

# Zur Reliabilität der Beschreibung morphologischer Merkmale in der traditionellen chinesischen Zungendiagnostik

Yanqing Li<sup>a</sup> Klaus Linde<sup>a</sup> Jingzhang Dai<sup>b</sup> Jisheng Zhang<sup>b</sup> Stefan Hager<sup>b</sup> Reinhard Saller<sup>c</sup>  
Dieter Melchart<sup>a,c</sup>

<sup>a</sup> Zentrum für naturheilkundliche Forschung, II. Medizinische Klinik und Poliklinik, Klinikum rechts der Isar, Technische Universität, München, Deutschland

<sup>b</sup> Klinik für Traditionelle Chinesische Medizin, Kötzing, Deutschland

<sup>c</sup> Institut für Naturheilkunde, Department für Innere Medizin, Universitätsspital Zürich, Schweiz

## Schlüsselwörter

Traditionelle chinesische Medizin · Zungendiagnostik · Reliabilität

## Zusammenfassung

**Hintergrund und Fragestellung:** Die Zungenbeurteilung ist ein zentrales diagnostisches Verfahren der traditionellen chinesischen Medizin. In der vorliegenden Studienserie wurde untersucht, inwieweit unterschiedliche Beurteiler bei der morphologischen Beschreibung von Zungenmerkmalen übereinstimmen. **Methoden:** In zwei Pilotstudien (jeweils n = 15 Patienten) und einer größeren Untersuchung (Hauptstudie; n = 101) wurden Digitalfotografien der Zungenbefunde von Klinikpatienten durch zwei bzw. drei erfahrene Beurteiler mithilfe von Formularen verblindet bewertet. Das wichtigste Zielkriterium war die Übereinstimmung über die erwartete zufällige Übereinstimmung hinaus (Cohens Kappa).

**Ergebnisse:** In der ersten Pilotstudie wurden in hohem Maße variable Kappa-Werte beobachtet (−0,15 bis 0,76). Bei zahlreichen Merkmalen zeigten sich jedoch nur geringe bis mäßige Übereinstimmungen. In der zweiten, mit einer verbesserten Methode durchgeföhrten Pilotstudie zeigten sich zwar wiederum in hohem Maße variable Kappa-Werte (−0,10 bis 1,00), bei 7 der 18 erhobenen Merkmale war die überzufällige Übereinstimmung jedoch sehr gut (Kappa ≥ 0,75). In der Hauptstudie lagen die Kappa-Werte zwischen 0,15 und 0,83. **Diskussion:** In den vorliegenden Studien wurden erste Schritte hin zu einer systematischen Untersuchung der Reliabilität der Zungendiagnose im Rahmen der traditionellen chinesischen Medizin unternommen. Die Ergebnisse deuten darauf hin, dass die Beurteilung morphologischer Merkmale eine befriedigende Reliabilität haben dürfte.

## Key Words

Traditional Chinese medicine · Tongue diagnostics · Reliability

## Summary

*On the Reliability of the Description of Morphological Characteristics in Traditional Chinese Tongue Diagnostics*

**Background and Objective:** The assessment of the tongue is a crucial diagnostic tool of traditional Chinese medicine. In a series of studies we aimed to investigate to what extent independent raters agree in the description of morphological tongue characteristics. **Methods:** In two pilot studies (n = 15 each) and one larger study (n = 101) two to three physicians experienced in traditional Chinese medicine assessed morphological characteristics in digital photos of tongues by use of a rating form and under blind conditions. The primary outcome measure was agreement beyond chance (Cohen's kappa). **Results:** Kappa values varied strongly in the first pilot study (−0.15 to 0.76) and for many items, agreement was weak or moderate. In the second pilot study which used improved methods kappa values still varied considerably (−0.10 to 1.00), but for 7 out of the 18 items assessed there was an excellent agreement (kappa ≥ 0.75). In the confirmatory study, kappa values ranged between 0.15 and 0.83. **Discussion:** The performed studies have to be seen as a first attempt to develop adequate methods in order to systematically investigate the reliability of traditional Chinese tongue diagnostics. The study findings suggest that the description of morphological characteristics within traditional Chinese tongue diagnostics has acceptable reliability.

## Einleitung

Die Beurteilung der Zunge spielt in der traditionellen chinesischen Medizin (TCM) im Rahmen der Syndromdiagnostik eine wichtige Rolle. Eine kritische empirische Evaluation der Zungendiagnostik erfolgte bisher jedoch kaum. Eine Evaluation, die sowohl westlichen wissenschaftlichen Ansprüchen als auch der Theorie und Praxis der TCM gerecht wird, ist nicht einfach zu bewerkstelligen. Aus wissenschaftlicher Sicht muss ein diagnostisches Verfahren reliabel und valide sein. Als reliabel wäre die Zungendiagnose zu bewerten, wenn mehrere (angemessen ausgebildete und erfahrene) Ärzte bei der Beurteilung zur gleichen Bewertung kämen. Ein diagnostischer Test wird als valide bezeichnet, wenn er das Vorhandensein bzw. das Nicht-Vorhandensein einer Krankheit weitgehend sicher erkennt [1].

