2 **An investigation of reliability of the Sunderland Tracheosophageal Voice**
3 **Perceptual Scale**
4                                                                                              1
5
6
7

2 **Abstract**

3

4 Introduction: The consensus on how to effectively evaluate alaryngeal voice

5 outcomes remains limited. The Sunderland Tracheosophageal Voice Perceptual

6 scale (SToPS) was developed as a perceptual rating scale specifically for

7 tracheosophageal voice [1]. Currently, it is the only tracheosophageal voice specific

8 perceptual scale available and aims to address the limitations of previous scales.

9

10 Objective: To investigate inter rater reliability of the Sunderland

11 Tracheosophageal Voice Perceptual Scale when analysing alaryngeal voice across

12 a range of voice prostheses.

13

14 Methods: Prospective evaluation of inter rater reliability of the SToPS based on

15 audio recordings of 230 voice samples from 41 laryngectomy patients rated by 3

16 experts.  Interval data was analysed using Intraclass Correlation Coefficients (ICC)

17 while categorical data was analysed using Kappa.

18

19 Results:  ICC of above 0.6 was observed between raters for each prosthesis on a

20 majority of parameters demonstrating a good level of reliability.  Reliability was

21 fair (ICC of between 0.40-0.59) on Q11 (Articulatory precision) and Q12

22 (Paralinguistics).  Reliability was also fair (0.21-0.40) or slight (0.00-0.20) for Q2

23 (Tonicity), which was analysed using Kappa.  Kappa of above 0.61 signified a good

24 level of reliability.

25

2    Conclusions: This study demonstrates good rater <u>reliability</u> for the majority of

3    parameters on the <u>SToPS</u> scale, supporting the use of this tool within the clinical

4    realm.    <u>However further research is required to ascertain if any methods of</u>

5    <u>increasing inter rater reliability on those parameters which did not reach good</u>

6    <u>reliability can be identified</u>.

7    Level of evidence: 2b Individual cohort study

8

9

10

2   **Introduction**

3

4

5   Laryngectomy involves the removal of the larynx in its entirety, usually as a

6   treatment for advanced laryngeal cancer.   As a consequence, this surgery

7   profoundly affects the ability to communicate.    The gold standard for

8   communication rehabilitation after laryngectomy is surgical voice restoration

9   (SVR) [2] [3] also known as tracheosophageal voice.   This technique involves the

10   placement of a one way valved voice prosthesis in a puncture between the trachea

11   and oesophagus [4] [5]. The voice prosthesis shunts lung air from the oesophagus to

12   a vibratory segment within the reconstructed throat to produce tracheosophageal

13   voice.   The ultimate objective of SVR is to provide the patient with the optimal

14   voice possible without a larynx [1].   However, consensus on the most appropriate

15   measure of voice outcome post laryngectomy is lacking..

16

17   **Evaluation of post laryngectomy voice**

18

19   Although most of the empirical research concerning laryngeal voice has focused

20   on acoustic measures of frequency, intensity and duration, these measures do not

21   necessarily indicate how well an individual communicates in a social situation.

22   Auditory perceptual rating involves an expert listener judging a voice sample

23   according to different parameters [6] which may include intelligibility, voice quality

24   and acceptability [7].  Auditory perceptual evaluation of tracheosophageal voice

25   quality has been posited as the most valid measure of SVR outcome [1].   There are

26   a number of well-established voice quality rating scales which provide perceptual

27   parameters for the patients with a larynx including the Buffalo Voice Profile [8], the

28   Vocal Profile Analysis Scheme [9], Grade, Roughness, Breathiness, Asthenia, Strain

2    (GRBAS) scale [10], and Consensus Auditory Perceptual Evaluation of Voice [11]. Of

3    these, the strongest validity and reliability has been established for the GRBAS

4    [12,13]. The GRBAS has been used to assess auditory perceptual aspects of

5    tracheosophageal voice in several studies [14-16] [17,18]. However, use of the GRBAS to

6    measure perceptual aspects of tracheosophageal voice has been considered

7    suboptimal due to the fundamental differences in tracheosophageal and laryngeal

8    voice [1]. As the phonatory source of alaryngeal voice (vibratory segment) contrasts

9    significantly with that of laryngeal voice (vocal folds), the use of a rating scale

10   validated for the latter population poses limitations for post laryngectomy

11   patients. Additionally, some perceptual features of alaryngeal voice such as tone

12   and extraneous noise when covering the stoma to produce voice are unique and

13   central to tracheosophageal voice quality and are not included in the GRBAS scale.

