



[biblio.ugent.be](http://biblio.ugent.be)

The UGent Institutional Repository is the electronic archiving and dissemination platform for all UGent research publications. Ghent University has implemented a mandate stipulating that all academic publications of UGent researchers should be deposited and archived in this repository. Except for items where current copyright restrictions apply, these papers are available in Open Access.

This item is the archived peer-reviewed author-version of:

Title : Standardization of Free Thyroxine Measurements Allows the Adoption of a More Uniform Reference Interval

Authors: De Grande, Linde A. C., Katleen Van Uytfanghe, Dries Reynders, Barnali Das, James D. Faix, Finlay MacKenzie, Brigitte Decallonne, et al.

In: *Clinical Chemistry* 63 (10): 1642–1652, 2017

**To refer to or to cite this work, please use the citation to the published version:**

De Grande, Linde A. C., Katleen Van Uytfanghe, Dries Reynders, Barnali Das, James D. Faix, Finlay MacKenzie, Brigitte Decallonne, et al. 2017. "Standardization of Free Thyroxine Measurements Allows the Adoption of a More Uniform Reference Interval." *Clinical Chemistry* 63 (10): 1642–1652  
DOI10.1373/clinchem.2017.274407

## **Standardization of free thyroxine measurements allows the adoption of a more uniform reference interval**

**Running head:** Standardization of serum FT4 measurements

Linde A.C. De Grande<sup>1</sup>, Katleen Van Uytfanghe<sup>2</sup>, Dries Reynders<sup>3</sup>, Barnali Das<sup>4</sup>, James D. Faix<sup>5</sup>, Finlay MacKenzie<sup>6</sup>, Brigitte Decallonne<sup>7</sup>, Akira Hishinuma<sup>8</sup>, Bruno Lapauw<sup>9</sup>, Paul Taelman<sup>10</sup>, Paul Van Crombrugge<sup>11</sup>, Annick Van den Bruel<sup>12</sup>, Brigitte Velkeniers<sup>13</sup>, Paul Williams<sup>14</sup>, Linda M. Thienpont<sup>1,15\*</sup> on behalf for the IFCC Committee for Standardization of Thyroid Function Tests (C-STFT).

<sup>1</sup>Department of Pharmaceutical Analysis, Faculty of Pharmaceutical Sciences, Ghent University, Ghent, Belgium.

<sup>2</sup>Ref4U, Laboratory of Toxicology, Faculty of Pharmaceutical Sciences, Ghent University, Ghent, Belgium.

<sup>3</sup>Department of Applied Mathematics, Computer Science and Statistics, Faculty of Sciences, Ghent University

<sup>4</sup>Biochemistry and Immunology Laboratory, Kokilaben Dhirubhai Ambani Hospital and Medical Research Institute, Mumbai, India.

<sup>5</sup>Clinical Chemistry and Immunology, Montefiore Medical Center, and Department of Pathology, Albert Einstein School of Medicine, New York, NY, USA.

<sup>6</sup>Birmingham Quality/UK NEQAS, University Hospitals Birmingham NHS Foundation Trust, Birmingham, UK.

<sup>7</sup>Department of Endocrinology, University Hospitals Leuven, Leuven, Belgium.

<sup>8</sup>Department of Infection Control and Clinical Laboratory Medicine, Dokkyo Medical University, Tochigi, Japan.

<sup>9</sup>Department of Endocrinology, Ghent University Hospital, Ghent, Belgium.

<sup>10</sup>Laboratory of Endocrinology, Department of Laboratory Medicine, AZ Maria-Middelares Sint-Jozef, Campus Maria-Middelares, Ghent, Belgium.

<sup>11</sup>Department of Endocrinology, OLV Ziekenhuis Aalst-Asse-Ninove, Aalst, Belgium.

<sup>12</sup>Department of Endocrinology, General Hospital Sint Jan, Bruges, Belgium.

<sup>13</sup>Department of Endocrinology, Universitair Ziekenhuis Brussel, Brussels, Belgium.

<sup>14</sup>Department of Endocrinology, Royal Prince Alfred Hospital, Camperdown, Australia.

<sup>15</sup>Linda M. Thienpont is now at Thienpont & Stöckl Wissenschaftliches Consulting GbR, Rennertshofen (OT Bertoldsheim), Germany

\*Corresponding author: (current affiliation) Thienpont & Stöckl Wissenschaftliches Consulting GbR, Erlbacher Strasse 11, Rennertshofen (OT Bertoldsheim), Germany.

Tel. +49 8434 94 365 22, email: linda.thienpont@ugent.be

### **Key words**

Thyroid hormones; SI-traceability; metrological traceability; method comparison; reference interval; proof-of-concept.

**Words: 3474** (max 3500)

**Figures: 4**

**Tables: 1**

### **Non-standard abbreviations**

TSH, Thyroid stimulating hormone; FT4, free thyroxine; RIs, reference intervals; C-STFT, Committee for Standardization of Thyroid Function Tests; IVD, in vitro diagnostic; SI-units, Système International d'Unités/International System of Units; RMP, reference measurement procedure; ED-ID-LC-MS/MS, equilibrium dialysis-isotope dilution/liquid chromatography/tandem mass spectrometry; MC, method comparison; CI, confidence interval; TE, total error; A-D, Anderson-Darling test.

## **Abstract**

**Background:** The IFCC Committee for Standardization of Thyroid Function Tests intended to standardize free thyroxine (FT4) immunoassays. We developed a Système International d'Unités traceable conventional reference measurement procedure (RMP) based on equilibrium dialysis and mass spectrometry. We describe here the latest studies intended to recalibrate against the RMP and supply a proof-of-concept, which should allow continued standardization efforts.

**Methods:** We used the RMP to target the standardization and reference interval (RI) panels, which were also measured by 13 manufacturers. We validated the suitability of the recalibrated results to meet specifications for bias (3.3%) and total error (8.0%) determined from biological variation. However, since these specifications were very stringent, we expanded them to 10% and 13%, respectively. The results for the RI panel were reported as if the assays were recalibrated. We estimated all but one RI using parametric statistical procedures and hypothesized that the RI determined by the RMP was suitable for use by the recalibrated assays.

**Results:** Twelve of 13 recalibrated assays had a bias meeting the 10% specification with 95% confidence; for 7 assays this applied even for the 3.3% specification. Only 1 assay met the 13% total error specification. Recalibration reduced the CV of the assay means for the standardization panel from 13% to 5%. The proof-of-concept study confirmed our hypothesis regarding the RI but within constraints.

**Conclusion:** Recalibration to the RMP significantly reduced the FT4 immunoassay bias, so that the RI determined by the RMP was suitable for common use within a margin of 12.5%.

## Introduction

The diagnosis of metabolic thyroid disorders and/or monitoring of treatment is based on laboratory testing of serum thyroid-stimulating hormone (TSH) and free thyroxine (FT4). Provided the hypothalamic-pituitary-thyroid axis is intact, a first line TSH result may suggest a number of thyroid disorders that could be clarified by follow-up measurement of FT4; however, immediate combined measurement is indicated for the differential diagnosis between mild (subclinical) primary hyperthyroidism, and secondary (central) hypothyroidism. Furthermore combined measurement is warranted during the first days/weeks of the follow-up of patients with severe thyroid dysfunction, where TSH has not yet return to a euthyroid baseline concentration and thus not representative of the actual thyroid functional status (e.g., in patients with autoimmune Graves' disease and high titers of TSH receptor antibodies or with increased human chorionic gonadotropin concentrations). On the other hand, FT4 is the primary test for the titration of levothyroxine replacement in patients with central hypothyroidism and/or with high risk differentiated thyroid cancer with need for a suppressed TSH (1-5). For maximum effectiveness, current FT4 immunoassays would benefit from improved clinical and analytical consistency (6, 7). Additionally, the issue of substantial inter-method variability needs to be resolved for improved everyday patient care because it requires interpretation of laboratory results against assay-specific reference intervals (RIs) and prevents incorporation of common decision levels in evidence-based practice guidelines (7, 8). Therefore, the IFCC Committee for Standardization of Thyroid Function Testing was commissioned to standardize FT4 measurements globally (9). The committee's efforts have been endorsed by the clinical community, which also called for general standardization of hormonal assays in the 21<sup>st</sup> century (10).

The committee conducted the standardization activities of FT4 measurements in partnership with the same *in vitro* diagnostic (IVD) manufacturers (with one exception) that had been involved in the TSH harmonization (11). The committee pursued a process similar to that used for the TSH assays, except for FT4 for they developed and used a reference

measurement system with traceability to the Système International d'Unités (SI) (12, 13). The committee defined the measurand and developed/validated a conventional reference measurement procedure (RMP) based on equilibrium dialysis combined with isotope dilution-liquid chromatography-tandem mass spectrometry (ED-ID-LC-MS/MS) (14-16), and undertook several method comparisons (MCs) with single-donation and commutable serum samples (Phase I to III studies) according to the "step-up" approach (8, 17-19). Each of the studies had a different focus, including documentation of the assays' intrinsic quality and demonstration of the feasibility of standardization of assay results by recalibrating the immunoassays to the RMP.

Here we report, on behalf of the Committee for Standardization of Thyroid Function Testing, our latest activities in the standardization process. We performed a Phase IV MC study between 13 immunoassays and the RMP. There were two objectives. First, to establish calibration traceability of the participating assays to the SI-traceable RMP. Second, to validate the efficiency of the process to eliminate the assay-specific biases. Subsequently, we conducted a RI study with a new panel of samples to test the proof-of-concept that, after standardization, immunoassays might accord sufficiently with the RMP to enable adoption of a common RI for diagnosis and follow-up of patients with thyroid dysfunction.

## **Material and methods**

### *Panels of clinical samples and value assignment*

We collected standardization and RI panels. The standardization panel comprised 91 clinically relevant samples and was intended to facilitate the calibration adjustment/readjustment by the manufacturers to the IFCC RMP. The aim of the RI panel was to let manufacturers evaluate their recalibration, for which we used 120 samples donated by apparently healthy American volunteers. The sources, eligibility and exclusion criteria, conditions for sampling, processing and storage were those described before for the TSH harmonization effort (11). Approval from a Bioethic Committee and informed consent from the patients was obtained along with a short description of the clinical background of the donating patients. The target values (mean of minimum 3 independent measurements) were assigned with the IFCC conventional RMP performed at the reference laboratory of Ghent University. Both are listed in the Database of the Joint Committee for Traceability in Laboratory Medicine (20).

### *Study participants and measurement protocol*

Thirteen IVD manufacturers participated in the current studies, each with one assay (coding and further details on the platforms/assays in Table 1). We requested that the IVD manufacturers perform all measurements according to a proposed randomized sequence, in singleton on each of two days, and include their master calibrators for measurement in parallel with the panel samples. The individual results were reported. The samples for the RI study were measured in order of their ascending ID number, in singleton and within a single run. Of note, that the organization and interpretation of internal QC was left to the discretion of each manufacturer.

### *Recalibration of immunoassays*

After submitting the results for measurement of the standardization panel with the assays' current calibrators, the IVD manufacturers received from us a preliminary validation report, comprising the target concentrations determined by the RMP. These were intended for use in value reassignment of the master calibrators. The manufacturers were entitled to use their in-house mathematical procedure to determine the relationship of their assay results to those from the RMP (11). After the readjustment of the master calibrators, the manufacturers recalculated and reported back the results for the standardization panel as if they were obtained with the recalibrated assays. The results for measurement of the RI panel were similarly reported after transformation to the revised calibration.

#### *Data treatment*

For consolidation of the MC study data we used Microsoft EXCEL<sup>®</sup> 2010. We concentrated on demonstrating and validating the efficiency of the recalibration process. We calculated for each assay *i*) the pre- and post-recalibration median deviation (%) to the RMP in several FT4 concentration intervals, *ii*) the mean deviation (%) or bias (and one-sided 95% confidence interval (CI)) after recalibration, *iii*) the total error (TE, %) from the first replicate after recalibration, and *iv*) the differences between the replicates (in % of the mean). We also compared the pre- and post-recalibration CVs (%) of the assay means.

