

1 **EQ-5D-5L and SF-6D utility measures in symptomatic thyroid nodules: Acceptability**
2 **and psychometric evaluation**

3
4 Authors: Carlos K.H. Wong^{1*} PhD (ORCID: 0000-0002-6895-6071), Brian H.H. Lang^{2*} MS
5 (ORCID: 0000-0002-9362-0086), Hill M.S. Yu¹, Cindy L.K. Lam¹ MD

6 ¹ Department of Family Medicine and Primary Care, The University of Hong Kong

7 ² Division of Endocrine Surgery, Department of Surgery, The University of Hong Kong

8 * Joint corresponding author

9
10 Address for Correspondence:

11 Dr Carlos KH Wong

12 Department of Family Medicine and Primary Care, The University of Hong Kong

13 3/F, Ap Lei Chau Clinic, 161 Ap Lei Chau Main Street, Ap Lei Chau, Hong kong

14 Tel.: (852) 25185688, Fax No.: (852) 28147475, Email: carlosheo@hku.hk

15
16 Dr Brian HH Lang

17 Department of Surgery, The University of Hong Kong

18 Queen Mary Hospital, 102 Pokfulam Road, Pok Fu Lam, Hong kong

19 Tel.: (852) 22554232, Fax No.: (852) 28172291, Email: blang@hku.hk

20
21
22 **Keywords:** Quality-of-life; thyroid nodules; Psychometrics; Validity; Reliability;

23 **Running Title:** Psychometrics of EQ-5D-5L and SF-6D in thyroid nodules

24
25 **Informed consent:** Informed consent was obtained from all individual participants included
26 in the study.

27

28

1 **Abstract**

2

3 **Purpose:** To examine the acceptability, validity and reliability of EuroQoL 5-dimension 5-
4 level (EQ-5D-5L) and SF-6D health utility measures in patients with symptomatic benign
5 thyroid nodules.

6

7 **Methods:** Data from a randomized controlled trial (ClinicalTrials.gov Identifier:
8 NCT02398721) of 294 patients with symptomatic thyroid nodules were utilized for this
9 psychometric evaluation of HRQOL measurement. Three HRQOL questionnaires, generic
10 12-item Short Form Health Survey (SF-12v2) and EQ-5D-5L, SF-6D, were interviewer-
11 administered at baseline and 2 weeks afterwards. Responses to SF-6D were transformed to
12 SF-6D utility scores using Hong Kong population scoring algorithm derived by Standard
13 Gamble whereas response to EQ-5D-5L were mapped onto EQ-5D-3L response via interim
14 mapping algorithms and then converted to EQ-5D-5L utility scores using Chinese-specific
15 value set. Construct validity was determined by evaluating Spearman correlation between SF-
16 12v2 scores and utility scores. Two-week test-retest reliability was assessed using intra-class
17 correlation coefficient.

18

19 **Results:** No significant (>15%) floor and ceiling effects were observed for SF-6D utility
20 scores. The SF-6D utility scores had a moderate Spearman rank correlation with the SF-12v2
21 domain score providing evidence for adequate construct validity. The SF-6D utility scores
22 showed good test-retest reliability (0.794; range: 0.696-0.860). Better reliability was observed
23 in SF-6D utility score than in EQ-5D-5L utility score.

24

25 **Conclusions:** While the EQ-5D-5L instrument was less reproducible, the SF-6D instrument
26 appeared to be applicable, valid, and reliable measure in assessing the HRQOL of Chinese
27 patients with symptomatic benign thyroid nodules. Impact of utility score selection on the
28 effectiveness and cost-effectiveness of clinical interventions targeted to these patients needs
29 further exploration.

30

31 **Clinical trial number and registry:** NCT02398721, ClinicalTrials.gov

32

33 **Key points**

- 34
- 35 • The SF-6D instrument appeared to be applicable, valid, and reliable measure in
36 assessing the HRQOL of Chinese patients with symptomatic benign thyroid nodules.
 - 37 • The adoption of SF-6D utility score for health economic evaluation of clinical
38 management on patients with symptomatic benign thyroid nodule was supported.

Manuscript Text

INTRODUCTION

Thyroid nodules are common and can be discovered by clinical palpation in 5% of normal individuals and by high-resolution ultrasonography (USG) in 60% of the general population[1, 2]. Although most nodules are benign and do not grow, some nodules do become large and symptomatic [1, 2]. Some patients may experience pain in the neck, jaw or ear and cosmetic unsightliness that adversely impacts appearance and social functioning. If left untreated, some patients may even experience difficulty with breathing and swallowing and even hoarseness of voice leading to impaired physical and mental health-related quality-of-life (HRQOL).