14   Critically, studies which have used the GRBAS [14-18] or other perceptual scales [19,20]

15   have failed to specify an anchor baseline so it is unclear whether raters have

16   compared voice stimuli to that of normal laryngeal voice or optimal

17   tracheosophageal voice.

18

19   **STOPS**

20   The Sunderland Tracheosophageal Voice Perceptual scale (SToPS) was developed

21   as a perceptual rating scale specifically for tracheosophageal voice [1]. Currently, it

22   is the only tracheosophageal voice specific perceptual scale available. The SToPS

23   was developed as means of overcoming the major conceptual and methodological

24   problems inherent in other studies of tracheosophageal voice, such as poorly

25   defined terminology and impressionistic vocabulary [21]. The STOPS includes

26   specific and clear guidance to define terminology used for each parameter. In

2    addition, the SToPS crucially defines the anchor baseline for parameters as

3    optimal tracheosophageal voice rather than normal laryngeal voice.

4

5

6    **Reliability**

7

8    Measurement is a way of understanding, evaluating and differentiating

9    characteristics of people and objects [22] and forms the basis for making decisions

10   or drawing conclusions in scientific research. A crucial prerequisite for clinical

11   measurement is reliability. Reliability indicates the consistency and lack of errors

12   in a tool [22,23]. As the ability to simply produce voice with a prosthesis following

13   SVR is unlikely to be sufficient indication of functional ability to communicate in

14   everyday situations, it is of clinical relevance to investigate the reliability of the

15   SToPS. As intra rater reliability for expert raters had previously been established

16   as good or above for all parameters of the SToPS except for accent, reading ability

17   and articulatory precision [24], this study focuses on the investigation of inter rater

18   reliability.

19

20   **Aim**

21   To investigate inter rater reliability of the Sunderland Tracheosophageal Voice

22   Perceptual Scale

23   **Hypothesis**

24   Experts will not achieve a good level of inter rater reliability when they use the

25   SToPS to rate alaryngeal voice. Should a good level of inter rater reliability be

26   achieved, this will support the clinical relevance for the SToPS in identifying

27   functional tracheosophageal voice for patients post laryngectomy.

2

3

**Methods**

5
6  **SToPS**
7

The Sunderland Tracheosophageal Voice Perceptual Scale (SToPS) for professional raters was originally developed as a 14-item auditory perceptual scale divided into two domains: (i) Six Voice quality parameters (perceptual voice tonicity, strain, wetness, impairment of volume, impairment of social acceptability of voice and whisper), and (ii) seven parameters not related to voice quality (impression of intelligibility, stoma blast, impairment of fluency, impairment of articulatory precision, positive features of articulation, accent and poor reader) and an overall score voice rating.   The scale later underwent item reduction and now contains 10 parameters.  Ref

Tone relates to the amount of pressure used to produce tracheosophageal voice. The perceptual voice tonicity parameter is measured on an 11 point bipolar semantic scale reflecting the continuum of tone [25] from hypotonic (too little tone) to hypertonic (too much tone)[1]. As stenotic voice occurs only in the absence of tone it is measured with a separate arm to the tone scale [1].  As stenosis is either present or absent, it is not rated along a graded continuum.  For each individual voice sample, only one arm of the scale is chosen by a rater.  Each of the remaining 5 items in the voice quality parameters domain are measured on a 4 point equally appearing interval scale 0 (optimal tracheosophageal voice quality), 1 (mild), 2 (moderate) and 3 (severe).