We used CBstat (version 5.1) for statistical evaluation of the data from the RI study. This software evaluated normality of data distributions by the Anderson-Darling (A-D) test ( $P \geq 0.05$ ), did outlier testing on the basis of power-transformed values (limit 4SD), and supplied parametric (direct on the original data and/or after transformation) as well as non-parametric procedures to estimate the RI characteristics. For the normally distributed datasets we used the direct parametric procedure [RI estimated as  $\text{mean} \pm 1.96(1/(1-1/(4(n-1)))) * \text{SD}$ ]. For those datasets for which normality did not apply, we selected the procedure after a sequence of investigations, i.e., in addition to the detection of statistical outliers, we did a visual screening for aberrant differences (%) to the RMP targets. If after omission of the detected values the A-D test allowed acceptance of the hypothesis of normally distributed

data, we again selected the direct parametric procedure; if not, we verified the data for normality after log-transformation. If the A-D  $P$ -value was then  $\geq 0.05$ , we applied the parametric procedure. Finally, one dataset remained, which was submitted to the non-parametric bootstrap (500 replicates) procedure, to generate bootstrap estimates of the  $(2.5/100)N+0.5$  and  $(97.5/100)N+0.5$  ordered values (22). To test the hypothesis that after recalibration a common RI could be used by all manufacturers, we first investigated whether the probabilities that the 2.5 and 97.5 percentiles (further also referred to as lower and upper limit, respectively), estimated from the datasets of the immunoassays were located within the 90% CI from the RMP data percentiles (further referred to as reference percentiles), were reasonably large ( $>90\%$ ). We repeated the probability testing while using limits of 12.5% around the reference percentiles. Probability estimations were done in R 3.2.3 for all assays but assay K (Table 1), where the CIs were determined by CB-stat; for the latter, we used the R statistical software to perform a bootstrap procedure on the original RI dataset to simulate the distribution of the percentiles.

#### *Analytical specifications*

We demonstrated/validated the suitability of the recalibrated results to meet desirable specifications for bias and TE based on the biological variation, i.e., 3.3% and 8.0%, respectively (23). However, because of the extreme stringency of these values, we also used the empirical bias limit of 10% that was considered state-of-the-art in previous MC studies, and expanded the TE specification to 13% to account for any imprecision of the RMP (8, 16, 18, 19). The 12.5% limit used for testing the RI hypothesis was based on the state-of-the-art bias specification used above but would additionally account for the uncertainty of the location of the reference percentiles.

#### *Homogeneity and stability study*

We assessed the homogeneity and stability of the FT4 standardization panel in the same way as described for TSH (11).

## Results

### *Concentration range covered by the panels of clinical samples*

The FT4 standardization panel covered a concentration range from 4.5 pmol/L to 164 pmol/L (determined by the RMP). The expanded uncertainty of the targets (coverage factor  $k = 2$ ) was estimated to be on the order of 7.0% (16). The central 95% of the RI panel covered the range from 13.5 pmol/L ( $\pm 0.7$  pmol/L; 90% CI) to 24.3 pmol/L ( $\pm 0.7$  pmol/L) with the mean at 18.9 pmol/L.

### *Validation of the efficiency of recalibration*

The combined difference plots (Fig. 1) reflect the assays' calibration biases to the RMP before (Fig. 1A) and after recalibration (Fig. 1B). The effect of recalibration on the assay-specific median deviations (%) to the RMP targets in 4 concentration intervals is shown in Fig. 2A by a combined picture with indication of the 15<sup>th</sup>, 50<sup>th</sup> and 85<sup>th</sup> centiles, and in Fig. 2B by the individual deviations (more details in Supplemental Table 1). Before recalibration, deviations were negative across the FT4 measurement range for all but assay N (<10 pmol/L). Moreover, the deviations increased with increasing concentration. The highest median manufacturer deviations were -40.8% (assay J) (<10 pmol/L), -37.9% (assay F) ( $\geq 10$  and <25 pmol/L), -57.7% (assay B) ( $\geq 25$  and <100 pmol/L), and -72.7% (assay B) ( $\geq 100$  pmol/L). The lowest median manufacturer deviations were 7.4% (assay N), -13.7% (assay N), -25.6% (assay O) and -30.2% (assay G), respectively. Hence, the most discrepant assay pairs (assays J/N, F/N, B/O and B/G) deviated by 48.2%, 24.2%, 32.1% and 42.5%, respectively. After recalibration the ranges of the median deviations became -12.0% (assay O) to +8.2% (assay A) (<10 pmol/L), -8.9% (assay O) to +1.7% (assay H) ( $\geq 10$  and <25 pmol/L), -8.4% (assay H) to +9.5% (assay F) ( $\geq 25$  and <100 pmol/L), and -12.5% (assay O) to +11.9% (assay G) ( $\geq 100$  pmol/L), respectively. Fig. 3 shows the post-recalibration differences (%) and the assay biases (%) reflected against the used specifications. From the numbers in Supplemental Table 2, we can confidently assert that after recalibration the bias

(and one-sided 95% CI) of all assays but assay O complied with the empirical specification of 10% at a 95% probability; the bias of 7 assays (A, B, D, E, I, J and N) complied when assessed against the 3.3% specification (Supplemental Table 2) (24). With regard to the assays' TE after recalibration, only assay I met the expanded specification, i.e., had 95% of its differences within 13%, while for the other assays 8% to 35% of the differences violated it (Supplemental Fig. 1). The median differences between the replicates from 2 runs ranged from -1.5% (assay K) to 4.1% (assay F), and the  $SD_{diff}$  from 2.5% (assay H) to 5.9% (assay A) (Supplemental Table 3). Supplemental Fig. 2 shows that for several assays the differences (%) between replicates were concentration-dependent. After recalibration, the CV of the assay means (the latter calculated for each assay from all results) decreased from 13% to 5%.

#### *Reference interval study*

The RI characteristics from the ED-ID-LC-MS/MS measurements were obtained with the direct parametric procedure. This procedure was also used for the other normally distributed datasets, which excluded the assays A, G, H and K. In spite of a negative outlier test in CBstat for these 4 datasets, visual inspection of the plots of assays G and H (Supplemental Fig. 3) revealed aberrant differences (%) to the RMP targets (4 for assay G and 3 for assay H, respectively). After omission of these aberrant data, the A-D *P*-values became >0.26 and >0.25, respectively, which justified application of the direct parametric procedure to these assays. For the assay A, the hypothesis of normality was accepted after log-transformation of the data, again justifying the use of the parametric procedure; only for assay K did we have to use a non-parametric bootstrap procedure. Supplemental Table 4 lists the main characteristics of the respective RIs. The widths of the RIs by the immunoassays ranged from 9.4 pmol/L to 12.0 pmol/L versus 10.7 pmol/L for the RMP. The CIs for the respective percentiles ranged from 1.1 pmol/L and 2.4 pmol/L (at the 2.5 percentile) and 1.2 pmol/L to 2.4 pmol/L (at the 97.5 percentile) versus 1.4 pmol/L (for both percentiles of the RMP). The range of the means/medians of the RIs was from 17.2/17.0 pmol/L to 20.8/20.5 pmol/L

versus 18.9/18.8 pmol/L for the RMP. Supplemental Tables 5 and 6, plus Supplemental Figs. 4 and 5, demonstrate that none of the calculated probabilities for the assays met the minimum requirement of >90%. However, after expanding the reference percentile intervals to 12.5%, they did for assays E, F, G, H, I, J, L and N at the 2.5 percentile. For the 97.5 percentile, the >90% requirement was achieved by all but assay A. The graphical overview of the respective RIs (Fig. 4) shows that assays A and B had the most discrepant 2.5 percentiles (calculated to the mean of both percentile values, they were 28% apart), while this was the case for assays A and F for the 97.5 percentiles (21% apart).

#### *Homogeneity study*

Statistical testing confirmed that the hypothesis of homogeneity of the aliquots in the standardization panel ( $P > 0.05$ , Supplemental Table 7) could be accepted. The stability study is still ongoing.

## Discussion

The approach to the standardization of commercial FT4 immunoassays was similar to that previously described for TSH (11). The Phase I MC demonstrated that mathematical recalibration of measurement results for samples from presumably healthy volunteers was able to align the different immunoassays to the RMP. The Phase II and III MCs extended the findings for euthyroid individuals to patients with hypo- and hyperthyroidism, and provided proof-of-concept that manufacturers were also able to do the recalibration by adjusting their calibrators (8, 17-19). The current Phase IV MC was the natural next step in our standardization project, and the RI study was intended to assess whether recalibration would allow a uniform basis for the use of common RIs. The strengths of the FT4 standardization approach were the involvement during several years of the globally operating IVD industry and the use of a panel of commutable samples, collected to mimic clinical conditions. The concentrations of the samples spanned the measurement range of current assays, because they were sourced from euthyroid individuals and also from patients with overt hypothyroidism and hyperthyroidism.

The current study confirmed that establishing calibration traceability to the RMP significantly reduced the negative biases of the immunoassays, as well as the CV of the assay means. However, it is also important to appreciate the huge impact that standardization could have on future measurement results and reference intervals. After recalibration, 12 of 13 immunoassays had their bias (and CI) meet the empirical specification of 10% at a 95% probability, while 7 of them even passed the very stringent specification of 3.3% derived from the biological variation. Although this outcome is overall reasonable, it also points to the fact that the recalibration effectiveness was better for some assays than for others.

The fact that the standardization panel comprised sufficient native samples enabled us also to focus on the validation of the post-recalibration TE. This is a very important performance attribute because it reflects the accuracy of an assay for measurement of the

individual sample. Most assays violated the expanded TE limits in spite of reasonable recalibration. This might be due to the specification being too stringent, even after expansion. However, considering that in the previous MC studies we already highlighted the TE issue of many FT4 immunoassays due to their susceptibility to sample-related effects, it is more realistic to suggest that our current study confirms this limitation.

Finally, the results on the differences between replicates highlight the occasional high inter-run imprecision and lack of robustness of calibration (Supplemental Table 3, Supplemental Fig. 2). The importance of continual improvement of these performance attributes across all assays was discussed with the IVD manufacturers.

The aim of the current RI study was primarily to supply a proof-of-concept that after recalibration the use of a common RI may be feasible. We used the RI estimated from the measurement data by ED-ID-LC-MS/MS as reference and assessed whether the recalibrated assays could share it. We inferred the percentiles and mean of the central 95% of all but one RI by a parametric procedure applied to either the original or log-transformed data. Interestingly, the width of the interval by the RMP corresponded reasonably with that calculated from the FT4 biological variation, i.e., 10.7 pmol/L versus 9.6 pmol/L, as well as that estimated in another study using ED-ID-LC-MS/MS, i.e., 12.1 pmol/L (23, 25). However, it was most important to compare the derived immunoassay percentiles of the RIs with those of the RMP. In the employed statistical approach an immunoassay would be qualified to share the RI of the RMP if the probability that its percentiles were located within the CI around the reference was higher than 90%. None of the assays met this criterion. However, when an interval of  $\pm 12.5\%$  was adopted, the probabilities of eight assays met the  $>90\%$  requirement at the 2.5 percentiles, and of all but one assay also at the 97.5 percentiles. We present 3 reasons to justify the hypothesis of testing with the 12.5% margin around the reference percentiles. First is the observation that the magnitudes of the CIs around the reference percentiles were  $\leq 5\%$ , thus similar to or narrower than the assays' effective biases in the euthyroid range after recalibration (range 0% up to 9%). Second, we refer to the impact of the lot-to-lot variation on the RI study, which was performed with a time offset of at

least 6 months compared to the Phase IV MC. Third, we found it legitimate to account to a certain extent for the uncertainty of the location of the estimated reference percentiles due to the potential impact of an undetectable bias in the measurements with the RMP.

Nevertheless, even if the current margin of 12.5% accommodates the current state-of-the-art measurements, we advocate that in the future it should be decreased, particularly because of the low biological variation of serum FT4. We also recommended the IVD manufacturers of the assays that did not accord with the RMP to share its percentiles, in spite of adopting the 12.5% margin, to do root cause analysis.

In conclusion, the Phase IV MC study described here showed that, in general, the recalibration process was able to eliminate the considerable FT4 calibration biases to the RMP. In addition, the basic RI study provided the proof-of-concept since the percentiles of the RMP applied for most of the recalibrated assays within a margin of 12.5%. Although this result represents substantial progress in standardization of FT4 measurements, we recognize that it cannot be extrapolated to all clinical situations where FT4 testing is indicated, particularly when binding proteins are abnormal. Therefore, to better understand more-subtle assay differences in other patient cohorts such as pregnant females and patients with the non-thyroidal illness syndrome, we recommend that our approach serve as model for future studies. We also see surveillance of the sustainability of the recalibration basis as a final key component of our standardization approach. We propose that, after implementation of the recalibrated assays, the surveillance should be done under field conditions to account for the impact of variables like lot-to-lot changes and instrument instability. This could be done by using the Percentiler/Flagger applications described elsewhere as useful tools for continuous monitoring of the stability of performance/flagging frequency in laboratories grouped according to instrument/assay-specific peers (26). Another tool could be the organization of proficiency testing or external quality assessment surveys with commutable samples (27). We also recognize that we should expand the measurement capacity with the conventional RMP. Therefore, we are currently working on establishing a network of competent reference laboratories. Last but not least, from the perspective that

implementing the recalibrated FT4 assays will have a huge impact on future measurement results and RIs, we are committed to gaining broad consensus on this step (28).

