

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29

Comparing global judgments and specific judgments of teachers about students' knowledge: Is  
the whole the sum of its parts?

Karina Karst, Stefanie Dotzel, Oliver Dickhäuser  
University of Mannheim

Author Note:

Prof. Dr. Karina Karst, Assistant Professor of Quality of Instruction in Heterogeneous  
Contexts, Department of Psychology, University of Mannheim; Stefanie Dotzel, M.A., Department of  
Psychology, University of Mannheim; Prof. Dr. Oliver Dickhäuser, Full Professor of Educational  
Psychology, Department of Psychology, University of Mannheim

This research was supported in part by grants (PLI 3026A/ Kassel: Prof. Frank Lipowsky &  
PLI 3026B/ Bamberg: Prof. Gabriele Faust) from the German Federal Ministry of Education and  
Research (BMBF).

Correspondence concerning this article should be addressed to Prof. Dr. Karina Karst,  
Assistant Professor of Quality of Instruction in Heterogeneous Contexts, Department of Psychology,  
University of Mannheim, A5,6, 68131 Mannheim.

Email: [karst@uni-mannheim.de](mailto:karst@uni-mannheim.de)

This article may not exactly replicate the final version published in the journal. It is not  
the copy of record.

**1 Abstract**

2 Teachers' judgments about students' knowledge and skills can be global or specific depending  
3 on the diagnostic situation during teaching. We test the relationship between these judgments, their  
4 accuracy, and whether global judgment (GJ) accuracy can be measured by aggregating specific  
5 judgments (SJ). Judgments of 52 primary school teachers about their students' achievement in a  
6 standardized mathematics test were assessed. SJs and GJs correlated high. However, SJs were slightly  
7 more accurate than GJs. Additionally, teachers' GJ accuracy is not similar to the accuracy of  
8 aggregated SJs. We conclude that teachers use different judgment strategies for GJs and SJs.

9

10

**11 Highlights**

12 Teachers are quite consistent in their judgments about their students.

13 Judgment accuracy of teachers depends on the kind of judgment and on the used accuracy measure

14 Student-specific judgments (SJ) were slightly more accurate than student-global judgments (GJ).

15 The accuracy measures of aggregated SJs do not represent teachers' GJ accuracy.

16 The whole is not the sum of its parts

17

18

## 1 Introduction

2 Students' learning is more likely to be successful the more closely the learning process is  
3 related to their existing knowledge (Baumert, Lüdtke, Trautwein, & Brunner, 2009, Harwell et al.,  
4 2007; Weinert & Helmke, 1995). This is why it is helpful for teachers to have insight into students'  
5 prior knowledge in order to initiate adaptive learning opportunities for their students in the classroom  
6 (Brühwiler & Blatchford, 2011). Teachers receive this insight during various diagnostic situations in  
7 or outside of the classroom. Such situations require the teacher to make judgments about students'  
8 knowledge and skills. Thus, making diagnostic judgments about students seems to be a central and  
9 necessary teacher task (Burns, 1984; Herppich et al., in press).

10 In a recently suggested model of teachers' assessment competencies, Herppich et al. (in press)  
11 assumed that the kind of diagnostic judgments about the student *and* the accuracy of the judgments  
12 can differ with respect to the degree of specification of the judgment. This variation in specification  
13 arises from different aims of diagnostic situations during teaching that determine a particular  
14 perception of the teacher about the student. Some teaching situations may require a global judgment  
15 (e.g., when a teacher has to prepare test-tasks for an exam that are not too easy or difficult, and  
16 therefore reflects on the mean level of the students' competencies), whereas some situations may  
17 require a more specific judgment (e.g., when a teacher prepares specific individual support for a  
18 student).

19 The differentiation of either global or specific teacher judgments has only recently been made  
20 within the literature. A review of the research literature in 2012 suggested that a systematic analysis of  
21 similarities and differences between these two kinds of judgment is lacking (Südkamp, Kaiser, &  
22 Möller, 2012). Therefore, it seems relevant to investigate how student-specific judgments (SJs)  
23 correspond to student-global judgments (GJs), and whether SJ and GJ differ with regard to accuracy.  
24 Additionally, there is little educational research on processes underlying these kinds of judgment  
25 (Crisp, 2013; Herppich et al., in press). It is unclear whether teachers make GJ based on summing up  
26 many single SJs or not. If teachers in fact make GJ by aggregating many single SJs, the initial  
27 judgments and the judgment processes would be the same and we would neither assume any  
28 differences between the teacher's GJ and SJs, nor between the corresponding accuracy measures.

1 Thus, the whole (GJ) would be the sum of its parts (SJs). This would also reduce the time needed to  
2 investigate teachers' judgment accuracy. Assessing only one kind of teachers' judgment would suffice  
3 for gaining insight into overall teachers' judgment accuracy.

#### 4 **Diagnostic Situations During Teaching**

5 Diagnostic situations emerge during and beyond teaching and involve judgments about  
6 individual students regarding characteristics that are relevant for learning- (e.g., prior knowledge, self-  
7 concept, and emotional stability; Herppich et al., in press). The aim of the teacher within these  
8 situations is to gain information for an upcoming pedagogical decision (Author, 2001).

9 One can imagine a huge variability between those situations. Author (2001) postulated four  
10 independent dimensions that describe the variability of these situations. These dimensions are:  
11 Purpose of the judgment (assessment for learning/formative assessment or assessment of  
12 learning/summative assessment), the option that the situation can be planned in advance, the binding  
13 character and consequence of the judgment for the student, and the degree of specificity of the  
14 teachers' perception of the student. We will focus on the last dimension "specificity of the teachers'  
15 perception of the student". Concerning this dimension, it is possible to differentiate between two  
16 endpoints so that the perception of the student could be more global or more specific.

17 The following example illustrates the importance of student-global perceptions of a teacher.  
18 The teacher plans a lesson with the aim to strengthen and consolidate the students' newly acquired  
19 knowledge. For this aim, he/she realizes internal differentiation by adopting methods supporting an  
20 internal differentiation resulting in groups of learners with different achievement levels. In order to  
21 manage this situation, the teacher should consider the students' different competence levels and group  
22 their students according to their overall competence levels in the learning unit (Shavelson & Stern,  
23 1981). This is the step where GJ accuracy becomes relevant. The judgmental task for the teacher is to  
24 answer the question: "How good is my student at solving such exercises that I'd like to implement  
25 during collaborative learning?" This requires a rather global (comprehensive) judgment about a  
26 student (Author, 2012).

27 By contrast, SJs probably more strongly influence teachers' behavior in situations where  
28 individual support takes place (Helmke & Schrader, 1987). This assistance could be realized by, for

1 example, specific explanations or targeted use of teaching material or individual feedback. The  
2 underlying question the teacher has to answer is: "How good do I think my student would be at  
3 solving this specific exercise? Does he or she need a specific explanation?" Thus, the perception of the  
4 student and the judgmental task are more specific. It is a specific judgment of an individual student.  
5 The contrast to the student-global diagnostic situation is that the teacher focuses on an individual  
6 student's prerequisites of a single student that are needed to master certain exercises.

7 Because of these different situations, the perception of the teachers about their students varies,  
8 as do the kinds of judgments, which can change from more or less specific judgments to global  
9 judgments (Author, 2001). The question, therefore, is to identify whether there are differences in the  
10 judgments and in the accuracy of the judgments depending on the kinds of judgment.

### 11 **Judgment Processes**

12 It is not only important to consider the differences in the results of these kinds of judgment but  
13 one should also pay attention to the information processing that underlies the judgments. Theories of  
14 social judgments (Chaiken, 1980; Fiske & Neuberg, 1990) essentially state that people (e.g., teachers)  
15 can apply two different strategies to process information for judging other people (e.g., students).  
16 Individuals can use either a heuristic strategy or a controlled strategy. The heuristic or short-cut  
17 strategy is assumed to be relatively automatic and to require little cognitive effort. Using this strategy  
18 leads to a simplification of the information in order to handle the complexity of the judgment more  
19 easily. Teachers may use the strategy as a result of (adverse) external conditions during teaching such  
20 as less time or energy of the teacher for competing tasks. Using the controlled strategy, the individual  
21 collects more information and integrates it into a judgment.

