



Published in final edited form as:

Acad Radiol. 2008 October ; 15(10): 1322–1330. doi:10.1016/j.acra.2008.04.020.

Measurement Consistency from Magnetic Resonance Images

Dongjun Chung, M.A.^{*}, Moo K. Chung, Ph.D.[†], Reid B. Durtschi, M.S.[◇], R. Gentry Lindell, M.D.[‡], and Houri K. Vorperian, Ph.D.[◇]

^{*} Department of Statistics, University of Wisconsin-Madison, 1220 Medical Sciences Center, 1300 University Ave, Madison, WI 53706

[†] Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison, 1500 Highland Avenue # 437, Madison, Wisconsin 53705

[‡] Department of Radiology, University of Wisconsin Hospital and Clinics, 600 Highland Avenue, E1-311 Clinical Science Center, Madison, Wisconsin 53792

[◇] Waisman Center, University of Wisconsin-Madison, 1500 Highland Avenue # 430, Madison, Wisconsin 53705

Abstract

Rationale and Objectives—In quantifying medical images, length-based measurements are still obtained manually. Due to possible human error, a measurement protocol is required to guarantee the consistency of measurements. In this paper, we review various statistical techniques that can be used in determining measurement consistency. The focus is on detecting a possible measurement bias and determining the robustness of the procedures to outliers.

Materials and Methods—We review correlation analysis, linear regression, Bland-Altman method, paired *t*-test, and analysis of variance (ANOVA). These techniques were applied to measurements, obtained by two raters, of head and neck structures from magnetic resonance images (MRI).

Results—The correlation analysis and the linear regression were shown to be insufficient for detecting measurement inconsistency. They are also very sensitive to outliers. The widely used Bland-Altman method is a visualization technique so it lacks the numerical quantification. The paired *t*-test tends to be sensitive to small measurement bias. On the other hand, ANOVA performs well even under small measurement bias.

Conclusion—In almost all cases, using only one method is insufficient and it is recommended to use several methods simultaneously. In general, ANOVA performs the best.

Keywords for indexing

measurement consistency; bias; outlier; head; neck; Bland-Altman

Corresponding author contact information: Moo K. Chung, Ph.D., Associate Professor, Department of Biostatistics and Medical Informatics, Waisman Laboratory for Brain Imaging and Behavior, University of Wisconsin-Madison, 1500 Highland Ave #437, Madison, WI 53705, Office: (608) 217-2452, mkchung@wisc.edu, <http://www.stat.wisc.edu/~mchung/>.

Publisher's Disclaimer: This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

INTRODUCTION

This paper is motivated in part by the need to establish a reliable measurement protocol of head and neck structures involving both bony and soft tissue structures from magnetic resonance images (MRI) collected for the purpose of quantifying the growth pattern of various oral and pharyngeal structures or vocal tract structures (Vorperian et al. 2005, Vorperian et al, 2007). Figure 1 depicts a select set of such measurements obtained manually from MRI.

It is crucial to obtain accurate and reliable measurements particularly in developmental studies, and establish an accurate measurement protocol. Unfortunately, since the ground truth for manual measurements is never known, it is difficult to quantitatively determine if a given protocol produces consistent measurements. We have addressed this problem by placing reference landmarks, and obtaining repeated measures from MRIs by two trained raters. Next, using those paired measurements, we assessed the consistency of measurements of our measurement protocol. The purpose of this study is to determine the ideal analysis method to check for consistency of measurements. We will refer to this problem as the *measurement consistency problem*.

The measurement consistency problem occurs universally and it is of broad interest to researchers in diverse medical imaging disciplines. There are several major statistical approaches that have been used to check measurement consistency. The most widely used methods are correlation analysis, linear regression, paired *t*-test, and the Bland-Altman method (Krummenauer and Doll, 2000; Bland and Altman, 1986). A review on the measurement consistency problem can be found in Krummenauer and Doll (2000). Krummenauer and Doll (2000) state or conclude that using only one method is insufficient and several methods should be applied and compared. They also suggest making as many repeated measurements as time and cost permit for more accurate determination of measurement consistency.

In Bland and Altman (1986), the authors found that the correlation analysis, which is a popular method in establishing measurement consistency (Edvardsen et al, 2002; Liu et al, 2006; Van Oosterhout et al, 1995; Powell et al, 2000; Vallejo et al, 2000), is not appropriate. They proposed a visualization technique called the *Bland-Altman method* based on the difference between measurements. The detailed discussion on the Bland-Altman method can be found in Bland and Altman (1995) and Bland and Altman (2003). Braždžionytė and Macas (2007) claimed that the Bland-Altman method is more appropriate for assessing the measurement consistency, when compared to the correlation analysis and the linear regression. However, the shortcoming of the Bland-Altman's approach is that it is a visualization technique and lacks the numerical quantification.

