



Interobserver and intraobserver variability affecting the assessment of loss of autofluorescence of oral mucosal lesions

Monica Pentenero*, Daniela Todaro, Roberto Marino, Sergio Gandolfo

University of Turin, Department of Oncology, Unit of Oral Medicine and Oral Oncology, Italy

ARTICLE INFO

Keywords:

Autofluorescence
Oral mucosa
Inter-examiner variability
Intra-examiner variability
Observer-related variability

ABSTRACT

Objectives: The assessment of loss of tissue autofluorescence (LAF) has been proposed as an adjunct to comprehensive oral examination to enhance the detection of mucosal lesions harbouring dysplasia or carcinoma. The assessment of LAF is not based on completely objectified parameters therefore intraobserver and interobserver variability cannot be neglected alongside the issue of correct interpretation of LAF. The present study evaluated intraobserver and interobserver variability in the clinical assessment of LAF as performed by oral medicine practitioners (OMPs) or general dental practitioners (GDPs).

Materials and Methods: Couples of clinical pictures, acquired under white incandescent dental operatory light and during the assessment of LAF performed by VELscope were retrieved. Four OMPs and eight GDPs were asked to assess the pictures and to score the LAF. Kappa statistics allowed the assessment of intra- and inter-observer related variability.

Results: Pictures of 109 lesions representative of all oral mucosal sites and clinical appearances were selected. OMPs had a better intraobserver agreement than GDPs (substantial versus moderate). The moderate ($k = 0.506$) interobserver agreement observed among both OMPs and GDPs in a 2-score model (positive versus negative), lowered down to poor values only among GDPs when a 3-score or 4-score model (including uncertain judgments) was applied.

Conclusions: A good agreement ($k > 0.8$) was never observed and the present results are similar to previously reported data about conventional oral examination. Irrespective of the diagnostic accuracy, the assessment of AF seems not to be able to improve observer-related variability in the clinical assessment of oral mucosal lesions.

1. Introduction

To date, early detection of oral mucosal lesions is based on conventional oral examination (COE) aiming to characterize a lesion through the assessment of clinical aspects. The most recent recommendation by the American Dental Association Council on Scientific Affairs and the Centre for Evidence-Based Dentistry states that “for patients seeking care for suspicious lesions, immediate performance of a biopsy or referral to a specialist remains the single most important recommendation for clinical practice” [1]. Aiming at improving oral mucosal lesion detection, evaluation and the diagnostic accuracy of COE, several adjunctive tests have been proposed among which is the assessment of optical fluorescence imaging or tissue autofluorescence (AF) [2].

In the presence of oral carcinoma or oral potentially premalignant lesions a loss of tissue autofluorescence (LAF) can be observed due to specific autofluorescent patterns of tissues with different epithelial and

stromal architectures [3]. Commercial devices able to highlight LAF have been developed in order to assist the detection of oral mucosal abnormalities [4].

Notwithstanding great variation in sensitivity and specificity reported in various studies [5], an overall analysis of the literature shows that AF based tests have a moderate to high sensitivity but low specificity. Regardless of diagnostic reliability, variations in LAF assessment can be largely attributed to the inconsistent use of the devices, to the lack of incorporation of diascopy as recommended by the manufacturers [6] and to clinical features of oral mucosal lesions [7]. Nevertheless, observer related issues namely personal skills, inter-observer and intraobserver variability cannot be neglected as they also are potentially able to affect the consistency of LAF assessments. Even if devices available on the market have primarily been designed for use by general dentists, most of data from the literature come from secondary or tertiary settings, with only a few studies utilizing the device in general dental practice for detection of oral mucosal lesions [8–10].

* Corresponding Author at: Dipartimento di Oncologia, Regione Gonzole 10, 10043 Orbassano, TO, Italy.

E-mail address: monica.pentenero@unito.it (M. Pentenero).

<https://doi.org/10.1016/j.pdpdt.2019.09.007>

Received 29 July 2019; Received in revised form 16 September 2019; Accepted 27 September 2019

Available online 30 September 2019

1572-1000/ © 2019 Elsevier B.V. All rights reserved.

