PD the spikes were moved over the test site. In both tests, the patient was asked whether he felt one or two spikes. The smallest distance between two spikes for which the patient gave three correct answers out of four trials was noted in millimetres for that test site. In case the RDF-39-D was used, the threshold was set at 8 millimetres (aberrant: >8 millimetres), based on previously reported normative data.¹⁶ S2PD was performed on five test sites which correspond to the nerve distribution of the foot: 1) plantar hallux (medial plantar nerve, branch of the tibial nerve); 2) medial heel (calcaneal nerve, branch of the tibial nerve); 3) first dorsal web space (deep peroneal nerve); 4) lateral foot (sural nerve); and 5) plantar fifth toe (lateral plantar nerve, branch of the tibial nerve). M2PD was only tested on four of the five sites (plantar hallux, medial heel, first dorsal web space and lateral foot), as the area of the plantar fifth toe is too small. S1PD was examined with the 20 piece Semmes-Weinstein monofilament set, ranging from 0.008 to 300 grams. Monofilaments were applied on the skin and slightly pressed into a C-shape. Patients were asked to indicate whether and where they felt a stimulus. The thinnest monofilament in the test series that was felt was recorded in grams for that site. In case the RDF-39-D was used, the threshold was set at 10 grams (aberrant: >10 grams), based on current international guidelines.^{7,17} S1PD was assessed on the same five test locations as S2PD. Vibration sense was evaluated using a Rydel-Seiffer tuning fork on the dorsum of the interphalangeal joint of the hallux and on the medial malleolus. Patients were instructed to indicate the moment they no longer perceived the vibration. In case the RDF-39-D test battery was used, outcomes were dichotomised by comparing them with normal vibration threshold values.¹⁸ Sensation of cold was tested on the medial arch of the foot by applying the tuning fork at room temperature on the skin. Romberg's test was used to assess sense of balance. Patients were asked to stand in an erect position with their feet together and their arms held forward with the palms facing upwards. The test was scored as positive when the patient was unable to maintain his balance with his eyes closed. Information on complaints of numbness and history of foot ulceration or amputation was derived from the patient interview.

In case the RDF-39-C test battery was used, items were scored as a continuous variable for S2PD, M2PD, S1PD and vibration sense, as 'positive' or 'negative' for Romberg's test and cold perception, and as 'yes' or 'no' for complaints of numbress and history of ulcer or amputation.

For the RDF-39-D test battery, results of all 39 items were broken down into dichotomous variables, scored as either 0 (normal) or 1 (abnormal). The sum of the individual item scores yielded a total score. Patients with a total score of 24 points or more were classified as having a high risk of DFU development, based on previously reported data.¹⁹

Test results of participants in whom the continuous version of the RDF-39 was performed were dichotomised in order to enlarge the sample size of the RDF-39-D.

The 13-item Rotterdam Diabetic Foot Study test battery (dichotomous)

The 13-item RDF is a shorter version of the RDF-39 and contains 13 dichotomized items, which are listed in Table 1. It takes approximately three to four minutes to complete. The threshold values used for the RDF-39-D test battery were also used for the RDF-13-D in order to dichotomise the test results. All items were scored as either 0 (normal) or 1 (abnormal) and summed into a total score. Based on previously reported data, a cut-off value of 7 points was used to identify patients with and without a high risk of DFU development.¹⁹

In order to create a larger sample size, the results of the patients in whom the RDF-39-D was tested were converted to the 13 items of the RDF-13-D.

| RDF-39 | | RDF-13 | | |
|-----------------------|-----------------------|--|-----------------------|--|
| Left lower extremity | Right lower extremity | / Left lower extremity Right lower ext | | |
| S2PD hallux | S2PD hallux | S2PD hallux | | |
| S2PD medial heel | S2PD medial heel | | S2PD medial heel | |
| S2PD first dorsal web | S2PD first dorsal web | | S2PD first dorsal web | |
| S2PD lateral foot | S2PD lateral foot | | | |
| S2PD fifth toe | S2PD fifth toe | | | |
| M2PD hallux | M2PD hallux | | | |
| M2PD medial heel | M2PD medial heel | | | |
| M2PD first dorsal web | M2PD first dorsal web | M2PD first dorsal web | | |
| M2PD lateral foot | M2PD lateral foot | | | |
| S1PD hallux | S1PD hallux | S1PD hallux | S1PD hallux | |
| S1PD medial heel | S1PD medial heel | | | |
| S1PD first dorsal web | S1PD first dorsal web | | | |
| S1PD lateral foot | S1PD lateral foot | | | |
| S1PD fifth toe | S1PD fifth toe | | | |
| Vibration sense MM | Vibration sense MM | Vibration sense MM | Vibration sense MM | |
| Vibration sense IP | Vibration sense IP | Vibration sense IP | Vibration sense IP | |
| Cold sensation | Cold sensation | | | |
| Amputation | Amputation | Amputation | Amputation | |
| Romberg's test | | | | |
| Numbness | | | | |
| Prior ulcer | | Prior ulcer | | |

Table 1. The 39-item and 13-item Rotterdam Diabetic Foot Study test battery

Abbreviations: RDF-39, 39-item Rotterdam Diabetic Foot Study test battery; RDF-13, 13-item Rotterdam Diabetic Foot Study test battery; S2PD, static two-point discrimination; M2PD, moving two-point discrimination; S1PD, static one-point discrimination; MM, medial malleolus; IP, interphalangeal joint.

Statistical analysis

Descriptive statistics were used to describe baseline data. Continuous data were reported as means and standard deviations (SD), categorical data were reported as frequencies and percentages.