Abate et al (1994) used the Bland-Altman method to analyze the measurement consistency between MRI and dissection for measuring adipose tissue mass. Powell et al (2000) used both a linear regression and the Bland-Altman method to analyze the measurement consistency between ultrasonic flowmeter measurements and phase-velocity cine MRI. Edvardsen et al (2002) used a paired *t*-test and the Bland-Altman method to compare the measurements from Tissue Doppler echocardiography with the measurements from MRI. Liu et al (2006) used the correlation coefficient to analyze the measurement consistency between manual delineation and automated segmentation of thermal coagulation on 3-D elastographic image.

In this paper, we review various quantitative techniques for determining measurement consistency, and provide an MRI study that describes the strength and the weakness of each technique. When comparing techniques, our main focus is on detecting the measurement bias and determining robustness to outliers. We provide further guidelines for using each of technique.

MATERIALS AND METHODS

Description of Head and Neck Imaging Data

MRIs from 10 male subjects between 0 and 4 years of age were used for this study. The landmarks for making measurements were placed on the MRI slice independently by two trained raters referred to as CC and RD. All landmarks and measurements were taken from the midsagittal slice of the MRI images from the imaging database. To insure unbiased placement of landmarks, RD and CC each placed landmarks on the image after suppressing the landmarks placed by the other. Thus each rater landmarked and measured the selected image independently of the other. All landmarks and measurements were made using the Sigma Scan Pro version 5 (Systat) and data was recorded onto a hard copy measurement sheet and entered into a measurement database for statistical analysis. All measurements were made in the centimeter unit.

Both CC and RD obtained measurements from ten MRIs independently at three separate times, resulting in a total of 60 measurements. These measurements were classified into four different categories: consistent, less consistent, biased, and with outliers. Of the 38 variables measured in the head and neck region, the following 6 variables are used to illustrate each case: Head length (HL), lower anterior facial height (LFH), anterior tongue length (ATL), Hyoid vertical distance from PNS (HVP), vocal tract length (VTL) and soft palate (SP). The definitions of those six variables are as follows (Figure 1).

Head Length (bony tissue – linear measurement)—The maximum linear distance from the glabella to the opisthocranium.

Lower Anterior Facial Height (bony and soft tissue - linear measurement)—The distance from the stomion to the gnathion. If the subject has an open mouth posture, the stomion was taken as the point at the antero-superior edge of the mandibular lip.

Anterior Tongue Length (soft tissue – curvilinear measurement)—The curvilinear distance along the dorsal superior contour of the tongue from the tongue tip to the intersection with the line dividing the hard palate and soft palate.

Vocal Tract Length (bony and soft tissue - curvilinear measurement)—The curvilinear distance along the midline of the tract (i.e. the distance along the midpoints of lines drawn between the inferior and superior boundaries of the vocal tract wall) starting at the level of the true vocal fold to the intersection with a line drawn tangentially to the lips.

Hyoid vertical distance from PNS (bony tissue – linear measurement)—The vertical distance from the inferior and anterior aspect of the hyoid bone to the level of the PNS.

Soft Palate Length (bony and soft tissue - curvilinear measurement)—The curvilinear distance from the posterior edge of the hard palate to the inferior edge of the uvula -- a projection of variable length from the free inferior border of the soft palate. The criterion used to identify the end of the hard palate and the beginning of the soft palate is a line drawn at the beginning of the hard palate/soft palate overlap.

The measurement errors themselves are relatively small and measured by the *average relative error* defined as

$$ARE = \frac{1}{n} \sum_{i=1}^n \frac{|RD_i - CC_i|}{|RD_i + CC_i|/2}, \quad (1)$$

where RD_i and CC_i be the i -th measurement of RD and CC respectively, and $n = 30$ be the number of measurements obtained by each rater. The average relative error for HL, LFH, ATL, HVP, VTL and SP are 0.016, 0.036, 0.041, 0.070, 0.046 and 0.1 respectively. The fairly large ARE of SP is caused by an outlier (Figure 2).

Figure 2 shows the scatter plot of the measurements of each head and neck structure. There are 30 data points on each scatter plot (three repeated measurements for 10 MRIs). The solid line ($y = x$) indicates the perfect consistency between two raters. Two raters measured HL and LFH *consistently* and most points are placed near $y=x$ line. ATL and HVP measurements are *less consistent* than LFH. For VTL, most points are under $y=x$ line and the measurements obtained by RD are *biased* against the measurements obtained by CC. For SP, there is an *outlier* caused by RD.