Insbesondere die Prüfung der Validität ist im Falle der Zungendiagnostik (und vieler anderer komplementärer Diagnoseverfahren) problematisch. In der westlichen Medizin erfolgt die Prüfung der Validität durch einen Vergleich des jeweiligen Tests mit dem sogenannten Goldstandard, also der besten und sichersten Möglichkeit, eine Diagnose zu stellen [1]. So würde man z.B. ein einfaches Screeninginstrument für Depressionen mit einer ausführlichen psychiatrischen Exploration vergleichen oder ein Belastungs-EKG zur Untersuchung einer koronaren Herzerkrankung mit einer Koronarangiographie. Im Fall der Zungendiagnostik würde man untersuchen, inwieweit eine alleinige Beurteilung der Zunge bereits die vollständige Syndromdiagnose vorhersagt bzw. inwieweit die Bewertung der Zunge zur Erstellung einer sicheren Syndromdiagnose beiträgt (im Vergleich zu einem diagnostischen Vorgehen ohne Zungendiagnose). Die chinesische Syndromdiagnose, die eher eine funktionale Zustandsbeschreibung des Organismus unter Berücksichtigung konstitutioneller Aspekte ist als eine «feste» Diagnose im westlichen Sinne, eignet sich für eine solche Vorgehensweise nur eingeschränkt. Dementsprechend fehlt eine eindeutige diagnostische Klassifikation, die die terminologischen Voraussetzungen für eine angemessene Validitätsprüfung schaffen würde.

Es ist daher forschungstechnisch einfacher, sich zunächst auf die Prüfung der Reliabilität zu konzentrieren. Dabei sind zwei Ebenen denkbar. Zum einen kann man überprüfen, ob verschiedene Beurteiler die einzelnen morphologischen Merkmale von Zungen übereinstimmend *beschreiben* (z.B. ob mehrere Beurteiler die Zungenfarbe als auffällig rot bewerten). Eine solche Prüfung scheint unproblematisch, ist jedoch nur von eingeschränkter klinischer Relevanz. Zum anderen kann man auf einer komplexeren Ebene untersuchen, ob die *Interpretation* der einzelnen Merkmale oder auch des Gesamtbildes der Zunge übereinstimmt. Diese Ebene hat eine deutlich größere klinische Relevanz. Für die notwendige terminologische Standardisierung sind jedoch umfangreiche Vorarbeiten und Vorprüfungen erforderlich.

In China scheinen sich die Forschungsbemühungen im Zusammenhang mit der Zungendiagnostik auf die Digitalisierung und Automatisierung der Beurteilung von Zungenbildern zu konzentrieren (z.B. [2–4]). Recherchen in PubMed fanden keinerlei Untersuchungen zur Reliabilität oder Validität. Den Autoren sind zwei lediglich als Abstracts verfügbare westliche Untersuchungen zu Aspekten der Reliabilität bekannt [5, 6]. Es wurde daher beschlossen, im Rahmen der Forschungsbegleitung der TCM-Klinik Kötzing – eines seit 1991 bestehenden Krankenhauses mit Akutversorgungsvertrag, in dem chinesische und deutsche Ärzte vorwiegend chronisch kranke Patienten gemeinsam betreuen – zu versuchen, Methoden für eine angemessene wissenschaftliche Überprüfung der Zungendiagnostik zu entwickeln. Das vorliegende Manuskript beschreibt die ersten im Rahmen einer Dissertation erarbeiteten Studien. Ziel der vorliegenden Studienreihe war eine exploratorische Überprüfung der Übereinstimmung mehrerer Beurteiler (Inter-Rater-Reliabilität) bezüglich der Beschreibung morphologischer Merkmale im Rahmen der Zungendiagnostik.

## Methodik

**Rahmenbedingungen.** Es wurden drei Studien durchgeführt: zwei Pilotstudien mit jeweils 15 Patienten bzw. Zungenfotografien und eine Hauptserie mit 101 Zungenfotografien. Die Durchführung der Studien erfolgte in der TCM-Klinik Kötzing. In dieser Klinik mit 75 Betten werden pro Jahr rund 1000 Patienten mit meist chronischen Erkrankungen stationär (mittlere Liegedauer 26 Tage) von 10 chinesischen und 6 deutschen Ärzten primär mit den Methoden der TCM (Akupunktur, chinesische Arzneimitteltherapie, Tuina und Qi Gong) behandelt. Zwischen 2003 und 2005 wurde routinemäßig von allen stationären Patienten der Klinik am Morgen nach der Aufnahme vor dem Frühstück im gleichen Raum unter gleichen Lichtbedingungen von einer deutschen Laborkraft mit einer Digitalkamera (Fuji Film Smart Media, Digital Camera, Fine Pix 2600 Zoom 3V) eine Fotografie der Zunge erstellt. Die Patienten saßen dabei immer auf dem gleichen Stuhl, um den Winkel, in dem fotografiert wurde, möglichst konstant zu halten. Die Fotografien wurden digital und kodiert elektronisch abgelegt. Für die Untersuchung wurden von einer an der Studie unbeteiligten Verwaltungskraft aus dem Pool von Fotodateien zufällig Fotografien gezogen.

**Beurteiler.** Die Zungenfotos wurden von zwei bis drei Beurteiern streng unabhängig voneinander (Blindbedingungen) mithilfe von Formularen bewertet. Die Bewertung erfolgte ausschließlich auf Basis der Zungenfotos ohne weitere Informationen. Insgesamt waren vier Beurteiler beteiligt. Beurteiler 1 (JD) und Beurteiler 2 (JZ) wirkten an allen drei Teilstudien mit, Beurteiler 3 (YL) ausschließlich an der ersten Pilotstudie und Beurteiler 4 (SH) ausschließlich bei der Hauptstudie. Die Beurteiler 1, 2 und 3 sind chinesische Ärzte mit abgeschlossenem TCM-Studium und langjähriger praktischer Erfahrung. Beurteiler 4 ist der deutsche Chefarzt der Klinik, der über 10 Jahre TCM-Erfahrung verfügt.