22 Regarding the two kinds of judgments (GJ vs. SJ), it seems likely that for a GJ, the teacher  
23 uses a heuristic strategy, whereas a controlled strategy is more likely to be applied when making the  
24 SJ. A SJ about students requires and facilitates the consideration of individualizing information about  
25 the student and about characteristics of the exercise (task characteristics). This, in turn, possibly  
26 requires more elaborated cognitive processes than a GJ because the teacher has to compare the ability  
27 of the students with the concrete task characteristics. GJs, by contrast, are probably easier to make  
28 because the teacher can use heuristics as a base for judgments. The idea that GJs and SJs are

1 associated with different types of processing is based on the reasoning developed by Sherman, Beike  
2 and Ryalls (1999). The authors argue, that general judgments are more likely to be influenced by  
3 heuristic processing as heuristics are more likely to exist for general cases. Consistent with this idea it  
4 has been found that abstract entities are more likely to be judged based on stereotypes (Fiske &  
5 Neuberg, 1990).

6         However, an assumed lower difficulty in making GJs does not necessarily mean that these  
7 judgments are more accurate than SJs. As SJs more likely require a systematic and controlled  
8 approach, one may assume that the accuracy for SJ is higher than for GJ. However, there are also  
9 reasons to assume higher accuracy of GJ compared with SJ. Brighton and Gigerenzer (2012) could  
10 show that heuristics do not necessarily lead to lower judgment accuracy than more complex  
11 procedures of information integration. Thus, information resulting from heuristic strategies might be a  
12 better predictor (Brighton & Gigerenzer, 2015) and thus lead to higher judgment accuracy. With  
13 regard to teachers, this could mean that they only consider information that is actually relevant for a  
14 global judgment, which in turn leads to a higher level of judgment accuracy.

15         Thus, we cannot make a clear prediction regarding if GJs or SJs lead to higher judgment  
16 accuracy. Nevertheless, we assume that different judgment processes underlie two different kinds of  
17 judgments. If teachers are asked to judge a student more globally, they will not do this based on many  
18 specific judgments. Low or no correspondence between GJ and aggregated SJs and its analogous  
19 accuracy measures would indicate those different processes. Thus, the sum of many SJs would not  
20 result in a GJ.

### 21 **Kinds of Judgment and Judgment Accuracy – State of Research**

22         As mentioned before, the differentiation between GJ and SJ is relatively new. However, this  
23 does not mean that previous research did not analyze differences in judgment accuracy for different  
24 kinds of judgment. The following section outlines the state of research concerning the accuracy of  
25 different kinds of judgments.

26         The level-, differentiation-, and rank component (Schrader, 1989, Helmke & Schrader, 1987)  
27 are the most widely used accuracy measures of teachers' judgments about their students' skills and  
28 knowledge. The level component shows whether the teacher is a good judge of the students' mean

1 achievement level. It derives from the difference between the average teacher judgment and the  
2 average criterion value (e.g., average student achievement in a standardized test). Values above 0  
3 show that the teacher overestimated the students' achievement, whereas levels below 0 indicate  
4 underestimation. The differentiation component is the quotient of the variance of teacher judgments  
5 and the variance of the student's test scores. Values above 1 indicate that the variance of the students'  
6 achievement was overestimated by the teacher, whereas values below 1 indicate underestimation. The  
7 rank component is the product-moment correlation between teacher judgments and students' test  
8 scores (criteria), and illustrates whether the teacher is aware of the rank order of his/her student  
9 concerning the judged characteristic/criterion.

10 Overall, it can be summarized that teachers' judgment accuracy is moderate. In a meta-  
11 analysis, Südkamp et al. (2012) reported a mean rank-component of .63 with a high variability  
12 (-0.03 up to 1.18; Fisher  $z$ -transformed correlations). The results are even more heterogeneous for the  
13 level and differentiation component. Whereas some studies reported an overestimation of students'  
14 achievement level (Bates & Nettelback, 2001; Südkamp, Möller, & Pohlmann, 2008), some studies  
15 reported an underestimation or over- and underestimation within one study for two measurement  
16 points (Feinberg & Shapiro, 2009; Stang & Urhahne, 2016). The same heterogeneity concerning the  
17 results emerges from the differentiation component. Furthermore, there were only weak correlations  
18 between the three different judgment accuracy measures (Lorenz & Artelt, 2009; Praetorius, Karst,  
19 Dickhäuser, & Lipowsky, 2011; Südkamp, Möller, & Pohlmann, 2008; Spinath, 2005; Schrader,  
20 1989). These weak correlations indicate independency between the judgment accuracy measures and  
21 make a separate analysis of each accuracy measure necessary to gain a comprehensive picture about  
22 teacher's judgment accuracy.

23 Within their meta-analysis and by conceptualizing a heuristic model of judgment accuracy,  
24 Südkamp et al. (2012) attempted to find reasons for this heterogeneity of results. They postulated that  
25 different kinds of judgments would be associated with different levels of judgment accuracy.  
26 According to Südkamp et al. (2012) the following factors (a-d) characterize the kind of judgment: (a)  
27 informed versus uninformed; (b) congruence of domain specificity, (c) direct versus indirect, and (d)  
28 specificity of the judgment.

1 An (a) informed teacher judgment means that the teacher judges their students' test  
2 performance in an achievement test, which is also simultaneously administered to the students. In  
3 contrast to this, an uninformed judgment means that the teacher judges the competence of their  
4 students in any domain using a Likert-type rating scale. Here, teachers are uninformed about the  
5 criterion their judgment will be compared with (e.g., "How do you assess the language skills of  
6 student x?"). The latter kind of judgment is usually more comprehensive and global about the student.  
7 Südkamp et al. (2012) found that the judgment accuracy is significantly higher for informed  
8 judgments than for uninformed judgments ( $\beta = .15$ ;  $p = .045$ ).

9 The factor (b) congruence of domain specificity refers to the accordance of teacher judgments  
10 and the achievement test. If the teacher judgment is specific (judgment of an academic ability in one  
11 subject) *and* the achievement test is specific (covering a single subject), Südkamp et al. (2012) called  
12 it a high correspondence of domain specificity. In their meta-analysis, they found higher judgment  
13 accuracy for teachers within studies that measured judgment accuracy with a high congruence of  
14 domain specificity ( $\beta = -.13$ ;  $p = .009$ ).

15 A (c) direct judgment is an informed judgment of the teacher with more or less  
16 correspondence of domain specificity, whereas an indirect judgment is defined as an uninformed  
17 judgment with low correspondence. For this factor, Südkamp et al. (2012) found no significant  
18 differences for judgment accuracy. Additionally, Hoge and Colardarci (1989)—authors of an earlier  
19 meta-analysis—found slightly higher judgment accuracy (rank components) for direct judgments, but  
20 the differences were not significant neither. In addition, in a narrative review, Harlen (2005) reported a  
21 higher accuracy for direct judgments as compared to indirect judgments.

22 Lastly, teachers' judgments about students' knowledge and skills can be differentiated  
23 according to the factor (d) specificity. Based on suggestions by Hoge and Colardarci (1989), Südkamp  
24 et al. (2012) differentiated between five levels of specification. The most specific judgment is needed  
25 when a teacher is asked "Will your student x solve the exercise y?" (Level 5). Less specific and thus  
26 more global is the judgment when a teacher is asked "How many of these seven exercises will your  
27 student x solve?" (Level 4). The lowest degree of specificity (Level 1) is realized by ratings when the  
28 teacher rates their students' competence on a rating-scale (e.g., poor to excellent). Judgment accuracy



1 was significantly higher for judgments with a higher level of specificity in the Hoge and Coladarci  
2 meta-analysis (1989), but not in the meta-analysis of Südkamp et al. (2012).