2  Each of the parameters not related to voice quality, with the exception of positive

3  features of articulation is measured on a 4 point equally appearing interval scale

4  0 (optimal tracheosophageal voice quality), 1 (mild), 2 (moderate) and 3 (severe).

5  Positive features of articulation are measured on an alternatively worded 4 point

6  equally appearing interval scale 0 (neutral), 1 (good), 3 (excellent), and 4

7  (outstanding).

8

9  The parameter 'overall grade" is measured using a four point interval scale 0 =

10  Excellent; 1 = Good; 2 = Adequate; 3 = Poor.  This design is similar to the GRBAS

11  scale [10] except that the value 0 represents optimal tracheosophageal voice quality

12  as opposed to "normal" laryngeal voice quality.

13

14

15  **Raters**

16  Three  Speech and Language Therapy raters were chosen.  Each rater had at least

17  five years experience specialising in the rehabilitation of communication post

18  laryngectomy and other head and neck cancer patients and had completed

19  advanced training in the field.

20

21  **Training of raters**
22
23  Each rater participated in three hours of training with the investigator in the use

24  of the SToPS.  This training took place during two conference calls of 90 minutes

25  length and included practice ratings of ten anonymised audio samples of

26  laryngectomy participants reading the Rainbow Passage. During training queries

2     about individual items on the <u>SToPS</u> scale were raised.  These parameters were

3     discussed with the main author of the <u>SToPS</u>.  Clarifications provided were passed

4     onto all three raters regardless of how many raters had initially raised a query.

5

2 **Voice stimuli**
3
4 230 voice samples were elicited from 41 post laryngectomy participants. <u>Please</u>

5 <u>see table 1 for demographic details.</u> Participants were recruited from the

6 outpatient caseload of Head and Neck cancer patients at a large centre in __.

7 Exclusion criteria included participants without a voice prosthesis, less than 3

8 months post surgery or post operative oncological treatment. Each participant

9 trialled up to 6 randomised voice prostheses over 2 appointments within a 72

10 hour period. <u>Participants were blinded to prosthesis type</u> and a voice sample was

11 provided for each for each prosthesis. This data was used in a subsequent study

12 investigating the differences between voice prostheses in terms of voice outcome.

13

14 For each prosthesis trial, participants had a Speedlink SL-8691-SBK spes clip on

15 metal microphone (Speedlink, Weertzen, Germany) attached to their clothing 10

16 cm lateral to the stoma on the opposite side to the hand used to occlude the stoma

17 during voicing. All subjects produced voice by occluding their stoma rather than

18 depressing a humidification exchange device or using a hands free attachment.

19 Subjects read a short version of the Rainbow passage [26], (see appendix). This was

20 recorded onto a Sony ICD-PX820 Digital Voice Recorder with flash 2 GB (Sony,

21 Weybridge, UK) in MP3 format to be rated later by experts.

22

23 **Data analysis**

24 Recordings of voice samples with individual voice prostheses were extracted in

25 MP3 format and transferred to Final Cut Pro (Apple, California, USA) to allow titles

26 to be added to indicate anonymised subject number and anonymised voice

27 prosthesis letter. Voice samples were then exported to 3 Verbatim 4GB pinstripe

2    USB memory sticks (Verbatim, Surrey, UK). Raters were blinded to subject,

3    prosthesis type, gender, type of laryngectomy surgery (extended laryngectomy or

4    standard total laryngectomy) and history of radiotherapy and chemotherapy.

5    Voice samples were posted to 3 expert Speech and Language Therapy raters along

6    with blank numbered and lettered SToPS forms which corresponded to each voice

7    sample for each subject.

8

9    **Statistical analysis**

10

11    Data was entered and analysed in IBM SPSS (Statistical Product and Service

12    Solutions) version 23 (IBM Armonk, New York). The SToPS consists of 14

13    parameters, 13 of which (Q1, Q-Q14), are rated from 0-3 on an interval scale. A

14    further parameter, Q2 of the SToPS is rated on an 11 point bipolar semantic scale

15    which yielded categorical data.