The interrater reliability was determined using the interrater agreement (IRA), the intraclass correlation coefficient (ICC) for continuous variables and the Cohen's kappa coefficient (κ) for dichotomous variables. The level of IRA was divided into four different categories: 'total agreement' (i.e. raters agreed on all outcomes), 'rater difference of ≤1 point' (i.e. raters agreed on all ratings or ratings between raters differed by one point), 'rater difference of ≤2 points' (i.e. ratings between raters differed by two or less points) and 'rater difference of ≤3 points' (i.e. ratings between raters differed by three or less points). IRA values above 75% were considered as an acceptable degree of agreement, values above 90% indicated a high level of agreement.²⁰ ICC values and their 95% confidence intervals (CI) were determined using a single measurement, absolute agreement, two-way random effects model. Values less than .50 indicated poor reliability, between .50 and .75 moderate reliability, between .75 and .90 good reliability and values greater than .90 were indicative of excellent reliability.²¹ For κ coefficients, cut-off values were based on the Landis and Koch classification, considering κ below .00 as poor agreement, .00 to .20 slight agreement, .21 to .40 fair agreement, .41 to .60 moderate agreement, .61 to .80 substantial agreement and .81 to 1.00 almost perfect agreement.²²

Sub-analyses of the different tests included in the test batteries were conducted as well to provide more insight in the interrater reliability of each test. Since information about complaints of numbness, history of ulceration and previous amputation were fixed data, the results of these items did not differ between raters. Hence, only data of variable tests were analysed separately, which included S2PD, M2PD, S1PD, vibration sense, cold sensation and Romberg's test for the RDF-39 test batteries and S2PD, M2PD, S1PD and vibration sense for the RDF-13-D test battery. For the continuous version of the RDF-39, the number of times raters agreed on *ratings of individual test items* was added up and divided by the total number of ratings, in order to calculate the IRA of the whole test battery and the IRA of the different tests. For the dichotomous test batteries (RDF-39-D and RDF-13-D), the IRA of the whole test battery and the IRA of the different tests agreed on the *total score* of the whole test battery or of a particular test, and dividing that number by the total number of cases rated.

In the current study, an acceptable level of IRA was achieved when raters completely agreed or their ratings differed by one point (i.e. 'rater difference of ≤ 1 point') in at least 75% of cases. However, for the S2PD, M2PD and S1PD tests of the RDF-13-D test battery only complete agreement was considered acceptable since these tests consisted of only one to three test items. Hence, the total agreement rate for these tests had to be at least 75% to achieve an acceptable level of agreement.

All analyses were performed using IBM SPSS Statistics, version 27 (IBM Corp). The threshold for statistical significance was set at $p \le .05$.

Results

A total of 65 diabetic patients (aged 45-76 years) were included in the study, whom all provided data for one or more versions of the test batteries (i.e. RDF-39-C, RDF-39-D, RDF-13-D), see Supplementary Figure 1. Overall, 48 (73.8%) patients were male and 54 (83.1%) had diabetes type 2. The mean duration of diabetes was 16.8 years (SD \pm 12.5 years). Further baseline characteristics of the participants are listed in Table 2.

| Characteristic | Participants (n = 65) |
|--|-----------------------|
| Sex, n (%) | |
| Male | 48 (73.8%) |
| Female | 17 (26.2%) |
| Age group, <i>n</i> (%) | |
| 18-25 years | 0 (0.0%) |
| 26-35 years | 0 (0.0%) |
| 36-45 years | 1 (1.5%) |
| 46-55 years | 8 (12.3%) |
| 56-65 years | 27 (41.5%) |
| 66-76 years | 29 (44.6%) |
| Ethnicity, n (%) | |
| Caucasian | 62 (95.4%) |
| Indo-Surinamese | 1 (1.5%) |
| African | 1 (1.5%) |
| Asian | 0 (0.0%) |
| Hindu | 0 (0.0%) |
| Other | 1 (1.5%) |
| Height (m), mean (SD) | 1.79 (0.1) |
| Weight (kg), mean (SD) | 91.3 (21.2) |
| BMI (kg/m ²), mean (SD) | 28.4 (6.1) |
| Type of diabetes, n (%) | |
| Туре 1 | 11 (16.9%) |
| Type 2 | 54 (83.1%) |
| Duration of diabetes (years), mean (SD) | 16.8 (12.5) |
| HbA1c (mmol/mol), mean (SD) | 59.0 (16.5) |
| Systolic blood pressure (mmHg), mean (SD) | 135.7 (14.6) |
| Diastolic blood pressure (mmHg), mean (SD) | 80.1 (8.5) |

Table 2. Baseline characteristics.

Abbreviations: n, number; SD, standard deviation; BMI, body mass index; HbA1c, glycated hemoglobin.

RDF-39-C test battery

The interrater reliability of the RDF-39-C test battery across rater 1 and rater 2 was assessed in a cohort of 11 patients. Table 3 summarises the cumulative IRA values of the RDF-39-C and Figure 1 displays the distribution of interrater differences of the test battery. Over all 11 patients (and 429 ratings), raters agreed in 69.0%, indicating a low level of agreement. However, a rater difference of \leq 1 points was achieved in 84.6%, which indicates an acceptable level of agreement. Interrater reliability analysis of the different tests showed low total agreement rates for S1PD and vibration sense (IRA =

28.2% and 52.3%, respectively), whereas an acceptable level of agreement for S2PD and M2PD (IRA = 86.4% and 79.5%, respectively), and even perfect agreement for cold perception and Romberg's balance test (IRA = 100% for both tests) was found.

IRA and ICC values of all individual test items of the test battery are listed in Supplementary Table 1.

Table 3. Interrater reliability of the RDF-39-C test battery reported as cumulative interrater agreement between rater 1 and rater 2.

Interrater reliability of RDF-39-C

| | Rater 1 vs. rater 2 |
|--|---------------------|
|--|---------------------|

| (n =) | 11) |
|--------|-----|
|--------|-----|

| | IRA (%) | | | | | |
|---------------------|-----------|---------------|---------------|---------------|--|--|
| | Total | Rater | Rater | Rater | | |
| | agreement | difference ≤1 | difference ≤2 | difference ≤3 | | |
| Test | | points | points | points | | |
| S2PD | 86.4 | 90.0 | 94.5 | 95.5 | | |
| M2PD | 79.5 | 86.4 | 88.6 | 92.0 | | |
| S1PD | 28.2 | 64.5 | 83.6 | 92.7 | | |
| Vibration sense | 52.3 | 90.9 | 100 | - | | |
| Cold perception | 100 | - | - | - | | |
| Romberg's test | 100 | - | - | - | | |
| Whole RDF-39-C test | 69.0 | 84.6 | 92.1 | 95.3 | | |
| batterv | | | | | | |

Abbreviations: RDF-39-C, continuous version of the 39-item Rotterdam Diabetic Foot Study test battery; n, number; IRA, interrater agreement; S2PD, static two-point discrimination; M2PD, moving two-point discrimination; S1PD, static one-point discrimination.