Patients' HRQOL are not only affected by thyroid nodules themselves, but also by the subsequent treatments. There are essentially three therapeutic options for symptomatic benign thyroid nodules[3]. First is close surveillance with lifetime regular follow-ups. Although it is the least costly option, more costly treatments may follow if the disease progresses. Second is surgical resection or surgery. This is a "one-off" treatment and is associated with a high initial cost and patients may still need regular visits to check on their thyroid function afterwards. Third is non-surgical minimal invasive option which includes any form thermal ablative therapies or ethanol injection. This third option is usually a "one-off" treatment in most instances but it has the highest initial cost of the three options and regular visits are necessary. Recent longitudinal study on 40 patients with symptomatic nodule [4] reported that the thermal ablation after a two-year period was associated with physical and mental aspects of HRQOL improvements, possibly due to the reduction in compressive and cosmetic symptoms[5]. Based on a recent study[6], the non-surgical, minimally invasive option may

1 lead to greater improvement in HRQOL as compared to surgery but is unlikely to be cost-
2 effective over time at its current price. Besides treatment issues, patients enrolled to
3 conservative management or close monitoring strategy may add uncertainty, anxiety and
4 worries about the development of thyroid cancer. Often, the cost-effectiveness analysis of
5 health interventions requires the input of health utility score, a composite preference-weight
6 of HRQOL theoretically ranging from zero for death to one for full health. Although there is
7 no general recommendation for a reference case and preferable health utility measure for
8 economic evaluations in Hong Kong, a valid and reliable utility instrument is desirable to
9 influence decision-making and cost-effectiveness analysis of treatment options for benign
10 thyroid nodules.

11

12 Therefore, measurement of HRQOL and utility scores are important in the subjective
13 outcome and health economic evaluation of management for patients with symptomatic
14 thyroid nodules. Indeed, to our knowledge, there are currently no disease-specific HRQOL
15 instruments specifically designed for symptomatic benign nodular goitre, suggesting the use
16 of existing generic HRQOL and utility score instruments might be acceptable. Although one
17 previous study[6] estimated utility scores of patients with thyroid nodules through EuroQoL
18 5-dimension (EQ-5D) instrument, the psychometric properties of the commonly used utility
19 instruments, EQ-5D and Short-form 6-dimension (SF-6D) and their pairwise comparisons in
20 patients with symptomatic benign thyroid nodules were not determined. Therefore, the
21 present study aimed to assess the acceptability, validity, and reliability of the EQ-5D and SF-
22 6D instruments in patients with symptomatic benign thyroid nodules by doing a head-to-head
23 comparison on psychometric properties among instruments.

1 **METHODS AND PATIENTS**

2 *Study design and Patients*

3 This was an interim analysis of data collected from an ongoing prospective trial (already
4 approved by the local institutional review board and registered with www.clinicaltrials.gov
5 (NCT02398721). Over a 9-month period, consecutive patients presenting for the first-time
6 with a thyroid swelling were evaluated. To be eligible, first, the thyroid swelling had to be
7 benign (i.e. Bethesda class II on fine needle aspiration cytology (FNAC) within 3 months of
8 recruitment and a low or very-low suspicion sonographic pattern on USG). Second, the
9 swelling (which could either be a solitary nodule or a dominant nodule in a multinodular
10 gland) had to be causing obstructive or pressure symptoms. A swelling simply causing non-
11 specific neck complaints or cosmetic concern was not included. Third, the index nodule had
12 to have all three orthogonal dimensions ≥ 10 mm on USG. Also, patients with severe existing
13 medical co-morbidities or terminal malignancy were not eligible. Patients with previous
14 history of thyroid surgery, indeterminate, suspicious of malignancy or malignant FNAC (i.e.
15 Bethesda class III or above) were excluded from current analysis.

16 After informed consent, eligible patients were interviewed and were asked to fill in a
17 structured HRQOL questionnaire consisting of the traditional Chinese (Hong Kong) version
18 of SF-12 Health Survey version 2, SF-6D, EQ-5D 5-level (EQ-5D-5L), and questions on
19 socio-demographics. As for consistency in administration modes, all patients were
20 interviewer-administered because elderly patients with low literacy level were incompetent to
21 self-complete questionnaires. The SF-12v2, SF-6D and EQ-5D-5L questionnaires were
22 interviewer-administered at baseline, whereas three mentioned questionnaires and global
23 rating of change scale were interview-administered at 2-week after baseline to assess test-
24 retest reliability.

25

26 *Study Instruments*

1 EuroQoL 5-dimension 5-level (EQ-5D-5L)

2 The EQ-5D-5L is a generic preference-based measure developed by the EuroQol Group for
3 measurement of health-related quality-of-life, providing descriptions of five dimensions of
4 health status. The EQ-5D-5L has five domain scales (mobility, self-care, usual activities, pain
5 and discomfort, and anxiety and depression) and five levels for each domain. Since the
6 Chinese-specific EQ-5D-5L value set / tariff is currently not available, we applied a two-step
7 indirect approach to estimate EQ-5D-5L scores applicable for Chinese population, as adopted
8 in previous studies[7-9]. The first step was the application of an indirect interim mapping
9 method[10]. The EQ-5D-5L health status was transformed to EQ-5D-3L health status
10 according to the transition probability matrix. Finally, EQ-5D-3L health status were scored
11 according to a recently developed Chinese-specific EQ-5D-3L value set[11], ranging from -
12 0.149 for the worst health status ('33333') to 1 for the full health ('11111'). A higher score in
13 EQ-5D-5L indicated better HRQOL.

