

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository:<https://orca.cardiff.ac.uk/id/eprint/171424/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Osinga, Joris A. J., Nelson, Scott M., Walsh, John P., Ashoor, Ghaliya, Palomaki, Glenn E., López-Bermejo, Abel, Bassols, Judit, Aminorroaya, Ashraf, Broeren, Maarten A. C., Chen, Liangmiao, Lu, Xuemian, Brown, Suzanne J., Veltri, Flora, Huang, Kun, Männistö, Tuija, Vafeiadi, Marina, Taylor, Peter N. , Tao, Fang-Biao, Chatzi, Lida, Kianpour, Maryam, Suvanto, Eila, Grineva, Elena N., Nicolaidis, Kypros H., D'Alton, Mary E., Poppe, Kris G., Alexander, Erik, Feldt-Rasmussen, Ulla, Bliddal, Sofie, Popova, Polina V., Chaker, Layal, Visser, W Edward, Peeters, Robin P., Derakhshan, Arash, Vrijkotte, Tanja G. M., Pop, Victor J. M. and Korevaar, Tim I. M. 2024. Defining gestational thyroid dysfunction through modified nonpregnancy reference intervals: an individual participant meta-analysis. *The Journal of Clinical Endocrinology & Metabolism* 109 (11) , e2151-e2158. 10.1210/clinem/dgae528

Publishers page: <http://dx.doi.org/10.1210/clinem/dgae528>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See <http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



1 **Title: Defining gestational thyroid dysfunction through modified non-pregnancy reference intervals: an**
 2 **individual participant meta-analysis**

3 **Keywords:** Thyroid Gland, Thyroid Function Tests, Reference values, Pregnancy, Thyrotropin, Thyroxine

4 Authors: Joris A.J. Osinga, MD^{1,2}, Scott M. Nelson, MRCOG, PhD³, John P. Walsh, MB, PhD^{4,5}, Ghalia
 5 Ashoor, MD⁶, Glenn E Palomaki, PhD⁷, Abel López-Bermejo, MD, PhD^{8,9}, Judit Bassols, MD¹⁰, Ashraf
 6 Aminorroaya, MD¹¹, Maarten A.C. Broeren, PhD¹², Liangmiao Chen, MD, PhD¹³, Xuemian Lu, MD, PhD¹³,
 7 Suzanne J. Brown⁴, Flora Veltri, MD, PhD¹⁴, Kun Huang, PhD¹⁵, Tuija Männistö, MD, PhD¹⁶, Marina
 8 Vafeiadi, PhD¹⁷, Peter N. Taylor, FRCP, PhD¹⁸, Fang-Biao Tao, MD, PhD¹⁵, Lida Chatzi, MD, PhD¹⁹, Maryam
 9 Kianpour, PhD¹¹, Eila Suvanto, MD, PhD²⁰, Elena N. Grineva, MD, PhD²¹, Kypros H Nicolaidis, MD, PhD²²,
 10 Mary E. D'Alton, MD²³, Kris G. Poppe, MD, PhD¹⁴, Erik Alexander, MD, PhD²⁴, Ulla Feldt-Rasmussen, MD,
 11 DMSc²⁵, Sofie Bliddal, MD, PhD²⁵, Polina V. Popova, MD, PhD²¹, Layal Chaker, MD, PhD^{1,2,26}, W. Edward
 12 Visser, MD, PhD^{1,2}, Robin P. Peeters, MD, PhD^{1,2}, Arash Derakhshan, MD^{1,2}, Tanja G.M. Vrijkotte, PhD²⁷,
 13 Victor J.M. Pop, MD, PhD²⁸, Tim I. M. Korevaar, MD, PhD^{1,2}

14 **Affiliations**

- 15 1. Department of Internal Medicine, Erasmus University Medical Center, 3000 CA Rotterdam, the
- 16 Netherlands.
- 17 2. Academic Center for Thyroid Diseases, Erasmus University Medical Center, 3000 CA, Rotterdam,
- 18 the Netherlands.
- 19 3. School of Medicine, Dentistry and Nursing, University of Glasgow, G12 8QQ Glasgow, UK.
- 20 4. Department of Endocrinology and Diabetes, Sir Charles Gairdner Hospital, Nedlands 6009,
- 21 Western Australia, Australia.
- 22 5. Medical School, University of Western Australia, Crawley 6009, Western Australia, Australia.
- 23 6. Harris Birthright Research Center for Fetal Medicine, King's College Hospital, SE5 9RS London,
- 24 United Kingdom.
- 25 7. Department of Pathology and Laboratory Medicine, Women & Infants Hospital and Alpert
- 26 Medical School at Brown University, RI 02903 Providence, Rhode Island.
- 27 8. Pediatric Endocrinology Research Group, Girona Biomedical Research Institute (IDIBGI), Dr.
- 28 Josep Trueta Hospital, 17007 Girona, Spain.
- 29 9. Departament de Ciències Mèdiques. Universitat de Girona, 17003, Girona, Spain.

- 1 10. Maternal-Fetal Metabolic Research Group, Girona Biomedical Research Institute (IDIBGI), Dr.
- 2 Josep Trueta Hospital, 17007 Girona, Spain.
- 3 11. Isfahan Endocrine and Metabolism Research Center, Isfahan University of Medical Sciences,
- 4 81745-33871 Isfahan, Iran.
- 5 12. Laboratory of Clinical Chemistry and Haematology, Máxima Medical Centre, 5504 DB Veldhoven,
- 6 Netherlands.
- 7 13. Department of Endocrinology and Rui'an Center of the Chinese-American Research Institute for
- 8 Diabetic Complications, Third Affiliated Hospital of Wenzhou Medical University, 325035 Wenzhou,
- 9 China.
- 10 14. Endocrine Unit, Centre Hospitalier Universitaire Saint-Pierre, Université Libre de Bruxelles (ULB),
- 11 Brussels, Belgium.
- 12 15. Department of Maternal, Child and Adolescent Health, Scientific Research Center in Preventive
- 13 Medicine; School of Public Health; Anhui Medical University, 230032 Anhui, China.
- 14 16. NordLab, Oulu and Translational Medicine Research Unit, University of Oulu, 90570 Oulu,
- 15 Finland.
- 16 17. Department of Social Medicine, School of Medicine, University of Crete, 710 03 Heraklion, Crete,
- 17 Greece.
- 18 18. Thyroid Research Group, Systems Immunity Research Institute, Cardiff University School of
- 19 Medicine, CF10 3EU Cardiff, UK.
- 20 19. Department of Preventive Medicine, Keck School of Medicine, University of Southern California,
- 21 90089 CA, USA.
- 22 20. Department of Obstetrics and Gynecology and Medical Research Center Oulu, University of
- 23 Oulu, 90570 Oulu, Finland.
- 24 21. Institute of Endocrinology, Almazov National Medical Research Centre, 197341 Saint Petersburg,
- 25 Russia.
- 26 22. Department of Women and Children's Health, Faculty of Life Sciences and Medicine King's
- 27 College London, SE5 9RS London, United Kingdom.
- 28 23. Department of Obstetrics and Gynecology, Columbia University Irving Medical Center, NY 10032
- 29 New York, USA.
- 30 24. Division of Endocrinology, Hypertension and Diabetes, Brigham and Women's Hospital, Harvard
- 31 Medical School, Boston, 02115 MA, USA.
- 32 25. Department of Medical Endocrinology and Metabolism, Copenhagen University Hospital,
- 33 Rigshospitalet, and Department of Clinical Medicine, Faculty of Health and Clinical Sciences,
- 34 Copenhagen University, 2100 Copenhagen, Denmark.
- 35 26. Department of Epidemiology, Erasmus University Medical Center, 3000 CA Rotterdam, the
- 36 Netherlands.
- 37 27. Department of Public and Occupational Health, Amsterdam UMC, University of Amsterdam,
- 38 Amsterdam Public Health Research Institute, 1081 HV Amsterdam, the Netherlands
- 39 28. Department of Medical and Clinical Psychology, Tilburg University, 5000 LE Tilburg, The
- 40 Netherlands.