Figure 1. Interrater agreement between rater 1 and rater 2 of the continuous 39-item Rotterdam Diabetic Foot (RDF-39-C) Study test battery. Each coloured bar represents a difference in ratings of either 0 (green), 1 (yellow), 2 (orange) or \geq 3 (red) points between the raters.

RDF-39-D test battery

The test results of 30 patients were used in the interrater reliability analysis of the RDF-39-D test battery. Cumulative IRA rates and ICC values are summarised in Table 4A and the distribution of rater differences is displayed in Figure 2. Regarding the total score of the RDF-39-D test battery, the level of complete agreement between rater 1 and rater 2 was found to be low (IRA = 60.0%). In 73.3% the total score differed by \leq 1 point between raters, still indicating low interrater agreement. Nevertheless, ICC analysis of the total score of the test battery demonstrated good reliability (ICC = .85, 95% CI .71-.93). Levels of total agreement of the different tests of the test battery ranged between 73.3% and 100%, with M2PD having the lowest and both cold perception and Romberg's test having the highest IRA rate. ICC values for each test ranged between .02 and .89, with S1PD and vibration sense both showing good reliability (ICC = .77, 95% CI .57-.88 and ICC = .89, 95% CI .78-.94, respectively).

IRA and κ values of all individual test items of the RDF-39-D test battery are displayed in Supplementary Table 2.

Regarding the risk categories of diabetic foot ulcer development (either low or high risk, based on the cut-off value of 24 points), the level of interrater agreement of risk classification was high (IRA = 93.3%). The κ value was significant (κ = .82, *p* < .001), indicating almost perfect agreement (Table 4B).

Table 4. Interrater reliability of the RDF-39-D test battery between rater 1 and 2, reported as cumulative interrater agreement and intraclass correlation coefficients (Table 4A), and interrater reliability of the diabetic foot risk classification of the RDF-39-D test battery between rater 1 and 2, reported as interrater agreement and Cohen's kappa coefficient (Table 4B).
A. Interrater reliability of RDF-39-D

| Rater 1 vs. rater 2 | | | | | |
|---------------------|-----------|------------|------------|------------|--------------|
| (n = 30) | | | | | |
| | | | | | |
| | Total | Rater | Rater | Rater | |
| | agreement | difference | difference | difference | |
| Test | | ≤1 points | ≤2 points | ≤3 points | ICC (95% CI) |
| S2PD | 86.7 | 90.0 | 96.7 | 100 | .02 (3538) |
| M2PD | 73.3 | 90.0 | 93.3 | 100 | .48 (.1571) |
| S1PD | 83.3 | 86.7 | 100 | - | .77 (.5788) |
| Vibration sense | 76.7 | 83.3 | 100 | - | .89 (.7894) |
| Cold perception | 100 | - | - | - | _ |
| Romberg's test | 100 | - | - | - | _ |
| Total score of RDF- | 60.0 | 73.3 | 90.0 | 93.3 | .85 (.7193) |
| 39-D test battery | | | | | |
| | • | | | | |

| D. L. (| - (DELL - | | |
|---------------------------|------------|------------------|----|
| B. Interrater reliability | voi dfu r | isk classificati | or |

Rater 1 vs. rater 2

(n = 30)

| | IRA (total agreement) (%) | Cohen's κ | <i>p</i> -value |
|-------------------------|---------------------------|-----------|-----------------|
| DFU risk classification | 93.3 | .82 | <.001**** |

Abbreviations: RDF-39-D, dichotomous version of the 39-item Rotterdam Diabetic Foot Study test battery; n, number; IRA, interrater agreement; ICC, intraclass correlation coefficient; CI, confidence interval; S2PD, static two-point discrimination; M2PD, moving two-point discrimination; S1PD, static one-point discrimination; DFU, diabetic foot ulcer.

**** $p \le .001$



Figure 2. Interrater agreement between rater 1 and rater 2 of the dichotomous 39-item Rotterdam Diabetic Foot (RDF-39-D) Study test battery. Each coloured bar represents a difference in ratings of either 0 (green), 1 (yellow), 2 (orange) or \geq 3 (red) points.

RDF-13-D test battery

Table 5A displays the cumulative IRA and the ICC for the assessment of interrater reliability of the RDF-13-D test battery among rater 1 versus 2 and rater 2 versus 3. Figure 3 provides a graphical representation of the distribution of interrater differences of the test battery for both pairs of raters. For reliability analysis among rater 1 and 2 the test results of 43 patients were used, whereas a cohort of 14 subjects was assessed between rater 2 and 3. For the total score of the RDF-13-D, a low degreee of total agreement but a good level of reliability was found for both pairs of raters (rater 1 versus 2: IRA = 41.9%, ICC = .86, 95% CI .75-.92; rater 2 versus 3: IRA = 42.9%, ICC = .81, 95% CI .51-.94). More detailed analysis of agreement between rater 1 and rater 2 showed a high level of agreement for M2PD (IRA = 95.3%), acceptable levels of agreement for S2PD and S1PD (IRA = 88.4% for both tests) and a low agreement level for vibration sense (IRA = 67.4%). ICC analysis of these tests demonstrated good reliability with respect to S1PD and vibration sense (ICC = .88, 95% CI .78-.93 and ICC = .85, 95% CI .73-.92, respectively). Poor reliability was found for both S2PD and M2PD testing (ICC = .16, 95% CI -.14-.43 and ICC = .49, 95% CI .23-.68, respectively).