**Beurteilung der Zungenfotografien.** Auf der Basis einer Probeauswertung der Routinedokumentation in der elektronischen Patientenakte der Klinik und einer Konsensusdiskussion im Studienteam wurde ein erstes Formular erstellt, auf dem folgende Merkmale als vorhanden (ja) oder nicht vorhanden (nein) bewertet und dokumentiert werden sollten: Zunge rot, Zunge blass, Zunge bläulich, Zunge rötlich-purpur; Vorhandensein von Flecken, Zahneindrücken, Rissen; dicker Belag, klebriger Belag, sich auflösender Belag, dünner Belag. Diese Aufteilung erfolgte, um in allen Fällen eine klare Ja-Nein-Bewertung zu ermöglichen (z.B. Nichtvorhandensein eines dicken Belags kann sowohl kein Belag als auch dünner

Die einfachste Maßzahl der Übereinstimmung der Bewertung eines diagnostischen Merkmals von zwei Beurteilern als vorhanden oder nicht vorhanden ist die prozentuale Übereinstimmung. Diese berücksichtigt jedoch nicht, dass ein Teil dieser Übereinstimmung zufällig ist bzw. sein kann. Werden z.B. 2 Münzen geworfen, erwartet man in 50% der Würfe, dass beide übereinstimmend Kopf oder Zahl zeigen. Ein Maß zur Quantifizierung der Übereinstimmung jenseits des Zufalls ist der Kappa-Index nach Cohen [7]. Das Maß berücksichtigt außerdem auch, ob ein Merkmal (z.B. eine auffällige rote Zunge) sehr häufig oder sehr selten auftritt. Die Bestimmung erfolgt in 3 Schritten. Im 1. Schritt wird die tatsächlich beobachtete Übereinstimmung ermittelt. Im folgenden Beispiel beträgt diese 0,9 oder 90% (in 5 von 100 Fällen sagen beide Beurteiler, die Zunge sei rot, in weiteren 85 sagen beide, sie sei nicht rot).

Beurteiler 1	Beurteiler 2		Summe
	Rot	Nicht rot	
Rot	5	5	10
Nicht rot	5	85	90
Summe	10	90	100

Im 2. Schritt werden für die vier Kombinationsmöglichkeiten die erwarteten Häufigkeiten berechnet, indem die jeweilige Zeilensumme mit der Spaltensumme multipliziert und dann durch die Gesamtzahl beurteilter Fälle geteilt wird. Für die Beurteilung rot durch beide Beurteiler ist die Zeilensumme 10, die Spaltensumme ebenso 10; die bei nur zufälliger Übereinstimmung erwartete Häufigkeit der Bewertung rot durch beide Beurteiler wäre demnach  $1 (10 \times 10 / 100)$ . Entsprechend ergibt sich für die beiden Nichtübereinstimmungsmöglichkeiten jeweils die erwartete Häufigkeit  $9 (10 \times 90 / 100)$  und für das Urteil nicht rot durch beide Beurteiler die erwartete Häufigkeit  $81 (90 \times 90 / 100)$ . Die erwartete Gesamtübereinstimmung beträgt somit 0,82 bzw. 82%.

#### *Beispiel mit den erwarteten Häufigkeiten in Klammern*

Beurteiler 1	Beurteiler 2		Summe
	Rot	Nicht rot	
Rot	5 (1)	5 (9)	10
Nicht rot	5 (9)	85 (81)	90
Summe	10	90	100

Im 3. Schritt wird der eigentliche Kappa-Index berechnet, indem die erwartete von der beobachteten Übereinstimmung abgezogen ( $0,9 - 0,82 = 0,08$ ) und dann durch  $1 - \text{erwartete Beobachtung}$  ( $1 - 0,82 = 0,18$ ) geteilt wird. Das errechnete Kappa für das Beispiel ist somit  $0,08 / 0,18 = 0,44$ . Kappa-Werte  $< 0,40$  werden als geringe überzufällige Übereinstimmung klassifiziert, Werte zwischen 0,40 und 0,59 als moderate Übereinstimmung, Werte zwischen 0,60 und 0,74 als gute und Werte  $\geq 0,75$  als sehr gute Übereinstimmung [7].

**Abb. 1.** Beispiel zur Bewertung von Übereinstimmungen mit Cohens Kappa.

Belag bedeuten). Die Beurteilung der Zungenfotos ( $n = 15$ ) erfolgte ohne weitere Absprache (Rater-Training) zwischen den Beurteilern am gleichen Computerbildschirm zu unterschiedlichen Zeitpunkten.

Aufgrund der Ergebnisse der ersten Pilotstudie wurde das Vorgehen bei der Beurteilung der Zungenmerkmale für die zweite Pilotstudie modifiziert. Um die starke Beeinflussung des Farbeindrucks auf dem Computerbildschirm durch den Blickwinkel des Beurteilers zu vermindern (die Farbe erscheint auf dem Bildschirm umso dunkler, je höher die Augen des Beobachters liegen, d.h. einem größeren oder höher sitzenden Betrachter erscheint die Zunge dunkler als einem kleineren bzw. tiefer sitzenden), wurde ein Gestell gebaut, auf das alle Beurteiler bei der Bewertung der Bilder ihr Kinn legten. Damit wurde erreicht, dass der Blickwinkel für alle Beurteiler weitgehend identisch war. Als weiteres Problem bei der Bewertung der Farbe erwies sich die Tatsache, dass diese in unterschiedlichen Regionen der Zunge variieren kann. Für die Beurteilung aus Sicht der TCM ist es von Bedeutung, ob ein auffälliger Farbbefund z.B. die Zungenspitze oder deren Ränder betrifft. Das Formular wurde daher entsprechend überarbeitet. Ein drittes Problem der ersten Pilotstudie war, dass unklare Befunde unterschiedlich beurteilt wurden. Während ein Beurteiler z.B. nur pathologisch relevante Zahneindrücke als «Zahneindrücke vorhanden» bewertete, zählte ein anderer Beurteiler auch weniger deutliche, aber sichtbare Zahneindrücke. Vor der Bewertung der 15 Fotografien der zweiten Pilotstudie erfolgte daher eine Abstimmung zwischen den Beurteilern auf der Basis einer gemeinsamen Durchsicht der 15 Bilder der ersten Pilotstudie.

