

A COMPARISON OF MORAL REASONING STAGES USING A MODEL OF HIERARCHICAL COMPLEXITY

TERRI LEE ROBINETT
Scottsdale, Arizona, USA

Empirical evidence has demonstrated the validity and reliability of moral development instruments such as the Defining Issues Test (DIT) and Moral Judgment Test (MJT). Rasch item reliability for each was .95. A newer instrument generated using the Model of Hierarchical Complexity had item reliability of .97. Rasch scores of responses to each instruments' items correlated well with the items' measure of hierarchical complexity, $r = .286, .372, .557, .767$. Test items used to measure moral reasoning were significantly correlated. In general, stage of moral development did not predict political affiliation or voting, supporting Kohlberg's claim that structure, not content, underlies stage.

KEYWORDS: Defining Issues Test, Model of Hierarchical Complexity, moral development, Moral Judgment Test, Rasch.

This study had several different but related purposes. First, it was an important developmental stage validation study as it utilized a mathematical model of hierarchical complexity of tasks to predict Rasch scored stages of moral reasoning. This is in contrast to the cognitive approach of interpreting stage from a participant's performance on a test. Second, this study utilized items from multiple measures of moral reasoning: items from existing tests and from new tests based on the Model of Hierarchical Complexity. All of the items were scored using the Hierarchical Complexity Scoring System (Commons, Miller, Goodheart, and Danaher-Gilpin, 2005) rather than traditional standard scoring systems.

Empirical evidence has repeatedly demonstrated the validity and reliability of moral development instruments such as James Rest's (1975) Defining Issues Test (DIT) and Georg Lind's (1985) Moral Judgment Test (MJT) in measuring Kohlberg's (1984) construct of moral reasoning. Proponents of the DIT cite numerous studies by Crowson (2002), Davidson (1979), Rest, Thoma, and Edwards (1997), Thoma (2002), Thoma, Narvaez, Rest, and Derryberry (1999), and Walker (2002) among others that support the reliability and validity of the DIT as a measure of moral development. Similarly, based on important validation criteria and numerous worldwide validation studies, Lind (2004) has concluded that the MJT is a valid measure of moral judgment competence and moral attitude.

Address correspondence to Terri Lee Robinett, SAP America, 4343 N. Scottsdale Road, Scottsdale, AZ 85251, USA. E-mail: tleerobinett@yahoo.com

The use of a variety of dilemma-based moral measures has consistently found correlations between moral reasoning stages and political beliefs and attitudes. For example, Fishkin, Keniston, and MacKinnon (1973) conducted a study with 75 undergraduate students in which the participants were given Kohlberg's Moral Judgment Interview (MJI), as well as measures of political ideology. Participants who reasoned at the conventional stage of moral reasoning (stages 9 or 10) tended to be politically conservative, whereas those who scored at the postformal, that is, postconventional, stage tended to reject conservative ideas in favor of liberal views. Their results demonstrated that Kohlberg's theory of moral development does identify a cognitive-developmental dimension of personality with a high correlation with political ideology. Alker and Poppen (1973) also examined the results of Kohlberg's MJI, which they had administered to 192 students at Cornell University. They found a strong correlation between liberal ideology and the choice of principled moral thinking. Likewise, a closed belief system correlated with lower-stage moral reasoning.

These findings might cast doubt on cognitive-developmental theory, suggesting that the stages represented in the instruments are not content free, but reflect a bias for liberal democratic norms (Gross, 1996). According to Emler, Renwick, and Malone (1983), the moral dilemmas used in these tests are typically composed of liberal values, such as civil rights, abortion, the right to die, and the death penalty, as well as the conflict between individual conscience and authority. They pointed out that if liberalism is in fact highly correlated with postformal's higher-principled thinking, then it makes sense that relationships will occur between the moral reasoning stage of the individual and his or her political leanings.

Another reason studies found differences in moral reasoning between liberals and conservatives may have to do with the tests themselves. Although the moral reasoning measures used in these studies rest on Kohlberg's (1984) cognitive-developmental moral stage theory, they use different testing methods and scoring procedures.

METHODOLOGY

A total of 163 participants, largely from a California community college, took 5 online instruments, including items from moral reasoning instruments. These included the Defining Issues Test-2, the Moral Judgment Items, the Politician-Voter Problem, and the Right to Bear Arms test. A demographic questionnaire was completed that was designed to collect data regarding the independent and sample variables.

These assessments of moral reasoning have been plagued by fundamental problems in measuring stages of difficulty because they are based on comparisons of performances that depend on content and context (see Day, this issue). This has led to problems with reliability and validity when comparing performances on tests. Therefore, the Hierarchical Complexity Scoring System was used to score all of the test items, including those that were not developed using the Model of Hierarchical Complexity. This was done in order to counter objections mentioned above to the subjective arbitrariness of existing measures of moral

stage. Each item was assigned an order of hierarchical complexity by examining the complexity of the tasks required to answer the item. The reliability of doing this varies depending on what method was used to construct the items. The items that were a short paragraph long that followed all the item construction rules described in the scoring manual were the most reliable. The shorter the items, the less reliable they were. The items' Rasch scores were regressed against item order of hierarchical complexity to determine if the items were measuring the same underlying construct. Rasch scores of a participant's successful performance on an item of a given order of complexity represented the person's stage of performance when completing the test item.