**Acknowledgments:** The Chair of the IFCC C-STFT (L.M. Thienpont) is grateful to (companies in alphabetical order; individual names also): D. Flanagan J. Reid and S. Ruetten (Abbott Diagnostics, USA); ; A. Adelman and J. Sackrison (Beckman Coulter Inc., USA); J.-M. Barbeaud (bioMérieux SA, France); I. Kutschera and G. Markowitz (DiaSorin S.p.A, Italy); K. Aoyagi, C. Hall and T. Niwa (Fujirebio Inc., Japan); S. Tashiro, and T. Ono (LSI Medience Corporation, Japan); P. Hosimer, M.-M. Patru and C. Thomas (Ortho-Clinical Diagnostics, UK); A. Hoppe and M. Rottmann (Roche Diagnostics GmbH, Germany); ZD. Chen, H.Xu, JY.Yuan and W.Li (Snibe Co.,Ltd, China); Y. Tao, L. Wan Ju and Y. De Qian (Sichuan Maccura Biotechnology Co., Ltd., China); R. Janzen, P. Sibley, R. Payne and V. Bitcon (Siemens Healthineers, USA); T. Sakata, M. Yamasaki, T. Kagawa, and K. Kishi (Sysmex Corporation, Japan); M. Kasai, S. Marivoet, S. Narayanan, H. Tsukamoto and M. Tsuura (TOSOH Corp., Japan; acting as representative on behalf of their organizations. Their efforts to review and provide comments on the manuscript are highly appreciated. The 14 organizations sponsored (all contributed equally) the study in terms of sample procurement and funding for the assignment with the reference method values.

**Role of Sponsor:** The funding organizations played no role in the study design, choice of enrolled patients, review and interpretation of data, or preparation or approval of manuscript.

## References (max 40)

1. Biondi B, Bartalena L, Cooper DS, Hegedüs L, Laurberg P, Kahaly GJ. The 2015 European Thyroid Association Guidelines on Diagnosis and Treatment of Endogenous Subclinical Hyperthyroidism. *Eur Thyroid J* 2015;4:149–63.
2. Ross DS, Burch HB, Cooper DS, Greenlee MC, Laurberg P, Maia AL, et al. 2016 American Thyroid Association Guidelines for Diagnosis and Management of Hyperthyroidism and other causes of Thyrotoxicosis. *Thyroid*. 2016;26:1343-1421.
3. Garber JR, Cobin RH, Gharib H, Hennessey JV, Klein I, Mechanick JI, et al. American Association of Clinical Endocrinologists and American Thyroid Association Taskforce on Hypothyroidism in Adults. Clinical practice guidelines for hypothyroidism in adults: co-sponsored by American association of clinical endocrinologists and the American thyroid association. *Endocrine Practice* 2012;6:988–1028.
4. Koulouri O, Auldin MA, Agarwal R, Kieffer V, Robertson C, Falconer Smith J, et al. Diagnosis and treatment of hypothyroidism in TSH deficiency compared to primary thyroid disease: pituitary patients are at risk of under-replacement with levothyroxine. *Clin Endocrinol (Oxf)* 2011;74:744-9
5. Demers LM, Spencer CA. The National Academy of Clinical Biochemistry Presents Laboratory medicine practice guidelines. Laboratory support for the diagnosis of thyroid disease. 13/2002.
6. Thyroid Disease Manager. Guidelines for diagnosis and management of thyroid disease. <http://www.thyroidmanager.org/> (accessed March 2017).
7. Thienpont LM, Van Uytvanghe K, Poppe K, Velkeniers B. Determination of free thyroid hormones. *Best Pract Res Clin Endocrinol Metab* 2013;27:689–700.
8. Thienpont LM, Van Uytvanghe K, Beastall G, Faix JD, Ieiri T, Miller WG, et al.; for the IFCC Working Group on Standardization of Thyroid Function Tests. Report of the IFCC Working Group for Standardization of Thyroid Function Tests, part 2: Free thyroxine and free triiodothyronine. *Clin Chem* 2010;56:912-20.

9. Committee for Standardization of Thyroid Function Tests (C-STFT). IFCC - Scientific Division (SD). SD Committees. <http://www.ifcc.org/ifcc-scientific-division/sd-committees/c-stft/> (accessed March 2017).
10. Wartofsky L, Handelsman DJ. Standardization of hormonal assays for the 21<sup>st</sup> century. *J Clin Endocrinol Metab* 2010;95:5141-3.
11. Thienpont LM, Van Uytfanghe K, De Grande LAC, Reynders D, Das B, Faix JD, MacKenzie F, Decallonne B, Hishinuma A, Lapauw B, Taelman P, Van Crombrugge P, Van den Bruel A, Velkeniers B, Williams P on behalf for the IFCC Committee for Standardization of Thyroid Function Tests (C-STFT). Harmonization of serum thyroid-stimulating hormone measurements paves the way for the adoption of a common reference interval – A technical report. *Clin Chem* 2017;63;Epub ahead of print.
12. Vesper HW, Thienpont LM. Traceability in laboratory medicine. *Clin Chem* 2009;55:1067-75.
13. ISO 17511 International Organization for Standardization (ISO). In vitro diagnostic medical devices—measurement of quantities in biological samples —metrological traceability of values assigned to calibrators and control materials. ISO 17511:2003. Geneva: ISO; 2003.
14. International Federation of Clinical Chemistry and Laboratory Medicine (IFCC), IFCC Scientific Division Working Group for Standardization of Thyroid Function Tests (WG-STFT), Thienpont LM, Beastall G, Christofides ND, Faix JD, Ieiri T, Miller WG, et al. Measurement of free thyroxine in laboratory medicine – proposal of measurand definition. *Clin Chem Lab Med* 2007;45:563-4.
15. International Federation of Clinical Chemistry and Laboratory Medicine IFCC, IFCC Scientific Division Working Group for Standardization of Thyroid Function Tests (WG-STFT), Thienpont LM, Beastall G, Christofides ND, Faix JD, Ieiri T, Jarrige V, et al. Proposal of a candidate international conventional reference measurement procedure for free thyroxine in serum. *Clin Chem Lab Med* 2007;45:934-6.
16. International Federation of Clinical Chemistry; Laboratory Medicine Working Group for Standardization of Thyroid Function Tests. Van Houcke SK, Van Uytfanghe K, Shimizu E,

Tani W, Umemoto M, Thienpont LM. IFCC international conventional reference procedure for the measurement of free thyroxine in serum. International Federation of Clinical Chemistry and Laboratory Medicine (IFCC) Working Group for Standardization of Thyroid Function Tests (WG-STFT). Clin Chem Lab Med 2011;49:1275-81.

17. Van Uytfanghe K, De Grande LA, Thienpont LM. A "Step-Up" approach for harmonization. Clin Chim Acta 2014;432:62-7.

18. Thienpont LM, Van Uytfanghe K, Van Houcke S. IFCC Working Group for Standardization of Thyroid Function Tests (WG-STFT). Standardization activities in the field of thyroid function tests: a status report. Clin Chem Lab Med 2010;48:1577-83.

19. Thienpont LM, Van Uytfanghe K, Van Houcke S, Das B, Faix JD, MacKenzie F, et al. IFCC Committee for Standardization of Thyroid Function Tests (C-STFT). A Progress report of the IFCC Committee for Standardization of Thyroid Function Tests. Eur Thyroid J 2014;3:109-16.

20. JCTLM Database of higher-order reference materials, measurement methods/procedures and services. <http://www.bipm.org/jctlm/> (accessed March 2017).

21. CLSI. Evaluation of Detection Capability for Clinical Laboratory Measurement Procedures; Approved Guideline—Second Edition. CLSI document EP17-A2. Wayne, PA: Clinical and Laboratory Standards Institute; 2012.

22. Linnet K. Nonparametric estimation of reference intervals by simple and bootstrap-based procedures. Clin Chem 2000;46:867-9.

23. Westgard QC. Desirable biological variation database specifications. Desirable specifications for total error, imprecision, and bias, derived from intra- and inter-individual biologic variation. <https://www.westgard.com/biodatabase1.htm> (Accessed March 2017).

24. Stöckl D, Rodríguez Cabaleiro D, Van Uytfanghe K, Thienpont LM. Interpreting method comparison studies by use of the bland-altman plot: reflecting the importance of sample size by incorporating confidence limits and predefined error limits in the graphic. Clin Chem 2004;50:2216-8.

25. Yue B, Rockwood AL, Sandrock T, La'ulu SL, Kushnir MM, Meikle AW. Free thyroid hormones in serum by direct equilibrium dialysis and online solid-phase extraction-liquid chromatography/tandem mass spectrometry. *Clin Chem* 2008;54:642-51.
26. De Grande LAC, Goossens K, Van Uytfanghe K, Das B, MacKenzie F, Patru MM, Thienpont LM; IFCC Committee for Standardization of Thyroid Function Tests (C-STFT). Monitoring the stability of the standardization status of FT4 and TSH assays by use of daily outpatient medians and flagging frequencies. *Clin Chim Acta* 2017;467:8-14.
27. Miller WG, Jones GR, Horowitz GL, Weykamp C. Proficiency testing/external quality assessment: current challenges and future directions. *Clin Chem* 2011;57:1670-80.
28. Thienpont LM, Faix JD, Beastall G. Standardization of free thyroxine and harmonization of thyrotropin measurements: A request for input from endocrinologists and other physicians. *Thyroid* 2015;25:1379-80.

**Table 1. Study participants (ordered by code), inclusive the platforms/FT4 assays examined for standardization. The listed reference and measurement intervals are those stated in the kit inserts.**

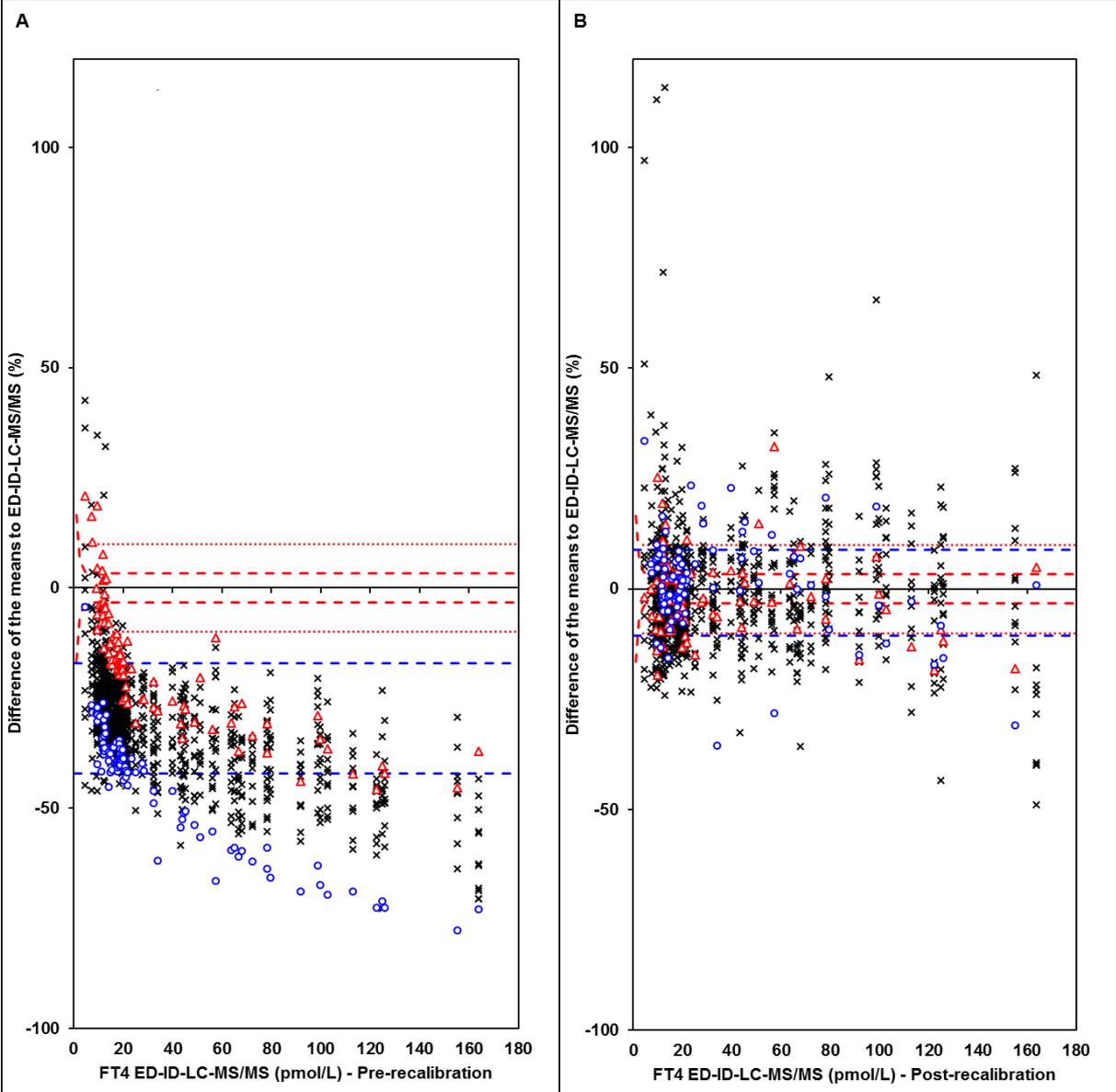
<b>IVD manufacturer Platform/Immunoassay</b>	<b>Code</b>	<b>Reference Interval (pmol/L)</b>	<b>Measurement Interval (pmol/L)<sup>c-g</sup></b>
Siemens Healthineers (Tarrytown, NY); <i>Advia Centaur XP</i>	A	11.5 - 22.7 (n = 388)	1.3 <sup>c</sup> - 155
Abbott Diagnostics (Abbott Park, IL); <i>Architect i2000</i>	B	9.0 - 19.1 (99%, n = 411)	5.2 <sup>d</sup> - 77
Ortho-Clinical Diagnostics (Buckinghamshire, UK); <i>Vitros ECI</i>	D	10.0 - 28.2 (98%, n = 535)	0.9 <sup>c</sup> - 90
bioMérieux SA (Marcy-l'Etoile, France) ; <i>Vidas</i>	E	10.6 – 19.4 (95%, n = 623)	1.1 <sup>e</sup> - 100
Beckman Coulter Inc. (Brea, CA); <i>Access 2</i>	F	7.9 – 14.4 (95%, n = 316)	3.2 <sup>e</sup> - 77
DiaSorin S.p.A (Saluggia, Italy); <i>Liaison® Analyser</i>	G	10.3 - 21.9 (95%, n = 517)	1.3 <sup>c</sup> - 129
<sup>a</sup> Sichuan Maccura Biotechnology Co., Ltd (Chengdu, China); <i>IS1200</i>	H	12.2 - 21.2 (95%, n = 175)	2.0 <sup>e</sup> - 100
Roche Diagnostics GmbH (Mannheim, Germany); <i>Elecsys (Cobas e 601)</i>	I	12.0 – 22.0 (95%, n = 801)	3.0 <sup>f</sup> - 100
Tosoh Corporation (Tokyo, Japan); <i>AIA-2000</i>	J	10.6 – 21.0 (95%, n = 618)	1.3 <sup>e</sup> - 103
<sup>a</sup> Snibe Co.,Ltd, (Shenzhen, China); <i>Maglumi 2000</i>	K	11.5 - 22.1 (95%)	1.3 <sup>g</sup> - 154