Correlation Analysis and Linear Regression

The correlation coefficient r measures the linear relationship between two variables, and ranges between -1 and 1 . If measurements are consistent, we expect to have a strong linear relationship and, in turn, correlation value close to 1 . On the other hand, if the measurements are less consistent, correlation value close to 0 is expected. Under the null hypothesis of $r=0$ (not consistent), the significance of correlation can be tested using a t -statistic with $n - 2$ degrees of freedom:

$$T = \frac{r \sqrt{n-2}}{\sqrt{1-r^2}}. \quad (2)$$

The correlation analysis has been previously used in measurement consistency (Edvardsen et al, 2002; Liu et al, 2006; Van Oosterhout et al, 1995; Powell et al, 2000; Vallejo et al, 2000). However, as we will show in the result section, it is not a proper procedure.

Alternately a linear regression can be used to determine the measurement consistency (Braždžionytė J and Macas A., 2007; Powell et al, 2000). The following regression model is used to fit measurements:

$$RD_i = \beta_0 + \beta_1 \times CC_i + \varepsilon_i.$$

When RD and CC are consistent, we expect the regression slope β_1 to be close to one. By testing if the slope is equal to one, we can quantitatively determine the consistency. The regression fit is given in Figure 2. Since the slope is proportional to the correlation coefficient, both the correlation analysis and the linear regression are equivalent approaches although this equivalence is not exploited previously (Chatterjee et al, 2000). Similarly one can test if the intercept β_0 is close to zero for testing a bias of if one rater is systematically obtaining larger or smaller measurements compared to the other rater.

Bland-Altman method and paired t-test

Although the Bland-Altman method has been discussed in various literatures (Bland and Altman, 1984; Bland and Altman, 1995; Bland and Altman, 2003; Krummenauer and Doll, 2000; Braždžionytė and Macas, 2007; Abate et al, 1994; Powell et al, 2000; Edvardsen et al,

2002), we briefly explain here for the completeness of the paper. Let d_i be the measurement difference, i.e. $d_i = CC_i - RD_i$. The measurement difference is the estimated bias of measurements between the two raters. Let \bar{d} and S_d^2 be the mean and the variance of the difference. Bland and Altman plotted d_i versus the average of measurements of two raters, with the reference lines, \bar{d} , $\bar{d} - 1.96S_d$ and $\bar{d} + 1.96S_d$ (Bland and Altman, 1984). The range between $\bar{d} - 1.96S_d$ and $\bar{d} + 1.96S_d$ provides the “limit of agreement” (Figure 3).

The weakness of the Bland-Altman method is that the measurement consistency is mainly determined visually without statistical significance attached to the plot. To give the statistical significance to the Bland-Altman method procedure, a paired t -test can be used. We test if the measurement difference is statistically small enough using the test statistic

$$T = \frac{\bar{d}}{\sqrt{S_d^2/n}}, \quad (3)$$

which is distributed as the t -distribution with $n - 1$ degrees of freedom.

ANOVA and within-rater consistency

All the previous methods can determine consistency between a set of paired measurements. When there are more than two raters the previous methods cannot be applied directly without significant modification. We propose to use ANOVA approach for more general cases. The strength of ANOVA is that it can be used to determine both between- and within-rater measurement consistency. If we have information about how each rater measures the same MRI consistently, we can determine who is more consistent. This additional information can be used to train less consistent raters further.

Let X_{ijk} be the k -th measurement on the j -th MRI by the i -th rater. Then, the two-way ANOVA model is given as

$$X_{ijk} = \mu + \alpha_i + \beta_j + \alpha\beta_{ij} + \varepsilon_{ijk}.$$

The usual measurement consistency between CC and RD can be determined by testing $\alpha_{CC} = \alpha_{RD}$. The interaction term $(\alpha\beta)_{ij}$ is used to determine the within-rater consistency for 10 MRIs. The within-rater consistency can be determined by simultaneously testing $\alpha\beta_{CC,1} = \dots = \alpha\beta_{CC,10} = \alpha\beta_{RD,1} = \dots = \alpha\beta_{RD,10}$.

We can also visualize the within-rater consistency patterns using *the box plot* (Tukey, 1977). The Box plot is one of popular data visualization methods and is drawn in the following way (Martinez and Martinez, 2005). First, we obtain the value corresponding to 25%, 50%, and 75% of the sorted observations. They are called the lower quantile q_1 , the median q_2 and the upper quantile q_3 respectively. The median q_2 provides the information about the center such that the half of the data is smaller than q_2 and the other half is larger than q_2 . Then, we draw “the box” from q_1 to q_3 with the line of q_2 within the box. This box provides the range containing 50% of the data around q_2 . Finally, we draw one line from q_1 to $q_1 - 1.5(q_3 - q_1)$ and another line from q_3 to $q_3 + 1.5(q_3 - q_1)$, which are called as “the whisker.” In box plot, the observations outside $q_1 - 1.5(q_3 - q_1)$ and $q_3 + 1.5(q_3 - q_1)$ are determined as potential outliers.

Let $d_{j,k}$ be the difference between k -th measurement of j -th MRI and the average measurements of j -th MRI by one fixed rater. The box plot of $d_{j,k}$ shows the diversity of measurements for each MRI. We can see how consistent each MRI is measured by a specific rater using the box plot of $d_{j,k}$. We can visually compare within-rater consistency by comparing the box plots between CC and RD (Figure 4).