14

15 Short-form 12-item Health Survey (SF-12)

16 The Chinese (Hong Kong) SF-12 Health Survey version 2 (SF-12v2) has been validated and
17 normed on the general Chinese population in Hong Kong. It measures eight domains of
18 HRQOL on physical functioning (PF), role physical (RP), bodily pain (BP), general health
19 (GH), vitality (VT), social functioning (SF), role emotional (RE) and mental health (MH) on
20 a scale with theoretical range from 0 to 100. A higher score indicates better HRQOL. The
21 eight domain scores were aggregated based on population-specific weights to calculate two
22 summary scores, the physical (PCS) and mental component summary (MCS) scores.

23

24 Short-form 6-dimension (SF-6D)

25 In this study, the SF-6D utility score was directly converted from raw responses from SF-6D
26 instrument. Specifically, estimation of SF-6D utility score is possible with the available item

1 responses of either SF-6D or SF-12v2[12] instrument. Previous research show that SF-6D
2 utility score measured by SF-6D instrument was more sensitive and preferable than that
3 derived from SF-12v2 instrument[13]. The Hong Kong SF-6D value set[14, 15] was derived
4 by standard gamble valuation method[12]. The theoretical range of SF-6D utility score
5 ranged from 1 for full health to 0.315 for the worse possible health state according to Chinese
6 Hong Kong population-specific scoring algorithm [14, 15]. The SF-6D utility score serves as
7 preference weighting input to quality-adjusted life-year outcomes in economic evaluation.

8

9 Global Rating of Change Scale

10 The global rating of change scale is a single item “Compared to the first visit, how would you
11 rate your overall health now?” with 7-point Likert scale rating from ‘extremely worse’ (rating
12 of -3) to ‘extremely better’ (rating of +3) health condition compared to previous assessment at
13 baseline. It has widely used as an external criterion of change in health condition[16, 17].
14 Patients self-rated ‘0’ as indication of ‘stable’ in health condition change were selected for
15 test-retest reliability assessment.

16

17 *Statistical Analysis*

18 Descriptive statistics including mean, standard deviation (SD) and percentage of floor and
19 ceiling of each subscale and each summary scale were calculated.

20

21 The acceptability of instruments[18] was assessed using the proportion of missing values, and
22 proportion of patients giving highest possible and lowest possible responses, denoted as
23 ceiling effect and floor effect. At least 15% of subjects achieved the lowest or highest
24 possible score was considered as the presence of floor or ceiling effect, respectively. Based
25 on patients with stable in health condition over a 2-week duration, 2-week test-retest
26 reliability was assessed by intra-class correlation (ICC) coefficient using a value ≥ 0.7 to

1 indicate adequate reproducibility. A weighted Kappa[19] of <0.2 was interpreted as poor
2 agreement of individual domain responses between baseline and at 2-week assessments, 0.21-
3 0.4 as fair, 0.41-0.6 as moderate, 0.61-0.8 as good and >0.8 as very good.

4

5 The construct validity[20] of the SF-6D and EQ-5D-5L utility score was assessed using
6 correlation test against the SF-12v2 subscale scores holding similar constructs. Concerning
7 the possible violation of normality assumption for HRQOL and health utility scores,
8 Spearman's rank correlation test was used. We hypothesized that the utility scores were
9 moderately correlated with SF-12v2 summary scores, since those scores were composite
10 scores of different important HRQOL domain scores. Data analyses were conducted using
11 SPSS Windows 23.0 (IBM SPSS Inc., Chicago, IL, USA). P-value < 0.05 was statistically
12 significant.

13

14

1 **Results**

2 *Baseline patient characteristics*

3 A total of 314 patients were enrolled and of these, 20 patients (6.5%) with previous history of
4 thyroid surgery were excluded. Therefore, a total of 294 patients was included for analysis.

5 Table 1 shows the characteristics of patients included in this study. The male:female ratio was
6 1:6.4. The mean (\pm SD) age at presentation was 56.59 ± 11.67 years. The mean body mass
7 index was 23.33 ± 3.29 kg/m². All of them were euthyroid on presentation. The mean serum
8 TSH level was 1.33 ± 0.99 mIU/L (normal: 0.05 – 4.2mIU/L). The majority (n=225, 71.7%)
9 had multiple benign-looking nodules in their gland (i.e. multinodular goitre). The mean size
10 of the dominant (largest) nodule was 2.4 ± 2.3 cm (range: 1.6 – 3.8).