41 **Corresponding author:** Joris Osinga, Generation R, Postbus 2040, 3000 CA Rotterdam,

42 j.osinga@erasmusmc.nl, ORCID-ID: 0000-0003-2527-5150

43 **Funding:** Netherlands Organization for Scientific Research (grant 401.16.020) and a Vidi grant

44 (016.176.331) from the Netherlands Organization for Scientific Research to R.P.P.

1 **Disclosures:** P.T. reports a travel grant from Society for Endocrinology (leadership development award).
2 E.N.G. received speaker's fees and payment for expert testimony from Merck and consulting fees from
3 Brunel Rus. T.G.M.V. reports grants from the Netherlands Organization for Health Research and
4 Development. L.C. received travel support by Pfizer. S.M.N. has received consultancy, speakers' fees, or
5 travel support from Access Fertility, Beckman Coulter, Ferring Pharmaceuticals, Merck, Modern Fertility,
6 Roche Diagnostics, and The Fertility Partnership. S.M.N. also reports payments for medical–legal work
7 and investment in The Fertility Partnership. T.I.M.K. reports lectureship fees from Berlin–Chemie,
8 Goodlife Healthcare, Institut Biochimique SA, Merck, and Quidel. U.F.R.'s research salary was sponsored
9 by an unrestricted grant from Kirsten and Freddy Johansen's Fund and reports lecture fee from Merck,
10 Darmstadt. S.B.'s research salary was sponsored by the Capital Region of Denmark's Research
11 Foundation and the Novo Nordisk Foundation (ID 0077221). S.B. received a lecture fee from Merck and
12 Novo Nordisk. All other authors declare no competing interests.

14 **Abstract**

15 Background: Establishing local trimester-specific reference intervals for gestational TSH and FT4 is often
16 not feasible, necessitating alternative strategies. We aimed to systematically quantify the diagnostic
17 performance of standardized modifications of center-specific non-pregnancy reference intervals as
18 compared to trimester-specific reference intervals.

19 Methods: We included prospective cohorts participating in the Consortium on Thyroid and Pregnancy.
20 After relevant exclusions, reference intervals were calculated per cohort in thyroperoxidase antibody-
21 negative women. Modifications to the non-pregnancy reference intervals included an absolute
22 modification (per 0.1 mU/L TSH or 1 pmol/L FT4), relative modification (in steps of 5%) and fixed limits
23 (upper TSH limit between 3.0 to 4.5 mU/L and lower FT4 limit 5-15 pmol/L). We compared (sub)clinical

1 hypothyroidism prevalence, sensitivity and positive predictive value (PPV) of aforementioned
2 methodologies with population-based trimester-specific reference intervals.

3 Results: The final study population comprised 52,496 participants in 18 cohorts. Optimal modifications
4 of standard reference intervals to diagnose gestational overt hypothyroidism were -5% for the upper
5 limit of TSH and +5% for the lower limit of FT4 (sensitivity 0.70, confidence interval [CI] 0.47-0.86; PPV
6 0.64, CI 0.54-0.74). For subclinical hypothyroidism, these were -20% for the upper limit of TSH and -15%
7 for the lower limit of FT4 (sensitivity 0.91, CI 0.67-0.98; PPV 0.71, CI 0.58-0.80). Absolute and fixed
8 modifications yielded similar results. Confidence intervals were wide, limiting generalizability.

9 Conclusion: We could not identify modifications of non-pregnancy TSH and FT4 reference intervals that
10 would enable centers to adequately approximate trimester-specific reference intervals. Future efforts
11 should be turned towards studying the meaningfulness of trimester-specific reference intervals and risk-
12 based decision limits.

13

14 **Introduction**

15 Thyroid dysfunction during pregnancy is associated with a higher risk of miscarriage, preeclampsia,
16 preterm birth, aberrant birthweight and lower offspring IQ¹⁻⁶. Current international guidelines
17 recommend defining gestational thyroid dysfunction according to population and pregnancy-specific
18 TSH and FT4 reference intervals, to take into account thyroid physiology during pregnancy, as well as
19 differences in TSH and FT4 determinants between populations and the use of different laboratory assays
20⁷⁻⁹. However, calculating such local reference intervals is generally not feasible for most centers^{10,11}. In
21 addition to the practical hurdles, most of the published reference intervals for TSH and FT4 are not in
22 accordance with the current ATA guidelines, as we recently showed by providing an overview of

1 published TSH and FT4 reference intervals and methodologies, showing that most studies included used
2 additional exclusion criteria based on health status, did not exclude TPOAb positive participants or used
3 different percentile cutoffs ⁸. This is in part because of changing guidelines and in part because many
4 centers use additional exclusion criteria or apply different reference limit cut-offs ⁸. These varying
5 methodologies hamper the adoption of reference intervals from other centers, and as such, the vast
6 majority of centers rely on non-pregnancy reference intervals for TSH with either a fixed limit approach
7 (upper limit of 4.0 mU/L for TSH) or a subtraction approach (subtraction of 0.5 mU/L of the upper limit
8 of TSH), while for FT4 varying local approaches are used including non-pregnancy reference intervals ¹²⁻
9 ¹⁴. These second-tier strategies are considered inferior compared to locally defined reference intervals
10 ¹⁵⁻¹⁷. In a follow-up study, we showed that the use of a fixed upper TSH limit or the subtraction approach
11 results in poor detection rates and high false positive rates for (subclinical) hypothyroidism in early
12 pregnancy with highly variable diagnostic performance between populations (sensitivity 0.63-0.82, false
13 discovery rate 0.11-0.35) ¹⁸.

14 In search of a method that is both easy to implement in clinical practice and would better identify
15 women with an abnormal thyroid function during pregnancy, we set out to investigate if it is possible to
16 modify the center-specific non-pregnancy TSH and FT4 reference intervals so that these are useful in
17 pregnancy. Such an approach could make the establishment of local pregnancy-specific reference
18 intervals obsolete while it takes account of the local assay and pre-existing laboratory harmonization
19 efforts ^{19,20}. A useful diagnostic approach would need to fulfill some conditions: 1) the diagnostic
20 performance should at least perform better than currently recommended alternative methods (TSH
21 upper limit of 4.0 mU/L or subtraction of 0.5 mU/L) ^{12,13}, and 2) the diagnostic performance should be
22 reasonably consistent between populations.

1 In this individual participant meta-analysis, we aimed to modify the center-specific non-pregnancy
2 reference intervals of TSH and FT4 in a standardized manner and study the sensitivity and the positive
3 predictive value (PPV) as compared to center-specific gestational reference intervals as calculated in
4 accordance with the current international guidelines.

5

6 **Methods**

7 The study inclusion and eligibility procedures are described in detail previously¹⁸. In short, eligible
8 studies were those participating in the Consortium on Thyroid and Pregnancy
9 (<https://www.consortiumthyroidpregnancy.org>). Exclusion criteria for participants were pre-pregnancy
10 thyroid disease, pregnancy through in-vitro fertilization/ intracytoplasmic sperm injection (IVF/ICSI), use
11 of thyroid (interfering) medication and multiple gestation. For this study, we followed the Preferred
12 Reporting Items for Systematic Reviews and Meta-Analyses guidelines for Individual Patient Data and
13 preregistered the study protocol (CRD42021270078), which can be found in the supplemental materials
14 along with an outline of protocol deviations²¹. Study quality and risk of bias were assessed using the
15 Newcastle-Ottawa scale (Supplemental materials²¹). All cohorts were approved by a local review board
16 and acquired participant informed consent or had been granted exemption from it by the local Ethics
17 Committee.

18 Defining gestational thyroid dysfunction

19 Non-pregnancy reference intervals were either published and/or provided by the principal investigator
20 of the included cohorts and are assay-specific. We defined the trimesters as 0 to 13 weeks, >13 to 27
21 weeks and >27 weeks of gestation. For cohorts containing participants with repeated measurements, we
22 used the first available sample for each trimester.