Agreement analysis among rater 2 and 3 revealed a high level of interrater agreement for S1PD (IRA = 92.9%), but low levels of total agreement for the other tests (S2PD: IRA = 42.9%; M2PD: IRA = 71.4%; vibration sense: IRA = 64.3%). In contrast, ICC analysis yielded poor reliability values for S2PD and M2PD (respectively ICC = .12, 95% CI - .45-.60 and ICC = .19, 95% CI - .35-.64), but demonstrated good reliability with regards to vibration sense (ICC = .87, 95% CI .65-.96).

The IRA and κ coefficients of all individual test items of the RDF-13-D test battery for both rater pairs are shown in Supplementary Table 3.

Regarding the DFU risk classification, an acceptable level of agreement was found for both pairs of raters (rater 1 versus 2: IRA = 83.7%, rater 2 versus 3: IRA = 78.6%). Cohen's κ analysis revealed substantial agreement between rater 1 and rater 2 (κ = .68, p < .001) and moderate agreement between rater 2 and 3 (κ = .57, p = .018) (Table 5B).

| Table 5. Interrater reliability of the RDF-13-D test battery between rater 1 and 2 and between rater 2 |
|---|
| and 3, reported as cumulative interrater agreement and intraclass correlation coefficients (Table 5A), |
| and interrater reliability of the diabetic foot risk classification of the RDF-13-D test battery between |
| rater 1 and 2 and between rater 2 and 3, reported as interrater agreement and Cohen's kappa |
| coefficient (Table 5B). |

| | | IRA (%) | | | | |
|---------------|----------------|-----------|------------|------------|------------|-------------|
| | | Total | Rater | Rater | Rater | |
| | | agreement | difference | difference | difference | ICC |
| Rater pair | Test | | ≤1 points | ≤2 points | ≤3 points | (95% CI) |
| Rater 1 vs. 2 | S2PD | 88.4 | 97.7 | 100 | - | .16 (1443) |
| (n = 43) | M2PD | 95.3 | 100 | - | - | .49 (.2368) |
| | S1PD | 88.4 | 100 | - | - | .88 (.7893) |
| | Vibration | 67.4 | 86.0 | 97.7 | 100 | .85 (.7392) |
| | sense | | | | | |
| | Total score of | 41.9 | 81.4 | 97.7 | 100 | .86 (.7592) |
| | RDF-13-D test | | | | | |
| | battery | | | | | |
| | | | | | | |
| Rater 2 vs. 3 | S2PD | 42.9 | 100 | - | - | .12 (4560) |
| (n = 14) | M2PD | 71.4 | 100 | - | - | .19 (3564) |
| | S1PD | 92.9 | 100 | - | - | .65 (.2287) |
| | Vibration | 64.3 | 85.7 | 100 | - | .87 (.6596) |
| | sense | | | | | |
| | Total score of | 42.9 | 64.3 | 92.9 | 100 | .81 (.5194) |
| | RDF-13-D test | | | | | |
| | battery | | | | | |

| A. Interrater reliabi | lity of RDF-13-D |
|-----------------------|------------------|

| B. Interrater reliability of DFU risk classification | | | | | | | |
|--|----------------|----------------|-----------|-----------------|--|--|--|
| | | IRA (total | | | | | |
| Rater pair | | agreement) (%) | Cohen's κ | <i>p</i> -value | | | |
| Rater 1 vs. 2 | DFU risk | 83.7 | .68 | <.001**** | | | |
| (n = 43) | classification | | | | | | |

| Rater 2 vs. 3 | DFU risk | 78.6 | .57 | .018* |
|---------------|----------------|------|-----|-------|
| (n = 14) | classification | | | |

Abbreviations: RDF-13-D, dichotomous version of the 13-item Rotterdam Diabetic Foot Study test battery; IRA, interrater agreement; ICC, intraclass correlation coefficient; CI, confidence interval; n, number; S2PD, static two-point discrimination; M2PD, moving two-point discrimination; S1PD, static one-point discrimination; DFU, diabetic foot ulcer.

 $^{*}p \leq .05;\,^{****}p \leq .001$



Figure 3. Interrater agreement of the dichotomous 13-item Rotterdam Diabetic Foot (RDF-13-D) Study test battery between rater 1 and 2 (Figure 3A) and between rater 2 and rater 3 (Figure 3B). Each coloured bar represents a difference in ratings of either 0 (green), 1 (yellow), 2 (orange) or \geq 3 (red) points.

Discussion

The purpose of the present study was to examine the interrater reliability of the RDF-39-C, RDF-39-D and RDF-13-D test batteries in patients with diabetes and symptomatic neuropathy. This study established that the continuous version of the 39item RDF test battery (i.e. RDF-39-C), which yielded over 80% agreement between raters, is a reliable tool to use in research and clinical settings. When measured dichotomously (i.e. RDF-39-D), however, the interrater agreement decreased to 73.3%, which is just below the threshold of 75% to be considered acceptable. The degree of correlation of the RDF-39-D test battery was found to be good (ICC = .85). The combination of the obtained low IRA and high ICC values likely indicate that rater 1 consistently assigned higher scores compared to rater 2 (data not shown), while the scores of both raters are highly correlated (see Supplementary Table 4 for a hypothetical example). The similarity between ratings rather than the correlation of raters' judgements is of interest in the current study, and hence the authors argue that the level of agreement provides the best representation of interrater reliability of the test battery. Hence, based solely on its low interrater agreement rate, the RDF-39-D does not seem reliable to use as a screening tool. Nevertheless, the high ICC value demonstrates that even though raters do not agree on the absolute total scores of the RDF-39-D test battery, their ratings are highly consistent. Concerning the shorter 13item test battery (i.e. RDF-13-D), discrepancy between agreement rates across rater pairs was found, with rater 1 versus 2 showing acceptable and rater 2 versus 3 showing low interrater agreement. These results assume that interrater reliability of this test battery varies widely depending on the pair of raters. A possible explanation for these conflicting results could be the difference in sample size between rater pairs. Rater 2 and rater 3 examined three times less subjects compared to rater 1 versus 2 (14 versus 43 patients). The small patient cohort of rater 2 and 3 may have prevented the identification of acceptable interrater reliability for the RDF-13-D test battery. ICC analysis of the RDF-13-D test battery revealed a good degree of correlation for both pairs of raters (ICCs > .80). These results indicate that raters highly agree on the relative ranking of the total score of the RDF-13-D test battery.