3 Looking at recently published studies that had not been included in the reported meta-  
4 analyses, the findings are heterogeneous. Karing, Matthäi, and Artelt (2011) analyzed direct (specific)  
5 and indirect less specific judgments and compared the rank components. For 64 primary school  
6 teachers, they reported a significantly higher rank component for the indirect judgment than for the  
7 direct judgment ( $t = 1.70, p < .05, d = 0.27$ ). They concluded that direct judgment accuracy was lower  
8 because direct judgments were more difficult for teachers to make due to their cognitive complexity.

9 In a study published by Zhu and Urhahne (2014) 16 primary school teachers were asked to  
10 make direct (multiple-item) judgments and indirect (single-item) judgments about their students'  
11 learning relevant motivation and emotions. The authors computed two different types of accuracy  
12 measures: the rank component (for the direct and indirect judgment) and percent agreement (only for  
13 the direct judgment). They reported no significant difference between the accuracy of indirect  
14 judgments and the accuracy of direct judgments. Thus, the authors regarded an indirect measure as a  
15 simple, efficient and valid method to evaluate teacher's judgment accuracy about students' motivation  
16 and emotions.

17 Considering the described state of the research, we conclude that the kinds of judgment  
18 teachers were asked to conduct in educational studies probably influences the judgment accuracy.  
19 However, it is not possible to decide whether teachers are better judges when they are prompted to  
20 make direct, or informed, or specific, or indirect, or uninformed, or less specific judgments.

21 The comparisons of direct and indirect judgments and their corresponding rank component (as  
22 one possible accuracy measure) dominated previous research. So far, no analyses exist that measure  
23 differences and relationships in judgment accuracy for the other two components, namely level and  
24 differentiation component. In most of the cases, such analyses were simply not possible because these  
25 components could not be calculated for single-item/indirect judgments. As these two components are  
26 also important elements of judgment accuracy, they should be integrated to investigate the question of  
27 whether judgment accuracy varies depending on the kind of judgment. However, a few studies exist,  
28 which compared various judgment accuracy measures (e.g., Zhu & Urhahne, 2014). Nevertheless,

1 those studies focused on the correlations between the accuracy measures and ignored to calculate the  
2 mean differences of the accuracy measures.

3 As seen from previous considerations, there is a clear overlap between the factors  
4 characterizing the kind of judgment and the present differentiation of GJ and SJ (Table 1).

5 Table 1

6 However, these factors characterizing the kinds of judgment are not entirely identical. An SJ is  
7 always a direct measure, but a GJ could be direct or indirect. In the manuscript, we analyze GJ on a  
8 direct level, which corresponds to the fourth category of specification according to Hoge and  
9 Coladarci (1989). Even though GJ and SJ are not identical to the factors of former research, they can,  
10 however, be integrated into the heuristic model of judgment accuracy (Südkamp et al., 2012). Both  
11 judgments (GJ and SJ) are informed judgments with a high congruence of domain specificity. In the  
12 present study, these two factors will be kept constant in order to avoid the confounding of more or less  
13 specific judgments with other judgment characteristics influencing judgment accuracy, such as direct  
14 or indirect measures. In almost all previous studies, the indirect measure was operationalized by a  
15 global judgment and the direct measure by a specific one. Thus, the effect of differences to the extent  
16 of judgment accuracy because of variations in the kind of judgment could not be ascribed to a single  
17 factor. Furthermore, the distinction between GJ and SJ takes the teaching activities and diagnostic  
18 tasks of teachers into account. Thus, the theoretical innovation and benefit of this distinction is the  
19 derived link to teaching situations. This theoretical derivation provides insight into when teachers need  
20 to make a GJ or SJ.

## 21 **Research Questions**

22 Based on this previously discussed research and the existing gaps, the present study intended  
23 to empirically investigate the following research questions:

- 24 (1) Does the accuracy of judgments differ with respect to the kind of judgments?
- 25 (2) Is it possible to express GJ-accuracy by aggregating single SJs of the teacher?

## 26 **Method**

### 27 **Participants**

1 To answer the above stated research questions empirically, we used the data set of a German  
2 primary school study. This study was conducted in three federal states of Germany: Saxony, Berlin,  
3 and Mecklenburg-Western Pomerania.

4 The teacher sample comprised 52 teachers working in 37 classes in 20 primary schools: the  
5 higher number of teachers compared to the number of classes is explained by a particular type of  
6 classroom organization in several schools, in which two teachers work in the same classroom with the  
7 same students for Mathematics and German (15 of the 37 classes in this case had two teachers). Of  
8 these teachers, 59% were in the age range 35 to 45, 98% (all but one) were female, and average time  
9 in teaching was 18.9 years (in general, on this basis, the teacher sample represented the primary  
10 teacher population in Germany (Tarelli, Lankes, Drossel, & Gegenfurtner 2012). The roughly gender-  
11 balanced student sample comprised 700 first graders (mean age 7 years, 4 months, 52% girls). The  
12 average class size included 18.9 students.

### 13 **Materials and Procedure**

14 Students completed a standardized math test at the end of grade 1. The *student's math test* was  
15 developed for this purpose and was based on a standardized achievement test (DEMAT 1 and 2;  
16 Krajewski, Küspert, & Schneider, 2002; Krajewski, Liehm, & Schneider, 2004). The test assessed  
17 students' arithmetic competence. The 45-minute test consisted of 49 test-tasks (items; example: Do the  
18 math of this addition!  $7+21 = \underline{\quad}$ ). All students received the same test (Author, 2012).

19 While the students filled out their test, the teacher judgments about a subset of the whole math  
20 test were collected via questionnaire. First, in order to assess teachers' *student-global judgment* (GJ),  
21 each teacher judged the number of each of their students' correct responses to items that represented  
22 two item-types (Example: "How many of these items will your student solve?"). The item-types were  
23 addition of two whole numbers (7 items) and reading numbers (8 items). For the *student-specific*  
24 *judgments* (SJ), the teachers had to judge 12 students of their class that were randomly chosen. The  
25 basis was six arithmetic items concerning basic arithmetic operations (e.g. addition with unknowns in  
26 all positions, basic word problems, and double a number). For each of these six items, the teachers  
27 reported whether each of the twelve students would be able to solve it. In comparison to the GJ, this

1 judgment asked the teacher to decide, for each item separately, whether a specific student would be  
2 able to solve the item (Author, 2012).

3 To achieve an informed test situation, the teachers received all test-items of the students' math  
4 test and the instruction of the administered math test. We used the same scale for the teacher  
5 judgments and student achievement. This enabled the direct comparison of the student responses with  
6 the two kinds of judgments and resulted in a high congruence of domain specificity.

7 The analyses for the GJ were conducted with 49 teachers and 639 students because three  
8 teachers did not completely fill in the GJ. For the SJ, the data were complete for all teachers.

9 The eight teacher judgments (two GJs and six SJs) had a high reliability, Cronbach's  $\alpha = .86$ .  
10 In addition, the subset of the students' math-test items (21 items: 7 addition items, 8 number reading  
11 items and 6 basic arithmetic operation items), which were the basis for the teachers' judgments, had a  
12 high reliability, Cronbach's  $\alpha = .88$ .

### 13 **Computing Accuracy Measures**

14 Before computing the accuracy measures, we transformed the teacher judgments and the  
15 students' test scores from absolute to relative numbers. This was necessary since we intended to  
16 compare the different kinds of judgments with the student achievement test scores directly.