The present study aims to evaluate observer-related issues able to affect the clinical assessment of LAF using a hand-held device (VELscope®; LED Medical Diagnostics Inc., Barnaby, BC, Canada). Particularly, interobserver and intraobserver variability will be assessed among both general dental practitioners (GDPs) and oral medicine practitioners (OMPs).

2. Materials and methods

2.1. Data collection

Digital clinical pictures were retrieved from the digital archives of the oral medicine section of the San Luigi Gonzaga University Hospital. The study was conducted according to Human Ethics Guidelines and was approved by the local Ethics Board (project number 6/2015). An investigator (S.G.) who did not participate as one of the reviewing examiners selected clinical pictures aiming to obtain a pool of lesions representative of all oral mucosal sites and all the clinical appearances. Each pair of pictures was formed by one clinical picture acquired during examination under white incandescent dental operatory light (WL picture) and one picture acquired during the assessment of LAF performed by VELscope (AF picture) according to the manufacturer's instructions (Fig. 1). Pairs of pictures of 109 oral mucosal lesions were selected from clinical pictures acquired in a 12-month period. Digital photographs were obtained using a Nikon D80 equipped with a Nikkor 105 mm Macro lens (Nikon Corporation, Japan). A Macro Ring Lite (Sigma EM-140 DG; Sigma Corporation of America, USA) was used for photographs under white light. The PhotoMed VELscope Photography Kit was used to obtain photographs through the VELscope, according to the manufacturer's instructions.

2.2. Study participants and protocol

A call for volunteers was sent out among dentists from the Turin area seeking for GDPs. At the same time, OMPs were recruited from the oral medicine section of the San Luigi Gonzaga university hospital. A total of eight GDPs and four OMPs participated, with each practitioner signing informed consent. Participants attended a 4-h workshop run by the authors (M.P. and S.G.) aiming to illustrate the study protocols and to inform participants about the assessment of LAF. In order to obtain a setting as similar as possible to a routine practice, participants were asked to apply the criteria reported by the manufacturer of the employed hand-held device and they had no calibration exercises.

Participants were asked to assess all together the same pairs of pictures during two assessment sessions 2 months apart (T1 and T2). In order to have the same setting, all participants were gathered in a darkened room where each pair of pictures was displayed on a LED panel, WL and AF pictures side by side. Looking at each pair of pictures, participants were asked to assess the presence of LAF; in the case of uncertain interpretation, participants still had to discriminate between a positive or negative judgment. Therefore, they had to choose among the following scores: “positive” (in the presence of LAF), “uncertain/positive”, “uncertain/negative” and “negative” (in the absence of LAF).

2.3. Statistical analyses

Data collection and management were conducted using Microsoft Office Excel 2013 (Microsoft, Washington, USA) in association with the SPSS 18.0 software package (Apache Software Foundation, Chicago, IL, USA) for the statistical analyses. Interobserver and intraobserver variability were tested using data from the two assessment sessions (T1 and T2). Variability parameters were assessed considering all the 4 scores “positive”, “uncertain/positive”, “uncertain/negative” and “negative” (4-score model), considering 3 scores obtained grouping the two “uncertain” scores (3-score model) and considering 2 scores obtained grouping the overall positive and negative scores (2-score model). The

intraobserver variability has been assessed by calculation of kappa statistic, a chance-corrected measure of agreement between two raters [11]. Being the score categories ordered we calculated the weighted kappa, which accounts for how far apart the two scores are. The interobserver agreement has been assessed by calculation of the Fleiss' extension of kappa, which allows a chance-corrected measure of agreement among three or more raters. The following grading of k values was used [11]: $k < 0.4$ corresponding to poor agreement, $k \geq 0.4$ and < 0.6 corresponding to moderate agreement, $k \geq 0.6$ and < 0.8 corresponding to substantial agreement, $k \geq 0.8$ corresponding to good agreement.

3. Results

The selected lesions were representative of all oral mucosal sites and all clinical appearance, as reported in Table 1. More than half of the lesions (62 out of 109; corresponding to 56.9%) had a non-homogeneous texture and about half of them (58 out of 109; corresponding to 53.2%) had some erosive/ulcerated or red aspect.

