

Karin Vogt • Jürgen Quetz

Assessment

im Englischunterricht

Kompetenzorientiert beurteilen und bewerten

HELBLING

Innsbruck • Esslingen • Bern-Belp

Inhalt

Vorwort Warum es sich lohnt, dieses Buch zu lesen	5	Kapitel 4 Hör- und Hör-/Sehverstehen	47
Das Autorenteam	6	Geeignete Hör- und Lesetexte finden	47
Kapitel 1 Auftakt: Zur Kunst der Beurteilung sprachlicher Schülerleistungen	7	Authentizität	48
Beurteilung und Bewertung	7	Hörtexte finden und bearbeiten	49
Test? Assessment? Klassenarbeit?		Von der Idee zum Hörtext	51
Prüfung? – Begriffe und ihre Geschichte	9	Vom Text zur Aufgabe	51
Qualitätsmerkmale von Tests	12	Aufgabenbeispiele	52
Funktionen von Assessment	13	Hör-/Sehverstehen	55
Nach dem Assessment ist vor dem Assessment	15	Kapitel 5 Leseverstehen: Zur Auswahl von Texten und Aufgabenformaten	59
Ausblick	20	Authentizität oder „The real thing“	59
Kapitel 2 Referenzrahmen, Bildungsstandards und andere tool boxes	21	Textquellen und -beispiele für Assessment	61
Wann lässt ein Assessment sinnvolle Schlüsse zu?	21	Aufgaben zur Überprüfung des Leseverstehens	65
Gemeinsamer europäischer Referenzrahmen für Sprachen: lernen, lehren, beurteilen	23	Spezifikationen zu Texten und Aufgaben	68
Wozu man die tool box GeR benutzen kann	27	Offene Antworten	68
Bildungsstandards und Vergleichsarbeiten	28	„Von der Idee zum Text und dann zur Aufgabe“ – dieses Mal zum Leseverstehen	68
Bildungsstandards für die Allgemeine Hochschulreife	30	Kapitel 6 Sprechen und Vortragen	75
Testspezifikationen	32	Eine Sprache sprechen = eine Sprache können?	75
Übungsaufgabe: Eine Aufgabe zum Leseverstehen beschreiben	34	Merkmale gesprochener Sprache	75
Wozu der ganze Aufwand?	35	Sprechen als komplexe Kompetenz	77
Kapitel 3 Texte verstehen	37	Sprechen beurteilen – wozu?	78
Das Verb „verstehen“	37	Formate für die Beurteilung mündlicher Sprachkompetenz	79
Ein komplizierter Fall: das Verstehen fiktionaler Texte	39	Arten von mündlichen Aufgaben	80
Fiktionale Texte im Referenzrahmen	41	Kapitel 7 Mündliche Kompetenzen beurteilen	85
Aufgaben für den Unterricht vs. Testaufgaben	43	Aufgaben pilotieren	85
Prozesse und Strategien	43	Bewertungsmöglichkeiten	86
Typen von Aufgaben zur Überprüfung von Hör- und Leseverstehen	45	Role play: A job interview	92
		Rückmeldung geben	94

Kapitel 8 Schreiben: von E-Mails zu formellen Texten	95	„It ain't necessarily so ...“ oder: Was ist ein Fehler?	130
Schreiben können in der Fremdsprache – was bedeutet das?	95	Wie verlässlich sind quantifizierende Bewertungen?	132
Was ist zu beurteilen? Schreiben als Konstrukt	96	Die Frage der Objektivität	133
Textsorten und Prüfungsaufgaben für das Schreiben	98	Unverzichtbar: das Beurteilertraining	134
Beurteilung von Schreibprozessen	98	Geschlossene Aufgabenformate und das Kriterium Reliabilität	134
		Fazit	136
Kapitel 9 Schreiben beurteilen	101	Kapitel 14 Vom Wiegen wird die Sau nicht fett	137
Geschlossene Aufgabenformate	101	PISA und die Folgen: Bildungsstandards	138
Bewerten von schriftlichen Leistungen mit offenen Aufgabenformaten	101	DESI: Bildungsmonitoring für den Englischunterricht	139
Kapitel 10 Sprachmittlung	109	„Is my BI your BI?“	140
Vier Fertigkeiten oder fünf?	109	Vergleichsarbeiten (VERA-8)	141
Sprachmittlung in Bildungsstandards und Abschlussprüfungen	110	Abschluss tests für Hauptschulen und den Mittleren Bildungsabschluss	141
Sprachmittlung im erweiterten Referenzrahmen	115	Das Zentralabitur	142
Kapitel 11 Was man alles nicht testen und prüfen kann (oder sollte)	119	Washback oder auch backwash	142
Dominanz von Kommunikativer Kompetenz in Prüfungen	119	Wird die Sau vom Wiegen doch fett?	143
Interkulturelle Kompetenz – was ist das eigentlich?	119	Kapitel 15 Es geht auch anders – Alternativen beim Assessment	145
Interkulturelle Kompetenz (IK) in der unterrichtsbasierten Leistungsbeurteilung	120	Alternativen beim Assessment	145
Lernerautonomie beurteilen	124	Welche Alternativen gibt es zur Beurteilung durch die Lehrkraft?	145
Kapitel 12 Prüfungen und Klassenarbeiten – ganz praktisch gesehen	125	Baustelle voraus: differenzierende Leistungsbeurteilung	151
Assessment ganz praktisch	125	Lösungsschlüssel	153
Vorbereitung von Prüfungen	125	Glossar	159
Nachteilsausgleich	126	Literaturverzeichnis	163
Nach der Prüfung	127	Texte und Bilder	166
Kapitel 13 Wie verlässlich sind Bewertungen?	129		
Verlässlichkeit (Reliabilität) als zentraler Aspekt von Bewertungen	129		
Offene und geschlossene Aufgaben – re-visited	129		

Kapitel I

Auftakt: Zur Kunst der Beurteilung sprachlicher Schülerleistungen

Im Titel dieses Buches wird das Wort „Assessment“ als übergeordneter Begriff für die unterschiedlichsten Dinge verwendet, die behandelt werden sollen. Assessment ist kein deutsches Wort, aber Sie kennen es sicher in Ausdrücken wie Assessment-Center. Wer in der Berufswelt ein Assessment-Center absolviert, nimmt an einem Testverfahren teil, in dem festgestellt wird, ob er oder

sie für die Position geeignet ist, die eine Firma ausgeschrieben hat. In der Pädagogik ist es an den Lehrerinnen und Lehrern¹, in Lerngruppen und deren Leistungen zu arbeiten und zu bewerten. Sowohl für die Schüler als auch für die Lehrkräfte spielt Assessment eine zentrale Rolle. Dies beginnt schon während des Unterrichts.

I Beurteilung und Bewertung

Ein ganz normaler Mittwochmorgen. Frau S. kommt noch einige Aufgaben im Workbook zu machen, zum Englischunterricht in ihre Klasse. Hausaufgabe war, Vokabeln zu lernen. Die Lehrerin fragt stumm die mündlich ab. Einzelne Schüler werden befragt, Frau S. markiert sich kurz eine Note in ihr Notenkonto: „Well done, Oliver, you've answered all my questions!“ – „Oh, Anne, the correct word is *Natives* and not *Indians*. You obviously didn't pay attention last week when we talked about this.“ Danach steht das Üben der *-ing form* auf dem Programm. Linus hat scheinbar noch Schwierigkeiten mit der Anwendung von Präpositionen. Frau S. nimmt sich vor, das noch einmal genauer zu erläutern. Nun ist das Textverständnis dran, es geht um Gruppenarbeit zum Lehrvideo. Frau S. gibt den Schülern knapp Feedback zu den Ergebnissen der Gruppenarbeit. „Group A has found out a lot of interesting details. Good presentation, however, you should have more contact with your audience.“ Eine Kollegin lässt die Schüler untereinander Feedback zu Lernergebnissen geben, das möchte Frau S. auch mal ausprobieren. Doch nun sind

noch einige Aufgaben im Workbook zu machen. Die Lehrerin fragt stumm die mündlich ab. Einzelne Schüler werden befragt, Frau S. markiert sich kurz eine Note in ihr Notenkonto: „Well done, Oliver, you've answered all my questions!“ – „Oh, Anne, the correct word is *Natives* and not *Indians*. You obviously didn't pay attention last week when we talked about this.“ Danach steht das Üben der *-ing form* auf dem Programm. Linus hat scheinbar noch Schwierigkeiten mit der Anwendung von Präpositionen. Frau S. nimmt sich vor, das noch einmal genauer zu erläutern. Nun ist das Textverständnis dran, es geht um Gruppenarbeit zum Lehrvideo. Frau S. gibt den Schülern knapp Feedback zu den Ergebnissen der Gruppenarbeit. „Group A has found out a lot of interesting details. Good presentation, however, you should have more contact with your audience.“ Eine Kollegin lässt die Schüler untereinander Feedback zu Lernergebnissen geben, das möchte Frau S. auch mal ausprobieren. Doch nun sind

Die Überprüfung, ob Vokabeln gelernt worden sind, liegt schon an der Grenze zu den formelleren Lernerfolgskontrollen. Hier erfolgt oftmals eine Benotung und somit eine *Bewertung*. Das wird bereits

¹ Um den Lesefluss nicht zu behindern, wird im Folgenden vornehmlich die maskuline Form von Schüler oder Lehrer verwendet, Schülerinnen und Lehrerinnen sind dabei selbstverständlich inkludiert. Wir bitten um Verständnis.

daran deutlich, dass das Notenbuch eine wichtige Rolle spielt: was dort nicht an Bewertungen, meist in Ziffernform, festgehalten ist, kann nicht rekonstruiert werden und spielt deshalb bei der Bewertung von Schülerleistungen langfristig keine Rolle. Terminologisch gilt es zwischen einer Bewertung

und der eher informellen Beurteilung zu unterscheiden. Beurteilungen können sich natürlich auch in schriftlicher Form finden, etwa wenn Schüler bei einer Übungsarbeit starke rote Kommentare erhalten, die über Stärken und Schwächen Auskunft geben. Hier ein Beispiel einer Klassenarbeit.

(1) Ganzheitliche Beurteilung

Couch Potatoes

Some couch potatoes haven't got friends. Their didn't know what they can do with the time. So many couch potatoes sit on the couch and watch TV or play computer or Playstation. But some children are playing on the computer for money. And for some children will it be a addiction. They come from the school and go on the computer or TV. A friend of me are a couch potato. All the day she sits on the couch and watch TV. She don't want to go out with me or her friends. But her mother did nothing. Her mother means that's okay. Sometimes I am a couch potato when I am ill. It's often the Sunday when I am don't know what I can do and then I am turn the TV on and watch all the soaps. But I didn't make it often. I think all the couch potatoes don't know what they can do with the time so they make the TV on and watch TV.

A nice little essay, Barbara! The reader learns from it that you are no couch potato, but you understand why others sometimes are. A good idea: parents should be interested in what their children do!

2. preparation for the next test: please study the tenses, especially the simple present and the negative constructions (do, don't, not: do'n't - doesn't) that generally well done!

(2) Analytische Beurteilung

Couch Potatoes

Some couch potatoes have no friends. Their didn't know what they can do with the time. So many couch potatoes sit on the couch and watch TV or play computer s or Playstation. But some children are playing on the computer for money. And for some children will it be a addiction. They come from the school and go on the computer or TV. A friend of me are a couch potato. All the day she sits on the couch and watch TV. She don't want to go out with me or her

pl. = potatoes / haven't
they don't / their

ξ games (expr.)
play (Gr) / ? (What do you mean?)
become / an (Gr) / word order
turn
mine is
watches
doesn't

friends. But her mother did nothing. Her mother means that is okay. Sometimes I am a couch potato when I am ill. It's often the Sunday when I am don't know what I can do and then I am turn the TV on and watch all the soaps. But I didn't make this often. I think all the couch potatoes do/nt know what they can do with the time so they make the TV on and watch TV.

*doesn't do anything
thinks (that this is)*

Gr!

*do
-oes / don't (love)
(same as second chance) / turn*

*Too many mistakes, Barbara!
I hope the next will be better!*

Eine ganzheitliche Beurteilung einer Schülerleistung erfordert von der Lehrkraft ein genaues Verständnis der Stärken und Schwächen einer sprachlichen Leistung und ist somit zeitaufwändiger und vielleicht auch schwieriger als der bloße Verweis auf einzelne Fehler. Andererseits ist sie sehr hilfreich für Schüler, weil konkrete Hinweise auf Fehlerschwerpunkte gegeben werden, auf die sie bei weiteren Leistungen ihre Aufmerksamkeit richten sollten.

Ähnlich wie ihre die meisten Lehrerinnen und Lehrer bei ihren Beiträgen. Um die Kommunikation nicht zu stören machen sie sich in Gedanken zu inhalieren oder typischen formalen Schwächen der Äußerung und thematisieren diese dann zu gegebenem Zeitpunkt.