1 Reference intervals, thyroid dysfunction (overt and subclinical hypothyroidism) and diagnostic test
2 properties were calculated separately for each cohort to account for inter-population differences. All
3 reference intervals were calculated as the 2.5th to 97.5th percentiles in TPOAb-negative participants. Our
4 primary aim was to optimize the diagnosis of thyroid dysfunction states for which treatment is indicated
5 or should be considered based on current guidelines, and thus we limited analyses to overt and
6 subclinical hypothyroidism¹³. A treatment indication was defined as either 1) overt hypothyroidism, 2)
7 subclinical hypothyroidism with TSH>10 mU/L or 3) subclinical hypothyroidism with TPOAb positivity. A
8 treatment consideration was defined as 1) TSH between 2.5 mU/L and the upper reference limit with
9 concomitant TPOAb positivity or 2) subclinical hypothyroidism without TPOAb positivity¹³. Treatment of
10 hyperthyroidism was outside the scope of this study, since gestational hyperthyroidism is often
11 considered physiological and we do not have data available to differentiate between gestational
12 transient thyrotoxicosis and Graves' hyperthyroidism¹³. The prevalence of thyroid dysfunction and
13 diagnostic performance measures were calculated according to several methods; 1) a relative
14 modification of the non-pregnancy upper limit of TSH varying from -5% to -40% in steps of 5%, with
15 modifications to the lower limit of FT4 varying from -20% to +20% in steps of 5% (relative modification
16 approach); 2) a subtraction from the non-pregnancy upper limit of TSH varying from -0.1 to -1.0 mU/L,
17 with modification of the non-pregnancy lower limit of FT4 varying from -5 to +5 pmol/L (-0.39 to +0.39
18 ng/dL; absolute modification approach) and 3) using fixed upper limits for TSH, varying from 3.0 to 4.5
19 mU/L, and fixed lower limits for FT4, varying from 5-15 pmol/L (0.39-1.17 ng/dL; fixed limit approach).
20 The choice for the range of modifications was based on previous recommendations (e.g. the fixed upper
21 limit of 4.0mU/L for TSH and 0.5 subtraction from this limit) and the optimal diagnostic performance in
22 this study, to keep the results organized. The results for each method were compared to the reference
23 standard (trimester-specific reference intervals), as is currently advised in international guidelines^{12,13}.

24

1 Diagnostic performance measures

2 The diagnostic performance of each assessed combination is described using the sensitivity (equivalent
3 to true positive rate, true positive rate among all with the disease according to the trimester-specific
4 method) and the PPV (equivalent to 1-false discovery rate, true positives among all with a positive test
5 result). Presenting the PPV, rather than the specificity, was preferred since the PPV is more informative
6 with regard to false positives for outcomes with a low prevalence ²². The aim was to maximize both
7 diagnostic performance markers, which poses a challenge, since maximizing sensitivity and the PPV is
8 often a trade-off.

9 The primary outcome was a single diagnostic performance measure, the F-score (also referred to as F1-
10 score), which is a combined measure of PPV (also referred to as 'precision') and sensitivity (also referred
11 to as 'recall') ²³. A higher F-score denotes a better overall diagnostic performance.

12 Prediction intervals and the I^2 statistic are presented to illustrate the expected inter-population variation
13 in diagnostic performance and between-study heterogeneity ^{21,24}. Prediction intervals are an attempt to
14 predict future individual values whereas confidence intervals give an indication of where the mean value
15 lies. To facilitate comparison of diagnostic performance markers between methods, interactive
16 heatmaps were constructed which can be found online ²⁵.

17 Statistical analyses

18 Diagnostic performance measures were calculated using 2x2 contingency tables (confusion matrices)
19 per cohort and pooled using random intercept logistic regression models utilizing maximum likelihood
20 for modeling between-study heterogeneity. This approach was chosen since it outperforms
21 conventional two-step inverse-variance approaches for sparse event datasets ^{26,27}. For each alternative
22 approach, the sensitivity, PPV and F-scores were calculated and compared with the trimester-specific

1 approach. All analyses were performed using R statistical software version 4.2.2²⁸, specifically using the
2 package 'meta'²⁹, 'ggplot2'³⁰ and 'heatmaply'³¹.

3

4 **Results**

5 After exclusions, the final study population comprised 52,496 participants included in 18 cohorts (Figure
6 1) of whom 8.6% were TPOAb positive (range across cohorts 5.7-17.1%; Supplemental table 1²¹). The
7 prevalence of thyroid function test abnormalities (in the first and second trimester, respectively)
8 according to the trimester-specific approach was 0.5% and 0.3% for overt hypothyroidism and 3.4% and
9 3.2% for subclinical hypothyroidism. The inclusion process and maternal demographics are described in
10 detail previously¹⁸. Cohort-specific prevalence of thyroid disease, reference limits, iodine status and
11 assay information can be found in Supplemental tables 2-6²¹. All figures are accompanied by
12 supplemental tables²¹ containing the diagnostic performance markers for each specific combination
13 (Figure 2 is an explanatory example of the diagnostic markers presented). To facilitate comparison of
14 diagnostic performance measures, an interactive version of the heatmaps including other diagnostic
15 performance measures can be found online and is also referred to throughout, as an alternative to the
16 supplemental tables²¹ (<https://www.consortiumthyroidpregnancy.org/heatmaps>²⁵).

17 **Diagnostic performance of alternative approaches**

18 Using the relative modification approach in the first trimester, the highest F-scores for overt
19 hypothyroidism were achieved with a relative subtraction of 5% for the upper reference limit of TSH and
20 a relative addition of 5% for the lower reference limit of FT4 (F-score 0.65; Figure 3A). The associated
21 sensitivity was 0.70 (95% confidence interval [CI] 0.47-0.86; 95% prediction interval [PI] 0.06-0.99; I²
22 64%), and the PPV was 0.64 (CI 0.54-0.74; PI 0.18-0.94; I² 45%; Figure 3A, Supplemental table 7²¹,
23 Interactive figures²⁵). For subclinical hypothyroidism the highest F-scores were achieved with a relative

1 subtraction of 20% for the upper reference limit of TSH and a relative subtraction of 15% for the lower
2 reference limit of FT4 (F-score 0.69; Figure 3B). Associated sensitivity was 0.91 (CI 0.67-0.98; PI 0.02-
3 1.00; I^2 95%) and PPV was 0.71 (CI 0.58-0.80; PI 0.20-0.96; I^2 95%; Supplemental table 8²¹, Interactive
4 figures²⁵).

5 Using the absolute modification approach in the first trimester, the highest F-scores for overt
6 hypothyroidism were achieved with a subtraction of either -0.1, -0.2 or -0.3 mU/L for the upper limit of
7 TSH and an addition of +1 pmol/L to the lower limit of FT4 and (F-score 0.62; Figure 3C). Associated
8 sensitivity (for upper limit TSH -0.2 mU/L) was 0.74 (CI 0.52-0.89; PI 0.08-0.99; I^2 66%) and PPV was 0.57
9 (CI 0.45-0.68; PI 0.24-0.84; I^2 39%; Supplemental table 9²¹, Interactive figures²⁵). For subclinical
10 hypothyroidism the highest F-scores were achieved with a subtraction of -0.8 mU/L from the upper limit
11 of TSH and a subtraction of either -1, -2, -3, -4 or -5 pmol/L from the lower limit of FT4 (F-score 0.64;
12 Figure 3D). Associated sensitivity (for lower limit FT4 -4 pmol/L) was 0.91 (CI 0.61-0.98; PI 0.01-1.00; I^2
13 95%) and PPV was 0.68 (CI 0.55-0.78; PI 0.20-0.95; I^2 95%; Supplemental table 10²¹, Interactive figures
14²⁵).