The intraobserver agreement is shown in Table 2. Most of raters had a moderate intraobserver agreement (weighted k value ranging from 0.4 to 0.6). When assessing GDPs or OMPs, the latter had slightly higher k values, reaching a mean substantial agreement irrespective of the adopted score model (2-score, 3-score or 4-score).

The interobserver agreement was assessed for both the two assessment sessions (T1 and T2) without significant variations (not shown data). A moderate agreement was always observed among OMPs irrespective of the score model, while among GDPs the moderate agreement observed in the 2-score model, lowered down to a poor one in the 3-score and 4-score models (Table 3).

4. Discussion

COE using incandescent light is the standard screening or case-finding test to diagnose oral cancer or PMDs. It has several limitations which can be summed up as follows: it does not identify all PMDs, nor does it accurately detect the small proportion of biologically relevant lesions that are likely to progress to cancer [12]. This is of interest not only for OMPs but also for GDPs: in fact, a quite high rate of the general population have oral mucosal abnormalities (up to 20%) [13] but the vast majority of them are clinically/biologically benign. Within such a frame GDPs would largely benefit from a non-invasive test able to identify lesions worth to be submitted to further specialist assessments. In addition, COE is definitely an observer-dependent test potentially affected by observer's skills and observer-related variability. Therefore, even irrespective of the diagnostic accuracy, the potential use of adjuncts as triage tools aiming to have more objective and consistent assessments in the evaluation of oral mucosal lesions is of paramount importance.

Alongside the issue of correct interpretation of clinical appearance, the presence of inter/intra observer variability cannot be neglected when a correct diagnosis relies on visual assessments of not completely objectified parameters. This issue has been addressed in other clinical settings implying subjective visual diagnosis as dermatology [14] and oral radiology [15]. Coming to stomatology, previous studies extensively reported the pitfall of the low agreement between pathologists in the assessment of epithelial dysplasia [16–18] and similarly a low agreement has been described when dealing with the clinical assessment of oral mucosal lesions [19,20].

At present commercially available systems designed to assess LAF rely on subjective interpretation of autofluorescence as demonstrated in the present study. When considering the intraobserver variability, the substantial agreement found among OMPs in LAF assessment is consistent with previously reported data for the clinical assessment of mucosal lesions [20]. Moreover, the better intraobserver agreement found in OMPs when compared to GDPs (substantial versus moderate as