Das sind aber alles noch Beurteilungen, die sich nicht unbedingt in einer Bewertung durch eine Note niederschlagen müssen.

2 Test? Assessment? Klassenarbeit? Prüfung? – Begriffe und ihre Geschichte

Wie kommt es eigentlich zu der Vielfalt an Begriffen für im Grunde gleiche oder ähnliche Dinge? Zu Beginn des vorigen Jahrhunderts gab es in Deutschland eigentlich zwei Begriffe: die „Klassenarbeit“ und die „(Abschließende Aufnahme...)-Prüfung“. Es folgt dann ein kurzer Abriss zur Entwicklung im Bereich von Tests bzw. Assessments, der sich zusammensetzt, Schneisen in den terminologischen Landschaften schlagen.

Traditionelle Assessments

In den Jahren vor dem 2. Weltkrieg war man der Meinung, dass sich die Beherrschung einer Sprache vor allem an der korrekten Lösung von Grammatikaufgaben feststellen lasse. Wirkliches Sprachkönnen unterstellte man aber erst, wenn die Lernenden in der Lage waren, Übersetzungen in die Zielsprache vorzunehmen oder Aufsätze zu

schreiben. Diese traditionsreichen Prüfungsformen waren allerdings schon vor dem 2. Weltkrieg als höchst subjektive Formen der Leistungsmessung in die Kritik geraten. Aufsätze kann man bekanntlich kaum objektiv beurteilen, wie zahlreiche Studien zeigen (vgl. u.a. Ingenkamp 1971).

Noten für den gleichen Aufsatz oder die gleiche Mathearbeit, das hat man immer wieder festgestellt, schwanken bei verschiedenen Korrektoren fast über die ganze Notenskala, und Schüler können sich nicht sicher sein, dass die gleiche Leistung von der gleichen Lehrkraft zu verschiedenen Zeitpunkten auch gleich beurteilt wird. Schon die Position eines Klassenarbeitshefts innerhalb eines Stapels von 30 anderen Heften kann die Note beeinflussen, wie man weiß: Wir ermüden bei der Korrektur, verlieren den Überblick und werden, je nach Naturell, gereizt oder gnädiger im Urteil, je mehr Hefte man korrigiert hat.



© Claire Bretécher, 1989

Schüler empfinden daher solche Bewertungen von Leistungen oft als subjektiv und ungerecht.

Da mit Klassenarbeiten in der Schule auf lange Sicht auch immer eine Auslesefunktion verbunden ist, kann man auf Abhilfe. Die Korrektur von Aufsätzen zu standardisieren und damit deren Benotung vergleichbarer zu machen, entpuppte sich jedoch als schwierig. Je akribischer ein Korrekturverfahren, desto unhandlicher ist es und damit ungeeignet für die tägliche Unterrichtspraxis. Wie lässt sich diesem Dilemma entkommen?

Neue psychometrische Testverfahren

In den 1960er Jahren wurde die behavioristische Lerntheorie und der strukturalistische Linguistik beruhende audiolinguale Methode immer populärer und man schuf den „richtigen“ Test. Diese Methode ging man davon aus, dass man einen begrenzten Vorrat an grammatischen Strukturen und Wörtern so einprägen kann, dass man sie ohne Nachdenken spontan abrufen kann und auf andere Situationen übertragen werden können. Voraussetzung dafür ist, dass man jede Struktur genau isolieren und gegen andere austauschen kann. Wie das Lehren kleinste Einheiten (*discrete points*) das formale Gerüst der Sprache beherrschen sollte, so wollte man auch die Fertigkeiten Hören, Sprechen, Lesen und Schreiben streng getrennt einüben. Tests des Hör- oder Leseverstehens zum Beispiel sollten so konstruiert sein, dass man bei der Lösung nicht auch noch schreiben musste. Der audiolinguale Unter-

richt war die Grundlage die Kunst des kleinschrittigen und wohlgeleiteten Lehrens, das dem Lernen die richtigen Wege bereiten sollte.

Die Behauptung, Sprachkompetenz resultiere aus der (nahe) mechanischen Addition isolierter Elemente, folgt dem Umkehrschluss, dass Lernende, die über die einzelnen Elemente verfügen auch Sprachkompetenz besitzen müssten. In den 1960er Jahren setzten sich folglich psychometrische Verfahren der Leistungsmessung durch, die bereits vor dem 2. Weltkrieg entwickelt und erprobt worden waren. Psychometrisch heißen diese Verfahren, weil sie in der Psychologie entwickelt wurden. Mit der Absicht Konstrukte wie Intelligenz, Motivation, Angst usw. zu messen, wurde von beobachtbaren Daten („Igitt, eine Spinne!“) auf Konstrukte geschlossen, die ihnen zu Grunde liegen (Angst, Ekel). Gleichermäßen wollte man von den beobachtbaren Daten (Zahl der richtig gelösten Aufgaben) auf die ihnen zu Grunde liegenden kommunikativen Fähigkeiten schließen: Kann ein Lernender einen Test lösen, so kann man wohl auch annehmen, dass er in der Zielsprache kommunizieren kann.

Man bemühte sich vor allem um die Objektivierung der Leistungsmessung, indem man die verwendeten Tests möglichst zuverlässig (engl. *reliable*) gestalten wollte. Zuverlässig heißt, dass Lernende für eine vergleichbare Leistung auch eine vergleichbare Punktzahl erhalten. Möglichst wenig in einem Test sollte der subjektiven Bewertung des Prüfers überlassen bleiben (siehe Cartoon oben). Des Weiteren

sollten die zu erreichenden Lernziele genau definiert sein und dann objektiv und verlässlich gemessen werden.

Geschlossene Testformate sollten störende Einflüsse bei der Lösung möglichst ausschließen. Die Beurteilung, ob die Ziele erreicht worden waren, sollte nach Möglichkeit kein subjektives Urteil mehr erfordern, sondern neutral und mechanisch erfolgen können. Die *Multiple-Choice*-Aufgabe hielt ihren Einzug in die Testdidaktik:

The singer ended the concert ... one of his most popular songs.

A by B with C in D as

Da die Fremdsprachendidaktik in den 1940er bis 1960er Jahren Sprache zu einem in kleinste Bausteine zerlegbaren Baukasten erklärt hatte, stieß die Isolation von messbaren Aufgaben kaum auf theoretische Bedenken. Die psychometrischen Messverfahren waren schließlich auf dem gleichen Pfad einer behavioristischen Psychologie entstanden wie die audiolinguale Methode, was sie für Fremdsprachendidaktiker ohne weiteres akzeptierbar machte. Testdidaktiker wie Robert Lado (im ersten epochalen Werk *Language Testing*, 1960) setzten psychometrische Tests als Standard durch.

Die USA, aber auch Großbritannien, brachten eine sehr große Zahl von Studierenden an, was die Universitäten wiederum vor erhebliche Probleme stellte, da sichergestellt werden musste, dass diese Studierenden hinreichend gut Englisch sprachen. Aufnahmeprüfungen wurden notwendig gemacht. Sie mussten vornehmlich maschinell auswertbar sein und zu zuverlässigen Ergebnissen führen. Die Entwicklung einer Testdidaktik im großen Stil ist folglich eng mit einer wissenschaftlichen Disziplin verknüpft – den Sprachprüfungen großer Organisationen.

In jedem Jahr absolvieren Millionen von Menschen Sprachprüfungen für Englisch ab. Die allermeisten davon wählen eins der großen, international anerkannten Examen der Universität Cambridge (heute *Cambridge English Language Assessment*, dessen bekanntester Test das *Certificate of Proficiency in*

English ist) oder den amerikanischen *Test of English as a Foreign Language (TOEFL)*, der von Ausländern für eine Immatrikulation an US-amerikanischen Universitäten vorzulegen. Man stelle sich vor, welche immense logistische Leistung es zu vollbringen ist! In *Measured Words* (1997) beschreibt Bernard Spolsky, wie die amerikanischen Universitäten Anfang der 60er Jahre TOEFL zu einem maschinell auswertbaren, objektiven, standardisierten und kostengünstigen Prüfungsinstrument entwickelten, um den Korrekturaufwand bei offenen Aufgabenformate in traditionellen Prüfungen zu reduzieren. Mittlerweile ist die maschinelle Auswertung Standard in allen großen Prüfungszentralen, und computerbasierte Sprachtests gibt es inzwischen längst.

Vom psychometrischen zum kommunikativen Test

Man hatte es bereits unter den Testdidaktikern der 1960er Jahre Auseinandersetzungen um den Sinn dieser psychometrischen Tests gegeben. In seinem 1967 erschienenen Buch *Fundamental considerations in testing for English language proficiency of foreign students* mahnte John B. Carroll (1961) beispielsweise an, dass Tests dieser Art bei komplexeren Fertigkeiten, wie Sprechen und Schreiben nicht sonderlich valide (also inhaltlich aussagekräftig) seien. Selbst wenn ein Kandidat die in kleinschrittige Aufgaben aufgebrochenen Teiltests richtig gelöst hätte, könnte dies nicht unbedingt Auskunft über den wirklichen Stand seiner Sprachkompetenz geben. Dies wurde zum zentralen Kritikpunkt an den Testverfahren der 1960er Jahre und führte nach und nach zu einer terminologischen Änderung: Statt von Tests sprach man jetzt lieber von Assessment und versuchte so, einen neuen Ansatz hervorzuheben.

Im Zuge der Hinwendung zum kommunikativen Fremdsprachenunterricht entwickelte die Testdidaktik in den 1970er Jahren auch kommunikative Sprachprüfungen. Man erkannte, dass kommunikatives Handeln ein komplexes Zusammenspiel verschiedener Teilfertigkeiten und sprachlicher Elemente erfordert, die verständnisvoll und strategisch

planvoll, oder auch kreativ und spielerisch eingesetzt werden. Zwar bemühte man sich auch in kommunikativen Curricula um detaillierte Lernzielkataloge (z. B. im wichtigsten Dokument dieser Zeit, dem *Threshold Level English* des Europarats, entwickelt von Jan van Ek 1975), begriff aber schnell die potenzielle Unendlichkeit der angestrebten Listen sowie ihre relative Beliebigkeit und benutzte sie daher eher als lockere Referenzmittel denn als strenge Lehrpläne. In kommunikativen Tests ging es wie im Unterricht darum, was man mit Wörtern tun kann, wie man sprachlich handelt, Absichten versprachlicht, Handlungsziele erreicht, sich mit anderen Menschen auseinandersetzt, Texte produziert, austauscht und rezipiert. Formales linguistisches Wissen wurde zwar als notwendige Grundlage für kommunikative Kompetenzen und kommunikatives Handeln akzeptiert, erhielt jedoch im didaktischen Denken eine

untergeordnete, quasi dienende Funktion. Als Problem der älteren psychometrischen Tests erachtete man die mangelnde *Validität* (Gültigkeit) der Aufgaben. Kommunikative Praktiker kritisierten, dass die objektivierten Testaufstellungen ihrer Vorgänger nicht die Lernziele repräsentierten, die einem modernen Fremdsprachenunterricht zugrunde liegen. Setzt man sich kommunikativ kompetentes Handeln zum Ziel, muss man sich fragen, ob die Lernenden auch wirklich kommunikativ handeln können. Dies lässt sich aber nur durch Testaufgaben (oder Formate) sicherstellen, die so offen angelegt sind wie menschliche Kommunikation. Um die entsprechende Schiebung zu markieren, wurde *testing* von *assessment* und somit zu einem *Language Assessment* (LAss) abgewandelt, in dem es um mehr und anderes gehen sollte.

You recently took part in a class discussion about choosing an interesting job. Your teacher has now asked you to write a composition, answering the following question and giving reasons for your choice. Would you rather be a politician, a teacher or a musician?

(ALTE: *Materials for the Guidance of Test Item Writers*, 1997, S. 139)

Natürlich hatte man dadurch schlafende Hunde geweckt, denn das Problem der Objektivität, das durch die verlässlichen (reliablen) psychometrischen Tests der vorherigen Jahre gebannt zu sein schien, stellte sich bei kommunikativen Tests erneut ein.

Der Weg vom behavioristischen zum kommunikativen Paradigma bedeutete gleichzeitig eine Verlagerung des Interesses vom Aspekt der Reliabilität auf den der Validität von Tests. Die beiden Begriffe spiegeln demnach unterschiedliche Konzepte von Sprachtests bzw. Assessment wider.