17 Following the approach of Schrader (1989), the accuracy components were calculated for each  
18 kind of judgment (student-global and student-specific) and for each teacher. The components are the  
19 level-, differentiation-, and ranking component. Additionally, we computed level- and differentiation  
20 errors for both kinds of judgments. The *level-component* derives from the difference between the  
21 average teacher judgment and the average students' criteria value. The *differentiation component* is  
22 defined by the quotient of the variance of the teacher judgment and the variance of the students'  
23 values. The optimum value of the level component is 0 and the optimum value of the differentiation  
24 component is 1. Deviations from those optimum values can be positive or negative. Values above the  
25 optimum value indicate an overestimation of the level or the variance of the students' knowledge and  
26 skills. Values beneath these optimums indicate an underestimation of the level or the variance. Both  
27 the level and the differentiation components lacked an absolute zero point. This makes it impossible to  
28 analyze directed relationships of the different kinds of judgments of each component. Therefore, we

1 also computed the corresponding errors of the level and the differentiation component. The level error  
 2 is the absolute value of the difference between the average judgment and the average students' score.  
 3 The differentiation error is the absolute value of the difference between the variance of the judgments  
 4 and the variance of the students' scores. The higher these absolute values are, the higher is the  
 5 corresponding error and the lower is the judgment accuracy of the teacher. The *rank component* is the  
 6 the product-moment correlation between teacher judgments and students' test scores (criteria).<sup>1</sup> To  
 7 report the mean of the rank components the teacher-specific correlations were transformed using a  
 8 Fisher-*Z* transformation. This converts correlations into an almost normally distributed measure and  
 9 allowed us to calculate the mean rank component. Then, this mean *Z*-value was recalculated to *r*. The  
 10 level of significance was set at 5% for all inferential statistical analyses. For effect sizes we refer to  
 11 Cohen (1992).

## 12 Results

13 First, we present the descriptive results for the teacher judgments and the pertaining student  
 14 test scores. Then we continue with the results for the two research questions.

15 Table 2 shows the descriptive results for the teacher judgments as well as the achievement test  
 16 scores.

### 17 Table 2

18 Overall, the teachers expected their students to perform better on the math items they assessed  
 19 for the specific judgments than for the items they assessed for the global judgments. This assessment  
 20 proves to be correct because students performed better in the six specific items than in the global  
 21 items. Moreover, the teacher judgments show a broader range than the students' test scores.

---

<sup>1</sup> The Pearson's correlation coefficient is rated as a stricter criterion in contrast to Spearman's rank correlation coefficient. The reason for using the Pearson's correlation coefficient is that the calculation not only refers to a teacher's capability to estimate the ranking position of individual pupils within the class, but also to a teacher's capability to estimate the relative performance gap between individual pupils (cf. Schrader, 1989).



1 The effect size with  $q = .38$  illustrates that this difference is a moderate effect (Cohen, 1992). Thus,  
 2 the rank component of the global judgment is significantly higher than the rank component of the  
 3 specific judgment.

4 **Research Question 2– Transferring Specific into Global Judgments.** To analyze the  
 5 second research question, we had to aggregate the six specific judgments by summing these judgments  
 6 up for each student. After that, we calculated the accuracy measures as described above.

7 At first, we had a look at the correlation between the two kinds of judgments. We computed  
 8 for each teacher semi-partial product-moment correlations between the GJ and aggregated GJ. Semi-  
 9 partial correlation was necessary to control for the two different item-types on which the teachers  
 10 judged in the GJ (addition vs. reading numbers). We also conducted class-specific product-moment  
 11 correlation between the averaged students' test scores underlying the GJs and the averaged test scores  
 12 underlying the SJs. Subsequently, we transformed all correlations using Fisher-Z transformation as  
 13 described above. This allowed us to compare the correlations to each other. Secondly, we calculated,  
 14 as before by a  $2 \times 4$  repeated measure ANOVA, if the accuracy measures of each kind of judgment  
 15 (original global and aggregated global) differ in their value. Thirdly, and finally, product-moment  
 16 correlations measure the relationship between the same accuracy measures of the original global and  
 17 the aggregated global judgment.

18 (1) The semi-partial product-moment correlation between original GJ and aggregated GJ  
 19 revealed a strong relationship. The mean correlation is  $r = .78$  ranging from .41 to .93 between the  
 20 teachers. The mean correlation between the students' test scores is significantly smaller with  $r = .55$   
 21 ( $p < .001$ ) than the correlation between original GJ and aggregated GJ ( $z = -10.22$ ;  $p < .001$ ). In  
 22 addition, the correlations between the student test scores show a larger range (from .01 to .80).

23 (2) The results of the  $2 \times 4$  ANOVA are as follows: The main effect for original versus  
 24 aggregated judgment,  $F(1,48) = 68.96$ ,  $p < .001$  indicates that the accuracy differs with respect to this  
 25 factor. The main effect for the accuracy measure,  $F(3,46) = 415.11$ ,  $p < .001$  indicates that the means  
 26 of the four accuracy measures (level and differentiation) are unequal. The interaction effect is also  
 27 significant,  $F(3,46) = 20.24$ ,  $p < .001$ . Figure 2 shows the descriptive statistics of the aggregated and  
 28 global accuracy measures in comparison with each other.

## Figure 2

1  
2 To test differences of the means between the GJ-accuracy measure and aggregated GJ-  
3 accuracy measures, we conducted pairwise comparison corrected by Bonferroni adjustment and the Z-  
4 test referring to Steiger (1980). The results are shown in Table 4. The results for the level component  
5 and the level error are equivalent to the results in Table 3 because the level component and the level  
6 error do not differ between original SJs and aggregated GJ.<sup>2</sup> As mentioned above the level component  
7 for the original GJ indicates underestimation of the students, whereas the level component of the  
8 aggregated GJ is highly accurate in the mean. However, there is no significant difference between the  
9 two level errors. The aggregated differentiation component is, on average, significantly higher than the  
10 original global differentiation component ( $d = 1.19$ ). This can also be shown for the differentiation  
11 error ( $d = 0.46$ ). Due to aggregating the specific judgments, there is no difference between the two  
12 rank components ( $p = .84$ ). Thus, measuring the rank component by global judgments shows  
13 approximately the same results as the mean rank component measured by aggregating specific  
14 judgments.

## Table 4

15  
16 (3) For the correlations between the accuracy measures, the findings suggest that there is no  
17 significant relationship between the original global level error and the aggregated one (Table 5). The  
18 relationship between the two differentiation errors is low and non-significant, but is negative. The  
19 relationship between the two rank components is slightly positive, but also not significant ( $r = .21$ ).  
20 These findings suggest that teachers' assessment accuracy is not stable over different kinds of  
21 judgments, and that it is not possible to make a statement about the global judgment accuracy of a  
22 teacher by summing up specific judgments.

---

<sup>2</sup> The aggregated level component and the level error are identical to their non-aggregated original components (specific judgment). The reason for this is that the level component and the level error are theoretically defined as aggregated components because the mean students' test scores are subtracted from the mean teacher judgments.



Table 5

**Discussion**

In the scope of this study, two research questions have been investigated that focused on the comparison of teachers' GJs and SJs about their students' knowledge and skills. These judgments are part of everyday school life when the teacher has to adapt instruction to learning requirements, and go together with a certain perception of the students, namely a global versus specific perception.

Looking at teachers' judgment accuracy, our data showed the following: On average, the teachers do very well for both kinds of judgment. For example, the level error of 0.09 for the GJs indicates that teachers, on average, only overestimated or underestimated the students' performance by approximately one and a half correctly solved items (test-task) out of a total of 15 items (test-tasks). Also, the global rank component ( $r = .64$ ) is fairly high and was found to be similar to the reported mean rank component ( $r = .63$ ) in the meta-analysis of Südkamp et al. (2012).