In our study, classification of risk of DFU development was acceptably reliable between raters, as agreement levels above 75% and significant κ coefficients were found. This finding demonstrates that patients are consistently being classified as having either low or high risk of foot ulcer development by different clinicians. Our study results are, especially for the use of the test batteries in screening settings, clinically meaningful, since shorter screening intervals are recommended in patients who are considered to be at higher risk of DFU development.⁷

Previous research has shown that repetitive nerve stimulation in diabetic patients resulted in abnormal responsiveness of mechanosensitive afferents due to fatigue of these nerves.²³ The data of their study indicate that repeated mechanical stimulation could lead to abnormal encoding of these stimuli. Despite the possibility of nerve irritation and fatiguability from repeated measures in the present study, it is interesting to note that the continuous and thereby most extensive version of the RDF-

39 yielded higher interrater reliability compared to the dichotomous version of this test battery. These results suggest that dichotomous measurement of the 39-item test battery affected reliability rather than duration of the testing sessions.

A closer inspection of the different tests included in the test batteries showed that S1PD seems to be reliable across raters when used in the dichotomous versions of the test batteries. This finding is in line with the results of previous studies, which all reported acceptable or even high levels of interrater reliability for S1PD testing in diabetic patients.¹²⁻¹⁴ In the present study S1PD was however less reliable in case the continuous version of the RDF-39 was used, as the extent of interrater agreement with a rating difference of ≤1 point was only 64.5%. Sensation of individual monofilaments could have been hampered due to the use of small inter-monofilament differences in the RDF-39-C test battery. This could have caused variability in ratings between raters and thus lower interrater agreement. Nevertheless, our study showed that agreement between raters for S1PD testing increased by almost twenty percent (to 83.6%) in case ratings did not differ more than two points. This finding shows that adequate interrater agreement levels can be achieved for S1PD testing, but only if a difference in ratings of two points or less is considered acceptable. Especially in clinical settings where more heterogeneity among clinicians and therefore less consistency in conducting the test batteries can be expected, accepting greater variability in ratings is something to consider.

The results of all three test batteries, except for RDF-13-D between rater 2 and 3, showed acceptable and high interrater agreement rates for the two-point discrimination tests (i.e. S2PD and M2PD). Contradictory results were however found regarding correlation coefficients, which were all less than .50, indicating poor reliability. Since the ICC partially depends on differences in ratings, its value could be misleadingly low in case ratings are tightly clustered. As in the present study most patients scored the worst result possible for S2PD and M2PD testing (data not shown), limited variation in ratings for these tests existed. A previous study has demonstrated that both static and moving two-point discrimination are the sensory functions that are lost first in patients with diabetic neuropathy.²⁴ This finding could explain the observed high scores on S2PD and M2PD testing in the current study, which have resulted in high agreement but low ICC values.

Perfect interrater agreement in the present study has shown both cold perception and Romberg's test to be reliable. These results implicate not only that both tests are conducted in a consistent fashion by different clinicians but also that proprioception and cold perception do not seem to be highly variable within patients with diabetic neuropathy. In contrast to the findings of our study, a previous study by Wasan *et al.* demonstrated low reliability across different examiners in testing cold perception.²⁵ In the study of Wasan and colleagues, whose study population consisted of patients with post herpetic neuralgia, it was however not tested whether patients could feel a cold stimulus but whether cold sensation at the affected site felt as more, less, different or the same in comparison to the control site. Their findings demonstrate that, in contrast to our method of testing cold perception, their method is not reliable across different

examiners. It is interesting to note that the results of their study implicate that patients with neuropathic symptoms may not rate the intensity of a cold stimulus consistently. The present study has some limitations which must be considered upon interpretation of the results. First, the level of agreement among raters could have been overestimated as percent agreement does not take chance agreement into account. However, it was not likely that raters guessed on scores in the current study, so the authors suggest that the interrater reliability could have been determined safely by calculating percent agreement. Second, the level of experience in conducting the tests differed between raters. A study by Marx et al. showed that clinicians with more expertise achieved higher ICC values, suggesting that level of experience plays an important role in reliability analysis.²⁶ Hence, the varying degrees of rater experience in the current study could have affected the results. A third limitation of this study is the small sample size, which limits the generalisability of our results. Though the test results of 30 and 43 subjects were used to assess the interrater reliability of respectively the RDF-39-D and RDF-13-D test battery among rater 1 and rater 2, only 11 patients were involved in reliability testing of the RDF-39-C. Moreover, the interrater reliability of the RDF-13-D among the other pair of raters (i.e. rater 2 versus rater 3) was assessed in a cohort consisting of only 14 patients. As a sample size of at least 30 subjects is used as a rule of thumb²⁷, studies with larger sample sizes are warranted.

The importance of estimating sensory loss in patients with diabetes has been clearly established. Batteries of sensory tests have been developed to assess sensibility in diabetic patients' feet. The current study provides insight into the interrater reliability of these test batteries. Our study shows that the continuous version of the RDF-39 is a reliable tool across examiners and can be used safely in research and clinical settings. The test battery is however less reliable when measured dichotomously. Concerning the shorter RDF-13-D test battery, the interrater reliability varies substantially depending on the pair of raters. Standardised training across healthcare professionals may be of importance to improve reliability of the three test batteries.