15 Using the fixed limit approach in the first trimester, the highest F-scores for overt hypothyroidism were
16 achieved with an upper limit of TSH of either 3.8, 3.9, 4.0, 4.1 and 4.4 mU/L and a lower limit of FT4 of
17 12 pmol/L (F-score 0.65; Figure 3E). Associated sensitivity (for upper limit TSH 4.0 mU/L) was 0.83 (CI
18 0.70-0.91; PI 0.41-0.97; I^2 0%) and PPV was 0.50 (CI 0.32-0.68; PI 0.05-0.95; I^2 70%; Supplemental table
19 11²¹, Interactive figures²⁵). For subclinical hypothyroidism the highest F- were achieved with an upper
20 limit of TSH of 3.2 mU/L and a lower limit of FT4 of either 5, 6, 7 or 8 pmol/L (F-score 0.70; Figure 3F).
21 Associated sensitivity (for lower limit FT4 8 pmol/L) was 0.99 (CI 0.88-1.00; PI 0.03-1.00; I^2 91%) and PPV
22 was 0.66 (CI 0.51-0.79; PI 0.11-0.97; I^2 96%; Supplemental table 12²¹, Interactive figures²⁵).

23 Additional analyses

1 In the second trimester, maximum F-scores were similar for the relative modification method, the
2 absolute modification approach and the fixed limit approach (Supplemental figure 1A-F²¹). However,
3 comparing the diagnostic performance measures of individual studies, the variability between studies
4 was very high, as reflected by overlapping confidence intervals for all methods, based on the wide
5 prediction intervals and based on high I^2 statistics for higher F-scores (Supplemental tables 13-18²¹). The
6 diagnostic performance of alternative methods to detect women for whom levothyroxine treatment is
7 indicated and those for whom treatment should be considered, according to ATA guidelines, in the first
8 trimester and second trimester were similar based on overlapping confidence intervals (Supplemental
9 figures 2, 3; Supplemental tables 19-30²¹).

10 **Discussion**

11 In this study, we systematically evaluated multiple standardized procedures to modify non-pregnancy
12 TSH and FT4 reference intervals with the aim of diagnosing the same individuals as having an abnormal
13 gestational thyroid function in line with the 'gold-standard' approach of center-specific and trimester-
14 specific reference intervals. Despite our efforts, we were unable to identify a standardized procedure
15 that achieved a satisfactory balance between sensitivity and PPV for gestational thyroid dysfunction
16 without considerable variability across different populations. These results underscore the inherent
17 challenge in balancing precise identification of gestational thyroid dysfunction with the practical
18 limitations of applying these diagnostic strategies universally in clinical settings, and indicate that
19 calculating local center and pregnancy-specific reference intervals for TSH and FT4 should still be
20 considered as current best practice.

21 Current recommendations on gestational reference interval definitions for TSH and FT4 are time and
22 resource consuming and are not feasible for most centers worldwide. The modification of non-
23 pregnancy reference intervals for the use in pregnancy could overcome feasibility problems. However, in

1 the current study we show that the variability in TSH and FT4 distributions leads to unacceptable
2 variation in diagnostic performance between cohorts. A possible explanation for this variation is that
3 even the non-pregnancy TSH and FT4 reference intervals are not an adequate reflection of the
4 distribution of thyroid function tests for a population if they are based on the manufacturer's
5 recommendation rather than local laboratory-specific establishment of the intervals. Methods for
6 determining reference intervals in pregnancy and outside pregnancy often differ, as current
7 recommendations on the establishment of reference limits in pregnancy include the local population
8 and are by definition a reflection of local TSH and FT4 distributions¹²⁻¹⁴, while reference limits outside
9 pregnancy are often supplied by the assay manufacturer, who mostly established reference intervals in
10 selected, non-pregnant populations^{32,33}. Global harmonization efforts for TSH and FT4 assays by the
11 International Federation of Clinical Chemistry and Laboratory Medicine (IFCC) Committee for
12 Standardization of Thyroid Function Tests (C-STFT) are ongoing to address this issue outside of
13 pregnancy, which could lead to an attenuation of this mismatch^{19,20}.

14 We also show that for overt hypothyroidism and for subclinical hypothyroidism, different and
15 sometimes opposing modifications of the reference limits of TSH and FT4 were needed to achieve
16 maximum diagnostic performance. For instance, when reviewing the relative modifications needed to
17 achieve the best diagnostic performance for overt hypothyroidism in the first trimester, we find that the
18 best F-score of 0.65 is achieved with the upper limit of TSH -5% and the lower limit of FT4 +5% (Figure
19 2A) while the best F-score for subclinical hypothyroidism of 0.69 is achieved with the upper limit of TSH -
20 20% and the lower limit of FT4 -15% (Figure 2B). We previously showed that the use of trimester-specific
21 reference intervals for FT4 are most important for the correct diagnosis of overt hypothyroidism while
22 for subclinical hypothyroidism the use of trimester-specific reference intervals for TSH are more
23 important¹⁸, which could explain the current results. This finding suggests that a uniform rule
24 established to diagnose both overt and subclinical disease, would be good at diagnosing one, at the cost

1 of incorrectly diagnosing the other. We also observe that the trends in diagnostic performance for a
2 treatment indication (Supplemental figure 2A, 2C, 2E ²¹) mostly overlap with the trend in diagnostic
3 performance for subclinical hypothyroidism (Figure 2B, 2D, 2F). This is because most women with a
4 treatment indication present with subclinical hypothyroidism with TPOAb positivity (73.6%) rather than
5 overt hypothyroidism (25.4%) or subclinical hypothyroidism with TSH>10 (1.1%; data not shown). Since
6 the prevalence of subclinical hypothyroidism is much higher than of overt hypothyroidism, it can be
7 expected that the best diagnostic performance of a test to detect a treatment indication is reached with
8 the same modifications as for subclinical hypothyroidism. This concept is important for future
9 recommendations on universal reference limits because diagnosing overt hypothyroidism, an entity with
10 an evident treatment indication, is generally prioritized in diagnostic strategies for gestational thyroid
11 dysfunction. However, failing to identify the more prevalent subclinical disease could also lead to
12 decreased benefits of (selective) screening. While we found no method with an agreeable trade-off in
13 terms of diagnostic performance, it is important to realize that the interpretation of diagnostic
14 performance of a test depends on the prior probability of disease ³⁴. This is a highly relevant concept
15 when thinking about differences between generalized population screening (with a low prior
16 probability) versus high-risk case-based screening (with higher prior probabilities). For example, for a
17 hypothetical diagnostic test with a sensitivity of 0.75 and a specificity of 0.99 (roughly equal to the tests
18 assessed in our study), a pre-test probability of 3% would result in a post-positive test probability (or
19 PPV) of 70% and a false discovery rate of 30%. Using the same sensitivity and specificity, a pre-test
20 probability of 10% would result in a post-positive test probability of 89% with a false discovery rate of
21 11%. The current study population consists of population-based cohort studies as a reflection of the
22 general population, which have a low prior probability of disease equal to the population prevalence
23 and similar to a universal screening approach. One option to improve how alternative reference interval
24 strategies could identify those with an abnormal thyroid function would be to increase the prior

1 probability of disease³⁴. This can be achieved by optimizing the identification of high-risk subgroups and
2 a risk-based screening approach, which could improve the accuracy of diagnostic strategies³⁵. Thus, the
3 implementation of universal screening will be inherently associated with the lowest prior probability of
4 disease and the highest rates of both over and underdiagnosis, especially if alternative strategies are
5 used to define thyroid function test abnormalities.