<sup>a</sup> Fujirebio Inc. (Tokyo, Japan); <i>Lumipulse G1200</i>	L	9.7 - 19.8 (95%, n = 141)	1.0 <sup>c</sup> - 129
<sup>b</sup> LSI Medience Corporation (Tokyo, Japan); <i>STACIA</i>	N	12.5 - 26.5	1.3 <sup>e</sup> - 103
<sup>b</sup> Sysmex Corporation (Kobe, Japan); <i>HISCL-5000</i>	O	9.9 - 20.5	3.2 <sup>e</sup> - 77

<sup>a,b</sup>Manufacturers who only joined in 2015<sup>a</sup> and/or 2016<sup>b</sup> for participation in the Phase IV method comparison study.

<sup>c-g</sup>The lower limit of the measurement intervals is: <sup>c</sup>limit of detection (according to the CLSI's EP-17 protocol); <sup>d</sup>functional sensitivity (CV 10%); <sup>e</sup>functional sensitivity (CV 20%); <sup>f</sup>limit of quantification at a total error of  $\pm 30\%$  (CLSI EP-17); <sup>g</sup>limit of quantification (CLSI EP-17) (21).

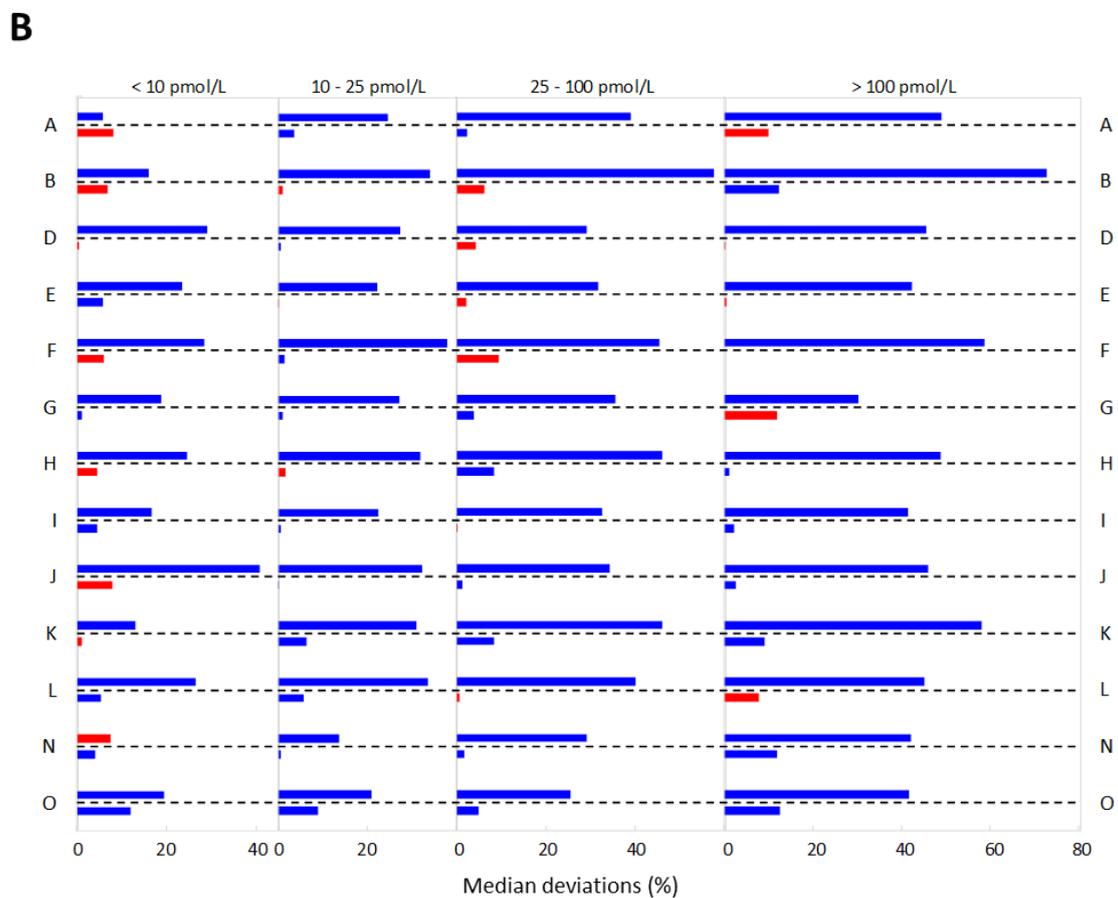
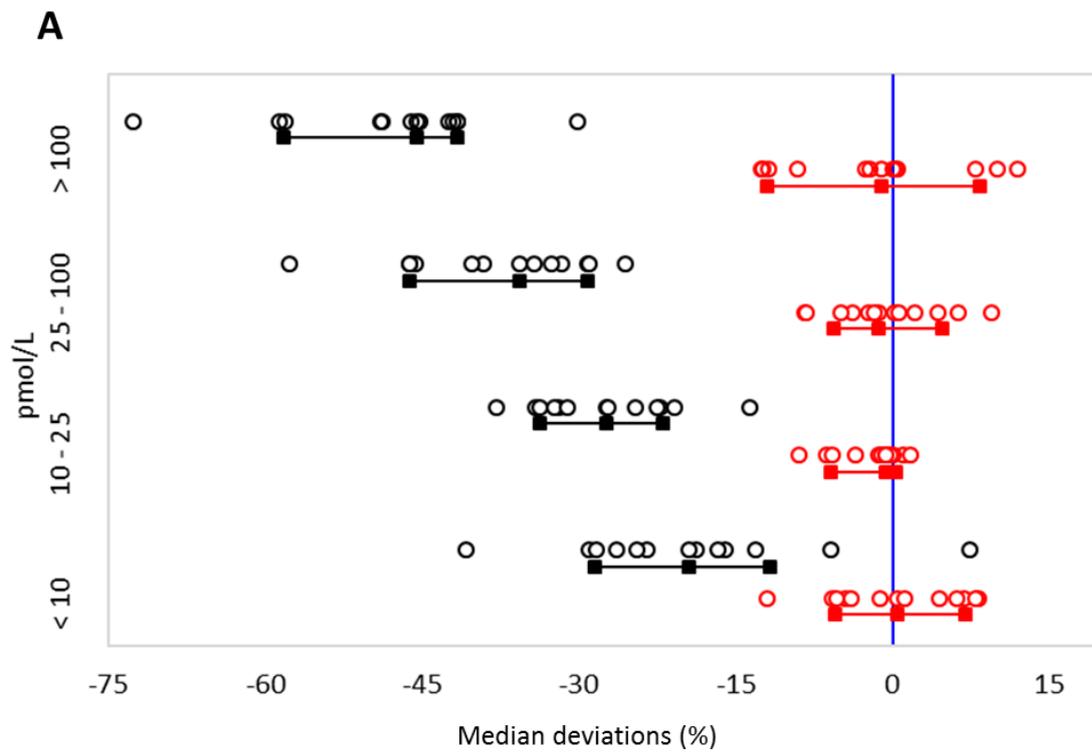
**Figure Legends**



**Figure 1. Combined difference (%) plots of the immunoassay results to those by ED-ID-LC-MS/MS, before (A) and after recalibration (B).**

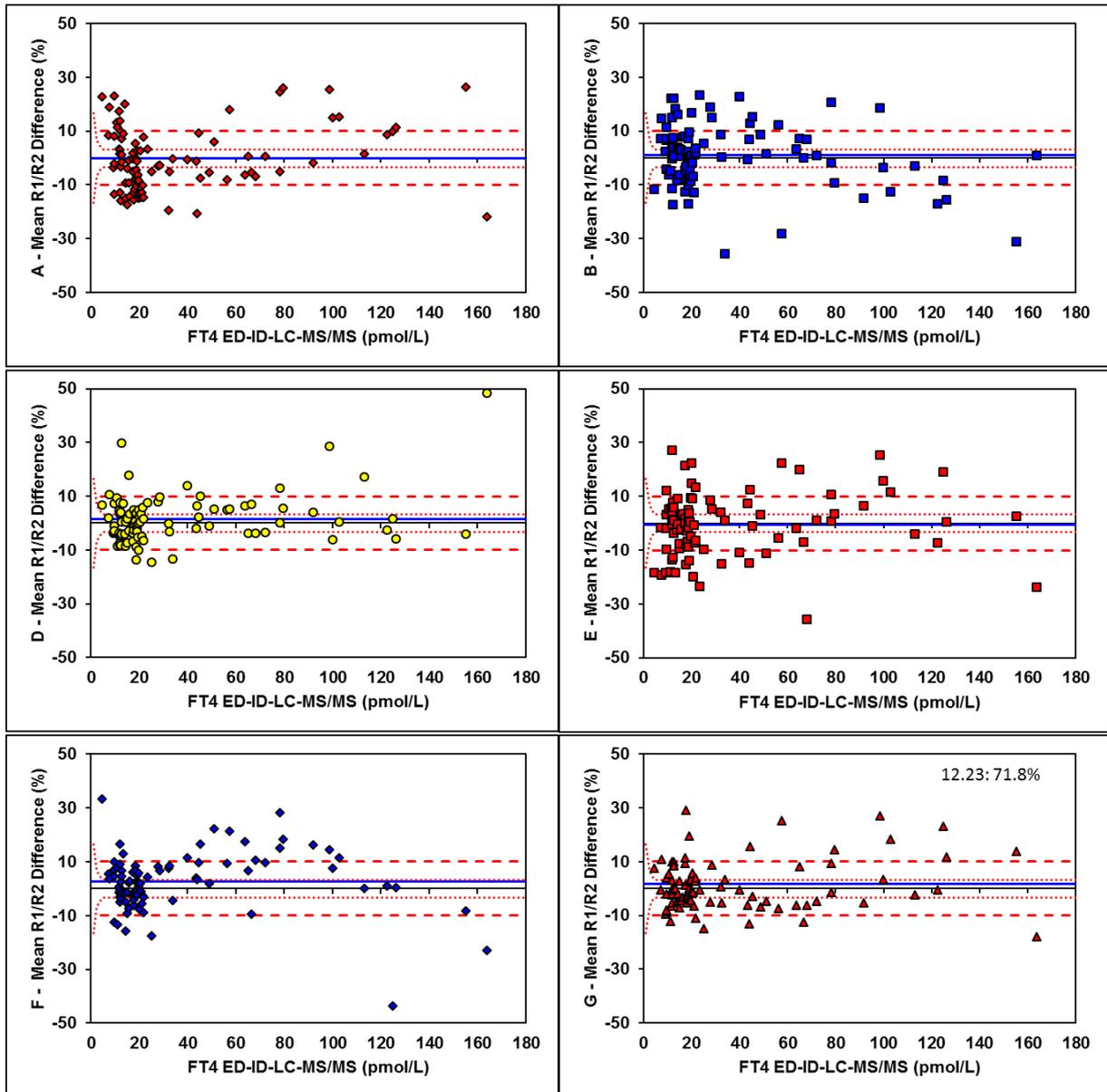
The most discrepant assays before recalibration are highlighted by colored symbols (blue circles for assay B (< 25 pmol/L) and assay F (> 25 pmol/L); red triangles for assay N), while all other assays are indicated with the symbol X. The red broken lines are the bias limits based on the biological variation concept:  $\pm 3.3\%$  (it is to note that we converted the percentage limit to 0.165 pmol/L for concentrations  $\leq 5$  pmol/L), while the red dotted lines are