References

- International Diabetes Federation. IDF Diabetes Atlas 10th edition. 2021.
 [Internet]. Available from: https://diabetesatlas.org/atlas/tenth-edition/.
- 2 Dyck PJ, Kratz KM, Karnes JL, Litchy WJ, Klein R, Pach JM *et al.* The prevalence by staged severity of various types of diabetic neuropathy, retinopathy, and nephropathy in a population-based cohort: The rochester diabetic neuropathy study. *Neurology* 1993; **43**: 817–824.
- 3 Edwards JL, Vincent AM, Cheng HT, Feldman EL. Diabetic neuropathy: Mechanisms to management. *Pharmacol Ther* 2008; **120**: 1–34.
- 4 Armstrong DG, Boulton AJM, Bus SA. Diabetic Foot Ulcers and Their Recurrence. *N Engl J Med* 2017; **376**: 2367–2375.
- 5 Pecoraro RE, Reiber GE, Burgess EM. Pathways to Diabetic Basis for Prevention. *Diabetes Care* 1990; **13**: 513–521.
- 6 Singh N, Armstrong DG, Lipsky BA. Preventing foot ulcers in patients with diabetes. *JAMA* 2005; **293**: 217–228.
- Schaper NC, Van Netten JJ, Apelqvist J, Hinchliffe RJ, Lipsky BA, Board IE.
 Practical Guidelines on the prevention and management of diabetic foot disease (IWGDF 2019 update). *Diabetes Metab Res Rev* 2020; 36: 1–10.
- 8 Young M, Breddy J, Veves A, Boulton A. The Prediction of Diabetic Neuropathic Foot Ulceration Using Vibration Perception Thresholds. A prospective Study. *Diabetes Care* 1994; **17**: 557–560.
- Kamei N, Yamane K, Nakanishi S, Yamashita Y, Tamura T, Ohshita K *et al.* Effectiveness of Semmes-Weinstein monofilament examination for diabetic peripheral neuropathy screening. *J Diabetes Complicat* 2005; **19**: 47–53.
- 10 Rinkel WD, Aziz MH, Van Neck JW, Cabezas MC, Van der Ark LA, Coert JH. Development of grading scales of pedal sensory loss using Mokken scale analysis on the Rotterdam Diabetic Foot Study Test Battery data. *Muscle Nerve* 2019; **60**: 520–527.
- Rinkel WD, Van der Oest MJ, Coert JH. Item reduction of the 39-item
 Rotterdam Diabetic Foot Study Test Battery using decision tree modelling.
 Diabetes Metab Res Rev 2020; 36: e3291.
- Lanting SM, Spink MJ, Tehan PE, Vickers S, Casey SL, Chuter VH. Non-invasive assessment of vibration perception and protective sensation in people with diabetes mellitus: Inter- And intra-rater reliability. *J Foot Ankle Res* 2020; 13: 1–7.
- 13 Bagherzadeh Cham M, Mohseni-Bandpei MA, Bahramizadeh M, Kalbasi S, Biglarian A. Reliability of semmes-weinstein monofilaments and tuning fork on pressure and vibration sensation measurements in diabetic patients. *Iran Rehabil J* 2019; **17**: 1–8.
- 14 Tentolouris A, Tentolouris N, Eleftheriadou I, Jude EB. The Performance and Interrater Agreement of Vibration Perception for the Diagnosis of Loss of Protective Sensation in People With Diabetes Mellitus. *Int J Low Extrem Wounds* 2021. doi:10.1177/1534734621994058.

- 15 McIllhatton A, Lanting S, Lambkin D, Leigh L, Casey S, Chuter V. Reliability of recommended non-invasive chairside screening tests for diabetes-related peripheral neuropathy: A systematic review with meta-analyses. *BMJ Open Diabetes Res Care* 2021; **9**: 1–11.
- 16 Rinkel WD, Aziz MH, Van Deelen MJM, Willemsen SP, Castro Cabezas M, Van Neck JW *et al.* Normative data for cutaneous threshold and spatial discrimination in the feet. *Muscle and Nerve* 2017; **56**: 399–407.
- 17 American Diabetes Association. Microvascular complications and foot care: Standards of medical care in Diabetes - 2018. *Diabetes Care* 2018; **41**: S105–S118.
- 18 Martina ISJ, Van Koningsveld R, Schmitz PIM, Van Der Meché FGA, Van Doorn PA. Measuring vibration threshold with a graduated tuning fork in normal aging and in patients with polyneuropathy. J Neurol Neurosurg Psychiatry 1998; 65: 743–747.
- Rinkel WD, Fakkel TM, Dijkstra DA, Castro Cabezas M, Coert JH. Mate van gevoelsverlies voorspelt het risico op de diabetische voet. *Podosophia* 2021; 29: 9–13.
- 20 Graham M, Milanowski A, Miller J. Measuring and Promoting Inter-Rater Agreement of Teacher and Principal Performance Ratings. *Online Submiss* 2012.
- 21 Koo TK, Li MY. A Guideline of Selecting and Reporting Intraclass Correlation Coefficients for Reliability Research. *J Chiropr Med* 2016; **15**: 155–163.
- 22 Landis JR, Koch GG. The measurement of observer agreement for categorical data. *Biometrics* 1977; **33**: 159–174.
- 23 Mackel R. Properties of cutaneous afferents in diabetic neuropathy. *Brain* 1989; : 1359–1376.
- Rinkel WD, Rizopoulos D, Aziz MH, Van Neck JW, Cabezas MC, Coert JH.
 Grading the loss of sensation in diabetic patients: A psychometric evaluation of the rotterdam diabetic foot study test battery. *Muscle and Nerve* 2018; 58: 559–565.
- 25 Wasan AD, Alter BJ, Edwards RR, Argoff CE, Sehgal N, Walk D *et al.* Testretest and inter-examiner reliability of a novel bedside quantitative sensory testing battery in postherpetic neuralgia patients. *J Pain* 2020; **21**: 858–868.
- 26 Marx RG, Hudak PL, Bombardier C, Graham B, Goldsmith C, Wright JG. The reliability of physical examination for carpal tunnel syndrome. *J Hand Surg Eur Vol* 1998; **23**: 499–502.
- 27 Browne RH. On the use of a pilot sample for sample size determination. *Stat Med* 1995; **14**: 1933–1940.