6 The heterogeneity between populations (as denoted by wide prediction intervals and high I^2 statistics)
7 underline that calculating local center and pregnancy specific reference intervals for TSH and FT4 should
8 still be considered as current best practice. However, other strategies for the improvement of the
9 diagnosis of gestational thyroid dysfunction might prove more effective. The trimester-specific approach
10 is currently accepted as the best diagnostic method for diagnosing thyroid dysfunction in pregnancy, but
11 the pragmatic division of the gestational period in trimesters does not necessarily reflect the
12 physiological changes of thyroid function tests during pregnancy³⁶⁻³⁸. Further studies are needed to
13 assess which gestational period reference intervals should be based upon to optimally identify the
14 women at increased risk of adverse events due to thyroid dysfunction, or if any form of standardization
15 to gestational age should be abandoned altogether. Current reference interval definitions are based on
16 outlying percentiles of TSH and FT4 distributions (2.5th and 97.5th percentiles), values above or below
17 those cutoffs were later shown to be associated with adverse pregnancy outcomes³⁹. With increasing
18 data availability in the literature, the ideal way to establish reference values would be to turn this
19 methodology around and base the cut-offs on the risk of adverse outcomes, similar to other fields^{40,41}.
20 Obvious adverse pregnancy events would be those associated with thyroid function tests in previous
21 studies such as preterm birth and offspring IQ scores^{3,4,6}. Since we did not identify an adequate or easily
22 implementable methodology to approach trimester-specific reference intervals in the current study, our
23 group will aim to establish risk-based decision limits.

1 In this study, we were able to leverage a large international dataset of multiple population-based
2 prospective cohort studies to assess novel strategies for diagnosing thyroid dysfunction in pregnancy.
3 The interpretation of the results of this study are limited to populations with sufficient or mild-to-
4 moderate iodine deficiency since studies with excessive status were excluded and no studies were
5 performed in an area of severe iodine deficiency. Additionally, multiple differences between the
6 included study populations, including differences in iodine supplementation, assays and determinants of
7 thyroid function tests, could have contributed to the variability in diagnostic performance of the non-
8 pregnancy reference interval adaptations assessed in this study. Adaptations of non-pregnancy
9 reference limits could be more accurate in specific populations, which we were not able to assess with
10 sufficient power. Nonetheless, this study reflects common practice, as these factors naturally vary
11 between populations. The results of the current study may not be optimally generalizable to present-
12 day populations since the inclusion periods for the majority of included cohorts were between the year
13 2000 and 2015. It is likely that determinants of thyroid function and assay calibrations standards have
14 changed over time ⁴². It can however be expected that large inter-population differences, as
15 demonstrated in this study, are still present to this day. Ongoing harmonization efforts by the IFCC could
16 improve the diagnostic performance of alternative strategies and future studies could assess if a
17 generalizable rule is more effective in cohorts established after the start of the harmonization efforts.

18 In conclusion, this is the first study to systematically quantify the diagnostic performance of
19 standardized modifications of non-pregnancy TSH and FT4 reference intervals in pregnancy. We show
20 that standardized modifications have poor overlap in diagnostic accuracy compared with cohort and
21 trimester-specific reference intervals, resulting in considerable variation in diagnostic performance
22 between populations. Future efforts should be turned towards studying the meaningfulness of
23 trimester-specific, pregnancy-specific reference intervals and the establishment of risk-based decision
24 limits.

1 Acknowledgements

2 The authors would like to gratefully acknowledge all participants, general practitioners, hospitals, and
3 midwives for their important contribution to the establishment of the cohorts and the resulting works.
4 Acknowledgements for individual cohorts are listed in the supplemental materials ²¹.

5 Data availability

6 The data that support the findings of this study are not publicly available due to local, national and
7 international restrictions aimed to protect the privacy of research participants.

9 References

- 10 1. Derakhshan A, Peeters RP, Taylor PN, Bliddal S, Carty DM, Meems M, Vaidya B, Chen L, Knight
11 BA, Ghafoor F, Popova PV, Mosso L, Oken E, Suvanto E, Hisada A, Yoshinaga J, Brown SJ, Bassols J,
12 Auvinen J, Bramer WM, Lopez-Bermejo A, Dayan CM, French R, Boucai L, Vafeiadi M, Grineva EN, Pop
13 VJM, Vrijkotte TG, Chatzi L, Sunyer J, Jimenez-Zabala A, Riano I, Rebagliato M, Lu X, Pizada A, Mannisto
14 T, Delles C, Feldt-Rasmussen U, Alexander EK, Nelson SM, Chaker L, Pearce EN, Guxens M, Steegers EAP,
15 Walsh JP, Korevaar TIM. Association of maternal thyroid function with birthweight: a systematic review
16 and individual-participant data meta-analysis. *Lancet Diabetes Endocrinol*. Jun 2020;8(6):501-510.
17 doi:10.1016/S2213-8587(20)30061-9
- 18 2. Toloza FJK, Derakhshan A, Mannisto T, Bliddal S, Popova PV, Carty DM, Chen L, Taylor P, Mosso
19 L, Oken E, Suvanto E, Itoh S, Kishi R, Bassols J, Auvinen J, Lopez-Bermejo A, Brown SJ, Boucai L, Hisada A,
20 Yoshinaga J, Shilova E, Grineva EN, Vrijkotte TGM, Sunyer J, Jimenez-Zabala A, Riano-Galan I, Lopez-
21 Espinosa MJ, Prokop LJ, Singh Ospina N, Brito JP, Rodriguez-Gutierrez R, Alexander EK, Chaker L, Pearce
22 EN, Peeters RP, Feldt-Rasmussen U, Guxens M, Chatzi L, Delles C, Roeters van Lennep JE, Pop VJM, Lu X,
23 Walsh JP, Nelson SM, Korevaar TIM, Maraka S. Association between maternal thyroid function and risk
24 of gestational hypertension and pre-eclampsia: a systematic review and individual-participant data
25 meta-analysis. *Lancet Diabetes Endocrinol*. Apr 2022;10(4):243-252. doi:10.1016/S2213-8587(22)00007-
26 9
- 27 3. Levie D, Korevaar TIM, Bath SC, Dalmau-Bueno A, Murcia M, Espada M, Dineva M, Ibarluzea JM,
28 Sunyer J, Tiemeier H, Rebagliato M, Rayman MP, Peeters RP, Guxens M. Thyroid Function in Early
29 Pregnancy, Child IQ, and Autistic Traits: A Meta-Analysis of Individual Participant Data. *J Clin Endocrinol*
30 *Metab*. Aug 1 2018;103(8):2967-2979. doi:10.1210/jc.2018-00224
- 31 4. Thompson W, Russell G, Baragwanath G, Matthews J, Vaidya B, Thompson-Coon J. Maternal
32 thyroid hormone insufficiency during pregnancy and risk of neurodevelopmental disorders in offspring:
33 A systematic review and meta-analysis. *Clin Endocrinol (Oxf)*. Apr 2018;88(4):575-584.
34 doi:10.1111/cen.13550