16

17    Intraclass correlation coefficients (ICC) were used to analyse reliability of interval

18    scale parameters. A 2 way mixed model was chosen as each subject was assessed

19    by the same set of raters who have been purposely and not randomly selected [27,28].

20    0.6 ICC has previously been indicated as signifying a useful [29] and good [30] level of

21    reliability. ICC of between 0.40 and 0.59 has been defined as signifying a fair level

22    of reliability [30]. This interpretation was used to benchmark inter rater reliability

23    interval level data.

24    Cohen's kappa was used to analyse reliability of categorical data extracted from

25    Q2 (Perceptual Tonicity – amount of pressure used to produce tracheosophageal

Reliability of the SToPS

2   voice) on the SToPS scale.   In order to examine inter rater reliability for Q2, data

3   were recoded into 4 categories as follows:

4     • Hypotonic 5, 4, 3, 2, 1 was recoded as 1

5     • Tonic 0 was recoded as 2

6     • Hypertonic 5, 4, 3, 2, 1 was recoded as 3

7     • Stenosis 5 was recoded as 4

8   Reliability was calculated using kappa to see whether raters agreed 2x2

9     • Rater 1x Rater 2

10     • Rater 1x Rater 3

11     • Rater 2 x Rater 3

12   Analysis was conducted for reliability by prosthesis type by splitting data by

13   prosthesis type and then using cross tabs for kappa analysis by rater 2x2.

14   The Landis and Koch [31] classification of 0.61 as a good level of reliability, 0.41-0.60

15   as moderate reliability 0.21-0.40 as fair reliability and 0.00-0.20 as slight

16   reliability was used to analyse categorical level data.

17   **Results**

18
19   **Reliability of interval scale data**
20

21   The majority of parameters (Q1,Q3,Q5,Q7,Q8,Q9,Q13,Q14) reached an ICC of 0.60

22   indicating a good level of reliability (table 2).  Parameters, which did not reach an

23   ICC of 0.60 are highlighted in greyscale.  While reliability was not observed on Q4

24   ("Wetness" of voice quality) for the Blom Singer Low pressure voice prosthesis

25   nor on Q10 (Impairment of fluency) for the Blom Singer Duckbill voice prosthesis,

26   the ICC for both prostheses on both parameters approached good reliability.

27   Reliability for Q11 (Impairment of articulatory precision) was fair as opposed to

2   good except for the low-pressure prosthesis.  Reliability was reached amongst

3   raters for only three of the voice prostheses (Blom Singer Duckbill, Blom Singer

4   Low pressure and Provox NID) but was fair for other prostheses on Q12 (Positive

5   features of articulation – paralinguistics/diction).

6
7
8   **Reliability of Q2 Bipolar Semantic Scale data from STOPS**
9
10  Results of this analysis are outlined in table 2.1, 2.2 and 2.3

11  Reliability between raters was therefore only fair or slight for Q2 Tonicity across

12  voice prostheses.

13
14
15  **Discussion**

16  ***Expert raters inter rater reliability on the SToPS***
17
18  Reliability was investigated to ascertain whether there was a good level of

19  agreement among all three raters when using the SToPS to perceptually judge

20  voice. Parameters with poor reliability were Q2 – Perceptual Voice Tonicity, Q11-

21  Impairment of articulatory precision and Q12 – Positive features of articulation

22  (paralinguistics/diction).  Q2 relates to tonicity of the vibratory segment or the

23  amount of pressure used to produce alaryngeal voice.  Clinically, a patient with a

24  tonic voice will be able to produce fluent sound of adequate intensity without

25  effort.  A tonic voice has been defined as the ability to sustain /a:/ for 10 seconds

26  and produce 10-15 syllables per breath [32] or to sustain /a:/ 8 seconds and count

27  from 1-15 on one breath [33]. A previous study [1] examined inter rater agreement

28  between 12 Speech and Language Therapists and 10 ENT surgeons for Q2 of the

29  STOPS.  While inter rater agreement was only moderate for the raters as a whole,

2    it was good for the subgroup of Speech and Language Therapists with specific