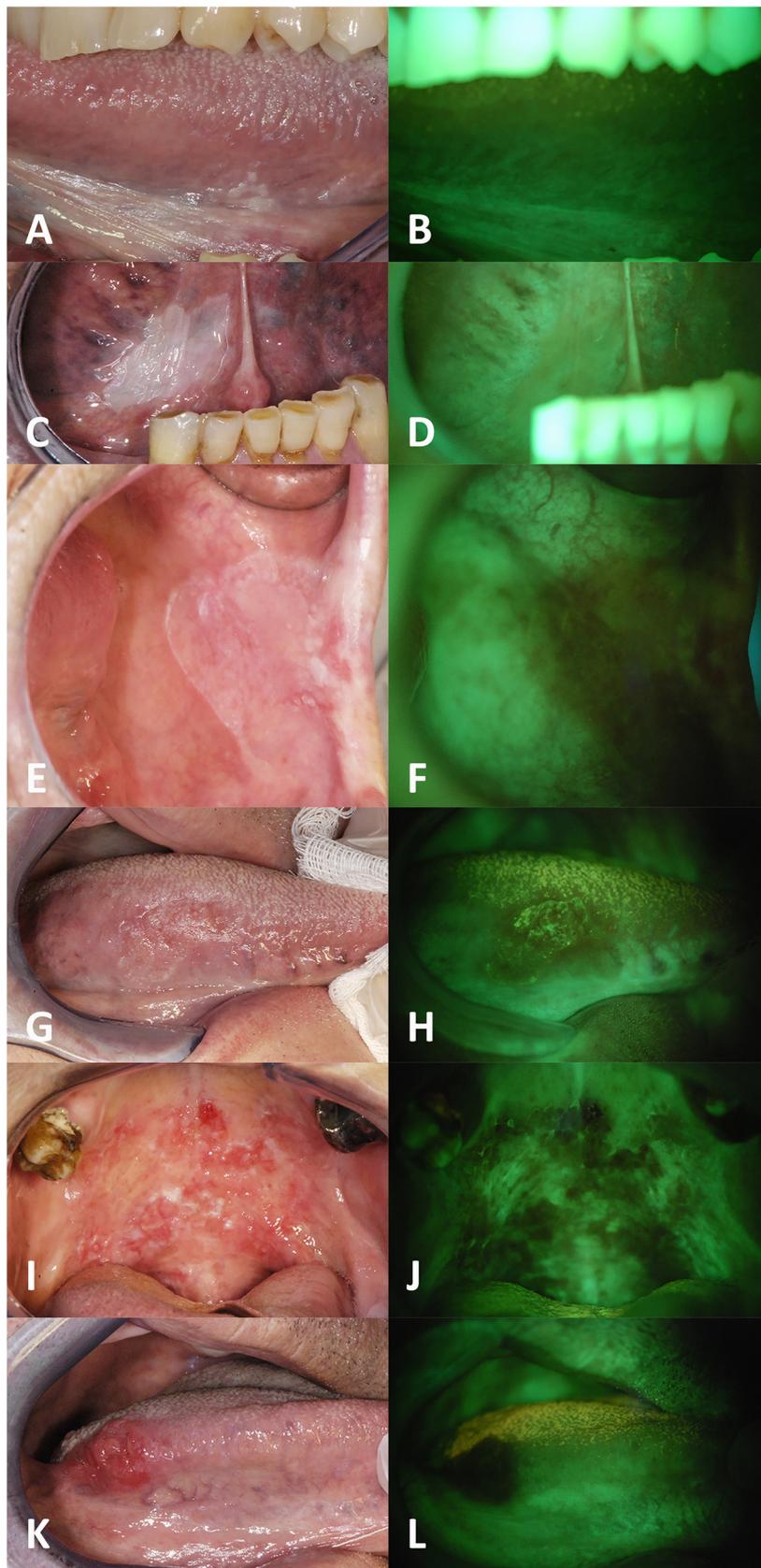


Fig. 1. Samples of the selected pairs of pictures showing on the left the clinical aspect of lesions as acquired during examination under white incandescent dental operator light (WL picture) and the corresponding AF pattern assessed by VELscope (AF picture) on the right. The selected lesions were characterized by different clinical features including patch-like (E, I), plaque-like (A, C, I), erosive (I), ulcerated (K) or verrucous surface (G) and showed maintenance of AF (B, D) or various degree of LAF (F, H, J, L).

Table 1
Clinical appearance of the selected lesions.

Site			Mucosal type			Main clinical aspect		
Buccal	30	27.5%	Thin	63	57.8%	Plaque	41	37.6%
Gingiva	4	3.7%	Vestibular	32	29.4%	Patch	54	49.5%
Tongue	44	40.4%	Masticatory	14	12.8%	Verrucous lesions	11	10.1%
Hard palate	9	8.3%				Erosion/Ulcer	3	2.8%
Floor of the mouth	9	8.3%						
Soft palate	11	10.1%						
Lip	2	1.8%						

Table 2
Intraobserver agreement as assessed by weighted kappa statistic.

Rater	K-value in a 2-score model ^a	K-value in 3-score model ^b	K-value in 4-score model ^c
GDP #01	0.662	0.541	0.607
GDP #02	0.229	0.117	0.153
GDP #03	0.522	0.504	0.510
GDP #04	0.479	0.523	0.506
GDP #05	0.470	0.511	0.496
GDP #06	0.721	0.739	0.733
GDP #07	0.495	0.500	0.498
GDP #08	0.440	0.405	0.417
OMP #09	0.506	0.462	0.476
OMP #10	0.657	0.658	0.658
OMP #11	0.960	0.857	0.892
OMP #12	0.730	0.651	0.677
Mean for GDPs	0.502	0.485	0.490
Mean for OMPs	0.713	0.657	0.676
Overall mean	0.573	0.542	0.552

^a “positive” OR “uncertain/positive” versus “uncertain/negative” OR “negative”.

^b “positive” versus “uncertain/positive” OR “uncertain/negative” versus “negative”.