- 1 5. Han Y, Gao X, Wang X, Zhang C, Gong B, Peng B, Li J, Liu A, Shan Z. A Systematic Review and
2 Meta-Analysis Examining the Risk of Adverse Pregnancy and Neonatal Outcomes in Women with
3 Isolated Hypothyroxinemia in Pregnancy. *Thyroid*. May 2023;33(5):603-614. doi:10.1089/thy.2022.0600
- 4 6. Korevaar TIM, Derakhshan A, Taylor PN, Meima M, Chen L, Bliddal S, Carty DM, Meems M,
5 Vaidya B, Shields B, Ghafoor F, Popova PV, Mosso L, Oken E, Suvanto E, Hisada A, Yoshinaga J, Brown SJ,
6 Bassols J, Auvinen J, Bramer WM, Lopez-Bermejo A, Dayan C, Boucai L, Vafeiadi M, Grineva EN, Tkachuk
7 AS, Pop VJM, Vrijkotte TG, Guxens M, Chatzi L, Sunyer J, Jimenez-Zabala A, Riano I, Murcia M, Lu X,
8 Mukhtar S, Delles C, Feldt-Rasmussen U, Nelson SM, Alexander EK, Chaker L, Mannisto T, Walsh JP,
9 Pearce EN, Steegers EAP, Peeters RP. Association of thyroid function test abnormalities and thyroid
10 autoimmunity with preterm birth: a systematic review and meta-analysis. *JAMA*. Aug 20
11 2019;322(7):632-641. doi:10.1001/jama.2019.10931
- 12 7. Krassas GE, Poppe K, Glinoe D. Thyroid function and human reproductive health. *Endocr Rev*.
13 Oct 2010;31(5):702-755. doi:10.1210/er.2009-0041
- 14 8. Osinga JAJ, Derakhshan A, Palomaki GE, Ashoor G, Mannisto T, Maraka S, Chen L, Bliddal S, Lu X,
15 Taylor PN, Vrijkotte TGM, Tao FB, Brown SJ, Ghafoor F, Poppe K, Veltri F, Chatzi L, Vaidya B, Broeren
16 MAC, Shields BM, Itoh S, Mosso L, Popova PV, Anopova AD, Kishi R, Aminorroaya A, Kianpour M, Lopez-
17 Bermejo A, Oken E, Pirzada A, Vafeiadi M, Bramer WM, Suvanto E, Yoshinaga J, Huang K, Bassols J,
18 Boucai L, Feldt-Rasmussen U, Grineva EN, Pearce EN, Alexander EK, Pop VJM, Nelson SM, Walsh JP,
19 Peeters RP, Chaker L, Nicolaidis KH, D'Alton ME, Korevaar TIM. TSH and FT4 reference intervals in
20 pregnancy: a systematic review and individual participant data meta-analysis. *J Clin Endocrinol Metab*.
21 Sep 28 2022;107(10):2925-2933. doi:10.1210/clinem/dgac425
- 22 9. Springer D, Bartos V, Zima T. Reference intervals for thyroid markers in early pregnancy
23 determined by 7 different analytical systems. *Scand J Clin Lab Invest*. Mar 2014;74(2):95-101.
24 doi:10.3109/00365513.2013.860617
- 25 10. Negro R, Attanasio R, Papini E, Guglielmi R, Grimaldi F, Toscano V, Niculescu DA, Paun DL, Poiana
26 C. A 2018 Italian and Romanian Survey on Subclinical Hypothyroidism in Pregnancy. *Eur Thyroid J*. Nov
27 2018;7(6):294-301. doi:etj-0007-0294 [pii]10.1159/000490944
- 28 11. Toloza FJK, Ospina NMS, Rodriguez-Gutierrez R, O'Keefe DT, Brito JP, Montori VM, Maraka S.
29 Practice Variation in the Care of Subclinical Hypothyroidism During Pregnancy: A National Survey of
30 Physicians in the United States. *J Endocr Soc*. Oct 2019;3(10):1892-1906. doi:10.1210/js.2019-00196
- 31 12. Lazarus J, Brown RS, Daumerie C, Hubalewska-Dydejczyk A, Negro R, Vaidya B. 2014 European
32 thyroid association guidelines for the management of subclinical hypothyroidism in pregnancy and in
33 children. *Eur Thyroid J*. Jun 2014;3(2):76-94. doi:etj-0003-0076 [pii]
- 34 13. Alexander EK, Pearce EN, Brent GA, Brown RS, Chen H, Dosiou C, Grobman WA, Laurberg P,
35 Lazarus JH, Mandel SJ, Peeters RP, Sullivan S. 2017 Guidelines of the American Thyroid Association for
36 the Diagnosis and Management of Thyroid Disease During Pregnancy and the Postpartum. *Thyroid*. Mar
37 2017;27(3):315-389. doi:10.1089/thy.2016.0457
- 38 14. Thyroid Disease in Pregnancy: ACOG Practice Bulletin, Number 223. *Obstet Gynecol*. Jun
39 2020;135(6):e261-e274. doi:10.1097/AOG.0000000000003893
- 40 15. Bliddal S, Feldt-Rasmussen U, Boas M, Faber J, Juul A, Larsen T, Precht DH. Gestational age-
41 specific reference ranges from different laboratories misclassify pregnant women's thyroid status:
42 comparison of two longitudinal prospective cohort studies. *European Journal of Endocrinology*. Feb
43 2014;170(2):329-39.
- 44 16. Liu J, Yu X, Xia M, Cai H, Cheng G, Wu L, Li Q, Zhang Y, Sheng M, Liu Y, Qin X. Development of
45 gestation-specific reference intervals for thyroid hormones in normal pregnant Northeast Chinese
46 women: What is the rational division of gestation stages for establishing reference intervals for
47 pregnancy women? *Clin Biochem*. Apr 2017;50(6):309-317. doi:S0009-9120(16)30630-0
48 [pii]10.1016/j.clinbiochem.2016.11.036

- 1 17. Mehran L, Amouzegar A, Delshad H, Askari S, Hedayati M, Amirshkari G, Azizi F. Trimester-
2 specific reference ranges for thyroid hormones in Iranian pregnant women. Article. *J Thyroid Res.*
3 2013;2013doi:10.1155/2013/651517
- 4 18. Osinga JAJ, Derakhshan A, Feldt-Rasmussen U, Huang K, Vrijkotte TGM, Mannisto T, Bassols J,
5 Lopez-Bermejo A, Aminorroaya A, Vafeiadi M, Broeren MAC, Palomaki GE, Ashoor G, Chen L, Lu X, Taylor
6 PN, Tao FB, Brown SJ, Sitoris G, Chatzi L, Vaidya B, Popova PV, Vasukova EA, Kianpour M, Suvanto E,
7 Grineva EN, Hattersley A, Pop VJM, Nelson SM, Walsh JP, Nicolaides KH, D'Alton ME, Poppe KG, Chaker
8 L, Bliddal S, Korevaar TIM. TSH and FT4 reference interval recommendations and prevalence of
9 gestational thyroid dysfunction: quantification of current diagnostic approaches. *J Clin Endocrinol*
10 *Metab.* Sep 22 2023;doi:10.1210/clinem/dgad564
- 11 19. Thienpont LM, Van Uytendaele K, De Grande LAC, Reynders D, Das B, Faix JD, MacKenzie F,
12 Decallonne B, Hishinuma A, Lapauw B, Taelman P, Van Crombrugge P, Van den Bruel A, Velkeniers B,
13 Williams P, Tests ICfSoTF. Harmonization of serum thyroid-stimulating hormone measurements paves
14 the way for the adoption of a more uniform reference interval. *Clin Chem.* Jul 2017;63(7):1248-1260.
15 doi:10.1373/clinchem.2016.269456
- 16 20. Thienpont LM, Van Uytendaele K, Van Houcke S, Das B, Faix JD, MacKenzie F, Quinn FA, Rottmann
17 M, Van den Bruel A, Tests ICfSoTF. A progress report of the IFCC committee for standardization of
18 thyroid function tests. *Eur Thyroid J.* Jun 2014;3(2):109-16. doi:10.1159/000358270
- 19 21. Osinga JAJ, Derakhshan A, Korevaar TIM. Data from: reference intervals. Consortium on thyroid
20 and pregnancy. Updated 18-07-2023. <https://www.consortiumthyroidpregnancy.org/referenceintervals>
- 21 22. Lutgendorf MA, Stoll KA. Why 99% may not be as good as you think it is: limitations of screening
22 for rare diseases. *J Matern Fetal Neonatal Med.* 2016;29(7):1187-9.
- 23 23. Goutte C, Gaussier E. A probabilistic interpretation of precision, recall and F-score, with
24 implication for evaluation. Springer; 2005:345-359.
- 25 24. Riley RD, Higgins JP, Deeks JJ. Interpretation of random effects meta-analyses. *BMJ.* Feb 10
26 2011;342:d549. doi:bmj.d549 [pii]10.1136/bmj.d549
- 27 25. Osinga JAJ, Derakhshan A, Peeters RP, Korevaar TIM. Data from: Heatmaps. Consortium on
28 Thyroid and Pregnancy. Updated 31-01-2024. <https://www.consortiumthyroidpregnancy.org/heatmaps>
- 29 26. Lin L, Chu H. Meta-analysis of proportions using generalized linear mixed models. *Epidemiology.*
30 Sep 2020;31(5):713-717. doi:00001648-202009000-00016 [pii]10.1097/EDE.0000000000001232
- 31 27. Stijnen T, Hamza TH, Ozdemir P. Random effects meta-analysis of event outcome in the
32 framework of the generalized linear mixed model with applications in sparse data. *Stat Med.* Dec 20
33 2010;29(29):3046-67. doi:10.1002/sim.4040
- 34 28. R Core Team R: A language and environment for statistical computing. R Foundation for
35 Statistical Computing, Vienna, Austria. 2022. <https://www.R-project.org/>.
- 36 29. Balduzzi S, Rucker G, Schwarzer G. How to perform a meta-analysis with R: a practical tutorial.
37 *Evid Based Ment Health.* Nov 2019;22(4):153-160. doi:ebmental-2019-300117 [pii]10.1136/ebmental-
38 2019-300117
- 39 30. Wickham H. ggplot2: elegant graphics for data analysis. Springer-Verlag New York.
- 40 31. Galili T, O'Callaghan A, Sidi J, Sievert C. heatmaply: an R package for creating interactive cluster
41 heatmaps for online publishing. *Bioinformatics.* May 1 2018;34(9):1600-1602. doi:4562328 [pii]btx657
42 [pii]10.1093/bioinformatics/btx657
- 43 32. Brochure. Roche Diagnostics GmbH. Reference Intervals for Children and Adults Elecsys Thyroid
44 Tests 2009.
- 45 33. Brochure. Abbott Diagnostic Division. Architect system TSH ref 7K62. 2010.
- 46 34. Bours MJ. Bayes' rule in diagnosis. *J Clin Epidemiol.* Mar 2021;131:158-160. doi:S0895-
47 4356(20)31225-7 [pii]10.1016/j.jclinepi.2020.12.021