11
12 Figure 1 plots the response distributions related to all dimensions of both the EQ-5D-5L and
13 SF-6D instruments. Table 2 entails the descriptive statistics of health utility scores and
14 detailed response distributions of EQ-5D-5L and SF-6D domains. The highest proportion of
15 ‘no problems’ response referred to social functioning domain in SF-6D and self-care domain
16 in EQ-5D-5L. The mean EQ-5D-5L and SF-6D utility scores were 0.901 ± 0.113 and $0.773 \pm$
17 0.139 , respectively, where distribution of EQ-5D-5L was heavily left-skewed. Ceiling effect
18 was observed for EQ-5D-5L. No significant (>15%) floor and ceiling effects were observed
19 for SF-6D utility score. Table 3 depicts the two-week test-retest reliability among HRQOL
20 and utility instruments. For two-week test-retest assessment, 130 patients completed the full
21 set of HRQOL questionnaires. The mean follow-up duration was 17.1 days, with range from
22 9 days to 51 days. Based on an external criterion indicating the change in health status
23 compared to baseline assessment, 104 patients perceived stable in health condition over the 2-
24 week duration, where the retest time frame is commonly adopted for test-retest reliability[21].
25 The SF-6D utility scores showed good test-retest reliability (0.794; range: 0.696-0.860) but
26 ICC of EQ-5D-5L score was less than 0.7.

1 Weighted kappa, indicating agreement between two assessments, was interpreted as fair-
2 moderate (0.223-0.445) for 6 dimensions of SF-6D, and poor-moderate (0.149-0.344) for 5
3 dimensions of EQ-5D-5L. Better reliability was observed in SF-6D utility score than in EQ-
4 5D utility score.

5
6 The SF-6D utility score had a moderate-strong Spearman rank correlation with SF-12v2
7 domain and summary scores (0.461-0.630) that conceptually measures the similar construct
8 providing evidence for adequate construct validity (Table 4). Only fair-to-moderate Spearman
9 rank correlations between EQ-5D-5L score (0.310-0.552), and SF-12v2 domain scores were
10 observed. Particularly Spearman correlations between EQ-5D-5L score and SF-12v2
11 summary scores were fair-to-moderate (0.257-0.457).

12 13 **Discussion**

14
15 This psychometric validation study is the first report to evaluate acceptability, reliability, and
16 validity of HRQOL and utility instruments in patients with symptomatic benign thyroid
17 nodules. In evaluating acceptability, about half of patients self-reported the EQ-5D-5L health
18 profile of '11111', suggesting an observed ceiling effect for the EQ-5D-5L score. As such,
19 EQ-5D-5L score may reflect the better health status condition of benign thyroid nodules and
20 the lacking in room for improvement in health utility score due to clinical interventions.

21 There was no observed floor and ceiling effects for the SF-6D score, demonstrating
22 substantial acceptability of two instruments. The reliability of instruments are essential
23 psychometric property of the HRQOL and utility assessment. It refers to the ability to
24 reproduce consistent responses within a short period of time. The test-retest reliability of SF-
25 6D score was satisfactory, whereas EQ-5D-5L score had an inadequate reproducibility. The
26 response test-retest reliability measured by weighted kappa was fair to moderate for SF-6D

1 dimensions but poor to moderate for EQ-5D-5L dimensions, in line with test-retest reliability
2 of their utility scores. The reliability of SF-6D instrument was superior to that of EQ-5D-5L
3 instrument for symptomatic benign thyroid nodules. To ascertain the construct validity of
4 HRQOL and utility instrument, the generic SF-12v2 instrument was used as an anchor
5 assessing whether dimension of SF-12v2 holds good correlation with dimensions in other
6 instruments holding similar construct. Correlation between SF-12v2 dimension scores and
7 utility scores reflected the significant strength of SF-6D over EQ-5D-5L. All subscale and
8 summary scores of SF-12v2 were more correlated with SF-6D score than EQ-5D-5L score, in
9 part due to the fact that items in SF-6D were extracted from SF-12 instrument. Moreover, SF-
10 6D score had moderate correlation to the physical ($r=0.630$) and mental ($r=0.461$) composite
11 summary scores of SF-12v2 but weak correlation ($r=0.257$) between the EQ-5D-5L score and
12 mental composite summary score of SF-12v2 was observed.