The low rank component of the specific judgment is rather surprising. A first conclusion could be that the specific judgments are more difficult than global judgments. In specific judgments, the focus is on individuating characteristics of the students (e.g. prerequisites needed to master an exercise) and on characteristics of the exercise. The teachers have to decide whether a student will solve the exercise or not. In this case, false positive and false negative deviations have a more substantial impact on the calculation of the product-moment correlation (Bortz, 2005, p. 227), than it is the case for the GJ. In line with this assumption is the meta-analysis of Machts, Kaiser, Schmidt, and Möller (2016), in which the authors compared, among others, the judgment accuracy between dichotomous to non-dichotomous judgment scales. The judgment accuracy (measured as rank component) for non-dichotomous judgments (more than two response categories) was significantly higher than for dichotomous response scales ( $r = .45$  vs.  $r = .36$ ).

Concerning the first research question, it is not possible to clearly state whether one kind of judgment is more accurate than the other. Whereas the value of the level error and the value of the differentiation error do not depend on the kind of judgment, the level and differentiation components show better values for SJ than for GJ. Furthermore, the level component of SJ does not significantly deviate from the optimum value zero. Thus, the teachers overall reach an accurate assessment of their

1 students' achievement level when they make specific judgments about their students. The same picture  
2 emerges concerning the differentiation component. This component does not significantly differ from  
3 the optimum value of one. Based on many specific judgments about their students, teachers reach an  
4 accurate assessment of overall heterogeneity in the class. The differences in accuracy between the  
5 kinds of judgments relate to the direction of the judgments in the sense of over- or underestimation of  
6 the students' knowledge and skills, and not to the absolute difference of students' tests scores and  
7 teacher judgments. By contrast, the rank component is significantly higher and thus better for GJ than  
8 for SJ. This leads to the conclusion that it might be easier for teachers to get an accurate estimation of  
9 the achievement level of their students and of the amount of heterogeneity of their students when they  
10 make SJs about the students. However, when it comes to estimating the rank order and the relative  
11 performance gap of their students, in terms of the rank component, the more accurate judgments are  
12 GJ.

13         The focus of the second research question was the comparability of original GJ and  
14 aggregated GJ and their corresponding accuracy measures.

15         At first, the results suggest a strong relationship between GJ and aggregated GJ. This means,  
16 that the teachers' assessments of the students is quite consistent and that the specific type of math  
17 exercise (e.g., addition of whole numbers; double a number) makes no difference to their assessment.  
18 The relationship between teachers' judgments for individual students is considerably closer than the  
19 relationship between the individual students' performances. The closer relationship may also indicate  
20 that the teachers expect more stable student achievement across the judged test-tasks than is actually  
21 present. This shows that teachers almost use an overall assessment benchmark to judge the  
22 performance of their students, and that the teachers do not differentiate between the difficulties of the  
23 exercises.

24         Concerning the judgment accuracy, we found mean differences between the aggregated and  
25 original global accuracy measures. The aggregated level component is closer to 0 and thus  
26 significantly better than the original global level component. The level errors do not differ  
27 significantly. The aggregated differentiation component is higher than the original one and  
28 significantly above the optimum value of 1. Whereas original GJ leads to an underestimation of the

1 amount of heterogeneity, the aggregated GJ leads to an overestimation of the amount of students'  
2 heterogeneity. This finding suggests that by the integration of individuating information—as could be  
3 assumed for SJ—differences between students were unconsciously accentuated. Hence, this  
4 accentuation does not become obvious until the SJs are aggregated to GJ. By aggregating and thereby  
5 averaging the test scores, intra-individual differences are eliminated, which would otherwise be part of  
6 the computation of heterogeneity and thus of the differentiation component. The misjudgment of  
7 heterogeneity by aggregating the SJ is also obvious for the comparison of the differentiation errors.  
8 The error is larger for the aggregated differentiation error than for the original global differentiation  
9 error. Finally, the results show that the aggregated specific rank component is almost as high as the  
10 original global rank component. Further, correlation analysis shows no substantial and significant  
11 relationships between the same accuracy measures of the aggregated global and original global  
12 variants. This means that a GJ is not the sum of many single SJs about one student, and indicates that  
13 the judgment processes differ between the kinds of judgments.

14 Theories of social information processing (Chaiken, 1980; Fiske & Neuberg, 1990) postulate  
15 that different (social) judgments depend of different types of information, different information  
16 integration processes, which may result in different outcomes. Applied to the SJ assessed in the  
17 present study, one may assume that an SJ fosters a controlled and systematic information processing.  
18 SJ means that a teacher judges a student with regard to the likelihood to master an exercise. In doing  
19 so, the focus of the teacher is on individuating information about the student *and* the exercise. The  
20 way of information integration is guided by these specific characteristics of the judgment. Concerning  
21 GJ, one may assume that other factors besides individuating information influence information  
22 processing and thus also influence the GJ. Such factors may be characteristics that result from social  
23 comparisons of the teacher between the students. Through these social comparisons, the teachers draw  
24 attention to general differences between their students' knowledge and skills and thus neglect specific  
25 characteristics of individual students, which might be necessary for a high level of judgment accuracy.  
26 These social comparisons might then result in lower accuracy of the level and differentiation  
27 component for the GJ than is the case for these accuracy measures resulting from SJ. We consider this  
28 finding as a hint that GJ is based in a rather heuristic processing of information. The results for the

1 rank component can be interpreted similarly, since the rank component indicates the knowledge of the  
2 teacher about the relative positions of his students in the class regarding achievement. The lower rank  
3 component emerges when the teachers make SJ. By aggregating these judgments, the rank component  
4 increased in the mean. There are no differences in the mean between original global and aggregated  
5 rank component. By aggregation, intra-individual differences are averaged and neglected. This  
6 indicates different judgment processes. If individuating information is ignorable, then the rank order of  
7 students can be better assessed.

### 8 **Implications**

9         The results show that the kind of judging has an impact on the accuracy of the judgment.  
10        Within the scope of teaching and teacher education, we can derive the following implications: GJ are  
11        more appropriate to assess the rank order of the students' skills and knowledge. On the contrary,  
12        global judgments (which may be assumed to be more strongly based on heuristic information  
13        processing, Scherman et al., 1999) would be less suitable to assess the level and differentiation of the  
14        students as accurately as possible. For the latter purpose, specific judgments (which are assumed to be  
15        subject to controlled information processing) seem particularly appropriate. As in previous studies,  
16        there is no clear statement possible about which judgment (in our case: specific vs. global judgment)  
17        leads to higher accuracy. One could say that it depends on the component of accuracy that is  
18        considered and perhaps on the teachers' information processing mode. Related to teachers'  
19        professional development and teaching practice, we can conclude that teachers should use *those*  
20        judgments which are more accurate for pedagogical decisions. If the teacher aims to gain information  
21        about the achievement level of the students and about the amount of heterogeneity, results of this  
22        study would advise the teacher to mostly judge the students specifically. If the teacher seeks  
23        information about the rank order of the students, global judgments would be more accurate. Following  
24        this assumption, this would mean that the global assessment of the students leads to a more  
25        appropriate internal differentiation in the class because the teacher is more able to determine the rank  
26        order of the students.

27        These findings also indicate that it will be necessary in further research to provide different  
28        kinds of judgments because it is impossible to draw inferences about one accuracy measure (e.g., GJ)

1 given the other measure (e.g., SJ). Teachers' GJ are not a result of aggregated single SJs. The  
2 differentiation between kinds of judgment, and thus between diagnostic situations during teaching  
3 with respect to the perspective of judges (teachers) about the judged target (student), seems necessary  
4 from an empirical and theoretical point of view. This means that for professional teacher development,  
5 teachers should be instructed to use such judgments that fit best in the underlying diagnostic situation.

### 6 **Limitations of the Present Study**

7         Based on the results of our study, we see initial evidence supporting the assumption that the  
8 information processing underlying global and specific judgments differs. However, further research is  
9 needed. Additionally, considering that judgment accuracy is not a uniform construct (Praetorius et al.,  
10 2011), it seems important to gain more insight into judgment processes that can lead to different  
11 outcomes in judgment accuracy.