3 Qualitätsmerkmale von Tests

Die großen Institutionen wie Cambridge English Language Assessment oder der Educational Testing Service (ETS) mit dem TOEFL/TOEIC haben es geschafft, die inhaltliche Gültigkeit, also die Validität ihrer Tests im Sinne moderner, kommunikativer Ansätze des Fremdsprachenunterrichts, zu aktualisieren. Außerdem tun sie alles, um durch immer raffiniertere statistische Auswertungsverfahren al-

les Zufällige bei der Bewertung in den Griff zu bekommen und dadurch die Verlässlichkeit von Tests zu verbessern (vgl. Kapitel 13). Von zentraler Bedeutung sind zudem die neueren Ansätze zur Verbesserung der Validität durch klare Orientierung an Dokumenten wie dem *Gemeinsamen europäischen Referenzrahmen für Sprachen: lernen, lehren, beurteilen* (Europarat 2001, vgl. Kapitel 2).

2 Gemeinsamer europäischer Referenzrahmen für Sprachen: lernen, lehren, beurteilen

Der *Gemeinsame europäische Referenzrahmen: lernen, lehren, beurteilen* (Europarat 2001, im Folgenden GeR abgekürzt) gilt als international wichtigstes Dokument für die Beschreibung von sprachlichen Kompetenzen. Es wäre jedoch zu kurz gegriffen, den GeR darauf zu reduzieren, dass er sprachliche Kompetenzen beschreibt und dokumentiert. Primär ist er ein sprachpolitisches Dokument, das Mehrsprachigkeit befürwortet und für die Förderung von weniger häufig gesprochenen Sprachen eintritt. Für den

Fremdsprachenunterricht beide interessant ist auch die Empfehlung für einen handlungsorientierten Ansatz.

Am bekanntesten geworden sind allerdings die Skalen 'Deskriptoren sprachlicher Kompetenzen'. Die nachhaltige Wirkung auf das sprachdidaktische Handlungsfeld entstand durch die im GeR verankerte Skala mit den Niveaustufen A1 bis C2. Die weltweit bekanntesten Sprachtests für Englisch, die auf dieses System basieren¹, sind z. B. die folgenden Tests:

A Elementare Sprachverwendung		B Selbständige Sprachverwendung		C Kompetente Sprachverwendung	
A1	A2	B1	B2	C1	C2
	Key English Test (KET)	Preliminary English Test (PET)	Certificate in English (FCE)	Certificate in Advanced English (CAE)	Cambridge Proficiency in English (CPE)

Die Buchstaben A, B und C stehen für die drei Stufen *Beginners* (Grundstufe), *Intermediate* (Mittelstufe) und *Advanced* (Oberstufe). Sie sind aus sprachlichen Gründen noch einmal unterteilt (A1, A2 usw.). Wenn man möchte, kann man im ersten noch weiter differenzieren (B1.1, B1.2 usw.).

Doch wie verhalten sich die Niveaus zueinander und wie leicht lassen sich die sechs oder mehr Stufen erklimmen? Die meisten internationalen Testanbieter sind sich darin einig, dass es von A1 zu A2 in etwa halb so lange dauert wie von A2 bis B1.

[Die Erfahrung zeigt, dass viele Lernende für den Weg von A1 bis A2 doppelt so viel Zeit benötigen wie für den Weg zu [A2]. Sie werden

folglich wohl auch mehr als doppelt so lange brauchen, um von [B1] aus [B2] zu erreichen – selbst, wenn die Niveaustufen auf der Skala den gleichen Abstand voneinander zu haben scheinen. Das ist darauf zurückzuführen, dass sich das Spektrum der Sprachaktivitäten, der Fertigkeiten und der sprachlichen Mittel nach oben hin notwendigerweise verbreitert. Diese Tatsache spiegelt sich darin, dass eine Skala von Niveaustufen oft als Diagramm „im Eistütenformat“ dargestellt wird, als ein dreidimensionaler Konus, der nach oben hin breiter wird. Was Aussagen zur durchschnittlichen Lernzeit für das Erreichen eines bestimmten Niveaus mit bestimmten Zielen betrifft, ist größte Vorsicht geboten.

(Europarat 2001: 29, Kap. 2.2)²

¹ Seit einigen Jahren ist eine Neubearbeitung auf dem Weg, die 2017 abgeschlossen wurde. Einige dieser Änderungen werden auch in Kapiteln unseres Buches aufgegriffen.

² Des Textflusses wegen zitieren wir nur aus der gedruckten deutschen Fassung, die 2001 im Langenscheidt-Verlag erschienen ist. Die „elektronische“ Version findet man auf der Website des Goethe-Instituts unter <http://www.goethe.de/Z/50/commeuro/deindex.htm> (letzter Zugriff: 01.12.2017).

Weil inzwischen immer mehr Menschen als politische Flüchtlinge oder aus wirtschaftlichen Gründen nach Europa kommen und die jeweilige Landessprache nie in der Schule lernen konnten, gibt es in der neuen Fassung des GeR (2018) eine Stufe *pre-A1* (dt. vor A1), also eine Art *survival level*. Aber auch für die Grundschule oder für dritte Fremdsprachen könnte dieses Niveau interessant sein.

Diese Erweiterung um pre-A1 wird wahrscheinlich zu einer Umstrukturierung des ganzen Systems der Referenzniveaus führen, die neu „geeicht“ werden müssten.

Wie die Beschreibung solcher Niveaus im Detail aussieht, zeigt die ursprüngliche Globalskala aus dem Jahr 2001, die hier auf anderen Skalen abgeleitet sind.

Kompetente Sprachverwendung	C2	<i>Kann praktisch alles, was er/sie liest und hört, mühelos verstehen. [...] Kann sich spontan, sehr flüssig und genau ausdrücken und auch bei komplexeren Sachverhalten feinere Bedeutungsnuancen deutlich machen.³</i>
	C1	<i>Kann ein breites Spektrum anspruchsvoller längerer Texte verstehen und auch implizite Bedeutungen erkennen. Kann die Sprache im gesellschaftlichen und beruflichen Leben [...] und flexibel gebrauchen. [...] Kann sich klar, strukturiert und ausführlich zu komplexen Sachverhalten äußern. [...]</i>
Selbstständige Sprachverwendung	B2	<i>Kann die Hauptinhalte komplexer Texte zu konkreten und abstrakten Themen verstehen; versteht im eigenen Spezialgebiet auch Fachdiskussionen. Kann sich so spontan und fließend verständigen, dass ein normales Gespräch mit Muttersprachlern ohne größere Anstrengung auf beiden Seiten gut möglich ist. Kann sich zu einer Vielzahl von Themen auf einer breiten Themenspektrum klar und detailliert ausdrücken, [...]</i>
	B1	<i>Kann die Hauptinhalte verstehen, wenn klare Standardsprache verwendet wird und wenn es um vertraute Dinge aus Arbeit, Schule, Freizeit usw. geht. Kann sich einfach und zusammenhängend über vertraute Themen und persönliche Interessengebiete äußern. [...]</i>
Elementare Sprachverwendung	A2	<i>Kann Sätze und kurze Gebrauchs-Ausdrücke verstehen, die mit Bereichen von ganz unmittelbarer Bedeutung zusammenhängen (z.B. Informationen zur Person und zur Familie, Einkaufen, Arbeit, nähere Umgebung). Kann sich in einfachen, routinemäßigen Situationen verständigen, in denen es um einen einfachen und direkten Austausch von Informationen über vertraute Themen und geläufige Dinge geht.</i>
	A1	<i>Kann vertraute, alltägliche Ausdrücke und ganz einfache Sätze verstehen und verwenden, die auf die Befriedigung konkreter Bedürfnisse zielen. [...] Kann sich auf einfache Art verständigen, wenn die Gesprächspartnerinnen oder Gesprächspartner langsam und deutlich sprechen [...].</i>

(Europarat 2001: 35, hier leicht abgeändert v. Verf.)

In der Globalskala sind sowohl rezeptive als auch produktive Kompetenzen aufgeführt. Um den Aufbau und die Struktur einzelner Deskriptoren zu verdeutlichen, wurden „Kann“ und das darauf folgende Verb kursiv hervorgehoben.

Diese Skala wird sich wegen des neuen Niveaus pre-A1 ein wenig verändern. Durch „vor A1“ verschieben sich auch zwangsläufig einige Elemente im gesamten Gefüge. Auf diesem neuen Niveau findet sich nun beispielsweise eine Variante des ursprüng-

³ Vor allem auf Niveau C2 enthält die Ergänzung von 2017 zum GeR einige Änderungen: wurden im GeR 2001 noch „Muttersprachler“ erwähnt, die man verstehen kann, heißt es jetzt „Gesprächspartner“, die flüssig und schnell sprechen.

lichen Deskriptors auf A1, der den Sprecher so einstuft, dass er Fragen über sich selbst stellen und beantworten kann, die unmittelbare Bedürfnisse und Routinen zum Inhalt haben und sich dabei formelhafter Redewendungen oder gar Gesten bedienen kann, um die Information zu unterstützen. Solche Bedeutungsunterschiede sind sehr schwer an konkreten Versprachlichungen festzumachen, weil die Unterschiede zwischen „vor A1“ und A1 oft minimal sind. Zumal sie für den Alltag des Assessments auch wenig bedeutsam sind, da sich Klassenarbei-

ten in der 5. Jahrgangsstufe eher am Lehrbuch orientieren werden als an Kompetenzen, über die ein Tourist oder Zuwanderer verfügen sollte.

Weitere Verästelungen innerhalb dieser Globalskala finden sich in unterschiedlichen Skalen zu den speziellen Fertigkeiten (z. B. Sprechen usw.). Diese ausführlichen Versionen sind präziser als die übergeordnete Globalskala. Exemplarisch soll dies an der Skala des GeR für mündliche Interaktion gezeigt werden und zwar beschränkt auf die Niveaus A1 bis B2.

Mündliche Interaktion allgemein	
B2	Kann die Sprache gebrauchen, um flüssig, korrekt und wirksam über ein breites Spektrum allgemeiner, wissenschaftlicher, beruflicher Themen oder über Freizeithemen zu sprechen und dabei Zusammenhänge zwischen Ideen deutlich zu machen.
B1	Kann Informationen austauschen, prüfen und bestätigen, mit wenigen routinemäßigen Situationen umgehen und erklären, warum etwas problematisch ist. Kann Gedanken zu eher abstrakten kulturellen Themen ausdrücken, wie z. B. zu Filmen, Büchern, Musik usw.
A2	Kann ohne übermäßige Mühe in einfachen Dialogen zu recht kommen; kann Fragen stellen und beantworten und in vorhersehbaren Alltagssituationen Gedanken und Informationen zu vertrauten Themen austauschen.
A1	Kann sich auf einfache Art verständigen, solange die Kommunikation völlig davon abhängig, dass etwas langsamer wiederholt, reformuliert oder korrigiert wird. Kann einfache Fragen stellen, einfache Feststellungen treffen oder auf solche reagieren, sofern es sich um unmittelbare Bedürfnisse oder um sehr vertraute Themen handelt.

(Europarat 2001: 114, Auszug A1 – B2)

Das System der Referenzniveaus seine hebelliche Dynamik entfaltet. Alle europäischen Testcenter in Europa (und viele weltweit) haben die Sprachprüfungen auf den GeR abgestimmt.⁴ Dieses *Common European Framework of Reference for Languages: Learning, Teaching, Assessment (CEFR)*, wie es im Original heißt (Council of Europe 2001), ist ein Hilfsmittel, mit dem man vor allem wichtige Prüfungen bewerten und „eichen“ kann. Ein großer Vorteil dieses Modells liegt darin, dass sich bei der Beurteilung unterschiedliche Profile in den Kompetenzen der Lernenden erfassen lassen: Einige von ihnen können z. B. flüssig lesen, andere hören und sprechen besser.

Eine ‚dritte Dimension‘

Im GeR fallen in erster Linie die Skalen mit den Deskriptoren für sprachliche Aktivitäten (Sprechen, Hören, Lesen, Schreiben usw.) ins Auge. Sie bestehen aus einer horizontalen Dimension (den Kompetenzbeschreibungen) und einer vertikalen Dimension (den Stufungen von „vor A1“ über A1 bis C2). Die Skalen basieren in erster Linie auf Textsorten, Themen und Inhalten und behandeln den jeweiligen Umgang damit. Hierfür braucht man allerdings vielfältige Mittel. Dazu zählen Wortschatz, Syntax, eine verständliche Aussprache, aber auch interkulturelles Wissen und das Verfügen über sozial angemessene Ausdrucksweisen. Diese Katego-

⁴ Eine Übersicht findet sich auf <http://www.alte.org/membership> (letzter Zugriff: 01.12.2017).

rien bilden eine Art ‚dritte Dimension‘, die zu der Erfassung von Sprachkompetenz beitragen.

Zu dieser ‚dritten Dimension‘ gibt es ebenfalls Skalen, mit deren Hilfe man sprachliche Kompeten-

zen grob beschreiben kann. Zur Illustration werden die Deskriptoren der Niveaus A1 und B1 im Hinblick auf die Verwendung von Grammatik angeführt.