13

14 One possible explanation for the differences in psychometric properties between SF-6D and
15 EQ-5D-5L utility scores was in part due to the differential scoring algorithms available for
16 the use in Hong Kong general and patient populations. Performance of acceptability and
17 psychometric properties may be influenced by the scoring algorithms used for conversion.
18 Unlike direct approach for calculation of SF-6D score, EQ-5D-5L utility score was derived
19 via two-step indirect approach, interim mapping algorithm plus EQ-5D-5L value set derived
20 by mainland China, which neither of them was not established based on Hong Kong Chinese
21 population. Moderate correlation ($r=0.556$) between two utility scores reflected the reality
22 that concepts in one utility score were not perfectly captured by another utility score, and thus
23 echoed importance of utility instrument selection. Impact of utility score selection, referring
24 to SF-6D score against EQ-5D-5L score, on quality-adjusted life-years, incremental cost-
25 effectiveness ratio and decision making of clinical interventions for thyroid nodule needs
26 further exploration.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

Limitations

Notably, four instruments used in current study may not be the best available HRQOL and utility instruments designed for benign thyroid nodules. A recent systematic review[22] appraising quality of 14 standardized instruments recommended the use of ThyPRO instrument, initially developed by Danish research group, for HRQOL assessment in patients with benign thyroid diseases. The cross-cultural validity of ThyPRO have been assessed at multiple clinical sites at multiple countries[23]. However, the Chinese version, irrespective of simplified or traditional version, has not been validated and was available at the time of study. Future works on psychometric properties of ThyPRO instrument or other existing HRQOL instruments specific to symptomatic thyroid nodules are warranted to enable the subjective evaluation of clinical interventions over time as well as integration of HRQOL measurement into routine care of thyroid nodules. Secondly, unlike psychometric testing on instruments used for moderate and severe disease, HRQOL scores in benign thyroid nodule were not hypothesized to have significant correlations with any clinical characteristics such as size of nodule, number of nodules and treatment modalities. Therefore, no analysis was conducted to assess correlations between HRQOL and clinical characteristics. Thirdly, the SF-12v2 scores were in principle more correlated with SF-6D scores than other scores, given the fact that SF-6D is the abbreviated form of SF-12v2. There was a possibility of favoring SF-6D in evaluation of construct validity. In current study, The SF-6D score was not calculated from seven items responses from SF-12v2 instrument, i.e. SF-12 derived SF-6D score. Both the SF-6D and SF-12v2 instruments were separately administered to each patient, gathering different set of responses from two instruments. The SF-12 was assumed to be ‘gold standard’ of HRQOL measurement testing against other HRQOL and utility scores. Nevertheless, this is presumably the best ‘gold standard’ among four instruments because this generic instrument has been validated in Hong Kong general population[24]. Since

1 instruments were not administered in a randomized order, there was also the possibility of
2 context effect (or ‘order effect’) on HRQOL and utility measurement. Fourth, study was the
3 secondary analyses of data collected within a RCT design. Since the RCT was not designed
4 for purpose of psychometric evaluation, this paper lacks important information for a
5 comprehensive acceptability investigation such as time for instrument completion, and
6 patients’ perceived views and feedback on each instrument. Finally, the patients were
7 sampled from endocrinology surgical outpatient clinic of one hospital in Hong Kong and so
8 our findings might be less generalizable in non-Chinese populations.

9

10 **Conclusions**

11

12 The SF-6D instrument demonstrated satisfactory acceptability and psychometric properties in
13 patients with symptomatic benign thyroid nodule, whilst the EQ-5D-5L instrument was less
14 reproducible than SF-6D instrument at the test-retest assessment. These findings supported
15 the use of the SF-6D instrument to evaluate the HRQOL in routine care as well as the
16 adoption of its utility score in the health economic evaluation of patients with symptomatic
17 benign thyroid nodule.

1 **Compliance with Ethical Standards**

2 Financial support for this study was provided by Health and Medical Research Fund
3 (HMRF#12132941) of Food and Health Bureau, HKSAR. The funding agreement ensured the
4 authors independence in designing the study, interpreting the data, writing, and publishing the
5 report.

6 Conflict of Interest: Carlos KH Wong, Brian HH Lang, Hill MS Yu, and Cindy LK Lam
7 declare that he / she has no conflict of interest.

8 Ethical approval: All procedures performed in studies involving human participants were in
9 accordance with the ethical standards of the institutional and/or national research committee
10 and with the 1964 Helsinki declaration and its later amendments or comparable ethical
11 standards.

12 Informed consent: Informed consent was obtained from all individual participants included
13 in the study.

14 **Author contributions:** CKHW wrote the manuscript, researched data, contributed to
15 statistical analysis and interpretation of results. BHHL contributed to study design,
16 acquisition of data and reviewed/edited the manuscript. HMSY contributed to acquisition of
17 data and reviewed/edited the manuscript. CLKL reviewed/edited the manuscript.