Die Beurteilung der Fotografien in der Hauptstudie erfolgte wie in der zweiten Pilotstudie. Es erfolgte lediglich eine Vereinfachung des Formulars beim Item «Zungenfarbe», da einzelne in der Pilotstudie abgefragte Merkmale extrem selten auftreten und daher von geringer praktischer Relevanz schienen.

**Statistik.** Die erhobenen Daten wurden in der begleitenden Forschungseinrichtung in eine SPSS-Datei eingegeben (Einfacheingabe mit Plausibilitäts- und Stichprobenechecks). Für jedes Beurteilerpaar wurde die prozentuale Übereinstimmung (Zahl der gleichen Bewertungen / Gesamtzahl der Bewertungen  $\times 100$ ) berechnet. Als Maß für die überzufällige Übereinstimmung zwischen einzelnen Paaren von Beurteilern wurde der Kappa-Index nach Cohen berechnet (Abb. 1). Für die Hauptstudie wurden mit StatsDirect statistical software Version 2.5.7 (07.08.2006) Kappa-Indizes mit dem zugehörigen 95%-Konfidenzintervall (95% CI) bestimmt. Zusätzlich wurde mit der Fleiss-Cusick-Erweiterung [7] ein Kappa (im Folgenden als Fleiss' Kappa bezeichnet) mit 95%-CI für die Übereinstimmung zwischen allen drei Beurteilern in der Hauptstudie berechnet.

## Ergebnisse

**Pilotstudie 1.** Die Übereinstimmung variierte zwischen den einzelnen Bewertungskriterien und zum Teil zwischen den einzelnen Beurteilern für ein einzelnes Kriterium deutlich (Tab. 1).

**Tab. 1.** Ergebnisse der ersten Pilotstudie (n = 15), Übereinstimmungen sind fettgedruckt

Item		Beurteiler 1 vs. 2				Beurteiler 1 vs. 3				Beurteiler 2 vs. 3			
		nein	ja	%	Kappa	nein	ja	%	Kappa	nein	ja	%	Kappa
Zunge rot	nein	<b>1</b>	<b>1</b>	47	-0,02	<b>2</b>	0	40	0,11	<b>7</b>	1	67	0,31
	ja	<b>7</b>	<b>6</b>			9	<b>4</b>			4	<b>3</b>		
Zunge blass	nein	<b>11</b>	2	73	-0,15	<b>7</b>	6	60	0,24	<b>7</b>	6	60	0,24
	ja	2	<b>0</b>			0	<b>2</b>			0	<b>2</b>		
Zunge bläulich	nein	<b>3</b>	5	53	0,08	<b>7</b>	1	47	-0,13	<b>5</b>	0	40	0,07
	ja	2	<b>5</b>			7	<b>0</b>			9	<b>1</b>		
Zunge rötlich-purpur	nein	<b>14</b>	1	93	n.b.	<b>15</b>	0	100	n.b.	<b>14</b>	0	93	n.b.
	ja	0	<b>0</b>			0	<b>0</b>			1	<b>0</b>		
Flecken	nein	<b>4</b>	11	27	n.b.	<b>14</b>	1	93	n.b.	<b>4</b>	0	33	0,05
	ja	0	<b>0</b>			0	<b>0</b>			10	<b>1</b>		
Gedunsen	nein	<b>7</b>	0	73	0,48	<b>5</b>	2	87	0,73	<b>5</b>	6	60	0,31
	ja	4	<b>4</b>			0	<b>8</b>			0	<b>4</b>		
Zahneindrücke	nein	<b>2</b>	0	40	0,11	<b>2</b>	0	67	0,30	<b>7</b>	4	73	0,48
	ja	9	<b>4</b>			5	<b>8</b>			0	<b>4</b>		
Risse	nein	<b>9</b>	0	67	0,19	<b>8</b>	1	73	0,41	<b>11</b>	3	80	0,33
	ja	5	<b>1</b>			3	<b>3</b>			0	<b>1</b>		
Dicker Belag	nein	<b>11</b>	0	73	n.b.	<b>9</b>	2	87	0,71	<b>9</b>	6	60	n.b.
	ja	4	<b>0</b>			0	<b>4</b>			0	<b>0</b>		
Klebriger Belag	nein	<b>5</b>	1	73	0,47	<b>5</b>	1	73	0,47	<b>6</b>	0	93	0,86
	ja	3	<b>6</b>			3	<b>6</b>			1	<b>8</b>		
Auflösender Belag	nein	<b>12</b>	0	93	0,76	<b>12</b>	0	87	0,44	<b>13</b>	0	93	0,63
	ja	1	<b>2</b>			2	<b>1</b>			1	<b>1</b>		
Dünner Belag	nein	<b>0</b>	4	73	n.b.	<b>0</b>	0	87	n.b.	<b>1</b>	3	73	0,19
	ja	0	<b>11</b>			2	<b>13</b>			1	<b>10</b>		

n.b. = Nicht berechenbar (bei  $\geq 2$  leeren Feldern in der 4-Felder-Tafel).