Grammatische Korrektheit	
B1	Kann sich in vertrauten Situationen ausreichend korrekt verständlich machen; im Allgemeinen gute Beherrschung der grammatischen Strukturen trotz deutlicher Fehler in der Muttersprache. Zwar kommen Fehler vor, aber es bleibt klar, was ausgedrückt werden soll.
A1	Zeigt nur eine begrenzte Beherrschung einiger weniger einfacher grammatischer Strukturen und Satzmuster in einem auswendig gelernten Repertoire.

(Europarat 2001: 114, Auszug A1 und B1, stark gekürzt)

Es gibt zwar keine Listen mit Wortschatz oder Syntax, mit deren Hilfe man Texte einstufen könnte, aber ein nachträglich entstandenes Werkzeug: *English Profile*⁵ hilft bei der Analyse von Texten und hat sich auch für das Erstellen von Tests als nützlich erwiesen.

Im Zusammenhang mit dieser ‚dritten Dimension‘ wird oft übersehen, dass beim Verstehen von Texten und bei anderen Produktion auch soziolinguistische, strategische und pragmatische Kompetenzen benötigt werden. Am einfachsten sind *soziolinguistische Kompetenzen* zu erläutern. Im Fall von Niveau B1 sieht es beispielsweise im GeR:

B1	Kann ein breites Spektrum von Sprachfunktionen realisieren und auf sie reagieren, indem er/sie die dafür gebräuchlichsten Redemittel in einem neutralen Register benutzt. Ist sich der wichtigsten Höflichkeitsformen bewusst und handelt entsprechend. Ist sich der wichtigsten Unterschiede zwischen den Sitten und Gebräuchen, den Einstellungen, Werten und Überzeugungen in der betreffenden Gesellschaft und in seiner eigenen bewusst und achtet auf entsprechende Signale.
-----------	--

(Europarat 2001: 122)

Zu den Kompetenzen in der ‚dritten Dimension‘ zählt außerdem *Handlungskompetenz (pragmatische Kompetenz)*, die sich zum Beispiel daran zeigt, wie man Gespräche organisiert, ein Wort ergreifen, zurückfragen etc. Wichtig sind weiterhin *strategische Kompetenzen*, zu denen auch Rezeptionsstrategien für das Lesen und Hören von Texten in einer Fremdsprache gehören.

Strategische Kompetenzen umfassen u. a. die Fähigkeit zum Umschreiben oder das Erbitten von Erklärung. Im GeR werden die strategischen Kompetenzen der Niveaus von A2 bis B2 folgendermaßen beschrieben:

⁵ Auf der Homepage <http://www.englishprofile.org/> (letzter Zugriff: 01.12.2017) werden Englischlehrkräften zahlreiche Hilfen und Materialien angeboten. Ein Clip auf der Internetseite sowie über <https://youtu.be/LI-HREEDa70> (letzter Zugriff: 01.12.2017) gibt einen guten Überblick über das Angebot.

B2	Kann eine Vielfalt von <i>Strategien</i> einsetzen, um das Verstehen zu sichern; dazu gehört, dass er/sie beim Zuhören auf Kernpunkte achtet sowie das Textverständnis anhand von Hinweisen aus dem Kontext überprüft.
B1	Kann in Texten mit Themen aus dem eigenen Fach- oder Interessengebiet bekannte Wörter aus dem Kontext erschließen. Kann die Bedeutung einzelner unbekannter Wörter aus dem Kontext erschließen und Satzbedeutung ableiten, sofern das behandelte Thema vertraut ist.
A2	Kann sich eine Vorstellung von der Gesamtaussage kurzer Texte und Äußerungen zu konkreten, alltäglichen Themen machen und die wahrscheinliche Bedeutung unbekannter Wörter aus dem Kontext erschließen.

(Europarat 2001: 78, Auszug A2–B2)

Solche Formulierungen unterstreichen noch einmal, dass der GeR ein *mehrdimensionales Modell von Kommunikation* ist, das Kompetenzen beschreibt und stuft, indem es eine erste und zweite Dimensi-

on, aber auch sprach- und soziolinguistische Kompetenzen miteinbezieht, die bei der Kommunikation eingesetzt werden.

3 Wozu man die tool box GeR benutzen kann

Nicht nur in Deutschland verbreitete sich der GeR wie ein Lauffeuer: Lehrbuchverlage erkannten blitzschnell den Wert für ihr Marketing, Volkshochschulen und auch das allgemeinbildende Schulwesen griffen vor allem die sechs Referenzniveaus und ihre inhaltlichen Füllungen, die Deskriptoren der Sprachkompetenzen, begierig auf. Ein Grund für die Attraktivität des GeR beim Massenmarkt war, dass mit der kommunikativen Wende des Fremdsprachenunterrichts seit den 1970er Jahren viele kompetenzorientierte Arbeitsformen im Unterricht gehalten haben. Klassenarbeiten und Abschlussprüfungen haben aber diese Entwicklung hinterher hinkicken lassen.

Defizite des deutschen Fremdsprachenunterrichts wurden durch Tests wie etwa dem *TEFL*⁶ aufgedeckt, und Studentinnen und Studenten mussten sich fragen, was das Abitur wirklich über ihre Sprachkenntnisse aussagte. Auch die Prüfungsämter für die Mittlere Reife verlangten immer stärker nach einer Objektivierung und Vergleichbarkeit ihrer Ergebnisse. Vor allem die *PISA*-Studien, in denen die Bundesrepu-

blik der Lesefähigkeit stets nur mittlere Plätze erreichte, schreckten die Bildungspolitik auf. Der GeR war daher ein willkommenes Mittel, um auch den Fremdsprachenunterricht in übersichtliche Systeme mit transparenten Abschlussprüfungen zu bringen (vgl. auch Kapitel 14).

Die sechs Referenzniveaus des GeR wurden von der vollständigen Konferenz der Kultusminister (KMK) als Basis für die Gestaltung der neuen Bildungsstandards benutzt. Der GeR ist also so etwas wie die Mutter aller Kompetenzskalen. Die folgenden Stufungen bzw. Vorstellungen davon, welche Kompetenzniveaus Lernende in der Regel bis zu einem bestimmten Zeitpunkt erreichen sollten, kristallisierten sich in den Köpfen der meisten Lehrerinnen und Lehrer – insbesondere aber bei den Schulverwaltungen – heraus:

Grundschule	A1 (bei 6 Jahren A2)
Hauptschule, 9. Klasse	A2
Mittlerer Schulabschluss	B1
Abitur	B2 bzw. C1

⁶ *Test of English as a Foreign Language*, vom Educational Testing Service in Princeton, USA, seit 1964 angeboten.

Kapitel 4

Hör- und Hör-/Sehverstehen

I Geeignete Hör- und Lesetexte finden

Dass Validität ein zentrales Merkmal von Assessment ist, wurde bereits in Kapitel 2 thematisiert. In diesem Zusammenhang wurden Hilfsmittel wie der GeR betrachtet, die vor allen Dingen zur Beschreibung inhaltlicher Validität nützlich sind. In Kapitel 3 stand dann die Konstruktvalidität im Zentrum. Zu deren Erläuterung wurden andere Skalen und Deskriptoren des GeR hinzugezogen, in denen ebenfalls Prozesse und Strategien beim Hör- und Leseverstehen beschrieben sind. In diesem Kapitel soll nun beides miteinander verknüpft werden. Dabei wird ein weiterer Begriff zur Beschreibung von Validität eine Rolle spielen: der der „Authentizität“.

Bei der Suche nach geeigneten Hör- und Lesetexten sind Bildungsstandards oder deren Anforderungen nicht sehr hilfreich, da sie im Unterschied zu

operationalen Curricula nur Listen mit Textsorten enthalten. Die Deskriptoren des GeR sind zudem nicht sehr systematisch formuliert. Sie erwähnen zum einen Textsorten und Situationen, die man verstehen soll, zum anderen die Art der Sprache, die dabei zum Nutzen und/oder die Bedingungen und Beschränkungen, unter denen das geschieht, die Relevanz dieser Merkmale variieren allerdings von Stufe zu Stufe.

Die folgenden Ausschnitte aus der Skala „Hör- und Leseverstehen Allgemein“ liefern einen Eindruck davon. Zunächst soll der Fokus auf den Themen des Hör- und Leseverstehens liegen, danach geht es um die Anforderungen, die beim Verstehen der Sprache gestellt werden. In beiden Fällen wurden die Deskriptoren durch Kursivsetzung hervorgehoben:

	Hörverstehen Allgemein – Themen	Hörverstehen Allgemein – Art der Sprache
CI	Kann genug verstehen, um längere Redebeiträgen über <i>nicht vertraute abstrakte und komplexe Themen</i> zu folgen, wenn auch gelegentlich Details bestätigt werden müssen, insbesondere <i>bei fremdem Akzent</i> . [...]	Kann genug verstehen, um längeren Redebeiträgen über nicht vertraute abstrakte und komplexe Themen zu folgen [...], <i>insbesondere bei fremdem Akzent</i> . Kann ein <i>breites Spektrum idiomatischer Wendungen und umgangssprachlicher Ausdrucksformen</i> verstehen und Registerwechsel richtig beurteilen. [...]
BI	Kann unkomplizierte <i>Sachinformationen über gewöhnliche alltags- oder berufsbezogene Themen</i> verstehen und dabei die <i>Hauptaussagen und Einzelinformationen</i> erkennen, sofern <i>klar artikuliert und mit vertrautem Akzent gesprochen</i> wird.	Kann unkomplizierte Sachinformationen über gewöhnliche alltags- oder berufsbezogene Themen verstehen und dabei die <i>Hauptaussagen und Einzelinformationen</i> erkennen, sofern <i>klar artikuliert und mit vertrautem Akzent gesprochen</i> wird.
AI	Kann [was gesagt wird] <i>die Verf.</i> verstehen, wenn <i>sehr langsam und sorgfältig gesprochen wird und wenn lange Pausen Zeit lassen, den Sinn zu erfassen</i> .	Kann verstehen, wenn <i>sehr langsam und sorgfältig gesprochen wird und wenn lange Pausen Zeit lassen, den Sinn zu erfassen</i> .

Ohne erkennbare Systematik wird mal dieses und mal jenes Merkmal betont. Das macht die Entscheidung nicht einfacher, ob ein gewählter Text für ein Assessment geeignet ist. Ähnliches lässt sich auch

bei den Bildungsstandards beobachten. Zur Diskussion der Konstruktvalidität soll daher das Konzept der Authentizität herangezogen werden.

2 Authentizität

Obgleich die meisten Menschen über die Aussprache des Begriffs stolpern, ist er in der Fremdsprachendidaktik unverzichtbar. Wörterbücher paraphrasieren den Begriff *authentic* oft mit *real*, *genuine* oder auch *known to be true*. Didaktisch greifen diese Umschreibungen jedoch zu kurz. Daher wird in der Fachdidaktik regelmäßig gefordert, dass die Sprache, die in Hörtexten Verwendung findet, der von realen Sprechsituationen oder Texten entsprechen soll. Damit wird allerdings ausgeblendet, dass auch Originaltexte, sobald sie für Lehrwerke oder Tests ausgewählt (und eventuell auch gekürzt oder bearbeitet) sind, nicht mehr völlig authentisch sind.

Aus der Forderung nach Authentizität ergeben sich eine Reihe von Qualitätskriterien: Erstens wäre wert wäre beispielsweise, dass Hör- oder Lesetexte nicht nur von Lehrbuchautoren, sondern auch von Testkonstrukteuren für didaktische Zwecke ausgewählt werden. Sie sollten vielmehr aus den üblichen Medien stammen, die für Muttersprachler*innen erstellt und von Muttersprachlern gelesen oder gesprochen werden. Finden lassen sich authentische Hör-/Sehtexte problemlos bei Funk und Fernsehen, in Bibliotheken, Mediatheken, bei Streaming-Diensten usw.

Texte, die für Hör- oder Sehtexte geeignet sind, haben zudem nicht nur die Eigenschaft authentisch zu sein, sie eröffnen auch neue Möglichkeiten des Umgangs mit ihnen: Einen Film sieht man, liest man ihn. Später diskutiert man eventuell mit Freunden darüber. Im Unterricht hingegen beantwortet man Lehrerfragen dazu oder macht Inhaltsangaben, bevor man mit anderen Lernenden über den Sinn diskutiert (natürlich unter Anleitung der Lehrkraft, die das steuert). Letzteres

wäre aber kein authentischer Umgang mit einem Text mehr, sondern ein „didaktischer“.