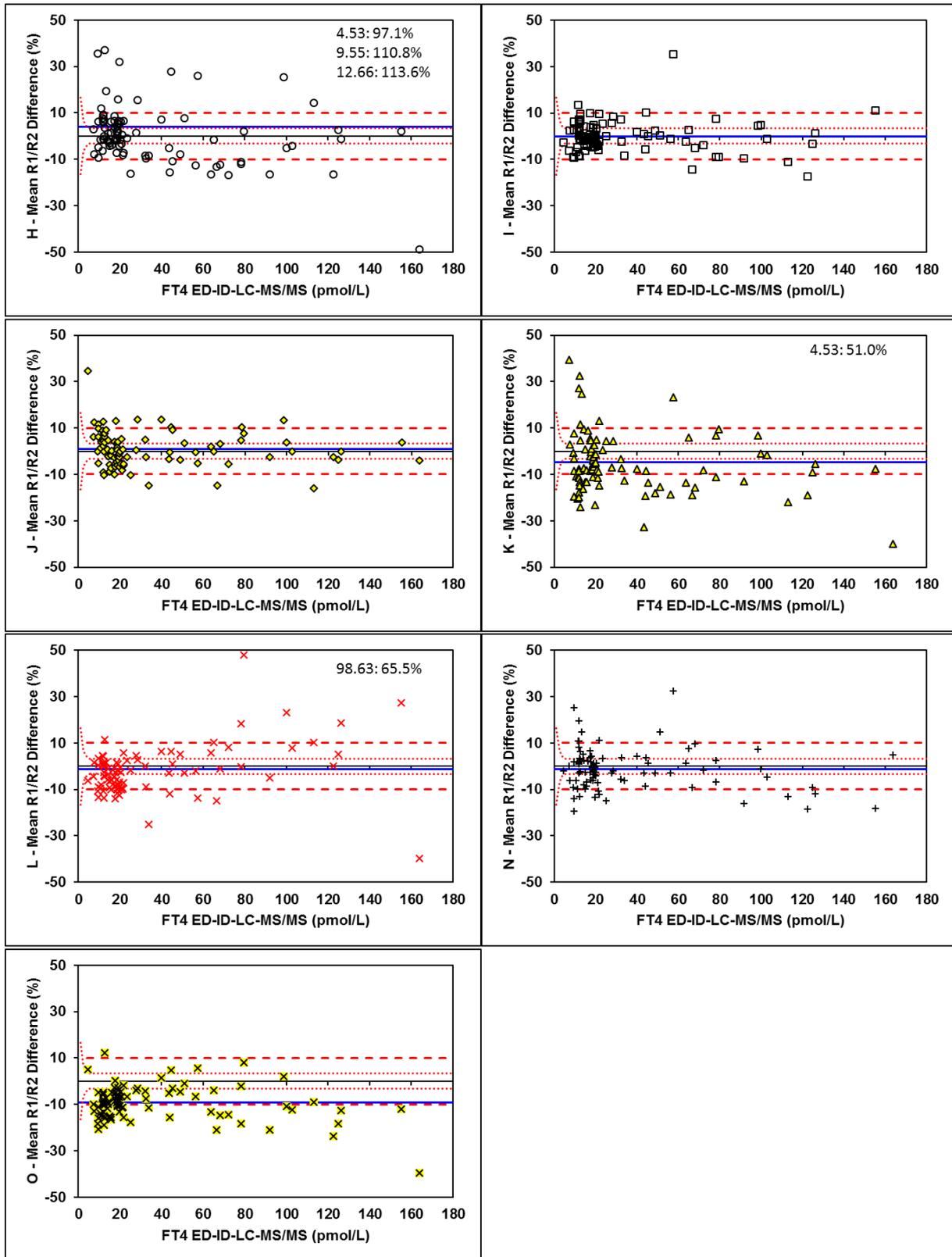
the empirical bias limits of 10% (8, 18,19). The blue broken lines represent the 15<sup>th</sup> and 85<sup>th</sup> centiles.



**Figure 2. Median deviations (%) of the immunoassays to ED-ID-LC-MS/MS before and after recalibration in 4 concentration intervals: <10 pmol/L, 10 – 25 pmol/L, 25 – 100 pmol/L and >100 pmol/L.**

(A) the overall improvement in terms of the median deviations (%) by recalibration. For each concentration interval, 2 pairs of data are shown; the black and red dots show the combined assay-specific median deviations before and after recalibration, respectively; the lines represent the 15<sup>th</sup>, 50<sup>th</sup> and 85<sup>th</sup> centiles; (B) the median deviations (%) of each assay by a pair of bars; the upper and lower bar shows the median deviation before and after recalibration, respectively. Note that the bars represent the unsigned magnitudes, while the colors refer to the signs (blue: negative, red: positive).

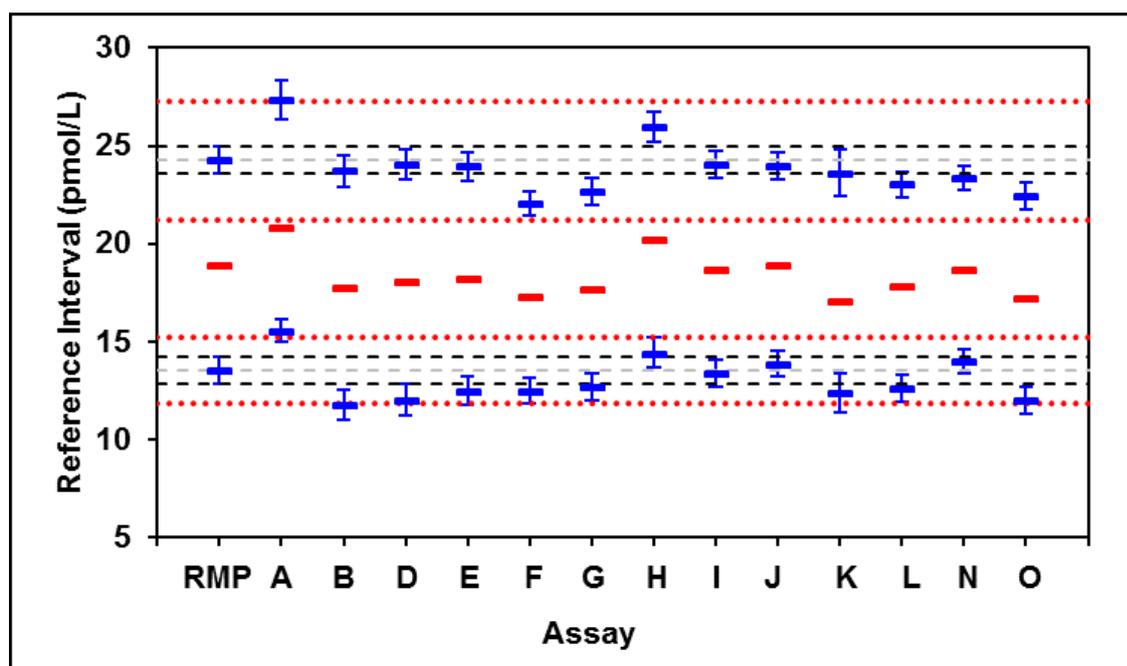




**Figure 3. Difference (%) plots after recalibration of the individual immunoassays.**

The red dotted lines are the 3.3% bias limits from the biological variation concept (converted to 0.165 pmol/L for concentrations  $\leq 5$  pmol/L), while the red broken lines stand for the

previously used empirical limits of 10% (8, 18,19). The blue line represents for each immunoassay the mean deviation or bias (%). The one-sided 95% CIs given in Supplemental Table 2 are not shown because of too little graphical resolution. To keep the Y-axis identical in all plots, certain % differences required omission (concentrations and % differences mentioned in the plots).



**Figure 4. Comparison of the reference interval percentiles of the individual immunoassays to those of ED-ID-LC-MS/MS (n = 120).**

The blue thick horizontal bars represent the respective 2.5 percentiles and 97.5 percentiles of each reference interval, while the blue vertical lines show the respective 90% CIs. The red thick horizontal bars for each assay stand for the mean (except for assay K, for which it shows the median). The grey and black broken horizontal lines represent the reference percentiles (from the data by the RMP) and the 90% CIs around them, respectively. The red dotted lines are the 12.5% limits of the interval around the reference percentiles.

**Supplement to “Standardization of free thyroxine measurements allows the adoption of a more uniform reference interval .”**

Linde A.C. De Grande, Katleen Van Uytfanghe, , Dries Reynders, Barnali Das, James D. Faix, Finlay MacKenzie, Brigitte Decallonne, Akira Hishinuma, Bruno Lapauw, Paul Taelman, Paul Van Crombrugge, Annick Van den Bruel, Brigitte Velkeniers, Paul Williams, Linda M. Thienpont on behalf for the IFCC Committee for Standardization of Thyroid Function Tests (C-STFT).

**Contents**

1	Assay-specific median deviations (%) (pre- and post-recalibration).....	33
2	Post-recalibration biases (%).....	35
3	Total error plots after recalibration .....	36
4	Differences (%) between the replicates .....	38
5	Reference interval study – Additional figures and tables .....	41
6	Statistical testing of the hypothesis that the percentiles of the reference interval by the RMP suits for common use .....	46
7	Summary of the results of the homogeneity study.....	50



## 1 Assay-specific median deviations (%) (pre- and post-recalibration)

**Table 1: Median deviation (%) of each of the immunoassays to ED-ID-LC-MS/MS before and after recalibration in 4 concentration intervals.**

Assay	Before recalibration				After recalibration			
	<10 pmol/L	≥10 <25 pmol/L	≥25 <100 pmol/L	≥100 pmol/L	<10 pmol/L	≥10 <25 pmol/L	≥25 <100 pmol/L	≥100 pmol/L
A	-5.9	-24.6	-39.1	-49.0	8.2	-3.6	-2.3	10.0
B	-16.0	-34.2	-57.7	-72.7	6.8	1.0	6.3	-12.4
D	-29.1	-27.4	-29.2	-45.5	0.4	-0.6	4.4	0.3
E	-23.5	-22.3	-31.7	-42.4	-5.8	0.1	2.1	0.5
F	-28.4	-37.9	-45.6	-58.7	6.1	-1.4	9.5	0.0
G	-18.8	-27.2	-35.7	-30.2	-1.2	-1.0	-3.9	11.9
H	-24.5	-32.0	-46.2	-48.9	4.5	1.7	-8.4	-1.1
I	-16.7	-22.5	-32.7	-41.6	-4.5	-0.6	0.2	-2.2
J	-40.8	-32.4	-34.3	-46.1	8.0	-0.2	-1.3	-2.6
K	-13.1	-31.1	-46.2	-58.1	1.1	-6.3	-8.3	-9.1
L	-26.4	-33.7	-40.2	-45.2	-5.4	-5.8	0.6	7.9
N	7.4	-13.7	-29.1	-42.1	-4.0	-0.6	-1.7	-11.9

O      -19.5      -20.9      -25.6      -41.7      -12.0      -8.9      -4.9      -12.5

---

## 2 Post-recalibration biases (%)

**Table 2: Assay biases/mean deviations (%) and one-sided 95% confidence interval (CI) after recalibration to the ED-ID-LC-MS/MS targets, and their assessment against 2 specifications: 3.3% inferred from the biological variation and 10% used as an empirical limit.**

Assay	Bias (%)	One-sided* 95% CI (%)	Upper bias limit (%) (Bias + CI)	Lower bias limit (%) (Bias - CI)
A	-0.1	2.0	1.9	-2.2
B	1.2	2.0	3.2	-0.8
D	1.5	1.6	3.1	-0.1
E	-0.5	2.0	1.6	-2.5
F	2.5	1.8	4.3	0.7
G	1.7	2.0	3.7	-0.4
H	<u>4.1</u>	4.0	8.1	0.1
I	-0.2	1.2	1.0	-1.4
J	0.9	1.3	2.2	-0.5
K	<u>-4.6</u>	2.5	-2.2	-7.1
L	-1.3	2.2	0.9	-3.6
N	-1.3	1.6	0.3	-2.9
O	<u>-9.2</u>	1.3	-7.9	-10.5

\*One-sided *t*-values (obtained from Excel with the function TINV(0.1, df)) were used for the calculation of the CI.