So war z.B. die Übereinstimmung bezüglich der Zungenfarbe (mit Ausnahme der Bewertung rötlich-purpur: ja/nein) zwischen den drei Gutachtern gering. Das Vorhandensein von Flecken wurde von den Beurteilern 1 und 3 sehr ähnlich bewertet, Beurteiler 2 kam zu deutlich abweichenden Ergebnissen. Bei den Merkmalen «Zahneindrücke» und «Risse» war die Übereinstimmung gering bis mäßig. Bezuglich der Beurteilung der Beläge stimmten die Bewertungen zum Teil besser überein.

*Pilotstudie 2.* In der zweiten Pilotstudie zeigten sich zum Teil sehr gute Übereinstimmungen (Tab. 2). Bei der Farbbeurteilung kam es jedoch immer noch zu deutlichen Diskrepanzen, insbesondere bei der Frage, ob Zungenränder und Zungenkörper blass seien. Da einzelne Merkmale wie eine blasse Zungenspitze oder Risse sehr selten vorkamen, sind die entsprechenden Bewertungen kaum zu interpretieren.

*Hauptstudie.* Die Patienten, deren Zungenfotos zur Auswertung herangezogen wurden, waren zu 63% weiblich, der Altersmedian lag bei 52 Jahren (Minimum 20, Maximum 87 Jahre). 52% der Patienten litten an chronischen Schmerzsyndromen, 12% an Erkrankungen des Nervensystems und 9% an Erkrankungen des Magen-Darm-Systems (27% sonstige Erkrankungen). Die Ergebnisse der Zungenbewertung sind in Tabelle 3 dargestellt. Bei der Beurteilung der Farbmerkmale war die überzufällige Übereinstimmung zwischen den Beurteilern gering bis moderat (Kappa-Werte 0,19–0,56), ebenso bei den Merkmalen «Flecken» (0,15–0,30), «dicker Belag»

(0,26–0,40), «klebriger Belag» (0,17–0,49) und «dünner Belag» (0,34–0,51). Gute bis sehr gute Übereinstimmungen ergaben sich insbesondere bei Zahneindrücken (0,65–0,83) und bei der Bewertung der Belagfarbe (0,64–0,83). Die Übereinstimmung bei den Merkmalen «gedunsene Zunge» (0,56–0,65), «Risse» (0,43–0,60), und «auflösender Belag» (0,55–0,78) schwankte zwischen moderat und sehr gut. Wurden die Angaben aller drei Beurteiler mit Fleiss' Kappa bewertet, ergab sich für fünf Merkmale eine geringe, für sechs eine moderate und für drei Merkmale eine gute Übereinstimmung.

## Diskussion

*Zusammenfassung der Hauptergebnisse.* Die vorliegende Untersuchungsreihe zeigt, dass eine Untersuchung der Reliabilität der morphologischen Beschreibung von Zungenmerkmalen möglich, aber mit erheblichen methodischen Schwierigkeiten verbunden ist. In den Voruntersuchungen wurde deutlich, dass ohne eine angemessene Standardisierung der Beobachtungsbedingungen und ohne Vorabsprachen zwischen den Beurteilern über die Bewertung unklarer Befunde keine sinnvolle Studiendurchführung möglich ist. Die Ergebnisse der Hauptstudie zeigen erwartungsgemäß mit wenigen Ausnahmen für alle Items statistisch signifikante Übereinstimmungen über den Zufall hinaus. Allerdings ist das Aus-

**Tab. 2.** Ergebnisse der zweiten Pilotstudie, Übereinstimmungen sind fettgedruckt

Item	Beurteiler 1 vs. 2			
	nein	ja	%	Kappa
Zungenspitze rot	nein <b>1</b>	1	93	0,63
	ja <b>0</b>	<b>13</b>		
Zungenspitze blass	nein <b>15</b>	0	100	n.b.
	ja <b>0</b>	<b>0</b>		
Zungenspitze bläulich	nein <b>15</b>	0	100	n.b.
	ja <b>0</b>	<b>0</b>		
Zungenränder rot	nein <b>12</b>	2	87	0,44
	ja <b>0</b>	<b>1</b>		
Zungenränder blass	nein <b>5</b>	8	47	0,14
	ja <b>0</b>	<b>2</b>		
Zungenränder bläulich	nein <b>14</b>	1	93	n.b.
	ja <b>0</b>	<b>0</b>		
Zungenkörper rot	nein <b>14</b>	0	100	1,00
	ja <b>0</b>	<b>1</b>		
Zungenkörper blass	nein <b>6</b>	6	60	0,29
	ja <b>0</b>	<b>3</b>		
Zungenkörper bläulich	nein <b>6</b>	0	67	0,39
	ja <b>5</b>	<b>4</b>		
Flecken	nein <b>12</b>	0	87	0,44
	ja <b>2</b>	<b>1</b>		
Zunge gedunsen	nein <b>9</b>	0	100	1,00
	ja <b>0</b>	<b>6</b>		
Zahneindrücke	nein <b>9</b>	0	100	1,00
	ja <b>0</b>	<b>6</b>		
Risse	nein <b>14</b>	0	100	1,00
	ja <b>0</b>	<b>1</b>		
Dicker Belag	nein <b>12</b>	2	80	-0,10
	ja <b>1</b>	<b>0</b>		
Klebriger Belag	nein <b>8</b>	0	100	1,00
	ja <b>0</b>	<b>7</b>		
Auflösender Belag	nein <b>14</b>	0	100	1,00
	ja <b>0</b>	<b>1</b>		
Dünner Belag	nein <b>5</b>	0	80	0,61
	ja <b>3</b>	<b>7</b>		
	gelb		weiß	
Belagfarbe	gelb <b>2</b>	0	100	1,00
	weiß <b>0</b>	<b>13</b>		

n.b. = Nicht berechenbar (bei  $\geq 2$  leeren Feldern in der 4-Felder-Tafel).

maß dieser Übereinstimmung (quantifiziert durch den Kappa-Index) bei der verwendeten Untersuchungsmethode für die einzelnen Items sehr variabel und in vielen Fällen nur moderat. Bei Merkmalen wie z.B. «Zahneindrücken» ergaben sich hohe Übereinstimmungen, während sich z.B. die Farbbeurteilung der Zungen nicht zuletzt aus technischen Gründen bei der Bewertung am Computerbildschirm als schwierig erwies.