^c “positive” versus “uncertain/positive” versus “uncertain/negative” versus “negative”.

reported in Table 2) is also in keeping with previous data from the literature showing that practitioners who regularly have to do with oral lesions are more consistent when compared to general dentists [19]. It has been seldom emphasized the importance of the experience of the observer for an accurate diagnosis [19,21,22] and a diagnostic aid able to reduce the gap between experienced and non-experienced observers could represent an effective benefit in the clinical assessment of oral mucosal lesions in a primary setting. Unfortunately, the present results are not able to support the use of LAF-based tests, in the case the VELscope device, in order to have similar intraobserver consistency in OMPs and GDPs.

The interobserver variability was assessed for both the assessment sessions (T1 and T2) with very similar even if not satisfactory results. Conceivably a binary scoring model is able to improve observers’ agreement as reported for example for the pathological grading of oral dysplasia [23], nevertheless even in a 2-score model a no more than moderate agreement was found. Also when dealing with interobserver

Table 3
Interobserver agreement among GDPs and OMPs as observed at the first assessment session (T1).

Rating model	GDPs			OMPs		
	K-value	95% CI	Agreement	K-value	95% CI	Agreement
2-score ^a	0.506	0.470–0.541	moderate	0.506	0.429–0.582	moderate
3-score ^b	0.378	0.352–0.404	poor	0.427	0.370–0.483	moderate
4-score ^c	0.354	0.330–0.377	poor	0.401	0.348–0.454	moderate

^a “positive” OR “uncertain/positive” versus “uncertain/negative” OR “negative”.

^b “positive” versus “uncertain/positive” OR “uncertain/negative” versus “negative”.

^c “positive” versus “uncertain/positive” versus “uncertain/negative” versus “negative”.

agreement OMPs showed better performance when compared to GDPs, but such differences should have a low clinical impact as they completely disappear in a 2-score model which better represents a potential clinical setting (Table 3).

Finally, also the significantly higher intraobserver than interobserver agreement was already reported for both clinical [20] and histopathological [24] assessments.

As previously said, variations in LAF assessment can be attributed to the lack of incorporation of diascopy. Diascopy is a test for blanchability performed by applying pressure and observing colour changes. Particularly, in the presence of LAF, the return of AF on blanching characterizes the so called diascopic fluorescence (DF). Vascular or inflammatory conditions are most often responsible for DF, due to higher content of light-absorbing haemoglobin. In such cases, blood dissipating intravascularly under compression results in a blanching effect. For this reason, diascopy has to be considered an essential part of LAF assessment, aiming to diminish the high rate of false-positive LAF findings to the point that it has been suggested to consider only LAF with no blanching as true LAF [10]. Nevertheless, this is not without potential pitfalls as DF has also been observed in the presence of oral epithelial dysplasia or squamous cell carcinoma [25,26]. Leaving out discussions about the improvement of the diagnostic accuracy of LAF assessment, performing diascopy has been found to be often very subjective (62% agreement on complete blanching after LAF) and difficult to be achieved even by experienced clinicians [25,26]. The assessment of DF requires skills to value LAF, to perform an adequate compression and to value the resulting blanching. Therefore, even in the presence of such specific skills, it would likely add other two observer-dependent factors potentially able to negatively affect the consistency of observations. Yet, the present study, based on the retrospective assessment of clinical pictures, is not able to investigate these issues in order to confirm or not such hypotheses. Even the assessment of brief videos could not be worth for that aim as diascopy should be performed by the observer him/herself.

Even if devices designed to assess LAF have been created and marketed to be used by GDPs, most of evidence still come from secondary or tertiary settings [1]. Aiming at achieving consistent results from AF assessment in general dental practice, two studies gave particular emphasis to the use of clinical interpretation and an established review protocol [8,9]. Nevertheless, the high level of interobserver disagreement observed in the present study highlights the need for new

devices possibly able to independently detect the LAF thanks to an automated image analysis. Such devices would allow consistent evaluations and would overcome the need for specific clinical training to perform the procedure and interpret the results [27]. Such devices would not only fill the potential gap between OMPs and GPDs, but could also be used by other health workers without specific clinical skills or experience.

Declaration of Competing Interest

None.