Interpretation: it can be confidently asserted that after recalibration the bias (and 95% CI) of all but assay O met the empirical specification of 10% with at 95% probability; for assay O (in spite of a bias below 10%) this statement does not apply as the lower 95% CI limit violated the specification (purple cell). When validating the biases against the 3.3% specification, those of 3 out of 13 assays (H, K, and obviously O; biases underlined) violated it; in addition, for 3 other assays (F, G and L), in spite of having biases less than  $\pm 3.3\%$ , it is not possible to state with 95% confidence that the specification was complied with (one of the 95% CI limits (orange cells) was outside either + or -3.3%) (Ref. 23 in the main text).

### 3 Total error plots after recalibration

**Figure 1:** FT4 total error (TE) plots. The TE at the level of the individual sample was estimated from the % difference to the RMP target of the first replicate after recalibration. For validation, we used the TE specification from the biological variation concept (red broken lines), but expanded it from 8.0% to 13% to account for the imprecision of the ED-ID-LC-MS/MS RMP. We also added the 95% limits of agreement (mean % difference  $\pm$  1.96  $CV_{diff}$  (%); blue broken lines) to emphasize on the fact that the magnitude of the scatter in the plots is different from assay to assay.

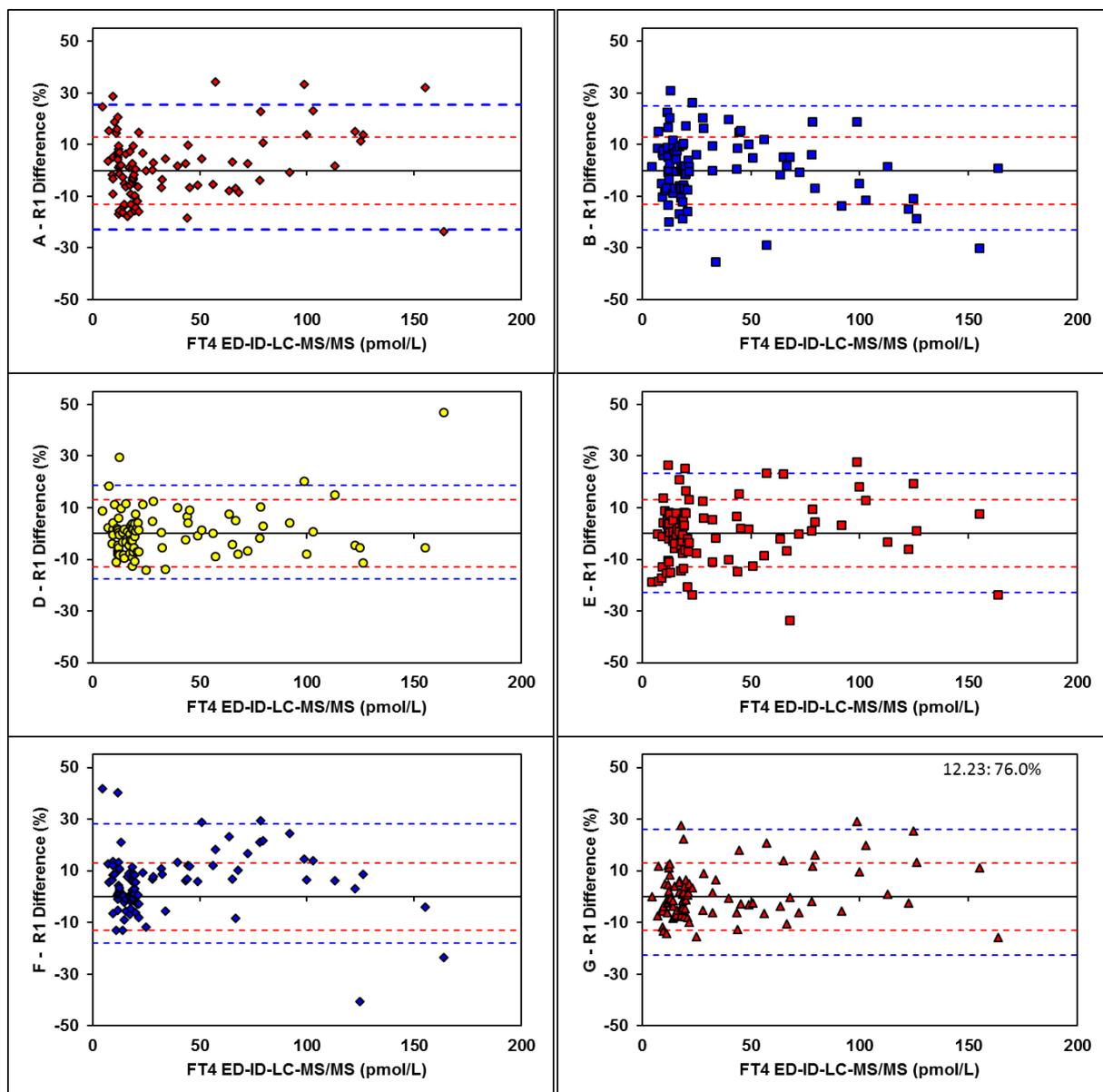
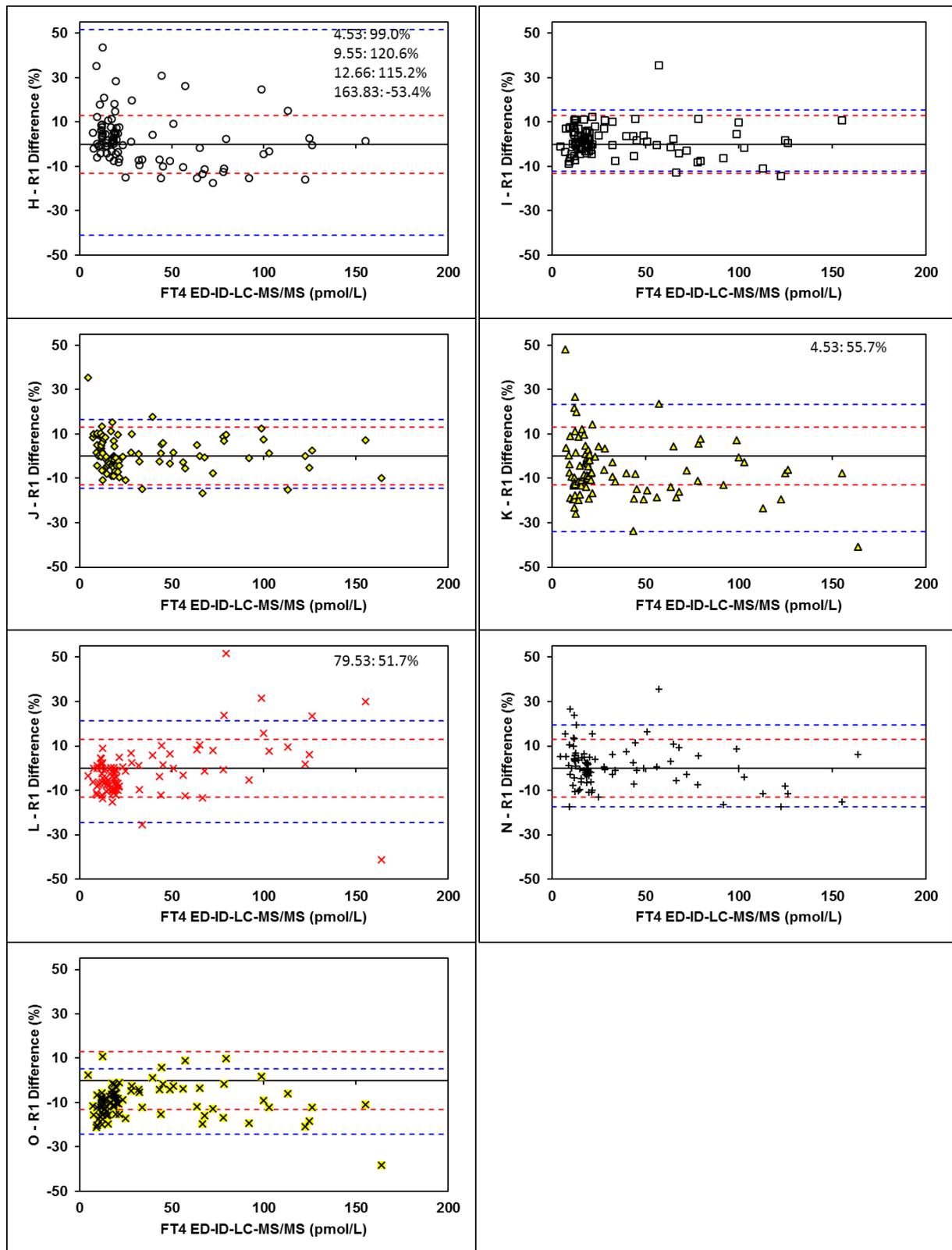


Figure 1: FT4 TE plots, continued.



### 3 Differences (%) between the replicates

**Table 3: Median differences (%) between the replicates from 2 runs (relative to the mean) and  $SD_{diff}$  (%). Note that the calculations were done from the reported results before recalibration.**

<b>Assay</b>	<b>Median difference (%)</b>	<b><math>SD_{diff}</math> (%)</b>
A	3.1	5.9
B	-0.4	4.1
D	-1.1	5.6
E	0.2	3.6
F	4.1	5.6
G	0.6	3.9
H	1.4	2.5
I	3.1	2.7
J	1.0	5.0
K	-1.5	3.5
L	0.9	4.4
N	2.4	4.5
O	-1.3	2.7

**Figure 2:** Difference (%) plots between the replicates obtained from different runs. Note that the samples for which the deviation was beyond 25% were not included in the plots; they are identified in the respective graphs by their concentration (according to ED-ID-LC-MS/MS) and difference (%).