*Stärken und Schwächen der verwendeten Methode.* Das methodische Vorgehen in den Pilotstudien und der Hauptstudie wurde in enger Absprache mit Ärzten entwickelt, die in der TCM sehr erfahren sind. Sowohl die beiden Pilotstudien

wie auch die Hauptstudie wurden unter strikten Blindbedingungen durchgeführt, d.h. die Beurteiler bewerteten die Zungenfotografien streng unabhängig voneinander und hatten keine weiteren Informationen zu den Patienten. Da die Bilder aus einem Pool von mehreren Tausend Fotografien ausgewählt wurden, ist es – vielleicht mit Ausnahme weniger sehr auffälliger Zungen – sehr unwahrscheinlich, dass sich die Beurteiler an die Patienten erinnerten. Dateneingabe und -auswertung wurden in einem unabhängigen Universitätsinstitut durchgeführt. Während in den Pilotstudien aufgrund der begrenzten Ressourcen nur 15 Fotografien verwendet wurden, ist die Hauptstudie mit 101 Fotografien angemessen umfangreich und erlaubt somit eine zuverlässige Aussage hinsichtlich der Übereinstimmung zwischen den teilnehmenden Beurteilern unter den gegebenen Studienbedingungen.

Trotz des systematischen Ansatzes trat bei der Untersuchungsreihe eine Vielzahl von Problemen auf, die zum Teil auch in zukünftigen Studien schwer zu lösen sein dürften. In den vorliegenden Studien erfolgte die Beurteilung aus Gründen der Durchführbarkeit und Standardisierbarkeit auf der Basis von digitalen Zungenfotografien und nicht direkt am Patienten. Die Verwendung von Fotografien und deren Beurteilung am Bildschirm erwies sich jedoch als problematisch. Obwohl die Zungenfotografien unter weitgehend standardisierten Bedingungen entstanden (gleiche Tageszeit, gleicher Raum, gleiche Sitzposition), war ihre Qualität nicht einheitlich. So führen z.B. bereits geringe Lageveränderungen oder die Art und das Ausmaß, wie weit der Patient die Zunge herausstreckt, zu deutlichen Unterschieden in der Ausleuchtung. Ein grundsätzliches Problem bei der Verwendung einzelner Zungenfotografien ist außerdem, dass im Gegensatz zur Beurteilung am realen Patienten, nur ein Eindruck aus einer einzigen Perspektive möglich ist. Die Entscheidung für eine Beurteilung am Bildschirm war getroffen worden, um eine mögliche weitere Verfälschung durch den Ausdruck der Fotografien zu vermeiden. Die Beurteilung am Bildschirm erwies sich jedoch durch den enormen Einfluss der Bildschirmqualität und des Blickwinkels als unbefriedigend. Um die Fehlerquellen zu minimieren, mussten daher alle Beurteilungen am gleichen Bildschirm und in Pilotstudie 2 und der Hauptstudie in einer extra angefertigten Vorrichtung zur weitgehenden Standardisierung des Blickwinkels erfolgen.

In der Praxis der Zungendiagnose spielt der Gesamteindruck der Zunge eine große Rolle. Eine detaillierte morphologische Beschreibung der Zunge erfolgt in der Routinedokumentation selten. Vielmehr erfolgt eine Interpretation der Zungenzeichen unter Einbezug der übrigen Informationen, Eindrücke und Befunde am Patienten. In der vorliegenden Studienserie wurde der Großteil dieses Vorgehens ausgebündet. Um eine einfache Auswertung der Daten zu gewährleisten, waren auf dem Beurteilungsformular einzelne Merkmale als vorhanden oder nicht vorhanden zu beurteilen. Da die Auswahl der Zungenfotografien zufällig erfolgte, ist davon auszugehen, dass eine erhebliche Zahl uneindeutiger