References

- [1] M.W. Lingen, M.P. Tampi, O. Urquhart, E. Abt, N. Agrawal, A.K. Chaturvedi, E. Cohen, G. D'Souza, J. Gurenlian, J.R. Kalmar, A.R. Kerr, P.M. Lambert, L.L. Patton, T.P. Sollecito, E. Truelove, L. Banfield, A. Carrasco-Labra, Adjuncts for the evaluation of potentially malignant disorders in the oral cavity: diagnostic test accuracy systematic review and meta-analysis—a report of the American Dental Association, *J. Am. Dent. Assoc.* 148 (11) (2017) 797–813 e52.
- [2] M.P. Rethman, W. Carpenter, E.E. Cohen, J. Epstein, C.A. Evans, C.M. Flaitz, F.J. Graham, P.P. Hujoel, J.R. Kalmar, W.M. Koch, P.M. Lambert, M.W. Lingen, B.W. Oettmeier Jr., L.L. Patton, D. Perkins, B.C. Reid, J.J. Sciubba, S.L. Tomar, A.D. Wyatt Jr., K. Aravamudan, J. Frantsve-Hawley, J.L. Cleveland, D.M. Meyer, American Dental Association Council on Scientific Affairs Expert Panel on screening for oral squamous cell C., Evidence-based clinical recommendations regarding screening for oral squamous cell carcinomas, *J. Am. Dent. Assoc.* 141 (5) (2010) 509–520.
- [3] I. Pavlova, M. Williams, A. El-Naggar, R. Richards-Kortum, A. Gillenwater, Understanding the biological basis of autofluorescence imaging for oral cancer detection: high-resolution fluorescence microscopy in viable tissue, *Clin. Cancer Res.* 14 (8) (2008) 2396–2404.
- [4] N. Bhatia, Y. Lalla, A.N. Vu, C.S. Farah, Advances in optical adjunctive aids for visualisation and detection of oral malignant and potentially malignant lesions, *Int. J. Dent.* 2013 (2013) 194029.
- [5] R. Macey, T. Walsh, P. Brocklehurst, A.R. Kerr, J.L. Liu, M.W. Lingen, G.R. Ogden, S. Warnakulasuriya, C. Scully, Diagnostic tests for oral cancer and potentially malignant disorders in patients presenting with clinically evident lesions, *Cochrane Database Syst. Rev.* (5) (2015) CD010276.
- [6] K.K. McNamara, B.D. Martin, E.W. Evans, J.R. Kalmar, The role of direct visual fluorescent examination (VELscope) in routine screening for potentially malignant oral mucosal lesions, *Oral Surg. Oral Med. Oral Pathol. Oral Radiol.* 114 (5) (2012) 636–643.
- [7] M. Ardore, G. Tempia Valenta, M. Pentenero, S. Gandolfo, Role of VELscope test in the assessment of oral mucosal lesions suspected to be oncologically relevant, *Dent. Cadmos* 80 (9) (2012) 538–546.
- [8] N. Bhatia, M.A. Matias, C.S. Farah, Assessment of a decision making protocol to improve the efficacy of VELscope in general dental practice: a prospective evaluation, *Oral Oncol.* 50 (10) (2014) 1012–1019.
- [9] D.M. Laronde, P.M. Williams, T.G. Hislop, C. Poh, S. Ng, C. Bajdik, L. Zhang, C. MacAulay, M.P. Rosin, Influence of fluorescence on screening decisions for oral mucosal lesions in community dental practices, *J. Oral Pathol. Med.* 43 (1) (2014) 7–13.
- [10] C.S. Farah, F. Dost, L. Do, Usefulness of optical fluorescence imaging in identification and triaging of oral potentially malignant disorders: a study of VELscope in the LESIONS programme, *J. Oral Pathol. Med.* (2019).
- [11] J.L. Fleiss, B. Levin, M.C. Paik, *The Measurement of Interrater Agreement, Statistical Methods for Rates and Proportions*, John Wiley & Sons, Inc., 2004, pp. 598–626.
- [12] M.W. Lingen, J.R. Kalmar, T. Karrison, P.M. Speight, Critical evaluation of diagnostic aids for the detection of oral cancer, *Oral Oncol.* 44 (1) (2008) 10–22.