Supplementary material



Supplementary Figure 1. Patient data flow diagram.

Abbreviations: n, number; RDF-39-C, continuous version of the 39-item Rotterdam Diabetic Foot Study test battery; RDF-39-D, dichotomous version of the 39-item Rotterdam Diabetic Foot Study test battery; RDF-13-D, dichotomous version of the 13-item Rotterdam Diabetic Foot Study test battery.

Supplementary Table 1. Interrater reliability of the individual test items of the RDF-39-C test battery, reported as interrater agreement and intraclass correlation coefficients. Interrater reliability of individual test items of RDF-39-C

| Rater 1 vs. rater 2 | | | | | | | |
|-----------------------|-----------|------------|------------|------------|--------------|--|--|
| (n = 11) | | | | | | | |
| | | | | | | | |
| | Total | Rater | Rater | Rater | | | |
| | agreement | difference | difference | difference | | | |
| Test | | ≤1 points | ≤2 points | ≤3 points | ICC (95% CI) | | |
| S2PD hallux | | | | | | | |
| Left foot | 72.7 | 81.8 | - | - | .61 (.0987) | | |
| Right foot | 72.7 | 81.8 | - | 90.9 | .66 (.1390) | | |
| S2PD medial heel | | | | | | | |
| Left foot | 100 | - | - | - | - | | |
| Right foot | 90.9 | - | 100 | - | - | | |
| S2PD first dorsal web | | | | | | | |
| Left foot | 72.7 | 90.9 | - | - | - | | |
| Right foot | 90.9 | - | 100 | - | .88 (.6497) | | |
| S2PD lateral foot | | | | | | | |
| Left foot | 100 | - | - | - | - | | |
| Right foot | 100 | - | - | - | - | | |
| S2PD fifth toe | | | | | | | |

| Left foot | 81.8 | - | 100 | - | .78 (.4094) |
|-----------------------|------|------|------|------|-----------------|
| Right foot | 81.8 | - | 90.9 | - | .67 (.1990) |
| M2PD hallux | | | | | |
| Left foot | 54.5 | - | 63.6 | 81.8 | .84 (.5196) |
| Right foot | 63.6 | - | - | 72.7 | .55 (0586) |
| M2PD medial heel | | | | | |
| Left foot | 90.9 | 100 | - | - | - |
| Right foot | 100 | - | - | - | - |
| M2PD first dorsal web | | | | | |
| Left foot | 81.8 | 90.9 | - | - | .47 (1082) |
| Right foot | 90.9 | 100 | - | - | .96 (.8799) |
| M2PD lateral foot | | | | | |
| Left foot | 81.8 | 90.9 | - | - | - |
| Right foot | 72.7 | 90.9 | 100 | - | .71 (.2091) |
| S1PD hallux | | | | | |
| Left foot | 36.4 | 63.6 | 90.9 | 100 | .96 (.8599) |
| Right foot | 36.4 | 72.7 | 90.9 | 100 | .91 (.7298) |
| S1PD medial heel | | | | | |
| Left foot | 36.4 | 63.6 | 90.9 | - | .68 (.1590) |
| Right foot | 18.2 | 45.5 | 81.8 | 100 | .27 (4374) |
| S1PD first dorsal web | | | | | |
| Left foot | 9.1 | 63.6 | - | 81.8 | .71 (.2491) |
| Right foot | 27.3 | 54.5 | - | 90.9 | .67 (.1890) |
| S1PD lateral foot | | | | | |
| Left foot | 27.3 | 63.6 | 90.9 | - | .72 (.2691) |
| Right foot | 36.4 | 72.7 | 90.9 | - | 1.00 (.98-1.00) |
| S1PD fifth toe | | | | | |
| Left foot | 27.3 | 72.7 | 90.9 | - | .95 (.8499) |
| Right foot | 27.3 | 72.7 | 90.9 | - | .57 (.0486) |
| Vibration sense MM | | | | | |
| Left foot | 45.5 | 90.9 | 100 | - | .94 (.8098) |
| Right foot | 54.5 | 100 | - | - | .96 (.8799) |
| Vibration sense IP | | | | | |
| Left foot | 45.5 | 100 | - | - | .96 (.8799) |
| Right foot | 63.6 | 72.7 | 100 | - | .92 (.6598) |
| Cold sensation | | | | | |
| Left foot | 100 | - | - | - | - |
| Right foot | 100 | - | - | - | - |
| Romberg's test | 100 | - | - | - | - |
| Numbness | 100 | - | - | - | - |
| Prior amputation | | | | | |
| Left foot | 100 | - | - | - | - |
| Right foot | 100 | - | - | - | - |
| Prior ulcer | 100 | - | - | - | - |

Abbreviations: RDF-39-C, continuous version of the 39-item Rotterdam Diabetic Foot Study test battery; n, number; IRA, interrater agreement; ICC, intraclass correlation coefficient; CI, confidence interval; S2PD, static two-point discrimination; M2PD, moving two-point discrimination; S1PD, static one-point discrimination; MM, medial malleolus; IP, interphalangeal joint.

Supplementary Table 2. Interrater reliability of the individual test items of the RDF-39-D test battery between rater 1 and rater 2, reported as interrater agreement and Cohen's kappa coefficient with its *p*-value.