12         Limits of our study are due to the different types of arithmetic items, which form the basis for  
13 the two kinds of judgment. Colardarci (1986) showed that the accuracy of teachers' judgments differs  
14 with respect to the judged item-types within one subject, namely math. In his study, the difference  
15 between the items, however, concerned the complexity of the cognitive processes, which were  
16 necessary to solve the test-items.). Such a difference is not applicable for the item-types used in our  
17 study. All test-items showed a rather similar level concerning the complexity of cognitive processes,  
18 which were relevant for solving. Moreover, all items, which were assessed by teachers, are scalable in  
19 accordance with the 1PL Rasch model (Wu, Adams, & Wilson, 1998). In line with Adams and Khoo  
20 (1996), the measured fit statistics of this scale are satisfactory (EAP/PV reliability = .95, WMNSQ:  
21 0.75-1.26). To avoid this limitation, it might be possible in future studies to prompt teachers to judge  
22 the same items for each kind of judgment. However, this might seem very strange and artificial to the  
23 teachers, so such an approach would need other research designs. Additionally, it would be positive to  
24 expand the thematic areas of the judged items. This would raise the validity and generalizability of the  
25 results.

26         The fact that each teacher was required only to offer SJs for each of 12 of their students and,  
27 even more importantly, for just six separate test items is another weak spot of our study. A more

1 robust basis for the SJ data – that is more students judged and more SJs provided per student – would  
2 strengthen the validity of the findings, and therefore in turn its interpretational value.

3 Further limit occurs with the operationalization of the GJs. Referring to the GJ the situation  
4 for the teachers in the study and the situation they experience in ordinary work are not directly  
5 comparable. Teachers usually will make even more global judgments about their students' knowledge  
6 and skills, for example about a larger domain. Thus, the global judgments in our study would be closer  
7 to the endpoint of specificity than to the global endpoint if one thinks about these kinds of judgments  
8 as a dimension with two endpoints (cf. Hoge & Coladarci, 1989). However, empirical research about  
9 the kinds of judgments that are part of teachers' daily work is still missing. Until evidence based  
10 research exists the specificity of judgments remains a question of definition.

11 Our teacher-sample is a non-representative convenience sample of primary school teachers in  
12 Germany. Although it is not a random sample, it is comparable to the whole population of primary  
13 school teachers in Germany, namely with regard to gender ratio and mean age (Tarelli et al.; 2012;  
14 [www.destatis.de](http://www.destatis.de)). However, the results are not transferable to the whole population. The analyses are  
15 limited to the subject of math in primary school classes. Additionally, a larger sample would be  
16 desirable to verify the results. A small sample size could lead to an underestimation of correlations and  
17 their significance because the standard errors of the correlation coefficients become greater with a  
18 decreasing sample size (Bortz, 2005).

19 The reported accuracy measures in this study are manifest measures. However, some recently  
20 published studies measured the accuracy of teachers' judgments via latent variable modeling in the  
21 scope of multilevel analyses (Johansson, Myrberg & Rosén, 2012; Leucht, Tiffin-Richards, Vock,  
22 Pant, & Köller, 2012; Author, 2017). This strategy would also be an elaborated and desirable  
23 opportunity to compare SJ and GJ in future research. However, it is not possible to measure the  
24 differentiation component and differentiation error via linear mixed models. Thus, we decided to use  
25 the manifest variant because we wanted to compare all accuracy measures.

## 26 **Conclusion**

27 Considering this state of knowledge, we can summarize as follows: SJs lead to higher  
28 judgment accuracy concerning the assessment of the achievement level and heterogeneity of the class,

1 whereas GJs lead to higher judgment accuracy concerning the assessment of the students' rank order  
2 of the students in the class. Thus, teachers and researchers should prefer a specific judgment to  
3 receiving accurate insight into the achievement level of the students and the amount of heterogeneity  
4 in a class. To get information about the rank order and relative gap between the students, teachers and  
5 researchers should refer to global judgments. However, if only specific judgments are available,  
6 results show that the aggregated specific judgments reveal almost the same mean rank component as is  
7 obvious for the global judgments.  
8

## References

- 1  
2 Adams, R. J., & Khoo, S. T. (1996). *Quest*. Melbourne: Australian Council for Educational Research.
- 3 Author, 2001. [details removed for peer review]
- 4 Author, 2012. [details removed for peer review]
- 5 Author, 2017. [details removed for peer review]
- 6 Bates, C., & Nettelbeck, T. (2001). Primary school teachers' judgements of reading achievement.  
7 *Educational Psychology, 21*, 177–187. doi:10.1080/01443410020043878
- 8 Baumert, J., Lüdtke, O., Trautwein, U., & Brunner, M. (2009). Large-scale student assessment studies  
9 measure the results of processes of knowledge acquisition: Evidence in support of the  
10 distinction between intelligence and student achievement. *Educational Research Review, 4*,  
11 165–176. doi:10.1016/j.edurev.2009.04.002.
- 12 Bortz, J. (2005). *Statistik für Human- und Sozialwissenschaftler* (Vol. 6). Heidelberg: Springer.
- 13 Brighton, H., & Gigerenzer, G. (2012). How heuristics handle uncertainty. In P. M. Todd,  
14 G. Gigerenzer, & ABC Research Group (Eds.), *Ecological rationality. Intelligence in the world*  
15 (pp. 33–60). New York: Oxford University Press.
- 16 Brighton, H., & Gigerenzer, G. (2015). The bias bias. *Journal of Business Research, 68*, 1772–1784.  
17 doi:10.1016/j.jbusres.2015.01.061
- 18 Brühwiler, C., & Blatchford, P. (2011). Effects of class size and adaptive teaching competency on  
19 classroom processes and academic outcome. *Learning and Instruction, 21*, 95–108.  
20 doi:10.1016/j.learninstruc.2009.11.004
- 21 Burns, R. B. (1984). The process and context of teaching: A conceptual framework. *Evaluation in*  
22 *Education, 8*, 95–112.
- 23 Chaiken, S. (1980). Heuristic versus systematic information processing and the use of source versus  
24 message cues in persuasion. *Journal of Personality and Social Psychology, 39*, 752–766.  
25 doi:10.1037/0022-3514.39.5.752
- 26 Cohen, J. (1992). A power primer. *Psychological Bulletin, 112*, 155–159. doi:10.1037/0033-  
27 2909.112.1.155



- 1 Coladarci, T. (1986). Accuracy of teacher judgements of student responses to standardized test items.  
2 *Journal of Educational Psychology*, 78, 141–146. doi:10.1037/0022-0663.78.2.141
- 3 Crisp, V. (2013). Criteria, comparison and past experiences: How do teachers make judgements when  
4 marking coursework? *Assessment in Education: Principles, Policy & Practice*, 20, 127–144.  
5 doi:10.1080/0969594X.2012.741059
- 6 Feinberg, A. B., & Shapiro, E. S. (2009). Teacher accuracy: An examination of teacher-based  
7 judgments of students' reading with differing achievement levels. *The Journal of Educational*  
8 *Research*, 102, 453–462. doi:10.3200/JOER.102.6.453-462
- 9 Fiske, S. T., & Neuberg, S. L. (1990). A continuum of impression formation, from category-based to  
10 individuating processes: Influences of information and motivation on attention and  
11 interpretation. In M. P. Zanna (Ed.), *Advances in experimental social psychology*  
12 (Vol. 23, pp. 1–74). New York: Academic Press.
- 13 Harlen, W. (2005). Trusting teachers' judgement: research evidence of the reliability and validity of  
14 teachers' assessment used for summative purposes. *Research Papers in Education*, 20, 245–270.  
15 doi:10.1080/02671520500193744
- 16 Harwell, M. R., Post, T. R., Maeda, Y., Davis, J. D., Cutler, A. L., Andersen, E., & Kahan, J. A.  
17 (2007). Standards-based mathematics curricula and secondary students' performance on  
18 standardized achievement tests. *Journal for Research in Mathematics Education*, 38, 71–101.
- 19 Helmke, A., & Schrader, F.-W. (1987). Interactional effects of instructional quality and teacher  
20 judgment accuracy on achievement. *Teaching and Teacher Education*, 3, 91–98.  
21 doi:10.1016/0742-051X(87)90010-2
- 22 Herppich, S., Praetorius, A.-K., Förster, N., Glogger-Frey, I., Karst, K., Leutner, D., ... Südkamp, A.  
23 (in press). Teachers' assessment competence: Integrating knowledge-, process-, and product-  
24 oriented approaches into a competence-oriented conceptual model. *Journal of Teaching and*  
25 *Teacher Education*. Advance online publication. <https://doi.org/10.1016/j.tate.2017.12.001>
- 26 Hoge, R. D., & Coladarci, T. (1989). Teacher-based judgments of academic achievement: A review of  
27 literature. *Review of Educational Research*, 59(3), 297–313. doi:10.2307/1170184