18

19

1

2 **References**

- 3 1. Gharib H, Papini E, Garber JR, Duick DS, Harrell RM, Hegedüs L et al. American
4 Association of Clinical Endocrinologists, American College of Endocrinology, and
5 Associazione Medici Endocrinologi Medical Guidelines for Clinical Practice for the
6 Diagnosis and Management of Thyroid Nodules – 2016 Update. *Endocrine Practice*.
7 2016;22(5):622-39. doi:10.4158/EP161208.GL.
- 8 2. Haugen BR, Alexander EK, Bible KC, Doherty GM, Mandel SJ, Nikiforov YE et al. 2015
9 American Thyroid Association Management Guidelines for Adult Patients with Thyroid
10 Nodules and Differentiated Thyroid Cancer: The American Thyroid Association Guidelines
11 Task Force on Thyroid Nodules and Differentiated Thyroid Cancer. *Thyroid*. 2016;26(1):1-
12 133. doi:10.1089/thy.2015.0020.
- 13 3. Che Y, Jin S, Shi C, Wang L, Zhang X, Li Y et al. Treatment of Benign Thyroid Nodules:
14 Comparison of Surgery with Radiofrequency Ablation. *American Journal of Neuroradiology*.
15 2015;36(7):1321-5. doi:10.3174/ajnr.A4276.
- 16 4. Valcavi R, Tsamatropoulos P. Health-related quality of life after percutaneous
17 radiofrequency ablation of cold, solid, benign thyroid nodules: a 2-year follow-up study in 40
18 patients. *Endocrine practice*. 2015;21(8):887-96. doi:10.4158/ep15676.or.
- 19 5. Spiezia S, Garberoglio R, Milone F, Ramundo V, Caiazzo C, Assanti AP et al. Thyroid
20 nodules and related symptoms are stably controlled two years after radiofrequency thermal
21 ablation. *Thyroid*. 2009;19(3):219-25. doi:10.1089/thy.2008.0202.
- 22 6. Yue W-W, Li X-L, Xu H-X, Lu F, Sun L-P, Guo L-H et al. Quality of Life and Cost-
23 Effectiveness of Radiofrequency Ablation versus Open Surgery for Benign Thyroid Nodules:
24 a retrospective cohort study. *Scientific Reports*. 2016;6:37838. doi:10.1038/srep37838.
- 25 7. Pan C-W, Sun H-P, Zhou H-J, Ma Q, Xu Y, Luo N et al. Valuing Health-Related Quality of
26 Life in Type 2 Diabetes Patients in China. *Medical Decision Making*. 2016;36(2):234-41.
27 doi:10.1177/0272989x15606903.
- 28 8. Pan C-W, Sun H-P, Wang X, Ma Q, Xu Y, Luo N et al. The EQ-5D-5L index score is more
29 discriminative than the EQ-5D-3L index score in diabetes patients. *Quality of Life Research*.
30 2015;24(7):1767-74. doi:10.1007/s11136-014-0902-6.
- 31 9. Cheung PHW, Wong CKH, Samartzis D, Luk KKD, Lam CLK, Cheung KMC et al.
32 Psychometric validation of the EuroQoL 5-Dimension 5-Level (EQ-5D-5L) in Chinese
33 patients with adolescent idiopathic scoliosis. *Scoliosis and Spinal Disorders*. 2016;11(1):19.
34 doi:10.1186/s13013-016-0083-x.
- 35 10. van Hout B, Janssen MF, Feng Y-S, Kohlmann T, Busschbach J, Golicki D et al. Interim
36 Scoring for the EQ-5D-5L: Mapping the EQ-5D-5L to EQ-5D-3L Value Sets. *Value in*
37 *Health*. 2012;15(5):708-15. doi:10.1016/j.jval.2012.02.008.