RESULTS

Correlation Analysis and Linear Regression

The linear regression fitting line for each head and neck structure appears as the dotted line in Figure 2. The measurements are more consistent when the dotted line is close to the solid line ($y=x$). Two lines were very close in HL, LFH, ATL, HVP and SP. In contrast, the dotted line was far from the solid line in VTL. The correlation coefficients of HL, LFH, ATL and HVP were 0.963, 0.987, 0.880 and 0.871 respectively (p-value < 0.001 in all cases). This implies the measurements are consistent for HL, LFH, ATL and HVP and this coincides with what we observe in Figure 2.

On the other hand, the correlation coefficient was 0.875 (p-value < 0.001) for VTL and this seems to contradict with Figure 2 because there was a clear systematic bias in VTL. We can infer from this that the correlation coefficient cannot detect the measurement inconsistency. Correlation coefficient of SP was 0.089 (p-value = 0.639). In spite of existing consistency between CC and RD, an outlier made the correlation coefficient close to 0. After removing the outlier, correlation coefficient of SP becomes 0.673 (p-value < 0.001). This implies that the correlation coefficient is very sensitive to outliers.

In summary, the correlation analysis has difficulty detecting the inconsistency between measurements. This is due to the fact that the correlation coefficient shows the degree of association not the degree of consistency. The correlation analysis is very sensitive to outliers. As a result, the correlation analysis is not appropriate as the measurement consistency analysis.

Bland-Altman method and paired t-test

Figure 3 shows the Bland-Altman plots for each head and neck structures. Even though these plots provide the degree of bias, it is not easy to infer about the measurement consistency based on these plots. This is because the Bland-Altman method lacks statistical significance attached to the plot. Moreover, in measuring SP, one outlier severely increases the limit of agreement. In summary, Bland-Altman method is not appropriate as a technique for determining measurement consistency.

The paired t -test indicates that there is significant inconsistency in measuring LFH (p-value = 0.008) and HVP (p-value = 0.038) although the scatter plots of LFH and HVP in the Figure 2 show measurement consistency. This contradiction can happen if one rater systematically measures either larger or smaller than the other rater. When this systematic bias becomes larger than the measurement variance, this contradiction will happen.

In summary, the paired t -test can detect measurement bias between raters fairly well in most cases. However, it may fail when one rater systematically measure either larger or smaller than the other rater.

ANOVA and within-rater consistency

ANOVA results show that measurements are consistent between raters in measuring HL (p-value = 0.110), LFH (p-value = 0.517), ATL (p-value = 0.576), HVP (p-value = 0.937) and SP (p-value = 0.279) but not in measuring VTL (p-value = 0.029). This finding exactly

coincides with what we found in Figure 2. The box plots in the Figure 4 and the interaction term in ANOVA show which rater performs better. RD is significantly more consistent than CC in measuring HL (the first row in the Figure 4; p -value < 0.001). CC is more consistent than RD in measuring LFH (the second row in the Figure 4) but the difference was not significant (p -value = 0.770). RD is significantly more consistent than CC in measuring ATL (the third row in the Figure 4; p -value = 0.008). CC is more consistent than RD in measuring HVP (the fourth row in the Figure 4) but the difference was not significant (p -value = 0.152). RD is significantly more consistent than CC in measuring VTL (the fifth row in the Figure 4; p -value = 0.016). CC is more consistent than RD in measuring SP (the sixth row in the Figure 4) but this difference was not significant (p -value = 0.115).

In summary, ANOVA extends the paired t -test method by considering the within-rater consistency. ANOVA analysis shows a good performance in detecting the measurement bias.

DISCUSSION

In this paper, we reviewed five techniques for determining measurement consistency of structures measured from head and neck MRI: the correlation analysis, the linear regression, the Bland-Altman method, the paired t -test and the ANOVA. We showed the strength and weakness of each technique in detecting the measurement bias and determining the robustness to outliers. Table 1 provides the summary of the strength and weakness of each technique.

A correlation analysis cannot detect the measurement inconsistency between raters and it is sensitive to outliers. So it is inappropriate to use the correlation analysis for determining measurement consistency. A linear regression should not be used either because it is equivalent to the correlation analysis.

It is not easy to make quantitative decision using the Bland-Altman method. This is mainly because the Bland-Altman plot does not have statistical significance attached to it. The paired t -test provides quantification for the Bland-Altman method and it shows a good performance in detecting measurement bias. However, when most of the measurements of one rater are consistently larger or smaller than the other rater, the paired t -test tends to fail.

ANOVA provides the best performance in all cases studied and showed accurate analysis results in determining the measurement consistency. In addition, it provides the additional information of within-rater consistency.