| Interrater reliability of i | ndividual test items of RDF-39-I |) | |
|-----------------------------|----------------------------------|-----------|-----------------|
| Rater 1 vs. rater 2 | | | |
| (n = 30) | | | |
| Test | IRA (total agreement) (%) | Cohen's κ | <i>p</i> -value |
| S2PD hallux | | | |
| Left foot | 96.7 | - | - |
| Right foot | 93.3 | - | - |
| S2PD medial heel | | | |
| Left foot | 96.7 | - | - |
| Right foot | 100 | - | - |
| S2PD first dorsal web | | | |
| Left foot | 100 | - | - |
| Right foot | 96.7 | - | - |
| S2PD lateral foot | | | |
| Left foot | 100 | - | - |
| Right foot | 93.3 | - | - |
| S2PD fifth toe | | | |
| Left foot | 93.3 | 03 | .850 |
| Right foot | 96.7 | - | - |
| M2PD hallux | | | |
| Left foot | 83.3 | .51 | .005*** |
| Right foot | 80.0 | .14 | .414 |
| M2PD medial heel | | | |
| Left foot | 93.3 | - | - |
| Right foot | 96.7 | - | - |
| M2PD first dorsal web | | | |
| Left foot | 100 | - | - |
| Right foot | 86.7 | .27 | .114 |
| M2PD lateral foot | | | |
| Left foot | 93.3 | 03 | .850 |
| Right foot | 96.7 | - | - |
| S1PD hallux | | | |
| Left foot | 96.7 | .78 | <.001**** |
| Right foot | 96.7 | .78 | <.001**** |
| S1PD medial heel | | | |
| Left foot | 100 | - | - |
| Right foot | 100 | - | - |
| S1PD first dorsal web | | | |
| Left foot | 96.7 | .65 | < .001**** |
| Right foot | 100 | _ | _ |
| S1PD lateral foot | | | |
| Left foot | 100 | _ | _ |
| Right foot | 96.7 | _ | _ |
| S1PD fifth toe | | | |
| Left foot | 86.7 | .43 | .014* |
| Right foot | 83.3 | .44 | .014* |
| Vibration sense MM | | | ·· * |
| | | | |

| Left foot | 93.3 | .87 | <.001**** |
|--------------------|------|-----|------------|
| Right foot | 80.0 | .60 | .001**** |
| Vibration sense IP | | | |
| Left foot | 90.0 | .77 | < .001**** |
| Right foot | 90.0 | .78 | <.001**** |
| Cold sensation | | | |
| Left foot | 100 | - | - |
| Right foot | 100 | - | - |
| Romberg's test | 100 | - | - |
| Numbness | 100 | - | - |
| Prior amputation | | | |
| Left foot | 100 | - | - |
| Right foot | 100 | - | - |
| Prior ulcer | 100 | - | - |

Abbreviations: RDF-39-D, dichotomous version of the 39-item Rotterdam Diabetic Foot Study test battery; n, number; IRA, interrater agreement; S2PD, static two-point discrimination; M2PD, moving two-point discrimination; S1PD, static one-point discrimination; MM, medial malleolus; IP, interphalangeal joint. * $p \le .05$; *** $p \le .005$; **** $p \le .001$ **Supplementary Table 3.** Interrater reliability of the individual test items of the RDF-13-D test battery between rater 1 and rater 2 and between rater 2 and rater 3, reported as interrater agreement and Cohen's kappa coefficient with its *p*-value.

| | Rater 1 vs. rater 2 (n = 43) | | | Rater 2 vs. rater 3 (n = 14) | | |
|--------------|---|-----------|-----------------|---------------------------------|-----------|-----------------|
| Tost | IRA (total agreement) | Cohen's ĸ | <i>p</i> -value | IRA (total agreement) | Cohen's ĸ | <i>p</i> -value |
| S2PD hallury | (70) | | | (/0) | | |
| J off foot | 00.7 | | | 02.0 | 62 | 011* |
| S2RD modial | 90.7 | - | - | 92.9 | .03 | .011 |
| 52PD mediai | | | | | | |
| Right foot | 95.3 | - 02 | 876 | 85.7 | 11 | 047* |
| S2PD first | 75.5 | 02 | .070 | 00.7 | .++ | .047 |
| dorsal web | | | | | | |
| Right foot | 95.3 | - 02 | 876 | 64.3 | _ | |
| M2PD first | 75.5 | 02 | .070 | 04.0 | | |
| dorsal web | | | | | | |
| Left foot | 95.3 | 48 | < 001**** | 71 4 | 18 | 469 |
| S1PD hallux | 70.0 | .10 | *.001 | / 1.1 | .10 | .107 |
| Left foot | 93.0 | 76 | < 001**** | 92.9 | 63 | 011* |
| Right foot | 93.0 | .73 | <.001**** | 100 | - | - |
| Vibration | ,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,, | | | 100 | | |
| sense MM | | | | | | |
| Left foot | 86.0 | .72 | < .001**** | 85.7 | .71 | .005*** |
| Right foot | 83.7 | .67 | <.001**** | 85.7 | .71 | .005*** |
| Vibration | | | | | | |
| sense IP | | | | | | |
| Left foot | 90.7 | .79 | <.001**** | 78.6 | .55 | .036* |
| Right foot | 81.4 | .61 | <.001**** | 100 | - | - |
| Prior | | | | | | |
| amputation | | | | | | |
| Left foot | 100 | - | - | 100 | - | - |
| Right foot | 100 | _ | _ | 100 | _ | _ |
| Prior ulcer | 100 | - | - | 100 | - | _ |

Interrater reliability of individual test items of RDF-13-D

Abbreviations: RDF-13-D, dichotomous version of the 13-item Rotterdam Diabetic Foot Study test battery; n, number; IRA, interrater agreement; S2PD, static two-point discrimination; M2PD, moving two-point discrimination; S1PD, static one-point discrimination; MM, medial malleolus; IP, interphalangeal joint. * $p \le .05$; *** $p \le .005$; **** $p \le .001$

| | Scenario 1: high | n IRA, high ICC | Scenario 2: low IRA, high ICC | |
|---------|------------------|-----------------|-------------------------------|---------|
| Patient | Rater X | Rater Y | Rater X | Rater Y |
| А | 1 | 1 | 1 | 3 |
| В | 2 | 2 | 1 | 3 |
| С | 3 | 3 | 3 | 5 |
| D | 4 | 4 | 3 | 5 |
| Е | 5 | 5 | 5 | 7 |
| F | 6 | 6 | 5 | 7 |

Supplementary Table 4. Hypothetical example of different levels of interrater agreement and intraclass correlation coefficient.

Abbreviations: IRA, interrater agreement; ICC, intraclass correlation coefficient.