- 1 35. Osinga JAJ, Liu Y, Männistö T, Vafeiadi M, Tao FB, Vaidya B, Vrijkotte TGM, Mosso L, Bassols J,
2 López-Bermejo A, Boucai L, Aminorroaya A, Feldt-Rasmussen U, Hisada A, Yoshinaga J, Broeren MAC,
3 Itoh S, Kishi R, Ashoor G, Chen L, Veltri F, Lu X, Taylor PN, Brown SJ, Chatzi L, Popova PV, Grineva EN,
4 Ghafoor F, Pirezada A, Kianpour M, Oken E, Suvanto E, Hattersley A, Rebagliato M, Riaño-Galán I, Irizar A,
5 Vrijheid M, Delgado-Saborit JM, Fernández-Somoano A, Santa-Marina L, Boelaert K, Brenta G, Dhillon-
6 Smith R, Dosiou C, Eaton JL, Guan H, Lee SY, Maraka S, Morris-Wiseman LF, Nguyen CT, Shan Z, Guxens
7 M, Pop VJM, Walsh JP, Nicolaidis KH, D'Alton ME, Visser WE, Carty DM, Delles C, Nelson SM, Alexander
8 EK, Chaker L, Palomaki GE, Peeters RP, Bliddal S, Huang K, Poppe KG, Pearce EN, Derakhshan A, Korevaar
9 TIM. Risk Factors for Thyroid Dysfunction in Pregnancy: An Individual Participant Data Meta-Analysis.
10 *Thyroid*. Mar 28 2024;
- 11 36. Korevaar TIM, Medici M, Visser TJ, Peeters RP. Thyroid disease in pregnancy: new insights in
12 diagnosis and clinical management. *Nat Rev Endocrinol*. Oct 2017;13(10):610-622.
13 doi:10.1038/nrendo.2017.93
- 14 37. Glinioer D, de Nayer P, Bourdoux P, Lemone M, Robyn C, van Steirteghem A, Kinthaert J, Lejeune
15 B. Regulation of maternal thyroid during pregnancy. *J Clin Endocrinol Metab*. Aug 1990;71(2):276-87.
16 doi:10.1210/jcem-71-2-276
- 17 38. Andersen SL, Andersen S, Carle A, Christensen PA, Handberg A, Karmisholt J, Knosgaard L,
18 Kristensen SR, Bulow Pedersen I, Vestergaard P. Pregnancy week-specific reference ranges for
19 thyrotropin and free thyroxine in the north Denmark region pregnancy cohort. *Thyroid*. Mar
20 2019;29(3):430-438. doi:10.1089/thy.2018.0628
- 21 39. van den Boogaard E, Vissenberg R, Land JA, van Wely M, van der Post JAM, Goddijn M, Bisschop
22 PH. Significance of (sub) clinical thyroid dysfunction and thyroid autoimmunity before conception and in
23 early pregnancy: a systematic review. *Hum Reprod Update*. Sep-Oct 2011;17(5):605-619.
24 doi:10.1093/humupd/dmr024
- 25 40. Group HSCR, Metzger BE, Lowe LP, Dyer AR, Trimble ER, Chaovarindr U, Coustan DR, Hadden DR,
26 McCance DR, Hod M, McIntyre HD, Oats JJ, Persson B, Rogers MS, Sacks DA. Hyperglycemia and adverse
27 pregnancy outcomes. *N Engl J Med*. May 8 2008;358(19):1991-2002. doi:358/19/1991
28 [pii]10.1056/NEJMoa0707943
- 29 41. Xu Y, Derakhshan A, Hysaj O, Wildisen L, Ittermann T, Pingitore A, Abolhassani N, Medici M,
30 Kiemeny L, Riksen NP, Dullaart RPF, Trompet S, Dorr M, Brown SJ, Schmidt B, Fuhrer-Sakel D,
31 Vanderpump MPJ, Muendlein A, Drexel H, Fink HA, Ikram MK, Kavousi M, Rhee CM, Bensenor IM, Azizi
32 F, Hankey GJ, Iacoviello M, Imaizumi M, Ceresini G, Ferrucci L, Sgarbi JA, Bauer DC, Wareham N, Boelaert
33 K, Bakker SJL, Jukema JW, Vaes B, Iervasi G, Yeap BB, Westendorp RGJ, Korevaar TIM, Volzke H, Razvi S,
34 Gussekloo J, Walsh JP, Cappola AR, Rodondi N, Peeters RP, Chaker L, Thyroid Studies C. The optimal
35 healthy ranges of thyroid function defined by the risk of cardiovascular disease and mortality: systematic
36 review and individual participant data meta-analysis. *Lancet Diabetes Endocrinol*. Oct 2023;11(10):743-
37 754. doi:S2213-8587(23)00227-9 [pii]10.1016/S2213-8587(23)00227-9
- 38 42. Van Uytendanghe K, Ehrenkranz J, Halsall D, Hoff K, Loh TP, Spencer CA, Kohrle J. Thyroid
39 Stimulating Hormone and Thyroid Hormones (Triiodothyronine and Thyroxine): An American Thyroid
40 Association-Commissioned Review of Current Clinical and Laboratory Status. *Thyroid*. Sep
41 2023;33(9):1013-1028. doi:10.1089/thy.2023.0169 [pii]
- 42 10.1089/thy.2023.0169
- 43
- 44

1 **Figure legends**

2 Figure 1 – Flowchart of included cohorts and participants

3 Figure 2 - Diagnostic performance of modified non-pregnancy reference intervals for overt hypothyroidism
4 using relative modification

5 Figure 2 legend: Diagnostic performance for relative modifications of non-pregnancy reference intervals
6 for the diagnosis of overt hypothyroidism, presented as F-scores. The zoomed in section presents
7 additional diagnostic performance markers for selected modifications, of which an interactive version
8 can be found online (<https://www.consortiumthyroidpregnancy.org/heatmaps>).

9 Figure 3 – Diagnostic performance of modified non-pregnancy reference intervals for overt and
10 subclinical hypothyroidism

11 Figure 3 legend – Diagnostic performance of modified non-pregnancy reference intervals are presented
12 using a relative modification (A, B), absolute modifications (C, D) and fixed limits (E, F) for overt and
13 subclinical hypothyroidism, respectively, of which an interactive version can be found online
14 (<https://www.consortiumthyroidpregnancy.org/heatmaps>).