As suggested by Krummenauer and Doll (2000), a good rule to follow is not to limit measurement consistency assessment on only one method, but to apply and compare several methods. We also recommend making as many repeated measurements as time and cost permit for more accurate determination of measurement consistency.

Acknowledgments

This work was supported in part by NIH Research Grants R03 DC4362 (Anatomic Development of the Vocal Tract: MRI Procedures), and R01 DC6282 (MRI and CT Studies of the Developing Vocal Tract), from the National Institute of Deafness and other Communicative Disorders (NIDCD). Also, by a core grant P-30 HD03352 to the Waisman Center from the National Institute of Child Health and Human Development (NICHD). We thank Celia Choih for assistance with placing the anatomic landmarks and making the necessary measurements.

References

1. Vorperian HK, Kent RD, Lindstrom MJ, Kalina CM, Gentry LR, Yandell BS. Development of vocal tract length during early childhood: A magnetic resonance imaging study. *The Journal of the Acoustical Society of America* 2005;117:338–350. [PubMed: 15704426]

2. Vorperian HK, Durtschi RB, Wang S, Chung MK, Ziegert AJ, Gentry LR. Estimating head circumference from pediatric imaging studies: an improved method. *Academic Radiology* 2007;14:1102–1107. [PubMed: 17707318]
3. Krummenauer F, Doll G. Statistical methods for the comparison of measurements derived from orthodontic imaging. *European Journal of Orthodontics* 2000;22:257–269. [PubMed: 10920558]
4. Bland JM, Altman DG. Statistical methods for assessing agreement between two methods of clinical measurement. *Lancet* 1986;1:307–310. [PubMed: 2868172]
5. Edvardsen T, Gerber BL, Garot J, Bluemke DA, Lima JAC, Smiseth OA. Quantitative assessment of intrinsic regional myocardial deformation by Doppler Strain Rate Echocardiography in Humans: Validation against three-dimensional tagged Magnetic Resonance Imaging. *Circulation* 2002;106:50–56. [PubMed: 12093769]
6. Liu W, Zagzebski JA, Varghese T, Dyer CR, Techavipoo U, Hall TJ. Segmentation of elastographic images using a coarse-to-fine active contour model. *Ultrasound in Medicine and Biology* 2006;32:397–408. [PubMed: 16530098]
7. Powell AJ, Maier SE, Chung T, Geva T. Phase-velocity cine Magnetic Resonance Imaging measurement of pulsatile blood flow in children and young adults: in vitro and in vivo validation. *Pediatric Cardiology* 2000;21:104–110. [PubMed: 10754076]
8. Vallejo E, Dione DP, Bruni WL, Constable RT, Borek PP, Soares JP, Carr JG, Condos SG, Wackers FJTh, Sinusas AJ. Reproducibility and accuracy of gated SPECT for determination of left ventricular volumes and ejection fraction: experimental validation using MRI. *The Journal of Nuclear Medicine* 2000;41:874–882.
9. Van Oosterhout MFM, Willigers HMM, Reneman RS, Prinzen FW. Fluorescent microspheres to measure organ perfusion: validation of a simplified sample processing technique. *American Journal of Physiology* 1995;269:H725–H733. [PubMed: 7653638]
10. Bland JM, Altman DG. Comparing methods of measurement: why plotting difference against standard method is misleading. *Lancet* 1995;346:1085–1087. [PubMed: 7564793]
11. Bland JM, Altman DG. Applying the right statistics: analyses of measurement studies. *Ultrasound in Obstetrics and Gynecology* 2003;22:85–93.
12. Brażdżionytė J, Macas A. Bland–Altman analysis as an alternative approach for statistical evaluation of agreement between two methods for measuring hemodynamics during acute myocardial infarction. *Medicina* 2007;43:208–214. [PubMed: 17413249]
13. Abate N, Burns D, Peshock RM, Garg A, Grundy SM. Estimation of adipose tissue mass by magnetic resonance imaging: validation against dissection in human cadavers. *Journal of Lipid Research* 1994;35:1490–1496. [PubMed: 7989873]
14. Chatterjee, S.; Hadi, AS.; Price, B. *Regression analysis by example*. 3. John Wiley & Sons, Inc; 2000.
15. Tukey, JW. *Exploratory data analysis*. New York: Addison-Wesley; 1977.
16. Martinez, WL.; Martinez, AR. *Exploratory data analysis with MATLAB*. Chapman & Hall/CRC; 2005.