- 1 11. Liu GG, Wu H, Li M, Gao C, Luo N. Chinese Time Trade-Off Values for EQ-5D Health
2 States. *Value in Health*. 2014;17(5):597-604. doi:10.1016/j.jval.2014.05.007.
- 3 12. Brazier JE, Roberts J. The estimation of a preference-based measure of health from the
4 SF-12. *Medical care*. 2004;42(9):851-9.
- 5 13. Wong CKH, Lam ETP, Lam CLK. Comparison of direct-measured and derived short form
6 six dimensions (SF-6D) health preference values among chronic hepatitis B patients. *Quality*
7 *of Life Research*. 2013;22(10):2973-81. doi:10.1007/s11136-013-0403-z.
- 8 14. Lam CLK, Brazier J, McGhee SM. Valuation of the SF-6D Health States Is Feasible,
9 Acceptable, Reliable, and Valid in a Chinese Population. *Value in Health*. 2008;11(2):295-
10 303. doi:10.1111/j.1524-4733.2007.00233.x.
- 11 15. McGhee SM, Brazier J, Lam CL, Wong LC, Chau J, Cheung A et al. Quality-adjusted life
12 years: population-specific measurement of the quality component. *Hong Kong medical*
13 *journal*. 2011;17 (Suppl 6):17-21.
- 14 16. Guyatt G, Walter S, Norman G. Measuring change over time: Assessing the usefulness of
15 evaluative instruments. *J Chronic Dis*. 1987;40(2):171-8. doi:10.1016/0021-9681(87)90069-
16 5.
- 17 17. Guyatt GH, Deyo RA, Charlson M, Levine MN, Mitchell A. Responsiveness and validity
18 in health status measurement: A clarification. *Journal of Clinical Epidemiology*.
19 1989;42(5):403-8. doi:10.1016/0895-4356(89)90128-5.
- 20 18. Fitzpatrick R, Davey C, Buxton M, Jones D. Evaluating patient-based outcome measures
21 for use in clinical trials: a review. *Health Technology Assessment*. 1998;2(14):1-74.
22 doi:10.3310/hta2140.
- 23 19. Cohen J. Weighted kappa: nominal scale agreement with provision for scaled
24 disagreement or partial credit. *Psychological bulletin*. 1968;70(4):213-20.
- 25 20. Terwee CB, Bot SDM, de Boer MR, van der Windt DAWM, Knol DL, Dekker J et al.
26 Quality criteria were proposed for measurement properties of health status questionnaires.
27 *Journal of Clinical Epidemiology*. 2007;60(1):34-42. doi:10.1016/j.jclinepi.2006.03.012.
- 28 21. Marx RG, Menezes A, Horovitz L, Jones EC, Warren RF. A comparison of two time
29 intervals for test-retest reliability of health status instruments. *Journal of Clinical*
30 *Epidemiology*. 2003;56(8):730-5. doi:10.1016/S0895-4356(03)00084-2.
- 31 22. Wong CKH, Lang BHH, Lam CLK. A systematic review of quality of thyroid-specific
32 health-related quality-of-life instruments recommends ThyPRO for patients with benign
33 thyroid diseases. *Journal of Clinical Epidemiology*. 2016;78:63-72.
34 doi:10.1016/j.jclinepi.2016.03.006.
- 35 23. Watt T, Barbesino G, Bjorner JB, Bonnema SJ, Bukvic B, Drummond R et al. Cross-
36 cultural validity of the thyroid-specific quality-of-life patient-reported outcome measure,
37 ThyPRO. *Quality of Life Research*. 2015;24(3):769-80. doi:10.1007/s11136-014-0798-1.
- 38 24. Lam CLK, Tse EYY, Gandek B. Is the Standard SF-12 Health Survey Valid and

1 Equivalent for a Chinese Population? *Quality of Life Research*. 2005;14(2):539-47.

2

1

2 **Figure Legend**

3 Figure 1. Response distribution of domains in SF-6D instrument (upper) and EQ-5D-5L
4 instrument (lower).

Tables

Table 1. Baseline Characteristics of Patients with symptomatic thyroid nodules

Characteristics	Total (N=294)		Characteristics	Total (N=294)	
	N	%		N	%
Gender			The largest nodule		
Male	40	13.61	Length (cm, Mean±SD)	1.52	0.37
Female	254	86.39	Width (cm, Mean±SD)	1.26	0.42
Age (years, Mean±SD)	56.59	11.67	Height (cm, Mean±SD)	1.21	0.42
Education			Volume (mL, Mean±SD)	1.47	1.17
No formal schooling	10	3.42	Position		
Primary	54	18.49	Upper	25	9.26
Secondary	149	51.03	Middle	107	39.63
Tertiary or above	79	27.05	Lower	138	51.11
Marital Status			Side		
Married	211	72.01	Right	145	50.88
Single	36	12.29	Left	123	43.16
Widow(er)	24	8.19	Isthmus	17	5.96
Separated or divorced	22	7.51	Number of nodules		
Currently Working			1	37	13.03
Yes	180	61.43	2	41	14.44
No	113	38.57	3	47	16.55
Unknown			4	54	19.01
Monthly income			5	29	10.21
≤HKD20,000	139	47.44	≥6-<10	28	9.86
>HKD20,000	154	52.56	≥10	48	16.90
BMI (kg/m ² , Mean±SD)	23.33	3.29	Treatment for thyroid disease		
serum TSH (mIU/L, Mean±SD)	1.33	0.99	RAI	3	1.02
FT4 levels (Mean±SD)	18.09	10.43	Radiation	7	2.39
			Operation	38	12.93
			Thyroid-related medication	3	1.02

Note:

SD=standard deviation; BMI=body mass index; RAI=radioactive ablation iodine; TSH=thyroid-stimulating hormone; FT4=free thyroxine; FNAC=fine-needle aspiration cytology

Table 2. Descriptive statistics of health utility scores and response distributions related to all dimensions of EQ-5D-5L and SF-6D instruments at baseline