3    voice experience.  Inter rater <u>reliability</u> was poor for three expert Speech and

4    Language Therapist raters in this study, each of whom had demonstrated a strong

5    understanding of tone within training sessions.  The experience of Speech and

6    Language Therapists in this study was primarily in head and neck cancer rather

7    than specifically with laryngeal voice.   This factor may account for the superior

8    agreement achieved on Q2 in a previous study [1] However, <u>the statistical</u>

9    <u>methodology</u> <u>which involved recoding data from Q2 from an 11 point equally</u>

10   <u>appearing interval scale into a four point categorical scale analysed with Kappa</u>

11   <u>may have</u> <u>been a further factor in the poor reliability found in this study.  Recoding</u>

12   <u>data in this manner changes tonicity from a continuum to a categorical scale and</u>

13   <u>thus may alter analysis.   The use of Cohen's Kappa for analysis is based on</u>

14   <u>absolute agreement.  In examining a parameter such as tonicity, it may not be</u>

15   <u>possible to attain absolute agreement within hypertonic and hypotonic aspects of</u>

16   <u>the continuum.   Both hypertonicity and hypotonicity contain a spectrum of</u>

17   <u>variety.</u>

18

19   Similarly the complexity of the scale used to measure Q2 may have influenced

20   levels of reliability achieved.

21

22   Q11- Impairment of articulatory precision demonstrated fair rater <u>reliability</u> only.

23   This parameter measures the degree of the lack of precision or "slurring" in

24   speech.  Lack of articulatory precision can be influenced by a number of factors

25   including fatigue and sometimes accent. During training of expert raters, Q11 was

26   not identified as one that needed further clarification. However, as the experience

2    of the expert raters involved in this study was predominantly with head and neck

3    cancer rather than with voice, it is possible that they were less familiar with the

4    defined baseline, which used the Vocal Profile Analysis scale as a reference.  This

5    factor may have accounted for the fair rater reliability on this parameter.  The final

6    parameter to demonstrate fair rater reliability was Q12 Positive features of

7    articulation (paralinguistics/diction).  Positive features of articulation refer to

8    diction, intonation or pause features that have an overall positive effect but are

9    not part of the voice signal.  Similarly to Q11, Q12 was not identified during

10   training as one that required further definition.  Fair rater reliability on this

11   parameter and on Q11 may simply reflect the difficulties of assessing articulation

12   and diction in laryngectomy patients, who present with an underlying disordered

13   voice.

14

15   This study examined the reliability of the SToPS across a range of voice prostheses

16   as part of the preparatory work for a later study examining differences between

17   prostheses in terms of voice quality.   Some voice prostheses notably differed in

18   levels of reliability achieved on parameters 4, 10, 11 and 12 of the STOPS.  The

19   attributes of different types of prostheses may affect tracheoesophageal voice and

20   therefore results of auditory perceptual analysis.  This is an area that may warrant

21   further research.

22

23   **Measurement of reliability**

24

25   This statistical methods used to analyse reliability in this study correspond with

26   those conventionally used for measurement of categorical data (Cohen's kappa) [34]

2      [22] and interval data (ICC) [27] [22]. As a previous study [1] [24] utilised weighted kappa to

3      evaluate reliability on all parameters of the STOPs, a possible limitation of this

4      study was the use of ICCs rather than kappa to measure interval data. ICCs have

5      been used extensively to measure reliability of pathological voice quality and for

6      this reason were utilised in this study. However, the use of ICC is largely based on

7      a framework of psychological testing. This framework substitutes listeners for test

8      items and voices for test subjects and implies that a new set of raters would

9      produce the same mean ratings for the same test voices. [35]. This approach has been

10     challenged as neither representing patterns of reliability nor overall agreement

11     for specific voice samples [35]. The alternative to ICC is weighted kappa. Weighted

12     kappa addresses the issue of Cohen's kappa failing to take into account the degree

13     of disagreement between raters by enabling greater weight to be assigned to some

14     rater disagreements than others [36]    However, kappa has been criticised as less

15     informative when used with more than 2 raters and analysing exact agreement

16     without accounting for "close" agreement [22]. In addition, with the use of kappa,

17     variance of subjects may be an issue, as a homogenous group of subjects is more

18     likely to show a high percentage of agreement, rather than a true reflection of

19     reliability [22]. The lack of consensus and limited evidence regarding the optimal

20     methodology to measure rater reliability in perceptual evaluation of both

21     laryngeal and alaryngeal voice supports the need for further research in this area.