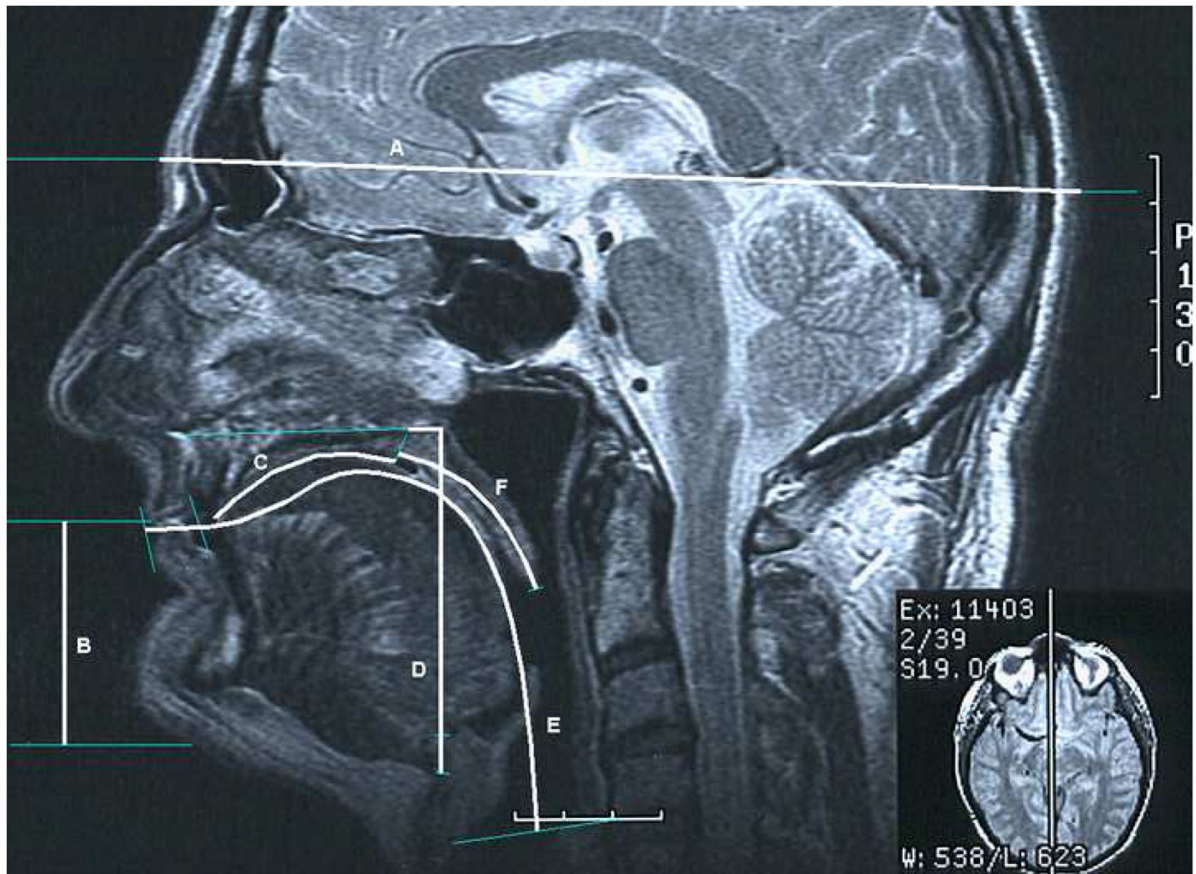


Figure 1. Midsagittal head and neck MRI with the six measurements used for measurement consistency comparison. A is Head Length (HL). B is Lower Anterior Face Height (LFH). C is Anterior Tongue Length (ATL). D is Hyoid vertical distance from PNS (HVP). E is Vocal Tract Length (VTL). F is Soft Palate (SP). See text for the definition of variables, and tissue type and measurement type of each variable.