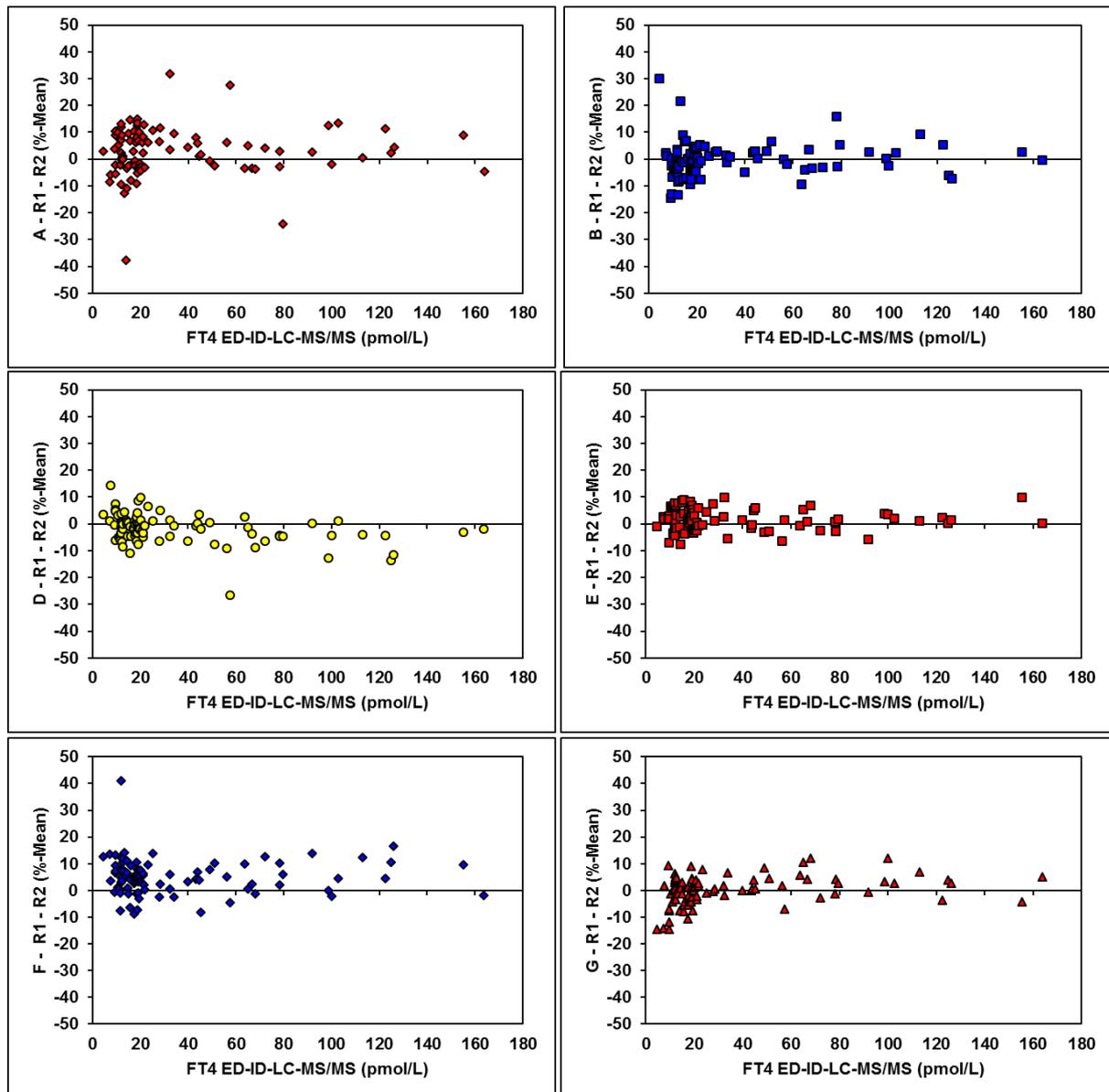
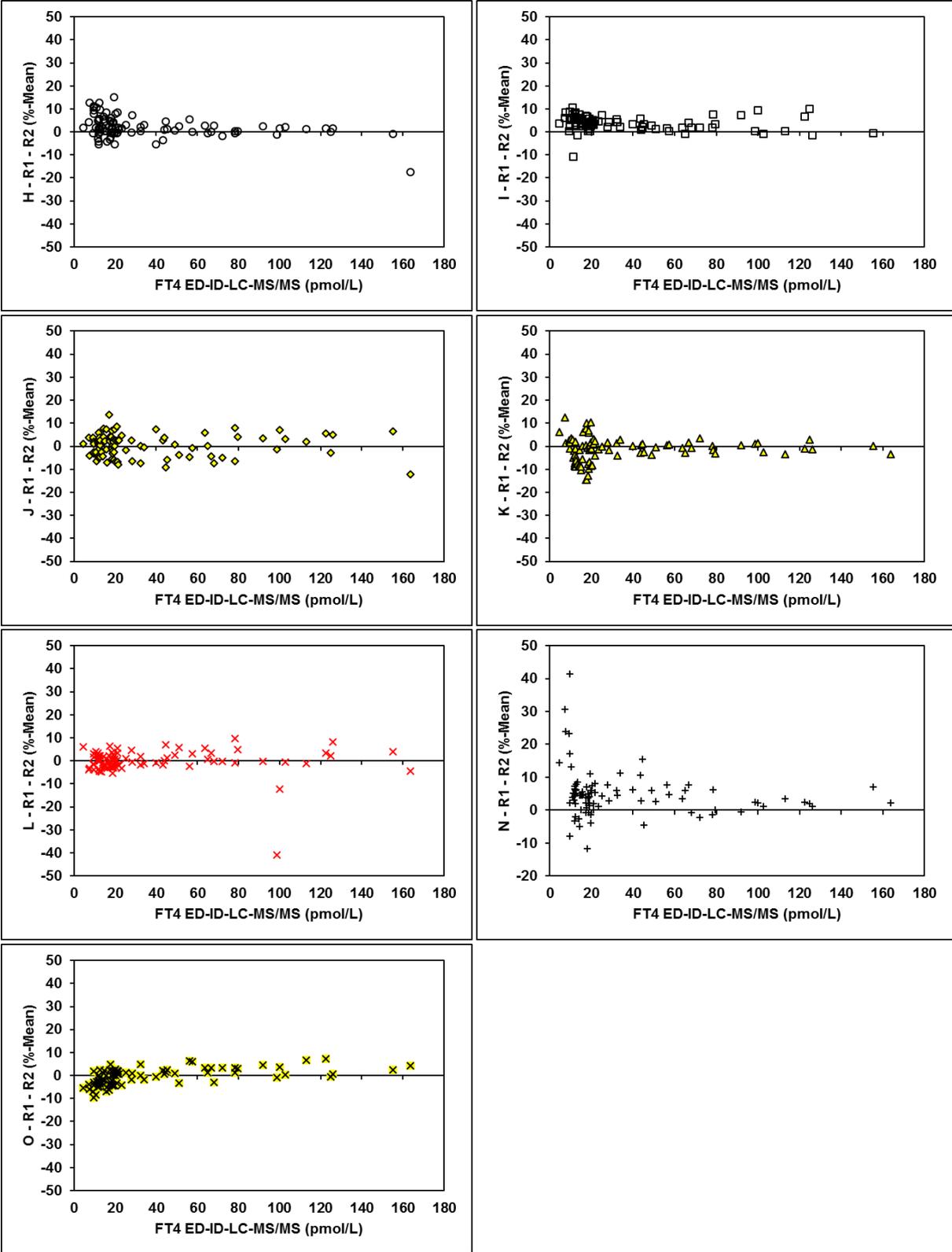
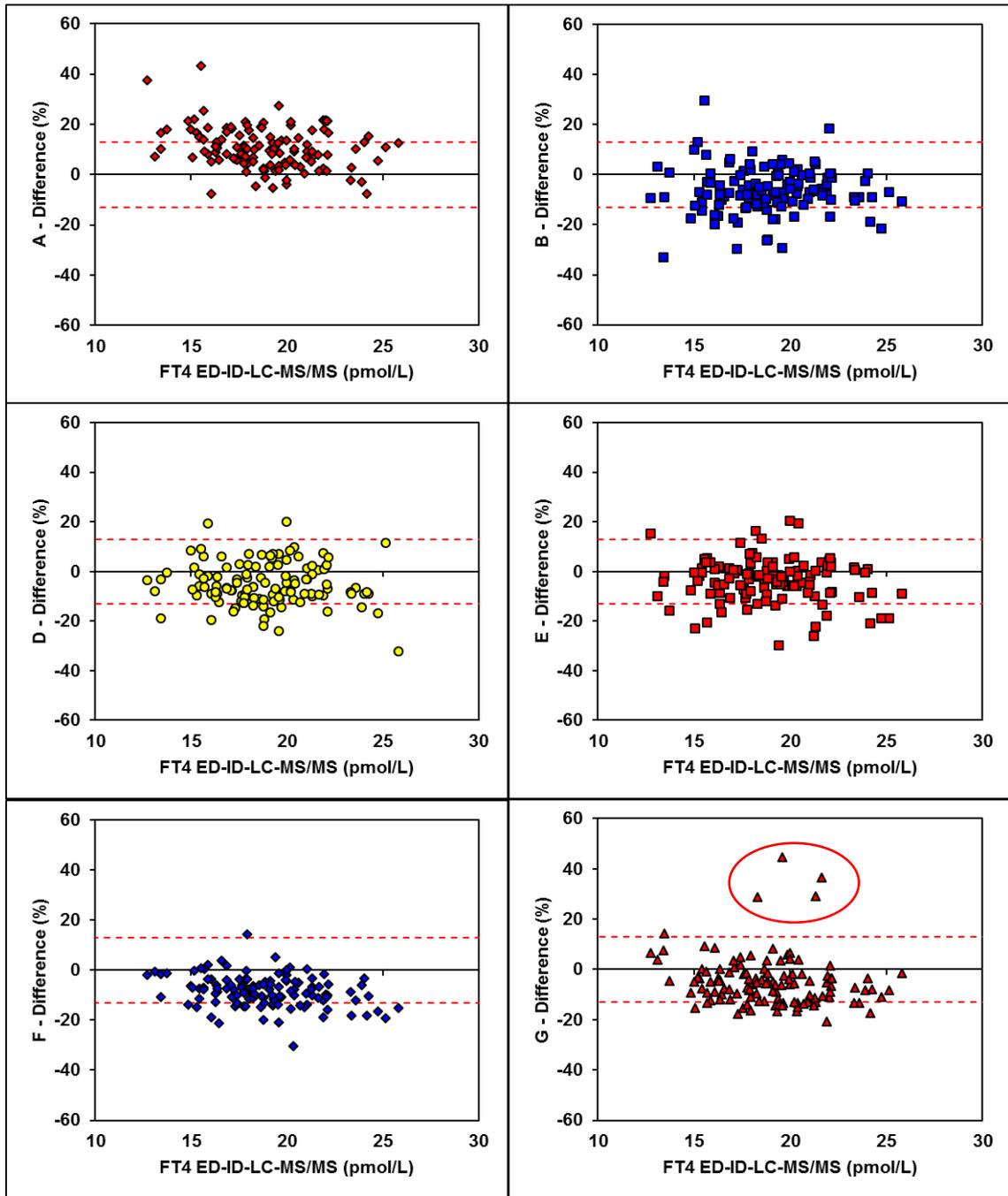


Figure 2: Difference (%) plots between the replicates, continued.

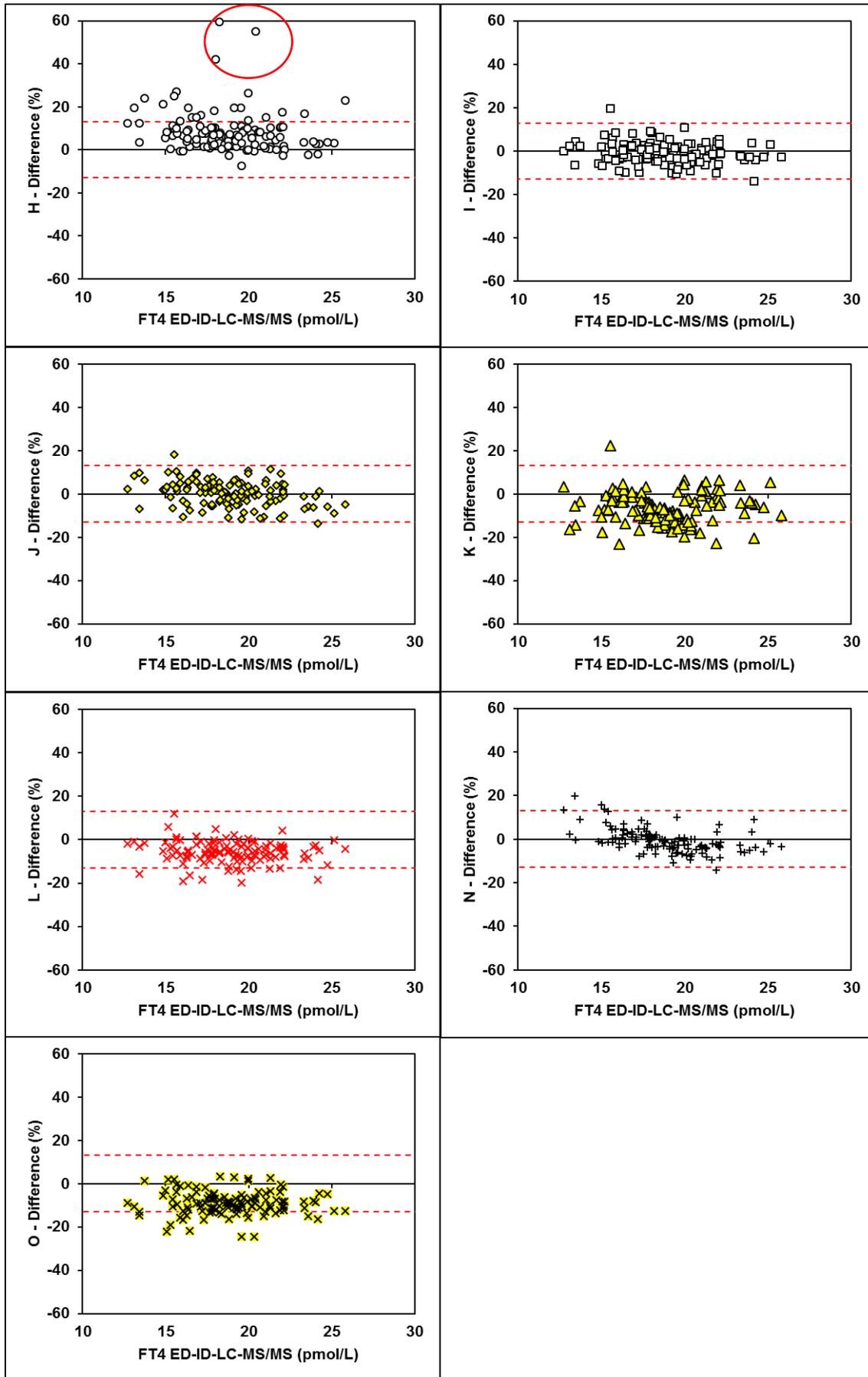


## 4 Reference interval study – Additional figures and tables

**Figure 3:** Difference (%) plots of the results by the immunoassays for the reference interval study against those by ED-ID-LC-MS/MS. The red broken line represents the expanded TE limit (13% for  $n = 1$ ). The red circles in the plots of assays G and H identify the differences (%) we visually identified as aberrant in comparison to the bulk of the data. As described in the main text, after removal of these data, the A-D test gave  $P > 0.05$ , which allowed use of the direct parametric procedure for estimating the RI characteristics of these 2 assays.