Utility score	Mean	Standard deviation	Observed Range	Theoretical Range	Floor (%)	Ceiling (%)
SF-6D score	0.773	0.139	0.41-1.00	0.32-1.00	0.0	1.3
EQ-5D-5L score	0.901	0.113	0.51-1.00	-0.15-1.00	0.0	47.1
Response Distribution						
SF-6D Dimension (%)	1	2	3	4	5	6
Physical functioning	24.5	32.7	28.6	12.2	1.7	0.3
Role functioning	65.0	12.9	7.1	15.0		
Social functioning	72.8	17.0	6.8	2.7	0.7	
Pain	29.6	24.5	22.1	13.9	6.8	3.1
Mental health	24.5	48.0	21.1	5.1	1.4	
Vitality	14.6	45.2	27.9	9.9	2.4	
EQ-5D-5L Dimension (%)	1 (No problems)	2 (Slight problems)	3 (Moderate problems)	4 (Severe problems)	5 (Unable to)	
Mobility	87.4	9.5	3.1	0.0	0.0	
Self-care	97.6	2.0	0.0	0.0	0.3	
Usual activities	91.5	6.8	1.4	0.3	0.0	
Pain/discomfort	50.3	34.7	13.9	0.3	0.7	
Depression/anxiety	72.8	19.0	7.5	0.7	0.0	

Table 3. Two-week test-retest reliability of health utility scores

Utility score	Intra-class correlation (n=104)		
	Estimate	95% CI	
SF-6D score	0.794	0.696	0.860
EQ-5D-5L score	0.575	0.374	0.712
	Weighted Kappa		
	Estimate	95% CI	
SF-6D			
Physical Functioning	0.278	0.159	0.396
Role Limitation	0.342	0.215	0.469
Social Functioning	0.340	0.178	0.501
Pain	0.223	0.112	0.334
Mental Health	0.445	0.319	0.570
Vitality	0.404	0.289	0.519
EQ-5D-5L			
Mobility	0.149	-0.007	0.306
Self-care	0.241	0.098	0.384
Usual activities	0.160	-0.033	0.352
Pain/discomfort	0.344	0.193	0.495
Depression/anxiety	0.231	0.072	0.390

Note: CI = confidence interval

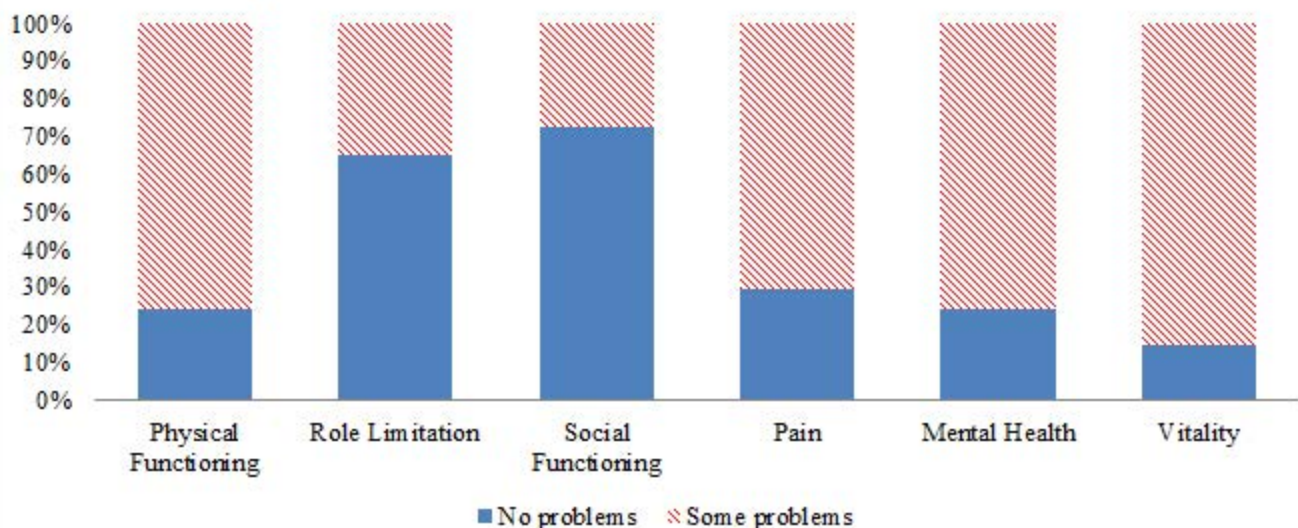
Table 4. Spearman Correlation Coefficients between SF-12v2 subscale and summary scores and the health utility scores

Utility score	SF-12v2 Subscale and Summary Scores									
	PF	RP	BP	GH	VT	SF	RE	MH	PCS	MCS
SF-6D score	0.581	0.589	0.610	0.495	0.537	0.610	0.501	0.517	0.630	0.461
EQ-5D-5L score	0.372	0.381	0.552	0.321	0.390	0.360	0.310	0.332	0.457	0.257

Note:

PF=physical functioning; RP=role physical; BP=bodily pain; GH=general health; VT=vitality; SF=social functioning; RE=role emotional; MH=mental health; PCS=Physical Composite Summary; MCS=Mental Composite Summary;

Response distribution of six domains in SF-6D



Response distribution of five domains in EQ-5D-5L