22

23     **Conclusions**
24
25     This study investigated inter rater reliability of the Sunderland Tracheosophageal

26     Voice Perceptual Scale.   The findings presented in this study supports the SToPS

27     as a reliable tool for the auditory perceptual rating of alaryngeal voice.  However,

2    it is acknowledged that further research may be required to improve levels of

3    agreement for parameters related to tonicity, articulatory precision  and positive

4    features of articulation.

5

6

2    **References**
3
4    1.    Hurren A, Hildreth A, Carding P. Can we perceptually rate alaryngeal
5           voice? Developing the Sunderland Tracheoesophageal Voice Perceptual
6           Scale. Clinical otolaryngology : official journal of ENT-UK ; official journal
7           of Netherlands Society for Oto-Rhino-Laryngology & Cervico-Facial
8           Surgery 2009; 34:533-538.
9    2.    Kazi R, Nutting C, Evans PR, Harrington K. A short perspective on the
10         surgical restoration of alaryngeal speech. Southern medical journal 2009;
11         102.
12   3.    Hancock K, Ward E, Lawson N, van As-Brooks CJ. A prospective,
13         randomized comparative study of patient perceptions and preferences of
14         two types of indwelling voice prostheses. Int J Lang Commun Disord
15         2012; 47:300-309.
16   4.    Singer M, Blom E. An endoscopic technique for restoration of voice after
17         laryngectomy. The Annals of otology, rhinology, and laryngology 1980;
18         89:529-533.
19   5.    Blom E, Singer M. Surgical-Prosthetic Approaches for Post-laryngectomy
20         Voice Restoration *Laryngectomy Rehabilitation*. Houston: College Hill
21         Press, 1979.
22   6.    Carding P, Carlson E, Epstein R, Mathieson L, Shewell C. Formal
23         perceptual evaluation of voice quality in the United Kingdom. Logoped
24         Phoniatr Vocol 2000; 25:133-138.
25   7.    Doyle P, Eadie T. The perceptual nature of alaryngeal voice and speech. In:
26         Doyle P, Keith R, eds. *Contemporary Considerations In The Treatment And
27         Rehabilitation Of Head And Neck Cancer: Voice, Speech, And Swallowing*.
28         Austin, Texas: Pro Ed, 2005:113-139.
29   8.    Wilson D. Voice problems of children. Baltimore: Williams and Wilkins,
30         1987.
31   9.    Laver J. The phonetic description of voice quality. Cambridge: Cambridge
32         University Press, 1980.
33   10.   Hirano M. Clinical Examination of Voice New York: Springer Verlag, 1981.
34   11.   Kempster GB, Gerratt BR, Verdolini Abbott K, Barkmeier-Kraemer J,
35         Hillman RE. Consensus auditory-perceptual evaluation of voice:
36         development of a standardized clinical protocol. Am J Speech Lang Pathol
37         2009; 18:124-132.
38   12.   De Bodt MS, Wuyts FL, Van de Heyning PH, Croux C. Test-retest study of
39         the GRBAS scale: influence of experience and professional background on
40         perceptual rating of voice quality. J Voice 1997; 11:74-80.
41   13.   Wuyts FL, De Bodt MS, Van de Heyning PH. Is the reliability of a visual
42         analog scale higher than an ordinal scale? An experiment with the GRBAS
43         scale for the perceptual evaluation of dysphonia. Journal of Voice 1999;
44         13:508-517.
45   14.   Omori K, Kojima H. Neoglottic vibration in tracheoesophageal shunt
46         phonation. European archives of oto-rhino-laryngology : official journal of
47         the European Federation of Oto-Rhino-Laryngological Societies (EUFOS) :
48         affiliated with the German Society for Oto-Rhino-Laryngology - Head and
49         Neck Surgery 1999; 256:501-505.