**Figure 3:** Difference (%) plots for the reference interval study (continued).



**Table 4: Characteristics of the different reference intervals.**

ID	Mean/Median concentration	Width RI	2.5 centile	90% CI	97.5 centile	90% CI	$\Delta^5$ 2.5 cent.	$\Delta^5$ 97.5 cent.
(pmol/L)							(%)	
RMP	18.9/18.8	10.7	13.5	12.8 - 14.2	24.3	23.6 - 25.0		
A <sup>1</sup>	20.8/20.5 <sup>4</sup>	11.8	15.5	15.0 - 16.1	27.3	26.3 - 28.3	14.9 <sup>6</sup>	12.6 <sup>6</sup>
B	17.7/17.8	12.0	11.7	11.0 - 12.5	23.7	22.9 - 24.5	-13.2 <sup>6</sup>	-2.3
D	18.0/17.7	12.0	12.0	11.2 - 12.8	24.0	23.3 - 24.8	-11.1	-0.9
E	18.2/18.5	11.5	12.5	11.7 - 13.2	23.9	23.2 - 24.7	-7.7	-1.3
F	17.3/17.1	9.6	12.5	11.9 - 13.1	22.1	21.4 - 22.7	-7.8	-9.1
G <sup>2</sup>	17.7/17.5	10.0	12.7	12.0 - 13.3	22.7	22.0 - 23.3	-6.1	-6.5
H <sup>2</sup>	20.2/19.8	11.6	14.4	13.6 - 15.2	26.0	25.2 - 26.8	6.5	7.1
I	18.7/18.6	10.7	13.4	12.7 - 14.1	24.0	23.3 - 24.7	-1.2	-0.9
J	18.9/18.7	10.1	13.8	13.2 - 14.5	24.0	23.3 - 24.6	2.4	-1.3
K <sup>3</sup>	17.7/17.1	11.2	12.4	11.3 - 13.8	23.6	22.4 - 24.8	-8.6	-2.7
L	17.8/17.6	10.4	12.6	11.9 - 13.3	23.0	22.3 - 23.7	-6.7	-5.2
N	18.7/18.4	9.4	14.0	13.3 - 14.6	23.4	22.8 - 24.0	3.4	-3.7
O	17.2/17.0	10.4	12.0	11.3 - 12.7	22.4	21.7 - 23.1	-11.5	-7.6

<sup>1</sup>Parametric after log-transformation.

<sup>2</sup>Parametric after removal of visually observed aberrant % differences to the target outliers (see above Figure 4S).

<sup>3</sup>Non-parametric bootstrap. Note, the increased width of the CI for assay K compared to the other assays is related to the statistical procedure (non-parametric bootstrap) used to estimate the RI.

<sup>4</sup>Geometric mean because the distribution of the dataset was not normal.

<sup>5</sup>Difference (%) to the reference 2.5 and 97.5 percentile, respectively.

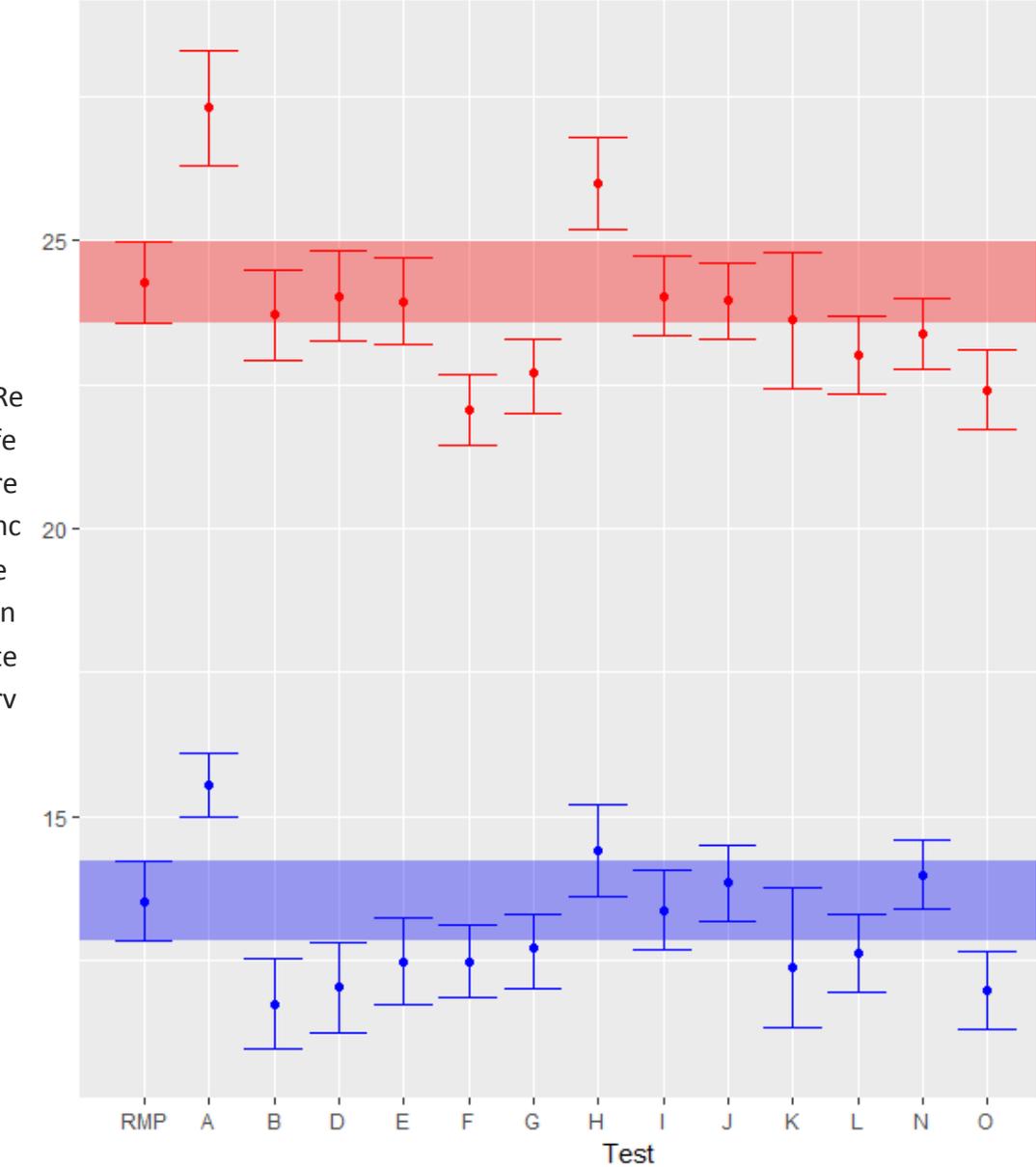
<sup>6</sup>Difference (%) to the reference percentile(s) exceeding the 12.5% limit.

## 5 Statistical testing of the hypothesis that the percentiles of the reference interval by the RMP suits for common use

**Table 5: Calculation of the probability that the 2.5- and 97.5-percentiles of the immunoassays are located in the interval flanked by the CI limits of the reference percentiles.**

<b>Assay</b>	<b>Probability 2.5-percentile (%)</b>	<b>Probability 97.5-percentile (%)</b>
A	0.0	0.0
B	1.1	61
D	4.7	81
E	23	78
F	18	0.0
G	38	1.4
H	36	1.7
I	88	85
J	82	83
K	20	25
L	31	9.1
N	74	31
O	2.1	0.3

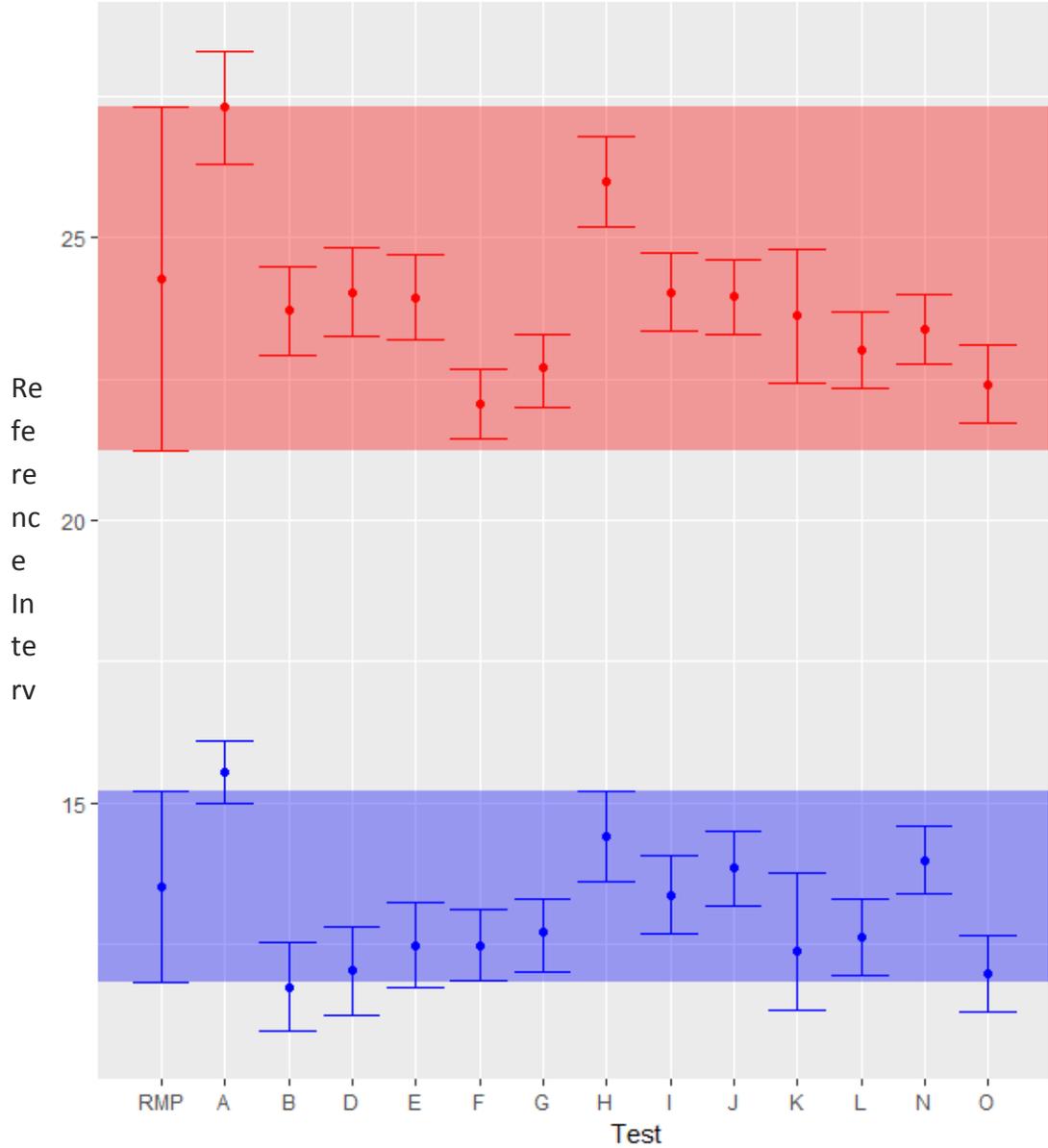
**Figure 4:** Visualization of the statistical test that calculates the probability that the 2.5- and 97.5-percentiles of the immunoassays are located in the interval flanked by the CI limits of the reference percentiles (indicated by the blue and red shaded zones). The blue and red vertical lines/horizontal bars represent the 90% CI of the RMP and the 13 assays for the 2.5- and 97.5-percentiles, respectively.



**Table 6: Calculation of the probability that the 2.5- and 97.5-percentiles of the immunoassays are located in the 12.5% interval around the reference percentiles.**

<b>Assay</b>	<b>Probability 2.5-percentile (%)</b>	<b>Probability 97.5-percentile (%)</b>
A	16	49
B	42	100
D	65	100
E	92	100
F	95	99
G	99	100
H	95	100
I	100	100
J	100	100
K	86	100
L	97	100
N	100	100
O	63	100

**Figure 5:** Visualization of the statistical test that calculates the probability that the 2.5- and 97.5-percentiles of the immunoassays are located in the 12.5% interval around the reference percentiles (indicated by the blue and red shaded zones). The blue and red vertical lines/horizontal bars represent the 12.5% interval around the reference percentiles, and the 90% CI of the 13 assays for the 2.5- and 97.5-percentiles, respectively.



## 6 Summary of the results of the homogeneity study

**Table 7: Summary of the results of the homogeneity study.**

Sample ID	Mean (pmol/L) (aliquots)	CV (%) (aliquots)	Mean (pmol/L) (pool)	CV (%) (pool)	<i>P</i> (F-test, 95% CL <sup>a</sup> )
1	18.5	0.7	18.5	1.1	0.2
2	17.9	0.6	18.0	0.2	0.5
3	7.7	0.7	7.7	1.0	0.3
4	9.0	0.3	9.1	0.8	0.5
5	31.3	0.6	31.3	0.8	0.4
6	83.4	1.3	83.0 <sup>b</sup>	1.0	0.5
7	11.1	0.6	11.1	0.8	0.5
8	28.2	0.6	28.3	1.0	0.2
9	17.9	0.7	17.8	1.0	0.4
10	14.1	0.6	14.1	0.5	0.6
11	18.0	0.5	18.0	0.7	0.7
12	9.0	0.8	9.0	0.6	0.5

<sup>a</sup>CL: confidence level

<sup>b</sup>: 1 Outlier identified with the Grubbs test