Nun lässt sich argumentieren, dass alle Unterrichtssituationen künstlich sind und „willing suspension of disbelief“ voraussetzen (der Begriff wurde schon 1817 von Samuel Taylor Coleridge geprägt). Hinzu kommt, dass authentische Hör- oder Lesetexte im Unterricht oft auf Ablehnung stoßen, was zum Teil auf schlechtes Englisch zu Folge zu schnell gesprochen wird, Hintergrundgeräusche und Dialekte das Verständnis erschweren usw. Didaktisch vereinfachte Tonaufnahmen der Schulbuchverlage dominieren

folglich den Unterrichtsalltag und führen dazu, dass viele Jugendliche erst in der Oberstufe authentisches Englisch zu hören bekommen. Dank des eigenen Medienkonsums haben sie bis dahin aber schon schon manche Varietät des Englischen zu hören bekommen und bestenfalls auch verstehen gelernt. Gesprochene Sprache zu verstehen ist ein wichtiges Ziel im kompetenzorientierten Unterricht. Authentische Sprache klingt aber ganz anders als die didaktisierte Sprache in bereinigten Tonaufnahmen, bei denen es durchaus passieren kann, dass der Regisseur einen Sprecher zur Wiederholung einer Passage auffordert: „*John, can you read that again? There was a stomach rumble.*“

Neben britischem und amerikanischem Englisch gibt es Dialekte und andere Varietäten, Auslassungen und Verschleifungen. Denken Sie beispielsweise an *fish'n'chips*, *bread'n'butter*, *lots of* (gesprochen /'lɒtsəv/) oder den „Peng!cake“ (= *pancake*).



© Pixabay/PublicDomainPictures

Bei mündlichen Äußerungen gibt es Brüche in der Syntax und Füllsel zur Überbrückung von Pausen (*erm, well*), in Dialogen gibt es Überlappungen und vieles mehr. Einmal ganz abgesehen davon, dass Menschen nicht im sorgfältig abgeschirmten Schallraum eines Studios sprechen, sondern ihre Stimme gegen Hintergrundgeräusche anheben, die Intonation auch emotional steuern etc. Wenn wir Lernende auf der

Dschungel gesprochener Sprache vorbereiten wollen, ist der Einsatz authentischer Hör- und Sehverstehenstexte im Unterricht und beim Assessment unumgänglich. Nur so können die Kompetenzen der Lernenden entwickelt und überprüft werden.

Leider ist Authentizität immer noch kein Qualitätskriterium im vollen Umfang akzeptiert wird, weder bei den Schülern noch bei Testanbietern. Man muss sich genau anhören und prüfen, wie stark sich Höraufnahmen zumindest um „Semi-Authentizität“ bemühen, und das auch schon auf niedrigeren Niveaustufen. Da gibt es durchaus graduelle Unterschiede. Im Grunde kann man aber schon zwischen zwei Typen unterscheiden, wenn man neben den Höraufnahmen zu den Werkzeugen für Assessment überhaupt Hörtexten findet, die in „natürlichem“ Sprechtempo und -melodie sind, mit altersgemäßen Stimmen und zumindest einige der oben genannten Merkmale authentischer Sprache widerspiegeln.

3 Hörtexte finden und bearbeiten

Für die Arbeit mit Hörtexten ist es zentral, eine Vorstellung davon zu haben, welche Arten und Sorten für die jeweilige Zielgruppe relevant sind und in Zukunft sein werden. Eine Steuerung des Schwierigkeitsgrads könnte dann mithilfe der Lernenden oder der Bildungsstandards erfolgen. Eine Schwierigkeit bleibt aber dennoch bestehen: Höraufnahmen für Lehrwerke und Tests sind in der Regel nicht authentisch. Wenn sie es in gewisser Weise sind, entsprechen sie meist nicht dem Niveau eines Kompetenzmodells wie dem CEFR/GeR.

Semi-authentische Hörmaterialien, das die Lernenden noch nicht kennen und das den Vorzug hat, auf bestimmte Aufgaben abzielen zu sein, findet man bei den großen Testanbietern:

- (1) Für einen Test auf dem Niveau A2, dem *Key English Test*¹ von Cambridge English Language Assessment. Auf dieser Website kann man sich auch zu allen anderen Niveaustufen durchklicken und sich eine Vielfalt von Aufgaben herunterladen, vor allem auch Hördateien im MP3-Format.
- (2) TELC = The European Language Certificates².
- (3) *DIALANG* (eine vom Europarat angebotene online Ressource)³.
- (4) *TOEFL/TOEIC*⁴.

Diese leicht zu erschließenden Hörmaterialien, die nicht nur auf *CEFR/GeR*-Niveaustufen angeboten werden, sondern dazu noch mit Aufgaben versehen

¹ <http://www.cambridge.org/ce/cambridgeenglish/catalog/cambridge-english-exams-ielts/key-schools/resources> (letzter Zugriff: 01.12.2017).

² <https://www.telc.net/pruefungsteilnehmende/sprachpruefungen.html> (letzter Zugriff: 01.12.2017).

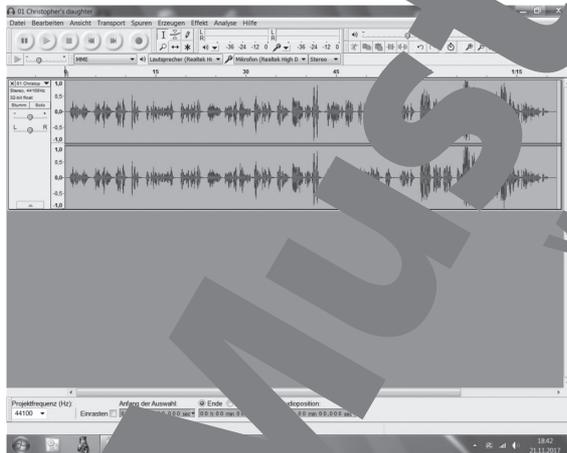
³ <http://www.lancaster.ac.uk/researchenterprise/dialang/about> (letzter Zugriff: 01.12.2017): Die Nutzung ist etwas kompliziert, unter der angegebenen Adresse gibt es aber eine ausführliche Anleitung.

⁴ http://www.ets.org/toefl/ibt/prepare/quick_prep/ (letzter Zugriff: 01.12.2017).

sind bilden eine ideale Quelle für eigenes Assessment im Bereich Hörverstehen.

Wenn man authentische Materialien für mittlere und obere Lernstufen sucht, sollte man die vielen Angebote von britischen oder US-amerikanischen Radio- und Fernsehsendern nutzen, bei denen man allerdings mit zusätzlichen Anforderungen konfrontiert wird. Es gilt nicht nur Hörtexte zu finden, sondern diese danach auch noch zu bearbeiten. Für beide Herausforderungen folgen nun einige Quellen und Hinweise für die Aufbereitung.

Es gibt eine außerordentlich nützliche Website der BBC mit einer großen Auswahl an Hörtexten,⁵ die sich für den Einsatz im Englischunterricht anbieten und besondere für das GeR-Niveau B1 und höher eignen. Auch bei YouTube gibt es Podcasts.⁶ Für ihre Nutzung sollte man sich allerdings mit einer Bearbeitungs-Software vertraut machen, da diese Texte von der Länge her nicht unbedingt für den Unterricht geeignet sind. Mit *Audacity*, einem kostenlosen Audio-Editor, der mit den bekanntesten Audio-Formaten zurechtkommt (MP3, WAV u.a.), geht das ganz einfach.



Screenshot einer in Audacity importierten Hördatei.

Sobald das Programm auf dem Rechner installiert wurde, ist es möglich, Audiodateien zu importieren und zu bearbeiten. So lassen sich beispielsweise einzelne Teile markieren und schneiden, die Klangeigenschaften können verbessert, Höhen und Tiefen beeinflusst oder aber entfernt werden. Die fertig bearbeitete Audiodatei kann schließlich auf einer CD, einem USB-Stick oder in der Cloud gespeichert und dann im Unterricht abgespielt werden.

Nützlich ist die Aufbereitung insbesondere bei vielen YouTube-Videos oder Mitschnitten von Podcasts. Bei YouTube-Videos benutzt man zunächst zum Herunterladen *Free YouTube Downloader* oder vergleichbare Freeware-Programme. Danach kann man in der Regel einen *MP4 to MP3 Converter*, der auch kostenlos im Internet erhältlich ist. Dieser ist notwendig, um die im MP4-Format gehaltenen Videos in Hördateien zu konvertieren. Wenn man also z. B. einen interessanten Vortrag findet, bei dem die visuelle Komponente nicht so wichtig ist, kann man das MP4-Video in eine MP3-Hördatei umwandeln und dann mit *Audacity* weiter bearbeiten (siehe Screenshot). Wenn das übersichtlich gespeichert und allen im Kollegium zugänglich gemacht wird, beteiligen sich wahrscheinlich auch immer mehr Personen am Suchen und Bearbeiten solcher Hörtexte. Auf diese Weise gelangt man schnell zu einem Pool an Höraufnahmen, den man sowohl für unterrichtliche Zwecke als auch bei der Erstellung von Klassenarbeiten nutzen kann.

⁵ www.bbc.co.uk/podcasts (letzter Zugriff: 01.12.2017).

⁶ <https://www.podcastchart.com/categories/top-200-podcasts?page=1#position-4> (letzter Zugriff: 01.12.2017) oder auch www.itunescharts.net/charts/podcasts/ (letzter Zugriff: 01.12.2017).

Kapitel 9

Schreiben beurteilen

1 Geschlossene Aufgabenformate

Bei der Beurteilung von schriftlichen Leistungen ist zunächst relevant, ob es sich um geschlossene oder offene Aufgabenformate handelt (vgl. Kapitel 3). Die Wahl der Aufgabenformate hängt wiederum von unterschiedlichen Faktoren ab. So spielt das Kompetenzniveau der Lernenden eine Rolle: Wieviel Sprache können sie zusammenhängend produzieren? Auch die Vorgaben von Lehrplänen müssen bedacht werden: Welche Kompetenzen sind besonders wichtig (z. B. schriftliche Produktion oder

Interaktion)? Schließlich müssen die Anforderungen an die Lernenden berücksichtigt werden: Welche Textsorten erwarten den Lernenden im späteren Leben? Welche Kompetenzen müssen die Lernenden in beruflichen Situationen bewältigen? Dies ist eine Frage der Zuverlässigkeit von Beurteilungsinstrumenten für welche Aufgabenformate man sich entscheidet. Geschlossene Aufgabenformate (*True/False, Multiple Choice* etc.) erlauben größere Zuverlässigkeit, also *Reliabilität* (vgl. Kapitel 13).

2 Bewerten von schriftlichen Leistungen mit offenen Aufgabenformaten

Bei der Bewertung offener Aufgabenformate steht eher die Validität im Vordergrund. Zentral ist die Zielgerichtetheit einer Aufgabe und die Leitfrage: Misst die Aufgabe die Kompetenz, die geprüft werden soll? Dadurch, dass Bewerter eine Entscheidung über die Qualität der sprachgelösten Aufgabe treffen müssen, sinkt die Reliabilität möglicherweise. Schließlich können Entscheidungen immer ein subjektives Element, sodass bei der Bewertung offener Aufgabenformate keine vollkommene Vergleichbarkeit gegeben ist. Man kann aber sehr wohl dafür sorgen, dass die Beurteilungsprozesse transparent sind.

Die Formen der Bewertung bei offenen Aufgaben sind denen aus Kapitel 7 vorgestellt wurden. Ebenso wie beim Sprechen können Checklisten oder Beurteilungsraster verwendet werden, um geschriebene Lernertexte zu beurteilen. Dabei eignen sich insbesondere für Lernende mit geringen Vorkenntnissen Checklisten: Sie können mit wenig

Aufwand auf eine bestimmte Aufgabenform abgestimmt werden und eignen sich auch zum Einsatz bei *self* und *peer assessment*. Die Checkliste (vgl. S. 102) bezieht sich auf eine per E-Mail ausgesprochene Einladung an einen Freund. Die Checkliste wird zuerst als Instrument für *self assessment* verwendet und dann für die unterrichtsbasierte Leistungsbeurteilung seitens der Lehrkraft.

Häufig kommen bei der Bewertung von produktiven und interaktiven schriftlichen Leistungen in der Fremdsprache auch Beurteilungsskalen zum Einsatz. In standardisierten Tests sind sie schon seit langem Usus, wenn es um die Bewertung schriftlicher Leistungen geht. Bei landesweiten Abschlusstests setzen sie sich ebenfalls durch. In der Funktion eines Erwartungshorizontes bzw. als Bewertungsraster sind sie darüber hinaus in manchen Bundesländern bei Klassenarbeiten Pflicht. Im Grunde unterscheiden sich die Beurteilungsskalen für Sprechen und

Important points	You		Your teacher	
	Yes ☺	No ☹	Yes ☺	No ☹
I have put in a greeting (Dear ...).				
I have asked how my friend is (How are you?).				
I have invited my friend to come to my hometown.				
I have written about interesting places or things to do in my hometown.				
I have asked him/her what s/he likes to do best.				
I have ended the e-mail with a greeting (Love, All the best etc.).				
I have written in a friendly tone.				
	You:		Your teacher:	
Comments/improvements:				

Writing an e-mail to a friend-Checkliste

Schreiben nur in einem Punkt, nämlich in ihrem Umfang. Beurteilungsskalen für das Schreiben können etwas detaillierter sein, weil der geschriebene Text im Gegensatz zum Sprechen vorliegt und von Beurteiler mehrmals gelesen werden kann.