**Tab. 3.** Ergebnisse Hauptstudie, Übereinstimmungen sind fettgedruckt

Item	Beurteiler 1 vs. 2				Beurteiler 1 vs. 4				Beurteiler 2 vs. 4				3 Beurteiler	
	nein	ja	%	Kappa (95%-CI)	nein	ja	%	Kappa (95%-CI)	nein	ja	%	Kappa (95%-CI)	Fleiss' Kappa (95%-CI)	
Zungenspitze rot	nein <b>7</b>	8	85	0,39	<b>12</b>	3	80	0,43	<b>12</b>	2	80	0,44	0,41	
	ja <b>7</b>	<b>78</b>		(0,15 bis 0,64)	17	<b>68</b>		(0,24 bis 0,63)	18	<b>69</b>		(0,25 bis 0,63)	(0,30 bis 0,52)	
Zungenränder rot	nein <b>63</b>	13	84	0,56	<b>72</b>	7	84	0,50	<b>64</b>	2	83	0,56	0,51	
	ja <b>2</b>	<b>19</b>		(0,39 bis 0,73)	9	<b>12</b>		(0,29 bis 0,71)	15	<b>17</b>		(0,38 bis 0,74)	(0,40 bis 0,62)	
Zungenkörper rot	nein <b>82</b>	14	85	0,24	<b>90</b>	5	93	0,33	<b>81</b>	1	87	0,42	0,32	
	ja <b>1</b>	<b>3</b>		(-0,01 bis 0,48)	2	<b>2</b>		(-0,04 bis 0,7)	12	<b>6</b>		(0,18 bis 0,67)	(0,20 bis 0,43)	
Zungenkörper	nein <b>50</b>	16	75	0,45	<b>60</b>	6	82	0,56	<b>53</b>	5	74	0,44	0,46	
blass	ja <b>8</b>	<b>23</b>		(0,27 bis 0,63)	12	<b>20</b>		(0,39 bis 0,74)	20	<b>20</b>		(0,27 bis 0,62)	(0,34 bis 0,57)	
Zungenkörper	nein <b>36</b>	18	76	0,52	<b>16</b>	36	57	0,19	<b>17</b>	24	70	0,36	0,30	
bläulich	ja <b>6</b>	<b>39</b>		(0,36 bis 0,68)	6	<b>39</b>		(0,03 bis 0,34)	6	<b>52</b>		(0,18 bis 0,53)	(0,18 bis 0,41)	
Flecken	nein <b>55</b>	2	70	0,3	<b>38</b>	16	62	0,15	<b>51</b>	25	66	0,18	0,14	
	ja <b>27</b>	<b>12</b>		(0,13 bis 0,46)	18	<b>17</b>		(0,04 bis 0,35)	6	<b>8</b>		(-0,01 bis 0,36)	(0,03 bis 0,25)	
Zunge gedunsen	nein <b>30</b>	2	83	0,65	<b>24</b>	8	81	0,56	<b>30</b>	14	81	0,61	0,60	
	ja <b>15</b>	<b>54</b>		(0,50 bis 0,80)	11	<b>56</b>		(0,39 bis 0,73)	5	<b>50</b>		(0,46 bis 0,77)	(0,49 bis 0,72)	
Zahneindrücke	nein <b>59</b>	4	92	0,83	<b>49</b>	12	83	0,65	<b>50</b>	11	85	0,69	0,70	
	ja <b>4</b>	<b>34</b>		(0,72 bis 0,94)	5	<b>33</b>		(0,51 bis 0,80)	4	<b>34</b>		(0,55 bis 0,84)	(0,59 bis 0,81)	
Risse	nein <b>79</b>	4	91	0,60	<b>68</b>	15	82	0,48	<b>68</b>	17	80	0,43	0,50	
	ja <b>5</b>	<b>11</b>		(0,39 bis 0,81)	3	<b>14</b>		(0,29 bis 0,68)	3	<b>12</b>		(0,24 bis 0,63)	(0,39 bis 0,61)	
Dicker Belag	nein <b>79</b>	17	83	0,32	<b>90</b>	6	91	0,26	<b>78</b>	1	84	0,40	0,32	
	ja <b>0</b>	<b>5</b>		(0,1 bis 0,53)	3	<b>2</b>		(-0,07 bis 0,6)	15	<b>7</b>		(0,17 bis 0,62)	(0,21 bis 0,43)	
Klebriger Belag	nein <b>47</b>	12	73	0,42	<b>45</b>	14	62	0,17	<b>53</b>	8	77	0,49	0,36	
	ja <b>15</b>	<b>26</b>		(0,24 bis 0,60)	24	<b>17</b>		(-0,02 bis 0,36)	15	<b>23</b>		(0,32 bis 0,67)	(0,25 bis 0,48)	
Auflösender Belag	nein <b>88</b>	1	97	0,78	<b>84</b>	4	93	0,55	<b>89</b>	4	94	0,59	0,50	
	ja <b>2</b>	<b>6</b>		(0,55 bis 1,0)	3	<b>5</b>		(0,25 bis 0,85)	2	<b>5</b>		(0,30 bis 0,89)	(0,39 bis 0,62)	
Dünner Belag	nein <b>26</b>	8	78	0,53	<b>20</b>	14	70	0,34	<b>24</b>	16	72	0,41	0,43	
	ja <b>14</b>	<b>53</b>		(0,36 bis 0,70)	16	<b>51</b>		(0,15 bis 0,54)	12	<b>49</b>		(0,23 bis 0,60)	(0,32 bis 0,54)	
	gelb weiß				gelb weiß				gelb weiß					
Belagfarbe	gelb <b>7</b>	1	95	0,71	<b>6</b>	2	94	0,64	<b>8</b>	1	97	0,83	0,71	
	weiß <b>2</b>	<b>88</b>		(0,48 bis 0,95)	2	<b>88</b>		(0,37 bis 0,90)	0	<b>90</b>		(0,64 bis 1,00)	(0,61 bis 0,82)	

n.b. = Nicht berechenbar (bei  $\geq 2$  leeren Feldern in der 4-Felder-Tafel).

Befunde selektiert wurde. In der ersten Pilotstudie erfolgten keinerlei Absprachen, wie «Grenzwertbefunde» zu bewerten seien. In der zweiten Pilotstudie und in eingeschränktem Ausmaß nochmals vor der Hauptstudie erfolgten mündliche Absprachen auf der Basis von Bildern aus den jeweiligen Vorstudien. Derartige «begrenzte» Absprachen scheinen für eine sinnvolle Studiendurchführung unumgänglich. Allerdings erfolgten die Festlegungen in den vorliegenden Studien nicht schriftlich.