- 1 Johansson, S., Myrberg, E., & Rosén, M. (2012). Teachers and tests: Assessing pupils' reading  
2 achievement in primary schools. *Educational Research and Evaluation*, 18, 693–711.  
3 doi:10.1080/13803611.2012.718491
- 4 Karing, C., Matthäi, J., & Artelt, C. (2011). Hängt die diagnostische Kompetenz von  
5 Sekundarstufenlehrkräften mit der Entwicklung der Lesekompetenz und der mathematischen  
6 Kompetenz ihrer Schülerinnen und Schüler zusammen? [Is there a relationship between lower  
7 secondary school teacher judgment accuracy and the development of students' reading and  
8 mathematical competence?] *Journal for Educational Research Online*, 3, 119–147.
- 9 Krajewski, K., Küspert, P., & Schneider, W. (2002). *Deutscher Mathematiktest für erste Klassen*  
10 *(DEMAT 1+)*. Göttingen: Hogrefe.
- 11 Krajewski, K., Liehm, S., & Schneider, W. (2004). *Deutscher Mathematiktest für zweite Klassen*  
12 *(DEMAT 2+)*. Göttingen: Hogrefe.
- 13 Leucht, M., Tiffin-Richards, S., Vock, M., Pant, H. A., & Köller, O. (2012). Diagnostische kompetenz  
14 von Englischlehrkräften bei der bewertung von schülerleistungen mit hilfe des Gemeinsamen  
15 Europäischen Referenzrahmens für Sprachen. [English teachers' diagnostic skills in judging  
16 their students' competencies on the basis of the Common European Framework of Reference.].  
17 *Zeitschrift Für Entwicklungspsychologie und Pädagogische Psychologie*, 44, 163–177.  
18 doi:10.1026/0049-8637/a000071
- 19 Lorenz, C., & Artelt, C. (2009). Fachspezifität und Stabilität Diagnostischer Kompetenz von  
20 Grundschullehrkräften in den Fächern Deutsch und Mathematik. [Domain specificity and  
21 stability of diagnostic competence among primary school teachers in the school subjects of  
22 German and mathematics.]. *Zeitschrift Für Pädagogische Psychologie*, 23, 211–222.  
23 doi:10.1024/1010-0652.23.34.211
- 24 Machts, N., Kaiser, J., Schmidt, F., T. C., & Möller, J. (2016). Accuracy of teachers' judgments of  
25 students' cognitive abilities: A meta-analysis. *Educational Research Preview*, 19, 85–103.  
26 doi:10.1016/j.edurev.2016.06.003
- 27 Praetorius, A., Karst, K., Dickhäuser, O., & Lipowsky, F. (2011). Wie gut schätzen Lehrer die  
28 Fähigkeitsselbstkonzepte ihrer Schüler ein? Zur diagnostischen Kompetenz von Lehrkräften.

- 1 [How teachers rate their students: On teachers' diagnostic competence regarding the academic  
 2 self-concept.]. *Psychologie In Erziehung Und Unterricht*, 58, 81–91.  
 3 doi:10.2378/peu2011.art30d
- 4 Schrader, F.-W. (1989). *Diagnostische Kompetenzen von Lehrern und ihre Bedeutung für die*  
 5 *Gestaltung und Effektivität des Unterrichts* [Diagnostic competencies of teachers and their  
 6 relevance for composition and effectiveness of teaching] (Vol. 289). Frankfurt am Main, Bern,  
 7 NY, Paris: Peter Lang.
- 8 Shavelson, R. J., & Stern, P. (1981). Research on teachers' pedagogical thoughts, judgments,  
 9 decisions, and behavior. *Review of Educational Research*, 51, 455–498. doi:10.2307/1170362
- 10 Sherman, S. J., Beike, D. R., & Ryalls, K. R. (1999). Dual-processing accounts of inconsistencies in  
 11 responses to general versus specific cases. In S. Chaiken & Y. Trope (Eds.), *Dual-process*  
 12 *theories in social psychology* (pp. 203-227). New York: Guilford Press.
- 13 Spinath, B. (2005). Akkuratheit der Einschätzung von Schülermerkmalen durch Lehrer und das  
 14 Konstrukt der diagnostischen Kompetenz. [Accuracy of Teacher Judgments on Student  
 15 Characteristics and the Construct of Diagnostic Competence.]. *Zeitschrift für Pädagogische*  
 16 *Psychologie*, 19, 85–95. doi:10.1024/1010-0652.19.12.85
- 17 Stang, J., & Urhahne, D. (2016). Wie gut schätzen Lehrkräfte Leistung, Konzentration, Arbeits- und  
 18 Sozialverhalten ihrer Schülerinnen und Schüler ein? Ein Beitrag zur diagnostischen Kompetenz  
 19 von Lehrkräften. [How teachers rate students achievement, attention, work habits and social  
 20 behavior? A contribution to the diagnostic competence of teachers.]. *Psychologie In Erziehung*  
 21 *Und Unterricht*, 63, 204–219. doi:10.2378/peu2016.art18d
- 22 Steiger, J. H. (1980). Tests for Comparing Elements of a Correlation Matrix. *Psychological Bulletin*,  
 23 87, 245–251. doi:10.1037/0033-2909.87.2.245
- 24 Südkamp, A., Kaiser, J., & Möller, J. (2012). Accuracy of teachers' judgements of students' academic  
 25 achievement: A meta-analysis. *Journal of Educational Psychology*, 104, 743–762.  
 26 doi:10.1037/a0027627
- 27 Südkamp, A., Möller, J., & Pohlmann, B. (2008). Der simulierte Klassenraum: Eine experimentelle  
 28 Untersuchung zur diagnostischen Kompetenz. [The simulated classroom: An experimental

- 1 study on diagnostic competence.]. *Zeitschrift für Pädagogische Psychologie*, 22, 261–276.  
 2 doi:10.1024/1010-0652.22.34.261
- 3 Tarelli, I., Lankes, E.M., Drossel, K., & Gegenfurtner, A. (2012). Lehr- und Lernbedingungen an  
 4 Grundschulen im internationalen Vergleich [Teaching and learning conditions at primary  
 5 schools in international comparison]. In W. Bos, I. Tarelli, A. Bremerich-Voss, & K.  
 6 Schwippert (Eds.), *IGLU 2011—Lesekompetenzen von Grundschulkindern in Deutschland im*  
 7 *internationalen Vergleich* [IGLU 2011—Reading competencies of primary schoolchildren in  
 8 Germany in international comparison] (pp. 137–173). Münster: Waxmann.
- 9 Weinert, F. E., & Helmke, A. (1995). Interclassroom differences in instructional quality and  
 10 interindividual differences in cognitive development. *Educational Psychologist*, 30, 15–20.  
 11 doi:10.1207/s15326985ep3001\_2
- 12 Wu, M. L., Adams, R. J., & Wilson, M. R. (1998). *ConQuest: Generalized item response modelling*  
 13 *software*. Melbourne: Australian Council for Educational Research.
- 14 Zhu, M., & Urhahne, D. (2014). Assessing teachers' judgements of students' academic motivation and  
 15 emotions across two rating methods. *Educational Research and Evaluation*, 20, 411–427.  
 16 doi:10.1080/13803611.2014.964261  
 17