15. Kazi R, Singh A, Mullan Get al. Can objective parameters derived from videofluroscopic assessment of post laryngectomy valved speech replace current subjective measures? An e tool based analysis. Clinical Otolaryngology 2006; 31:518.

16. Kazi R, Singh A, Venkitaraman R, Sayed S, RhysEvans P, Harrington K. Is electroglottography-based videostroboscopic assessment of post-laryngectomy prosthetic speech useful? Journal of Cancer Research and Therapeutics 2009; 5:85-92.

17. Kazi R, Kiverniti E, Prasad Vet al. Multidimensional assessment of female tracheoesophageal prosthetic speech. Clinical otolaryngology 2006; 31:511-517.

18. Schindler A, Mozzanica F, Ginocchio D, Invernizzi A, Peri A, Ottaviani F. Voice-related quality of life in patients after total and partial laryngectomy. Auris Nasus Larynx 2012; 39:77-83.

19. van As CJ, Hilgers FJ, Verdonck-de Leeuw IM, Koopmans-van Beinum F. Acoustical analysis and perceptual evaluation of tracheoesophageal prosthetic voice. Journal of voice 1998; 12:239.

20. Finizia C, Dotevall H, Lundstrom E, Lindstrom J. Acoustic and perceptual evaluation of voice and speech quality: a study of patients with laryngeal cancer treated with laryngectomy vs irradiation. Archives of otolaryngology--head & neck surgery 1999; 125:157-163.

21. vanAs C. Tracheoesophageal speech: A multidimensional assessment of voice quality. *Institute of Phonetic Sciences*. Amsterdam: University of Amsterdam, 2001:199.

22. Portney L, Watkins M. Foundations of Clinical Research - Applications to practice. New Jersey: Pearson Education International, 2009.

23. Lachin JM. The role of measurement reliability in clinical trials. Clinical trials (London, England) 2004; 1:553-566.

24. Hurren A. The development of a new rating scale for the perceptual assessment of tracheoesophageal voice quality outcome following total laryngectomy *Institute of Health and Society*: Univeristy of Newcastle, 2014:307.

25. Perry A. Vocal rehabilitation after total laryngectomy *Leicester School of Speech Pathology*. Leicester: De Montfort University, 1989:176.

26. Fairbanks D. Voice and Articulation Drill book. New York: Harper & Brothers, 1960.

27. Shrout P, Fleiss J. Intraclass correlation: Uses in assessing rater reliability. Psychological bulletin 1979; 86:420-428.

28. McGraw K, Wong S. Forming inferences about some intraclass correlation coefficients. . Psychol Methods 1996; 1:30-46.

29. Chinn S. Statistics in respiratory medicine. 2. Repeatability and method comparison. Thorax 1991; 46:454-456.

30. Cicchetti D. Guidelines, Criteria, and Rules of Thumb for Evaluating Normed and Standardized Assessment Instrument in Psychology. 1994.

31. Landis JR, Koch GG. The measurement of observer agreement for categorical data. Biometrics 1977; 33:159-174.

32. Lewin JS, Baugh RF, Baker SR. An objective method for prediction of tracheoesophageal speech production. J Speech Hear Disord 1987; 52:212-217.

2    33.    Blom ED, Singer MI, Hamaker R. Tracheoesophageal voice restoration
3           following total laryngectomy. Singular Pub. Group, 1998.
4    34.    Cohen J. A Coefficient of Agreement for Nominal Scales. Educational and
5           Psychological Measurement 1960; 20:37-46.
6    35.    Gerratt BR, Kreiman J. Theoretical and methodological development in
7           the study of pathological voice quality. Journal of Phonetics 2000; 28:335-
8           342.
9    36.    Cohen J. Weighted kappa: nominal scale agreement with provision for
10         scaled disagreement or partial credit. Psychological bulletin 1968;
11         70:213-220.
12