*Andere Untersuchungen zum Thema.* Wie bereits angekennert, scheint der Untersuchung der Reliabilität der Zungendiagnose bisher kaum Aufmerksamkeit geschenkt worden zu sein. Entsprechende Arbeiten aus China konnten auch mithilfe chinesischer Fachleute nicht identifiziert werden. Die beiden nur als Abstracts publizierten westlichen Untersuchungen [5, 6] können nur als Pilotstudien interpretiert werden. Biltstein und Bonzak [5] ließen drei erfahrene Akupunkteure 31 Merkmale bei 20 Zungenbildern beurteilen. Die

ermittelten Kappa-Werte lagen zwischen 0,21 und 0,6. Unbefriedigende Kappa-Werte betrafen besonders die Beurteilung der Zungenspitze. Gareus et al. [6] versuchten in einem ersten Schritt, möglichst optimale und standardisierte Bedingungen für die Zungenfotografien zu erreichen, z.B. um immer einen möglichst identischen Lichteinfall zu gewährleisten. Zwei Beurteiler bewerteten dann den Zungenkörper und den Belag direkt am Patienten sowie am Computerbild und nach einer längeren Zeitspanne nochmals das Computerbild. Die Kappa-Indizes beim Vergleich von Computerbild und Direktbeurteilung am Patienten waren 0,35 in Bezug auf den Zungenkörper und 0,34 in Bezug auf den Zungenbelag beim ersten und 0,65 bzw. 0,62 beim zweiten Beurteiler. Die Werte für die Test-Retest-Reliabilität betrugen 0,53 bzw. 0,48 für den ersten und 0,65 bzw. 0,62 für den zweiten Beurteiler. Über das Internet konnte ein Artikel gefunden werden, der eine Übersicht über den Stand der Zungendiagnoseforschung in China gibt [8]. Demnach untergliedert sich die Forschung in diesem Bereich

dort in folgende Bereiche: «Die Systematisierung der schriftlichen Informationen zur Zungendiagnose, die Untersuchung von gesunden Probanden zur Aufstellung eines Standards, die experimentelle Erforschung der Zungendiagnose, die Methodik der Diagnose des Zungengrundes, die Erforschung des Zungenfarbgrades, die Untersuchung der Zungenwärme mit Infrarotmessungen, die Entwicklung von Diagnosegeräten zur Zungendiagnose und schließlich die experimentelle und klinische Untersuchung der Zungenveränderung bei Krankheiten.» Untersuchungen zur Reliabilität werden in der Übersicht nicht genannt.

*Schlussfolgerung und Empfehlungen für zukünftige Forschungsvorhaben.* Die vorliegenden Ergebnisse deuten darauf hin, dass die Beurteilung morphologischer Merkmale im Rahmen der Zungendiagnose eine befriedigende Reliabilität haben dürfte. Eine weitergehende Schlussfolgerung ist angesichts der geringen Zahl vorliegender Untersuchungen, aufgrund der methodischen Schwierigkeiten und aufgrund der unklaren Übertragbarkeit der Ergebnisse auf andere Personen, die TCM praktizieren, nicht vertretbar. Angesichts der Verbreitung der TCM und der Relevanz der Zungendiagnostik im Rahmen dieses Therapiesystems ist eine weitere

Erforschung wünschenswert. Studien zur Reliabilität der morphologischen Beurteilung sind sinnvoll, wenn auch nur als erster Schritt hin zu Untersuchungen, die die Reliabilität der Interpretation der Zungenzeichen überprüfen, sowie Studien, die sich mit der Validität der Zungediagnose im Rahmen der TCM auseinandersetzen. Angesichts der offensichtlich umfangreichen Bemühungen zur Digitalisierung und Objektivierung der Dokumentation in China [1–3, 8] muss versucht werden, eine bessere Übersicht über den Stand der Forschung zu erlangen und Erkenntnisse für weitere Untersuchungen auch im Westen zu nutzen.

## Dank

Die Autoren danken Prof. Dr. S. Wagenpfeil (Institut für medizinische Statistik und Epidemiologie der Technischen Universität München) für die Berechnung der Kappa-Indizes für die Hauptstudie.

## Conflict of Interests

Keine.

## Literatur

- 1 Fletcher RW, Fletcher SW: Clinical Epidemiology, ed 4. Philadelphia, Lippincott, Williams and Wilkins, 2005.
- 2 Chiu CC: A novel approach based on computerized image analysis for traditional Chinese medical diagnosis of the tongue. *Comput Methods Prog Biomed* 2000;61:77–89.
- 3 Wei BC, Shen LS, Wang YQ, Wang YG, Wang AM, Zhao ZX: A digital tongue image analysis instrument for traditional Chinese medicine [in Chinese]. *Zhongguo Yi Liao Qi Xie Za Zhi* 2002; 26:164–9.
- 4 Zhou Y, Yang J, Shen L: Methodological study on digitalization of tongue image in traditional Chinese medical diagnosis [in Chinese]. *Sheng Wu Yi Xue Gong Cheng Xue Za Zhi* 2004;21:917–20.
- 5 Blitstein R, Bonzak S: Interobserver reliability in traditional Chinese medicine tongue observation. *J Altern Complement Med* 2006;12:206.
- 6 Gareus IK, Tan L, Lüdtke R, Niggemeier C, Li FM, Bäcker M, Lehmann C, Klose P, Spahn G, Dobos GJ: Introducing a computer-assisted, digital tongue-imaging device to standardise and evaluate traditional Chinese tongue diagnosis. *Focus Alter Complement Ther* 2005;10(suppl 1):20–21.
- 7 Fleiss JL, Levin B, Paik MC: Statistical methods for rates and proportions, ed 3. New York, Wiley, 2003.
- 8 Zhang BL: Erforschung der Systematisierung und Objektivierung der TCM-Zungendiagnose [ohne Jahreszahl]. [www.tcminter.net/Artikel/Zungendiagnose.htm](http://www.tcminter.net/Artikel/Zungendiagnose.htm).