1

2 Table 1

3 *Characteristics of SJ and GJ integrated into earlier definitions of kinds of judgments*

4

|                   | Direct vs.<br>Indirect | Informed vs.<br>Uninformed | Congruence of<br>Domain Specificity | Specificity <sup>a</sup> |
|-------------------|------------------------|----------------------------|-------------------------------------|--------------------------|
| Specific Judgment | Direct                 | Informed                   | High                                | Level 5                  |
| Global Judgment   | Direct                 | Informed                   | High                                | Level 4                  |

5 *Notes.* <sup>a</sup> Specificity according to the definition of Hoge and Coladarci (1989)

6

1 Table 2

2 *Descriptive results for teacher judgments and student achievement test scores (relative frequencies)*

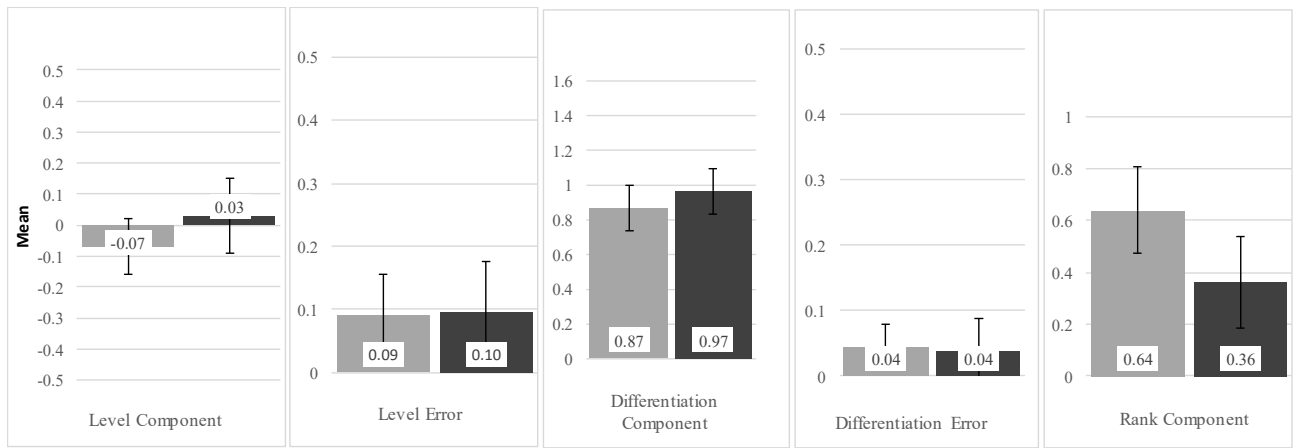
3

|                    | <i>n</i> | <i>M</i> | <i>SD</i> | Range   |
|--------------------|----------|----------|-----------|---------|
| Teacher judgment   |          |          |           |         |
| Global             | 49       | .57      | .26       | .11-.80 |
| Specific           | 52       | .69      | .29       | .38-.96 |
| Student test score |          |          |           |         |
| Global             | 700      | .64      | .29       | .35-.76 |
| Specific           | 434      | .66      | .25       | .41-.80 |

4

5

1



• global judgment accuracy | • specific judgment accuracy

2

3 *Figure 1. Descriptive statistics (M, SD) of the accuracy components depending on kinds of judgment.*

4

5

6

7

8

1 Table 3

2 *Pairwise comparisons for the accuracy measures depending on kind of judgment corrected using*

3 *Bonferroni adjustment*

4

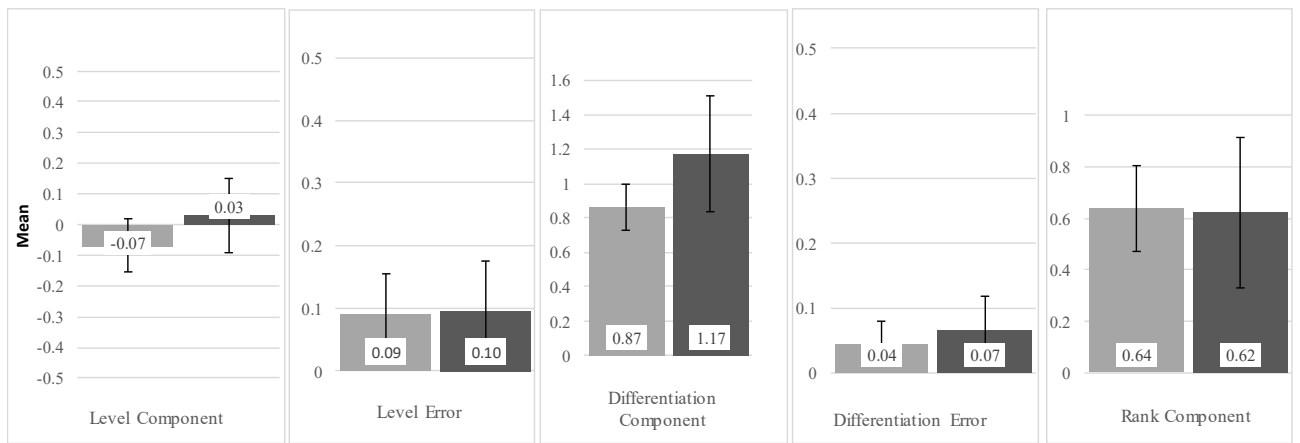
| Accuracy measure          | $\Delta$ | $p$   | $d$  |
|---------------------------|----------|-------|------|
| Level component           | -0.099   | <.001 | .926 |
| Level error               | -0.005   | .743  | .067 |
| Differentiation component | 0.100    | .001  | .748 |
| Differentiation error     | 0.007    | .472  | .157 |
|                           | $Z$      | $p$   | $q$  |
| Rank component            | 2.254    | .024  | .376 |

5 *Notes.*  $\Delta$  = difference between GJ-accuracy and SJ-accuracy;  $d$  = Cohen's  $d$ ;  $q$  = Cohen's  $q$ .

6



1



• global judgment accuracy | • aggregated global judgment accuracy

2

3 *Figure 2. Descriptive statistics (M, SD) of the accuracy measures depending on GJ and aggregated GJ.*

4

1 Table 4

2 *Pairwise comparisons for the accuracy measures depending on original GJ and aggregated GJ*

3 *corrected using Bonferroni adjustment*

4

| Accuracy measure          | $\Delta$ | $p$   | $d$   |
|---------------------------|----------|-------|-------|
| Level component           | -0.099   | <.001 | 0.926 |
| Level error               | -0.005   | .743  | 0.067 |
| Differentiation component | 0.308    | <.001 | 1.199 |
| Differentiation error     | -0.020   | .044  | 0.456 |
|                           | $Z$      | $p$   | $q$   |
| Rank component            | 0.201    | .841  | .026  |

5

6 *Notes.*  $\Delta$  = difference between GJ-accuracy and aggregated GJ-accuracy;  $d$  = Cohen's  $d$ ;

7  $q$  = Cohen's  $q$ .

8

9

1 Table 5

2 *Correlation between the accuracy of aggregated GJ and original GJ (N = 49)*

3

|                       | Original (Global) |                 |      |
|-----------------------|-------------------|-----------------|------|
|                       | Level             | Differentiation | Rank |
| Aggregated (Global)   |                   |                 |      |
| Level Error           | .00               |                 |      |
| Differentiation Error |                   | -.21            |      |
| Rank Component        |                   |                 | .21  |

4

5

6