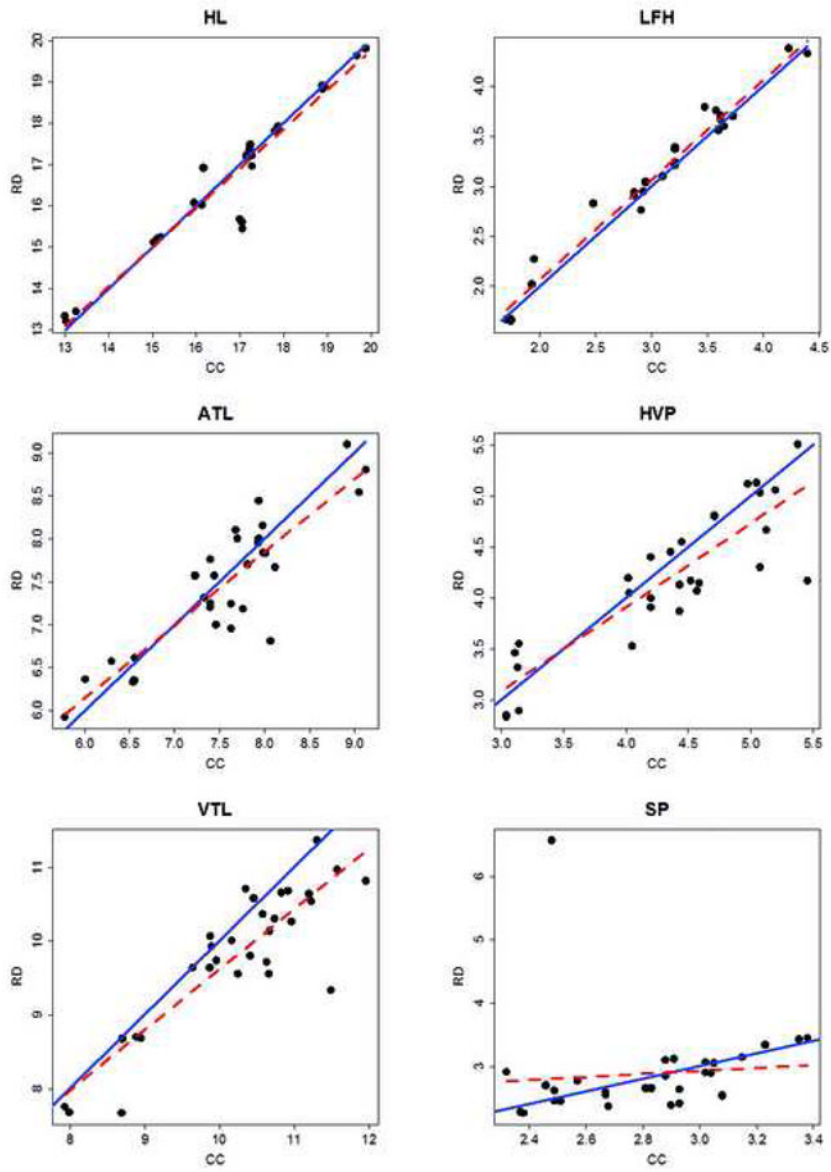


Figure 2. Scatter plots of HL, LFH, ATL, HVP, VTL, and SP. The solid lines ($y=x$) indicate the perfect consistency between two raters. The dotted lines are the linear regression fit.