Für die Verwendung von Skalen ist die Wahl zwischen *holistischen* Skalen (die eine Leistung ganzheitlich bewerten und keine Kriterien explizit ausweisen, aber dennoch welche beinhalten) und *analytischen* Skalen (die die Bewertungskriterien einzeln auswerten und mehrere Teilergebnisse in diesem dann zu einem Endergebnis zusammenfügen (vgl. auch Kapitel 7)). Es lässt sich eine fertige Skala für die Bewertung übernehmen oder anpassen müssen. Alternativ kann aber auch eine eigene Skala mit den Bewertungskriterien, die als wichtig erachtet werden und mit den jeweiligen Lernenden übereinstimmen, erstellt werden.

Das folgende Beispiel (vgl. S. 103) stammt von der *Association of Language Testers in Europe (ALTE)* (2007) und ist eine *holistische* Skala, die sich völlig aufgabenunabhängig verwenden lässt.

Der Vorteil dieser Bewertungsanleitung liegt darin, dass man sie für alle Aufgabentypen – auch über das Schreiben hinaus – verwenden kann. Darüber hinaus lässt sich neben sprachlichen Aspekten auch die kommunikative Wirkung der Äußerung bewerten, sodass die Lernenden angemessen beurteilt werden können. Deren Sprache weist zwar die eine oder andere Unebenheiten auf, sie trauen sich aber beispielsweise, komplexere Sprache zu verwenden, um eine Aufgabe auf kreative Weise zu lösen. Damit wird die Bewertungsanleitung interessant für aufgabenbasierte Beurteilung.

5	<p>Aufgabe voll erfüllt.</p> <ul style="list-style-type: none"> • Alle Inhaltspunkte werden in angemessenem Umfang behandelt. • Große Breite der Strukturen und des Wortschatzes gemäß der Aufgabe. • Minimale Fehler, vielleicht weil zu ehrgeizig; gut entwickelte Beherrschung der Sprache. • Effektiver Aufbau der Gedanken mit einer Vielfalt verbindender Elemente. • Register und Format durchgehend der Schreibabsicht und dem Leser angemessen. <p>Erzielt voll die gewünschte Wirkung auf den Leser.</p>
4	<p>Gute Erfüllung der Aufgabe.</p> <p>[...]</p>
3	<p>Ausreichende Erfüllung der Aufgabe.</p> <ul style="list-style-type: none"> • Alle hauptsächlichen Inhaltspunkte werden behandelt; es gibt keine Auslassungen. • Angemessene Breite der Strukturen und des Wortschatzes, die die Anforderungen der Aufgabe erfüllen. • Eine Reihe von Fehlern ist möglich, sie behindern nicht die Kommunikation. • Angemessener Aufbau der Gedanken mit einigen Verbindungen. • Ausreichender, wenn auch nicht immer erfolgreicher Versuch, ein der Schreibabsicht und dem Leser angemessenes Register und Format zu verwenden. <p>Erzielt im Großen und Ganzen die gewünschte Wirkung auf den Leser.</p>
2	<p>Versuch, die Aufgabe zu erfüllen, der nicht mit angemessenem gelingt.</p> <p>[...]</p>
1	<p>Schwacher Versuch, die Aufgabe zu erfüllen.</p> <ul style="list-style-type: none"> • Erhebliche inhaltliche Auslassungen und/oder häufig irrelevante Inhalte, möglicherweise wegen Nicht-Verstehen der Aufgabe oder Anweisungen. • Geringe Breite des Wortschatzes und der Strukturen. • Häufige Fehler, die die Kommunikation behindern; kaum entwickelte Beherrschung der Sprache. • Kein Aufbau der Gedanken mit verbindenden Elementen. • Kein oder nur geringes Bewusstsein von Register und Format. <p>Die Äußerung hat eine starke negative Wirkung auf den Leser.</p>

Holistische Skala zum Schreiben (Lorenz 1977: 12)

Für geschriebene Texte speziell auf mittlerem Niveau, beispielsweise für die 8. bis 10. Klasse, hat Charlotte Stolle in ihrer Arbeit für das hessische Staatsexamen

(2009) ein interessantes Beurteilungsraster entworfen, das in vereinfachter Form aufgegriffen werden soll.

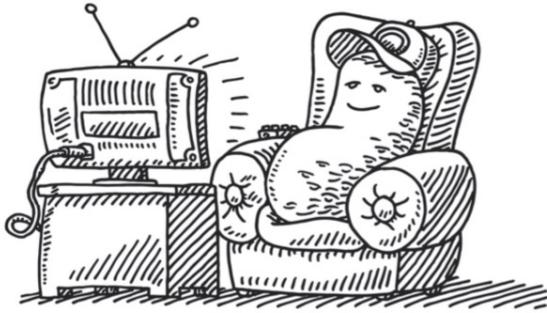
Kriterium	Unterteilung	Punkte	Gesamt
1. Gesamteindruck			
	Kohärenz: Form, Inhalt, Eigenständigkeit etc.	3/25	
2. Inhalt			
	• Quantität (evtl. laut Vorgabe)	3	7/25
	• Qualität	3	
	• sachliche Richtigkeit	1	
3. Adressatenbezug			
	• formale Aspekte (korrekte Anrede, Höflichkeit oder Vertrautheit etc.) • inhaltliche Aspekte (Vorwissen etc.)	4/25	
4. Sprache			
	• Korrektheit	3	7/25
	• Komplexität	3	
	• Angemessenheit Textsortenbezug	1	
5. Gliederung			
	• Einteilung in Absätze • Stimmigkeit	4/25	

Beurteilungsraster für die Mittelstufe (nach Stolle 2009).

Stolle integriert in ihrem Raster kommunikative und linguistische Aspekte. Beispielsweise sinnvolle Satzanfänge, englische Idiome, Satzdruckweise, die Verwendung von Konnektoren wie „und“ bzw. „aber“ und eine angemessene themenangemessene Wortwahl. Das Raster hat den Vorteil, dass es gleich ein Kriterium vornimmt, die man natürlich nach eigenen Schwerpunkten ändern oder erweitern kann.

Das folgende Beispiel stammt aus einer Real-schulabschlussprüfung:

Ausgehend von einem visuellen Impuls und drei Leitfragen (*What are the reasons for being a couch potato? Are you a couch potato or do you know a person who is a couch potato? Do you enjoy being a couch potato in your free time?*) sollen die Schüler etwa 150 Wörter zum Thema produzieren, sei es über sich selbst, eine/n Freund/in oder allgemein. Eine Textsorte ist nicht vorgegeben.



© iStock/Frank Ramspott

The Couch Potato Family (unkorrigiert)

My family is a idol for all the others couch potatoes in the world. I don't know why it's so. I expect that my father works the whole day and in the evening he is very tired. My mother doesn't want that the family go outside. She is a real couch potato because she doesn't work and she does nothing the whole day, apart from the housework. I am the only one in the family who do a sport. At the weekends my parents sometimes go to Frankfurt or in Rheingau. Then I meet my friends or my boyfriend. We go shopping in the cinema.

I do much sport with my friends. One example is video-clip dancing or athletics. I go back to my family. Perhaps I can understand the business by my father, because he works very hard. But I don't understand the laziness of my mother. She sits the whole day at home and she tells to her that we can go shopping. She doesn't want. That's my lazy family. (173 words)

Geht man sich in diesem Beispiels das Raster von Stolle (2007) an, sieht man, dass nicht alle Kri-

terien für diese etwas unspezifische Textsorte gleichermaßen hilfreich sind. Die Kategorie „Gesamteindruck“ ist unproblematisch, weil sie recht allgemein ist und textsortenspezifische Schwerpunktsetzungen zulässt. Der Umgang mit den Komponenten „Quantität“ und „Formale Richtigkeit“ ist ebenfalls handhabbar. Lediglich die Kategorie „Qualität“ würde vermehrt eine Diskussion unter den Bewertern im Lehrerbildungsjahr ergeben, etwa um sich darüber verständigen, ob die Leitfragen im Schülertext angedacht werden müssen. Das ist in diesem Beispiel nicht der Fall, obwohl der Text inhaltlich trotzdem kohärent ist.

Problematisch ist in diesem Fall das Kriterium „Angemessenheit“. Formale Aspekte wie die angelegene Rede oder ein passendes Register fallen weniger ins Gewicht. Ähnlich verhält es sich mit der inhaltlichen Adressatengemessenheit. Da der Leser für diesen Text kein Vorwissen benötigt und die Textsorte nicht spezifiziert war, lässt sich über dieses Kriterium wenig sagen.

Aufgrund der hohen Punktzahl fällt die Punkteverteilung schwer. Beim Kriterium Sprache stehen die Korrektheit und die Komplexität (im GeR in etwa unter „Spektrum der sprachlichen Mittel“ gefasst) sowie die Angemessenheit im Vordergrund. Folglich sind bereits viele Aspekte einer sprachlichen Leistung abgedeckt. Für die Beurteilung etwas umfangreicherer Texte in geschriebener Form ist zudem die Kategorie „Gliederung“ hilfreich, da sie den Bewertern helfen, die stimmige Struktur und den Aufbau des Texts sowie dessen Kohäsion und Kohärenz einzuschätzen.

Ohne die entsprechenden curricularen Vorgaben zu kennen, die dieser Leistung zugrunde liegen, ist es sicherlich schwierig, individuelle Leistungen zu bewerten. Ein Vorschlag zur Einschätzung der Schülerarbeit wäre:

Kriterium	Unterteilung	Gewichtung		Bewertung
1. Gesamteindruck				
	Kohärenz: Form, Inhalt, Eigenständigkeit etc.	3	5	3
Kommentar	Recht eigenständige Leistung mit geschickt gewähltem Anfang und Ende.			
2. Inhalt				
	<ul style="list-style-type: none"> Quantität (evtl. laut Vorgabe) Erfüllt die Zielvorgabe von 1500 Wörtern und ist umfangreich genug, um eine Gesamtaussage zu entwickeln.		7/25	3
	<ul style="list-style-type: none"> Qualität Sie beantwortet die Leitfragen zwar nicht vollständig, aber bringt eine pointierte, leicht ironische, persönliche Darstellung zur Meinung der Verfasserin inklusive eigenständig formulierter Vorwürfe an die Adresse der Verfasserin.			2
	<ul style="list-style-type: none"> sachliche Richtigkeit Bei der gewählten Textsorte (persönliche Meinungsäußerung) kann man eigentlich nichts falsch machen.	1		1
3. Adressatenbezug				
	<ul style="list-style-type: none"> formale Aspekte (korrekte Anrede, Höflichkeit, oder Vertrautheit etc.) inhaltliche Aspekte (Vorwissen etc.) 	3	3/25	3
Kommentar	Diese Kategorie ist schwierig, weil die Verfasserin der Adressatenbezug in Ordnung, weil die Verfasserin eine neutrale Adressierung wählt. Der Text erfordert der Text kein Vorwissen.			
4. Sprache				
	<ul style="list-style-type: none"> Korrektheit Unebenheiten in der Wortstellung (Einfluss des Deutschen auf den Satzbau) und einige kleinere grammatrische Fehler, die das Verständnis nicht beeinträchtigen:	3	7/25	2
	<ul style="list-style-type: none"> Komplexität Die Verfasserin verwendet einfache Strukturen und wenig komplexe Sätze, die sich auch wiederholt:	3		1
	<ul style="list-style-type: none"> Angemessenheit Textsortenbezug Die gewählten Worte der personalisierten Meinungsäußerung ist angemessen und schlägt sich in einem persönlichen, leicht gefärbten Stil nieder:	1		1
5. Gliederung				
	<ul style="list-style-type: none"> Abgrenzung von Absätzen Stimmigkeit 	4/4	4/25	4
Kommentar	Der Text ist trotz Redundanzen kohärent und hat einen geschickt gewählten Anfang und Ende.			

Bewertungsschlag: 20/25 Punkten.