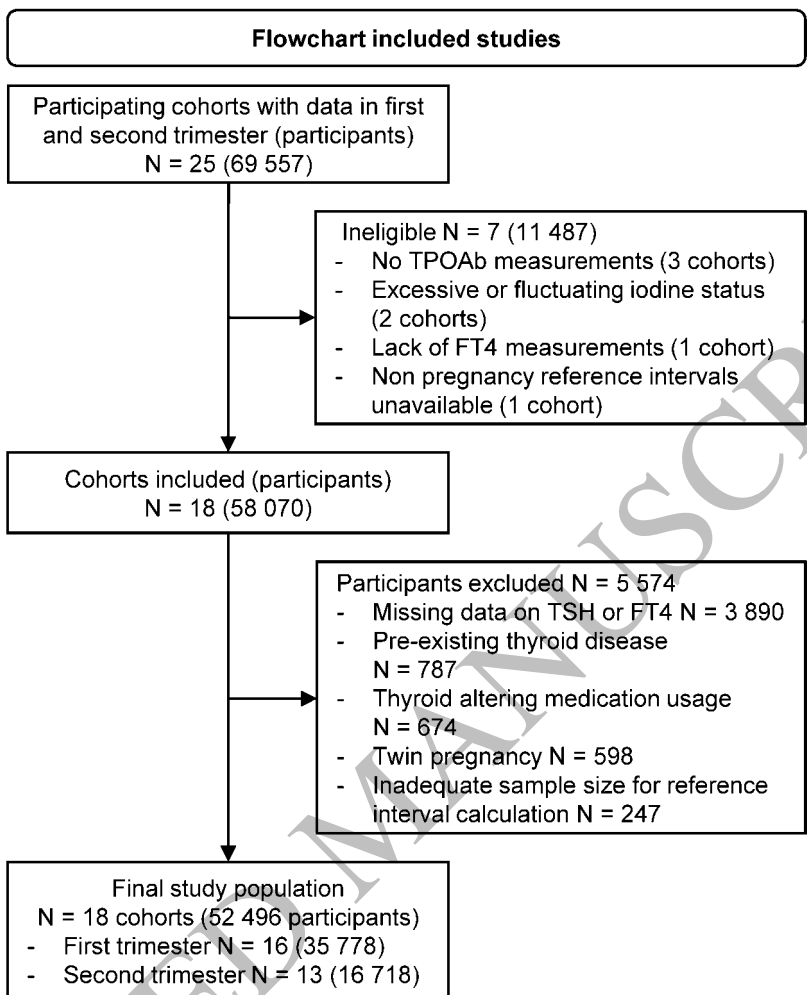


Figure 1
111x137 mm (x DPI)

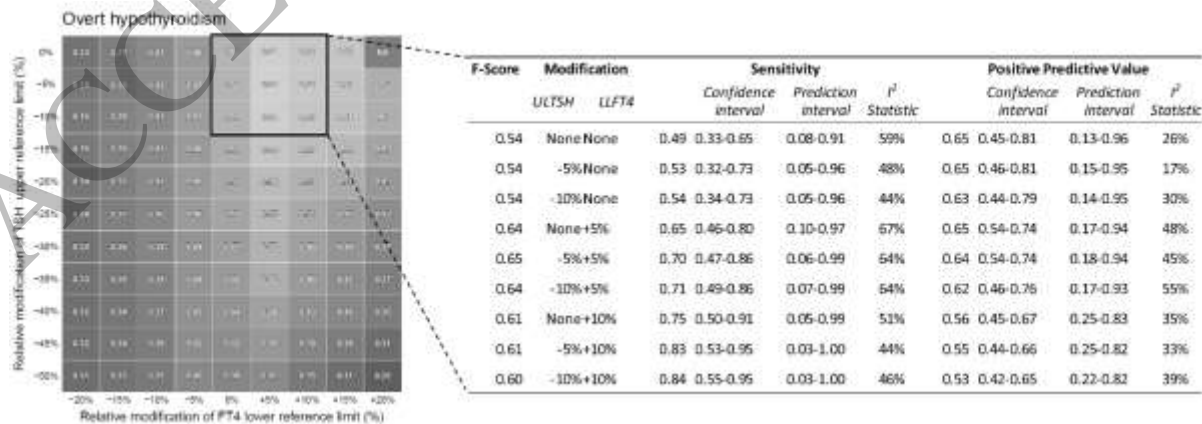


Figure 2
313x107 mm (x DPI)

1
2
3
4

5
6
7

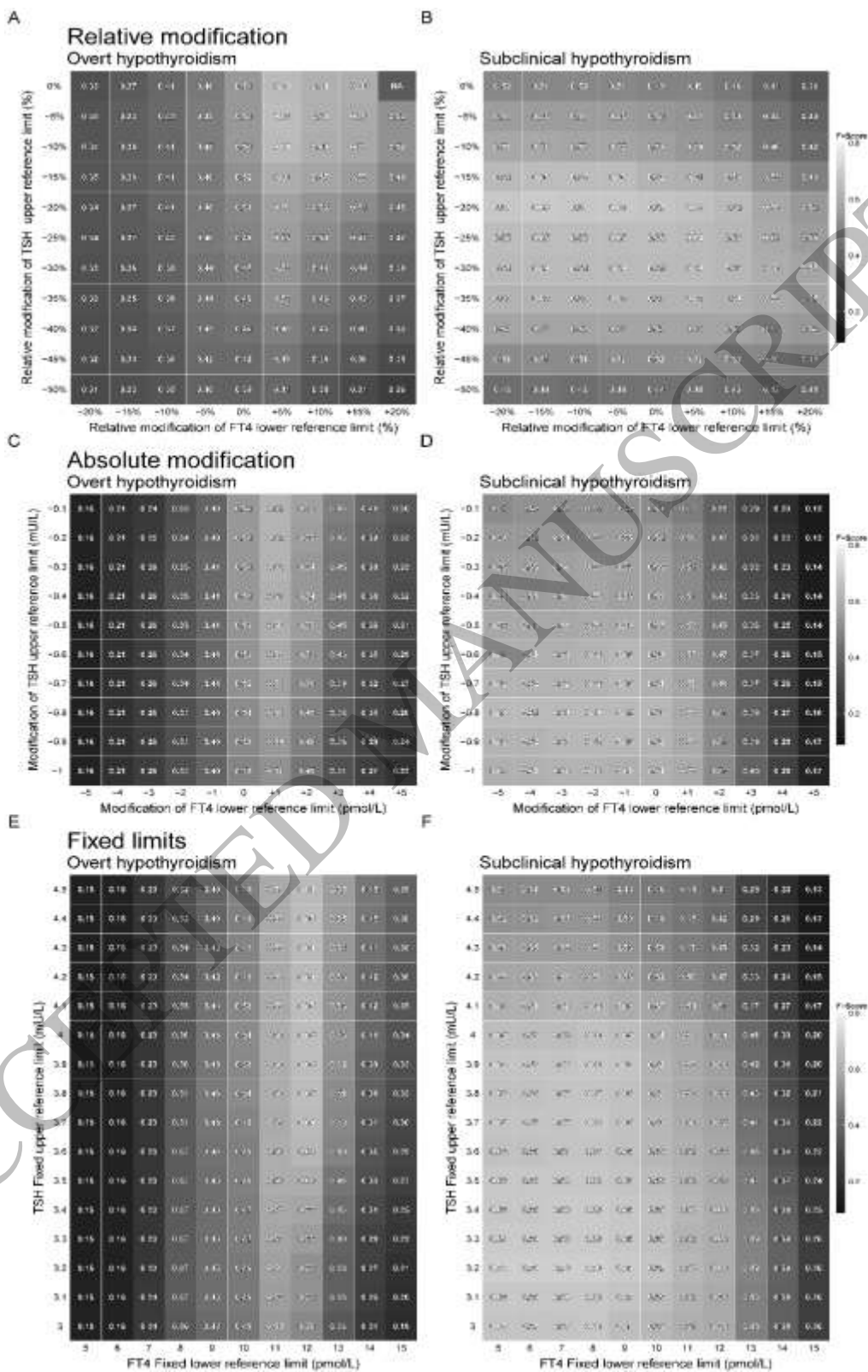


Figure 3
328x559 mm (x DPI)

1
2
3