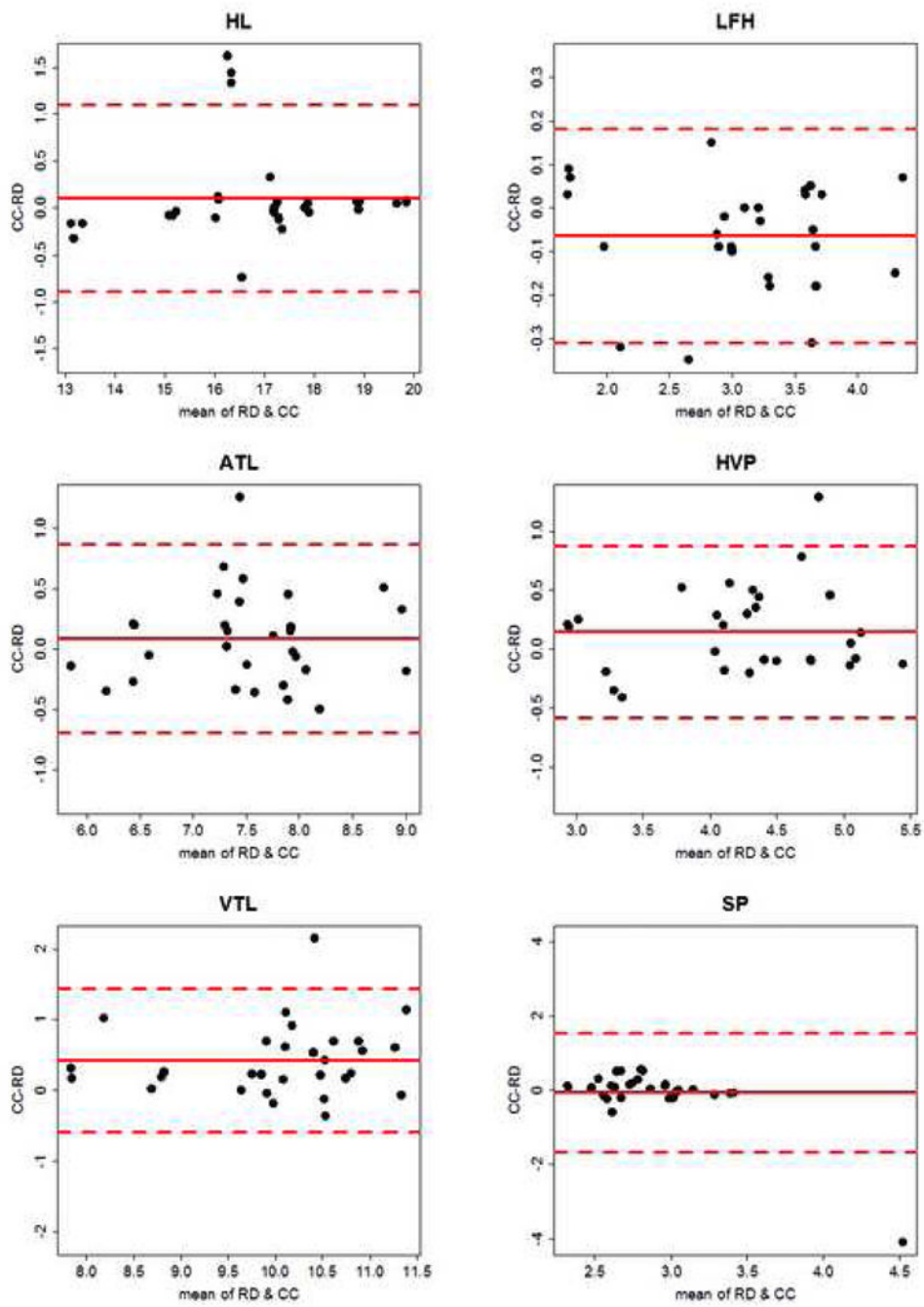


Figure 3. The Bland-Altman plots of HL, LFH, ATL, HVP, VTL and SP. The solid line is the mean difference \bar{d} , and the dotted lines are $\bar{d} - 1.96S_d$ (lower) and $\bar{d} + 1.96S_d$ (upper).

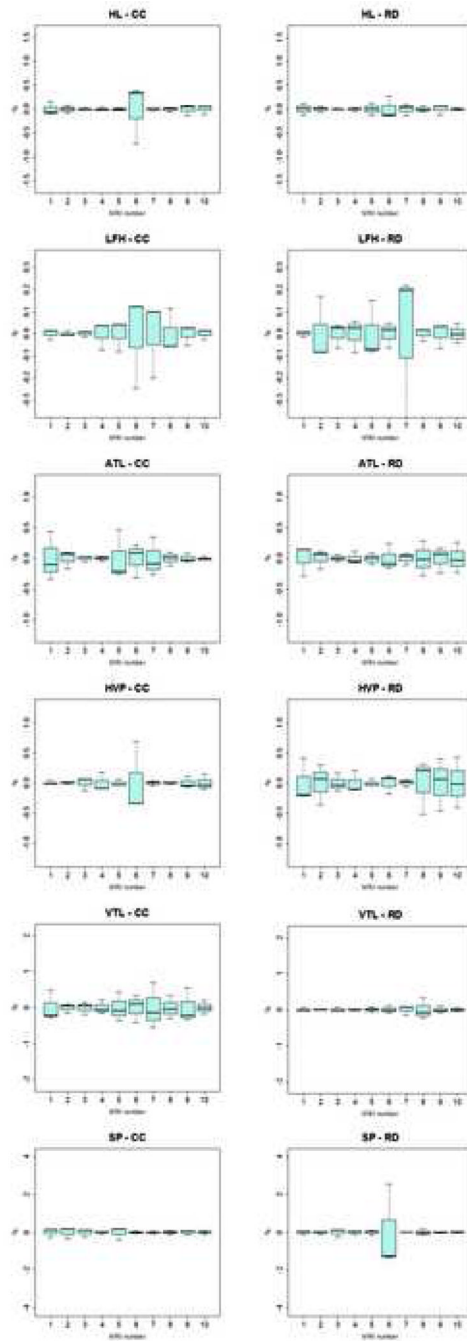


Figure 4. Within-rater consistency box plot of $d_{j,k}$ for ten MRIs of HL, LFH, ATL, HVP, VTL and SP for CC (left) and RD(right).

Table 1

Summary of statistical method used in determining the measurement consistency. The last two columns show if the method agrees with the ANOVA result for the six variables: Head Length (HL), Lower Anterior Face Height (L.FH), Anterior Tongue Length (ATL), Hyoid vertical distance from PNS (HVP), Vocal Tract Length (VTL), and Soft Palate (SP).

Method	Strength	Weakness	Agreement	Disagreement
Correlation & Regression	<ul style="list-style-type: none"> - Show degree of consistency - Simple procedure 	<ul style="list-style-type: none"> - Can not easily detect inconsistency - Sensitive to outliers 	HL, L.FH, ATL, HVP	VTL, SP
Bland-Altman method	<ul style="list-style-type: none"> - Visualization technique 	<ul style="list-style-type: none"> - Lacks statistical significance - Not easy to quantify the degree of consistency 	The method does not provide a decision.	
Paired t-test	<ul style="list-style-type: none"> - Detect bias fairly well - Simple procedure 	<ul style="list-style-type: none"> - Fails under systematic bias 	HL, ATL, VTL, SP	L.FH, HVP
ANOVA	<ul style="list-style-type: none"> - Best performance - Provide additional information of the within-rater consistency - Applicable for more than two raters. 	<ul style="list-style-type: none"> - Complicated procedure 	HL, L.FH, ATL, HVP, VTL, SP	