6 Unverzichtbar: das Beurteilertraining

Dass selbst die Benutzung von Beurteilungsbögen und Bewertungsrastern nicht immer eine hinreichende Objektivität garantiert, wurde bereits thematisiert. Deshalb empfiehlt es sich, auf Fachkonferenzen und mit Hilfe von gemeinsamen Korrekturen nach einem vorher diskutierten Beurteilungsraster ein größtmögliches Einvernehmen über Schülertexte zu erzielen. Abweichende Einschätzungen müssen dabei diskutiert werden, aber nicht mit dem Ziel, eine Kollegin oder einen Kollegen unbedingt zu

überzeugen, sondern besser zu verstehen, wo die Gründe für die Differenzen liegen: Vielleicht wurde das Korrekturraster von beiden Kollegen anders ausgelegt oder missverstanden. Der Nutzen von Beurteilertraining für die *inter-rater reliability* (Übereinstimmung zwischen verschiedenen Beurteilern) ist in zahlreichen Untersuchungen nachgewiesen worden. So lästig es sein mag, ein Beurteilertraining ist unverzichtbar, wenn man in einem gewissen Maße verantwortungsvoll mit Beurteilungsergebnissen umgehen will.

7 Geschlossene Aufgabenformate und das Kriterium der Reliabilität

Wenn man eine größtmögliche Zuverlässigkeit (Reliabilität) bei Bewertungen erreichen möchte, bieten sich die im ersten Kapitel erwähnten geschlossenen Aufgabenformate (in der Regel *multiple choice*) an. Eine nicht zu unterschätzende Schwierigkeit besteht dabei oft schon bei der Konstruktion solcher Aufgaben. Eingangs wurde bereits auf das Problem der Ratewahrscheinlichkeit bei solchen Mehrwahl-Aufgaben hingewiesen. Diese Wahrscheinlichkeit sollte man nicht auf die leichte Schulter nehmen.

Die Konstruktion von 4 Aufgaben zu Strukturen und Wortschatz ist es oft einfacher, für eine einfache binäre Entscheidung („X“ oder „Y“) eine richtige Lösung und vier Distraktoren zu beschreiben, die alle gleich plausibel sind und von denen nicht ein oder zwei von den Lernenden als offenkundig falsch durchschaut werden. Das würde nämlich das Ergebnis eines Assessments verfälschen.

In professionellen Testsystemen gibt es daher – statt nach der Erprobung oder dem eigentlichen

Aufgabenanalyse (nach Elisabeth Ingram 1968)¹

Anzahl der Schüler: 25 (: 3 =) 8 : 3 = 8

Item	oberes Drittel	mittleres Drittel	unteres Drittel	richtige Lösungen	% (1)	Differenz (2)	$\Sigma 1 - 3$ (3)
1	//// /	//// /	//// /	8	25	0	0
2	//// /	//// /	//// /	0	12	8	1
3	//// /	//// /	//// /	3	14	3	.38
4	//// /	//// /	//// /	6	22	2	.25
5	//// /	//// /	//// /	2	14	5	.63
6	//// /	//// /	//// /	7	11	- 4	- .50
7	//// /	//// /	//// /	6	24	2	.25

(1) Alle richtigen Lösungen geteilt durch Anzahl der Schüler $\times 100$ (in Item 3 also $14 : 25 \times 100 = 56\%$)

(2) Oberes minus unteres Drittel (in Item 3 also $1 : 6 - 3 = 3$)

(3) Differenz oberes/unteres Drittel geteilt durch die Anzahl der Schüler geteilt durch 3 (in Item 3 also $3 : 8 = .38$)

¹ Eine Vorlage der Aufgabenanalyse nach Ingram lässt sich auf der Homepage des Helbling Verlags zum Download finden.

Test – eine Aufgabenanalyse. Ein einfaches Verfahren schlägt Ingram vor. Nach einer Vorkorrektur aller Arbeiten (bei der alle Lösungen erfasst werden) ermittelt man für alle Schüler, hier 25, welche Arbeit wie gut ausgefallen ist. Daraufhin sortiert man die Hefte oder Lösungsbögen in absteigender Reihenfolge und unterteilt sie in ein oberes, mittleres bzw. unteres Drittel. Dann legt man für jedes Item eine Strichliste an und hält fest, wie oft Schüler des oberen Leistungsdrittels die Aufgabe richtig gelöst haben und wie oft das im mittleren bzw. unteren Drittel der Fall war. In der Tabelle (S. 134) ist solch eine Strichliste für 7 Aufgaben erfasst und ermöglicht durch die statistische Auswertung einige interessante Beobachtungen.

Interpretation der Werte

Richtige Lösungen/in %

Dieser Wert verrät den Schwierigkeitsgrad eines Items. Zu einfache (d.h. 90 – 100% richtige), und auf alle Fälle zu schwierige (= unter 20% richtige) Items sollte man aus der Bewertung herausnehmen, weil beide niedrige Trennschärfewerte haben und nicht viel über Leistungsunterschiede aussagen.

$\Sigma I - 3$ (= Trennschärfe): Idealwert hier ein Wert von + 1 (wie in Item 2), der zeigt,

Item genau zwischen den guten und den schwächeren Schülern unterscheidet. Inakzeptabel sind der Wert 0 (wie in Item 1, das von guten und von schwächeren Schülern gelöst wurde) oder negative Werte (wie in Item 4, das von den schwächeren Schülern besser gelöst wurde. Professionelle Testkonstrukteure vermeiden Items mit Werten unter .35 mit

Lehrkräfte sollten sich ein- oder zweimal im Schuljahr die Mühe machen ihre selbst erstellten Tests einer solchen Prüfung zu unterziehen, um ein Gefühl dafür zu entwickeln, welche Aufgaben zu leicht oder zu schwierig sind, und welche bei der Beurteilung von Schülerleistungen, sowohl für die Lehrkraft als auch für die Schüler, von angemessener Schwierigkeit sind. Dies gilt für Tests oder Lernerfolgskontrollen aus Lehrwerken, die auch nicht immer von Testexperten erstellt werden. Bei der professionellen Konstruktion von Testaufgaben sollte es selbstverständlich sein, dass diese erprobt werden und dass missglückte Items aus dem Test getilgt werden. Professionell erstellte Aufgaben-/Itemanalysen geben zusätzlich noch über weitere Dinge Auskunft.

Iteman

Item and Test Analysis Program – (AN) Version 3.50

Seq No.	Scale -Item	Prop. Correct	Discrim. Index	Point Bisser	Alternative	Prop. Total	Endorsing Low	High	Point Biser.	Key
8	2-1	.38	.59	.48	A	.00	.00	.00		
					B	.38	.13	.66	.48	*
					C	.12	.11	.12	-.01	
					D	.49	.74	.23	-.44	
					Other	.01	.00	.00	-.11	
9	2-2	.71	.42	.42	A	.07	.11	.01	-.16	
					B	.11	.18	.04	-.22	
					C	.10	.16	.00	-.22	
					D	.71	.53	.95	.42	*
					Other	.01	.00	.00	-.13	

Hier sind mithilfe eines Computerprogramms die Aufgaben/Items 8 und 9 ausgewertet worden. *Scale-Item* bedeutet die Nummer des Items in einem bestimmten Subtest (Prüfungsteil, wie z. B. Lese- oder Hörverstehen), und der Prozentsatz der korrekten Antworten (*Proportion Correct*) war bei 8 = 38% und bei 9 = 71%. Die Trennschärfe (*Discrimination Index*) lag für 8 bei .52, für 9 bei .42, war also in beiden Fällen zufriedenstellend.

Point Biserial ist eine andere Form der Berechnung der Trennschärfe, und auch hier sind bei bei-

den Items die Werte zufriedenstellend. Bei dieser gründlicheren Aufgabenanalyse wurde zusätzlich noch der Wert der vier Wahlmöglichkeiten A bis D berechnet (* = die richtige Antwort). Bei der Auswertung kann man so ohne Weiteres sehen, welche der Antwortmöglichkeiten am häufigsten oft gewählt wurde, was auf eine gute Qualität der Distraktoren schließen lässt. Bei Aufgabe 9 wurde zum Beispiel die Antwort A überhaupt nicht gewählt, sollte also bei einer erneuten Nutzung der Aufgabe neu formuliert werden.

8 Fazit

Während man früher vor allem in Prüfungen der renommierten Testanbieter allergrößten Wert auf Objektivität und Zuverlässigkeit (Reliabilität) von Tests legte, ist man heute eher bemüht, kompetenzorientiertes sprachliches Assessment sinnvoll (valid) und trotzdem auch verlässlich zu gestalten. Assessments in der Klasse sollte man sich ebenfalls um hohe Qualitätsstandards bemühen, wenn man aber auch einmal ein Auge zudrückt. Thomas von Aquin hat geschrieben, „Gerechtigkeit ohne Gnade Grausamkeit sein.“ Lehrenden daher nie vergessen, dass diejenigen Schüler, die in diesem Versetzungszeugnis noch mit einer Note 4 davongekommen sind, im nächsten Jahr die Fünferkandidaten sein werden. Und die Fünfer des laufenden Jahres die Klassenleser haben. Und da Lehrkräfte – im Gegensatz zu den großen kommerziellen Testanbietern – eingeschränkte Möglichkeiten zur Qualitätssicherung haben, empfiehlt es sich, ruhig auch ein wenig Barmherzigkeit im Sinne des Thomas von Aquin walten zu lassen.

McNamara, Tim (2000). *Language Testing*. Oxford University Press.
 McNamara, Tim (2007). *Materials for the Guidance of Test Item Writers*. *Handreichungen für Testautoren*. <http://www.alte.org/resources/filter> (letzter Zugriff: 12.02.2017).

McNamara, Tim (2013). *A Guide to Language Testing: Development, Evaluation, Research*. Boston, Mass.: Heinle & Heinle. Erste Auflage 1987.

McNamara, Tim (2017). *Einführung in statistische Verfahren bei Sprachtests*; die Neuauflage dokumentiert die Popularität und Aktualität dieses „Klassikers“.

McNamara, Tim (2000). *Language Testing*. Oxford University Press.

Erläutert neuere Verfahren der statistischen Prüfung bei Sprachtests wie *Item Response Theory*, das aktuell in der statistischen Auswertung von empirisch erhobenen Daten (nicht nur bei Tests!) favorisierte Verfahren. IRT ist keine einzelne Theorie, sondern eine Familie von formalen, mathematischen, probabilistischen Messmodellen, welche postulieren, dass dem beobachtbaren Testverhalten (manifeste Variable) eine bestimmte Fähigkeit / Eigenschaft bzw. Disposition (latente Variable) zugrunde liegt. Mehr dazu bei McNamara.

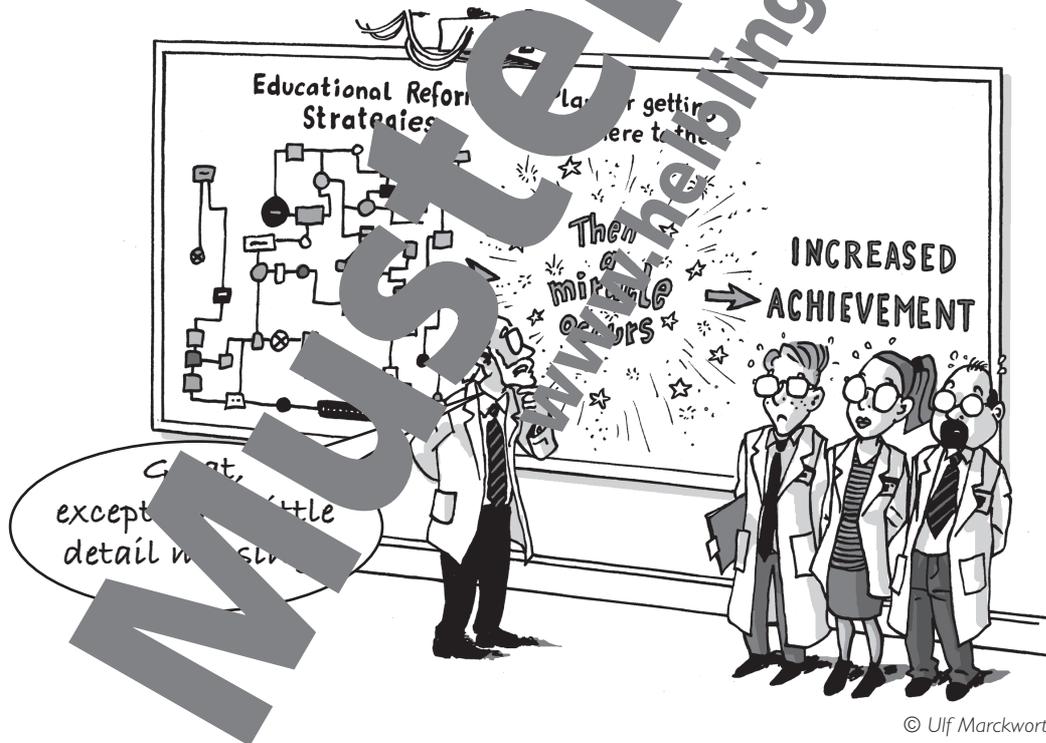
Kapitel 14

Vom Wiegen wird die Sau nicht fett

„Vom Wiegen wird die Sau nicht fett“ war ein beliebter Spruch von Lehrkräften, denen zentrale Lernstandserhebungen (von der Grundschule bis zum Abitur) und andere Eingriffe in ihren Unterricht den Nerv raubten. Das sei „Testeritis“! Statt den Erfolg von Unterricht messen zu wollen, solle man sich lieber um eine Verbesserung der Unterrichtsbedingungen kümmern.