- [13] M. Pentenero, R. Broccoletti, M. Carbone, D. Conrotto, S. Gandolfo, The prevalence of oral mucosal lesions in adults from the Turin area, *Oral Dis.* 14 (4) (2008) 356–366.
- [14] S. Chen, Q. Wang, T. Chu, M. Zheng, Inter-observer reliability in assessment of sensation of skin lesion and enlargement of peripheral nerves in leprosy patients, *Lepr. Rev.* 77 (4) (2006) 371–376.
- [15] B.G. Baksi, E. Alpoz, E. Sogur, A. Mert, Perception of anatomical structures in digitally filtered and conventional panoramic radiographs: a clinical evaluation, *Dentomaxillofac. Radiol.* 39 (7) (2010) 424–430.
- [16] A. Karabulut, J. Reibel, M.H. Therkildsen, F. Praetorius, H.W. Nielsen, E. Dabelsteen, Observer variability in the histologic assessment of oral premalignant lesions, *J. Oral Pathol. Med.* 24 (5) (1995) 198–200.
- [17] L.M. Abbey, G.E. Kaugars, J.C. Gunsolley, J.C. Burns, D.G. Page, J.A. Svirsky, E. Eisenberg, D.J. Krutchkoff, M. Cushing, Intraexaminer and interexaminer reliability in the diagnosis of oral epithelial dysplasia, *Oral Surg. Oral Med. Oral Pathol. Oral Radiol. Endod.* 80 (2) (1995) 188–191.
- [18] D.J. Fischer, J.B. Epstein, T.H. Morton, S.M. Schwartz, Interobserver reliability in the histopathologic diagnosis of oral pre-malignant and malignant lesions, *J. Oral Pathol. Med.* 33 (2) (2004) 65–70.
- [19] Y. Zadik, H. Orbach, A. Panzok, Y. Smith, R. Czerninski, Evaluation of oral mucosal diseases: inter- and intra-observer analyses, *J. Oral Pathol. Med.* 41 (1) (2012) 68–72.
- [20] E.H. van der Meij, K.P. Schepman, D.R. Plonait, T. Axell, I. van der Waal, Interobserver and intraobserver variability in the clinical assessment of oral lichen planus, *J. Oral Pathol. Med.* 31 (2) (2002) 95–98.
- [21] B. Mathew, R. Sankaranarayanan, K.B. Sunilkumar, B. Kuruvila, P. Pisani, M.K. Nair, Reproducibility and validity of oral visual inspection by trained health workers in the detection of oral precancer and cancer, *Br. J. Cancer* 76 (3) (1997) 390–394.
- [22] K.J. Patel, H.L. De Silva, D.C. Tong, R.M. Love, Concordance between clinical and histopathologic diagnoses of oral mucosal lesions, *J. Oral Maxillofac. Surg.* 69 (1) (2011) 125–133.
- [23] O. Kujan, R.J. Oliver, A. Khattab, S.A. Roberts, N. Thakker, P. Sloan, Evaluation of a new binary system of grading oral epithelial dysplasia for prediction of malignant transformation, *Oral Oncol.* 42 (10) (2006) 987–993.
- [24] E.H. van der Meij, J. Reibel, P.J. Slootweg, J.E. van der Wal, W.F. de Jong, I. van der Waal, Interobserver and intraobserver variability in the histologic assessment of oral lichen planus, *J. Oral Pathol. Med.* 28 (6) (1999) 274–277.
- [25] F. Kordbacheh, N. Bhatia, C.S. Farah, Patterns of differentially expressed genes in oral mucosal lesions visualised under autofluorescence (VELscope™), *Oral Dis.* 22 (4) (2016) 285–296.
- [26] C.S. Farah, L. McIntosh, A. Georgiou, M.J. McCullough, Efficacy of tissue autofluorescence imaging (VELscope) in the visualization of oral mucosal lesions, *Head Neck* 34 (6) (2012) 856–862.
- [27] E.C. Yang, M.T. Tan, R.A. Schwarz, R.R. Richards-Kortum, A.M. Gillenwater, N. Vigneswaran, Noninvasive diagnostic adjuncts for the evaluation of potentially premalignant oral epithelial lesions: current limitations and future directions, *Oral Surg. Oral Med. Oral Pathol. Oral Radiol. Endod.* 125 (6) (2018) 670–681.