Zentrale Prüfungen haben aber durchaus wichtige Funktionen jenseits des eigenen klasseninter-

nen Assessment, das im Folgenden eingegangen werden soll. In diesem Zusammenhang wird immer wieder der Begriff des „Bildungsmonitorings“ fallen. Mit diesem Begriff ist die Leistungsüberprüfung des deutschen Bildungssystems gemeint, auf das sich die Kultusministerkonferenz im Jahr 2006 geeinigt hat, um bestehende Defizite auszumachen und zu beheben. Als Instrumente wurden u. a. Vergleichsarbeiten (VERA) zur Überprüfung landesweiten und internationaler Leistungsfähigkeit beschlossen.



Bei „*Educational Reform Strategies*“, so der Text des Cartoons, möchte man sich nicht auf Wunder verlassen, sondern die Steuerung auf solide Daten stützen und dazu zählen natürlich möglichst ob-

jektive Messungen, was Schüler gelernt haben und können. Dies ist eine in angelsächsischen Ländern selbstverständliche Überlegung, in Deutschland war das bis PISA ungewohnt neu.

I PISA und die Folgen: Bildungsstandards

PISA wird eigentlich in Großbuchstaben geschrieben, da es für *Programme for International Student Assessment* steht und nicht für die Stadt mit dem schiefen Turm. Es ist ein Versuch der OECD, die Kenntnisse und Kompetenzen von Fünfzehnjährigen in Mathematik und Naturwissenschaften sowie ihre Lesefähigkeit in der jeweiligen Muttersprache zu ermitteln. Die OECD handelte im Auftrag von Regierungen, in Deutschland der Kultusministerkonferenz, um „Bildungsmonitoring“ zu betreiben, also nicht nur eine „Ist-Beschreibung“ der Ressourcen zu liefern, die einer Gesellschaft in Zukunft zur Verfügung stehen, sondern auch zu deren Verbesserung beizutragen. Warum also der Aufschrei in den Medien und unter den Lehrkräften?

Getroffene Hunde bellen, heißt es. Vorurteile und Weisheit. PISA entzauberte die deutsche Schule, die sich zuvor unter den besten Bildungssystemen der Welt gesehen hatte, denn die Bundesrepublik schnitt in dieser Studie anfangs nicht so gut ab. In der Folge wollte man zum einen genauer wissen, wo das deutsche Bildungssystem international steht und zum anderen, wie sich einzelne Bundesländer im Vergleich zueinander verhalten konnten. Dahinter steckten vor allem partei- und bildungspolitische Interessen, wie etwa die Frage nach der Struktur des Schulsystems (Stichwort: Gesamtschule) und der Förderung bestimmter sozialer Gruppen (Jungen, Mädchen, später auch Kinder von Migranten). Vor allem „katholische Mädchen aus Arbeiterfamilien auf dem Lande“ waren seit den 1960er Jahren, als Georg Picht in Deutschland die „Bildungskatastrophe“ ausrief, der

Prototyp für Förderungswürdige Schülergruppen. Dies sah man 30 Jahre später als behoben an. Doch dann kam die Kritik der PISA-Kritiker.

Man möchte sich zu fragen, was Finnland und andere skandinavische Länder in ihrem Schulwesen besser machen als die Bundesrepublik. Insbesondere im Bereich der Lesekompetenz lagen sie weit vor den meisten anderen Ländern an der Spitze der Vergleiche. Als Folge dieses „PISA-Schocks“ schuf die Kultusministerkonferenz „Bildungsstandards“ sowie zentrale Abschlussprüfungen, mit deren Hilfe die Analyse der Schulstandards auf den Prüfstand gestellt werden sollte.

Das rief die Kritiker noch stärker auf den Plan. Die PISA selbst: Bildung könne man nicht standardisieren. Kritete ein zentraler Vorwurf. Interessanterweise geschah dies oft unter Berufung auf den Humboldtschen Bildungsbegriff aus dem 19. Jahrhundert. Bildungsforscher wie Eckhard Klieme versuchten die Debatte zu versachlichen, indem sie begriffliche Präzisierungen vorschlugen:

Bildungsziele fallen nicht vom Himmel und sie haben nicht den Status unbefragbarer Gewissheiten, sondern verdanken sich historischen Kontexten und nationalen Traditionen. Wer „Allgemeinbildung“ sagt, der beansprucht eine – erkennbar deutsche – Tradition des „Bildungs“-denkens und der Interpretation von „Kultur“, [...] wer von „Basisfähigkeiten“ ausgeht, nimmt die Debatte über Standards und notwendige Erwartungen an Schule zur Kenntnis. Bildungsziele sind also, sichtbar an solchen Traditi-

¹ Die Diskussion um PISA ist in Wikipedia sehr gut dokumentiert: <https://de.wikipedia.org/wiki/PISA-Studien> (letzter Zugriff: 01.12.2017)

onen, in ihrer konkreten Gestalt immer Ergebnis gesellschaftlicher Entscheidungen und sozialer Machtlagen [...].

(Klieme 2003: 58)

In den Chor der Kritiker stimmten auf der „Frühjahrskonferenz zur Erforschung des Fremdsprachenunterrichts“ auch die meisten Fremdsprachendidaktiker der deutschen Hochschulen ein (Bausch et al. 2003). Mit „Standardisieren“ waren dabei nicht nur mögliche Stufungen gemeint, sondern die Profile von Abschlussprüfungen. Es waren nämlich just die Tests von Cambridge Assessment, die seit dem Erscheinen des *Gemeinsamen europäischen Referenzrahmens für Sprachen* (Europarat 2001) auf diesen „geeicht“ waren. Als Musteraufgaben wurden sie zunächst in den Bildungsstandards für die erste Fremdsprache (Englisch/Französisch) veröffentlicht: im Jahr 2003 für den Mittleren Bildung

abschluss und im Jahr darauf auch für den Hauptschulabschluss.

Im Grunde war klar, was auf Lehrkräfte (und Lernende) zukam. Begeisterung kostete dies allerdings nicht gerade aus. Die Cambridge-Prüfungen unterschieden sich erheblich von den gängigen Klassen- und Abschlussarbeiten. In dem, dass diese herkömmlichen Arbeiten einen höheren Bildungswert hatten, lag die Aufgabe nicht in standardisierten Tests benutzbar. Das übersah man großzügig. Viel schwerer fiel es Gewicht, dass Cambridge-, TOEFL- und andere Tests Alltagstexte und Alltagssituationen beinhalten, um Lese- und Hörverstehen zu überprüfen. Trotz der „kommunikativen Wende“ dominierte noch immer die Grammatik, die Hör- und Leseverstehen und Überprüfung von Hör- und Leseverstehen erst ganz langsam im Fremdsprachenunterricht Einzug hielt.

2 DESI: Bildungsmonitoring für den Englischunterricht

Parallel zu den Bildungsstandards gründete die KMK ein Konsortium von Fachdidaktikern und Bildungsforschern den Auftrag, ergänzend zu PISA eine Bestandsaufnahme in den Fächern Deutsch und Englisch zu versuchen: *DESI = Deutsch-Englisch-Schülerleistungen-International*.² Dies war der erste Anlauf zum Bildungsmonitoring für den Englischunterricht, mit dem sich Lehrkräfte konfrontieren mussten, nachdem Deutsch und Mathematik schon bei PISA auf dem Prüfstand waren. Die Ergebnisse dieses Bildungsmonitorings waren besorgniserregend. Am Ende der neunten Jahrgangsstufe erreichten nur zwei Drittel der Schüler im Mündlichen das GEN-Niveau A2, das die KMK für den Hauptschulabschluss erwartete. Ein Drittel erreichte B1, das für den Mittleren Schulabschluss erwartet wurde. Allerdings erreichten sogar neun Prozent B2, wie es eigentlich erst in der Ober-

stufe erwartet werden kann. Ähnliches gilt für andere Teilfertigkeiten.

DESI zeigt damit, dass wir im Englischunterricht – vor allem in den Gymnasien – eine sehr starke Leistungsspitze von 10 bis 15 Prozent der Schülerinnen und Schüler haben, deren Kompetenzen weit über das Anforderungsniveau der Lehrpläne hinausragen. Andererseits zeigt DESI vor allem in Hauptschulen, Integrierten Gesamtschulen und Schulen mit mehreren Bildungsgängen deutliche Defizite. Im Bildungsgang Hauptschule erreicht etwa nur ein Drittel der Schülerinnen und Schüler das Regelziel der Bildungsstandards.

(Klieme 2006: 2)

² Der Name war etwas hoch gehängt, denn „International“ war die Studie nur, weil auch einige Schulen in Österreich und die deutschsprachigen Schulen in Südtirol miteinbezogen wurden.

Des Weiteren zeigten die Ergebnisse, dass die fremdsprachlichen Leistungen von Mädchen denen von Jungen überlegen sind. Oft waren auch Kinder mit Migrationshintergrund denjenigen Kindern überlegen, die Deutsch als Muttersprache sprechen:

Das Aufwachsen in einer mehrsprachigen Familie ist unter sonst gleichen Lernbedingungen (sozialer Hintergrund, kognitive Grundfähigkeiten, Geschlecht, Bildungsgang) im Englischen mit einem Leistungsvorsprung verbunden, der den Gewinn mindestens eines halben Schuljahres ausmacht. Auch Schülerinnen und Schüler, die aus Migrationsfamilien mit ausschließlich nicht-deutschem Sprachhintergrund stammen, zeigen im Englischunterricht vergleichsweise gute Leistungen.

(DIPF 2006: 5)

3 „Is my BI your BI?“

Mit welchen Problemen sich Bildungsforscher auseinandersetzen müssen, kann man anhand von Studien beobachten, in denen es um die Ermittlung des Niveaus von Fremdsprachkenntnissen in Europa ging: dem *European Survey on Language Competences* (2012) und der Folgestudie *Study on Comparability of Language Testing in Europe* (2015). Obwohl Deutschland an der 2012er Studie nicht beteiligt war und auch die 2015er Studie nicht das deutsche Schulsystem betraf, lohnt sich ein kurzer Exkurs, um Einblicke in methodischen Schwierigkeiten zu gewinnen.

Beiden Untersuchungen lagen die Kompetenzniveaus des GeR zugrunde, die sich daran orientierten, sollten Sprachkompetenzen von Schülerinnen und Schülern überprüft.⁴ In beiden Fällen ging es aber eben nicht um eine Antwort auf J. Charles Aldersons bereits 2004 formulierte Fra-

Hinzu kam ein weiterer Befund, der nicht unbedingt anderen Studien entsprach: Sprachliche Kompetenzen (auch in Englisch) stehen laut *DESI* eher im Zusammenhang mit kulturellen Ressourcen der Familie und damit wie positiv die Familie Bildung allgemein gegenüber sieht und nicht unbedingt mit deren Einkommen.

Mit Aktionen des IASA und *DESI* werden Lehrkräfte immer häufiger sensibilisiert, sodass sich ein Gewöhnungsprozess eingestellt hat. Daher halten sich die abwehrenden Reaktionen darauf mittlerweile auch in Grenzen. Es gilt ebenfalls für die vom IQB vorgelegte Studie *Bildungstrend*³, die seit 2015 jährlich vorgelegt wird und die weiteren Bildungsberichte.

„Is my BI your BI?“ Denn nicht jedes am GeR angelegte, gerichtete Modell von Sprachkompetenz ist deckungsgleich mit anderen Definitionen und nicht jede Verteilung von Schülerleistungen auf Skalen entspricht der in anderen Schulsystemen. Selbst wenn die jeweiligen Befunde statistisch untermauert sind, stellen sich viele Fragen, wie das methodisch erfolgt ist.

Hinzu kommt, dass der GeR bzw. der *CEFR* überarbeitet wird. Aufgrund dieser Überarbeitung muss eventuell auch das Spektrum der Referenzniveaus neu „geeicht“ werden, da sowohl bei den niedrigen Niveaus „vor A1“, A1 und A2, als auch auf C2 einiges ins Rutschen gekommen ist. Das ist wichtig für die nächsten Abschnitte in diesem Kapitel, die sich mit Vergleichs- und Abschlussarbeiten beschäftigen werden – ein Thema, mit dem sich Lehrkräfte verstärkt auseinandersetzen müssen.

³ <https://www.iqb.hu-berlin.de/bt> (letzter Zugriff: 01.12.2017)

⁴ Das Design und die Ergebnisse der beiden Erhebungen sind in einem Beitrag von Neil Jones in der Zeitschrift für Fremdsprachenforschung (ZFF 2016, 27/1, 39–57) nachzulesen.