Essentially, hierarchical complexity replaces performance-based measures of moral reasoning with task-based measures. This allows for an objective measure of moral reasoning to which performance can be related (Commons and Robinett, 2006). Therefore, rather than measuring one's stage of development based on subjective criteria of test performance, stage of moral reasoning was determined objectively by scaling participants' performance in completing moral tasks correctly.

In this study, the Rasch analyses were performed using Winsteps, a Rasch model software package (Linacre and Wright, 2000). Winsteps provides reliability information for participants and items, as well as a unique scaled score for each individual and test item. Human cognitions and behaviors are extremely complex and can never really be satisfactorily expressed by one score on any test, or with any scoring system. However, the Rasch analysis provides key information regarding unidimensionality, construct validity, difficulty and ability estimation and error, and reliability. Construct validity information focuses on the idea that performances reflect a single underlying construct; in this case, order of hierarchical complexity.

The Rasch item fit statistics are indicators of how well each item fits within the underlying construct. Linacre and Wright (2000) developed a criterion for rejecting items with infit errors larger than 2.00. They suggested that it is possible that items with an infit score greater than 2.00 have characteristics that are sensitive to issues not reflective of the scale. They may not have fit because they are too extreme for the scale or because they lie on another dimension. In this study, those items that diverged significantly from the expected pattern as indicated by the infit and outfit statistics were removed from the data set.

Likert scale data provided interval stage estimates of participant responses. Standard dichotomous models for analyzing Likert data makes simple right and wrong distinctions for each response whereas the Rasch analysis is applied to polytomous data that establishes the relative difficulty of each item from the lowest to the highest stages of the items. Therefore, the Rasch model tests the hypothesis that each item reflects increasing stages of an attitude or trait, as intended (Bond and Fox, 2001). The rest of instruments had no significant correlation with affiliation or voting.

RESULTS

Overall results indicated that with a few specific exceptions, order of moral stage did not predict political affiliation. Only two of the instruments moderately

predicted political affiliation. These were the right to bear arms content ($r(102) = .118$, $F(1,102) = 4.182$, $p \leq .043$), and the anti-right to die items were found to be significant predictors of political identity ($r(103) = .106$, $F(1,103) = 8.540$, $p \leq .004$).

Predicting Rasch Scores from Hierarchical Complexity of the Items

The four figures that follow illustrate the relationship between the Rasch scaled scores and the hierarchical complexity of the items that predict them. The DIT and MJ item hierarchical complexity orders were not checked by the item makers. The range of person scores was attenuated because almost all the participants came from a community college. There were almost no higher stage participants.

Figure 1 shows an $r(58) = -.286$. The r may be small for a number of reasons. There may be an inadequate range of scores due to a limited sample.

Figure 2 shows an $r(58) = .372$. The items received a subject reliability of .81, and item reliability of .95, indicating very good test and item reliability.

Figure 3 depicts the scattergram results of the moral judgment items. Once again concrete responses cluster around $-.30$, with a range of $.10$ to $.60$, whereas the clustering of metasystematic items are around $.10$, with a very limited range from 0 to $-.20$. However, a pattern emerges indicating a clustering of scores from high to low; concrete to metasystematic. The moral judgment items received a person reliability index of .66, and an item reliability of .95, indicating low test reliability, but high item reliability.

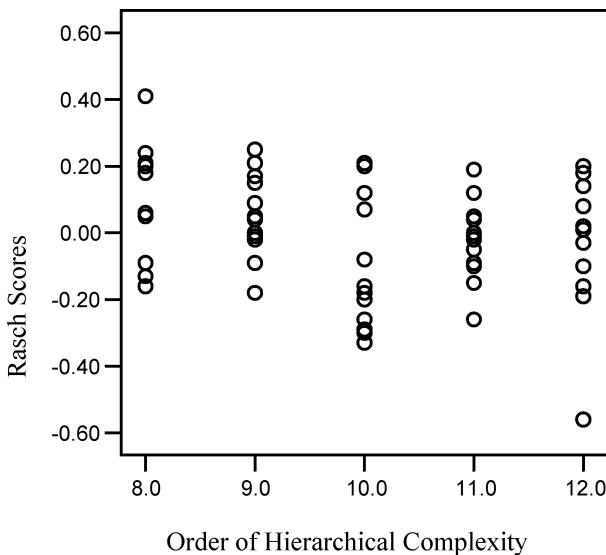


Figure 1. Regression of Rasch Scores on order of hierarchical complexity for the Right to Bear Arms items ($r(58) = -.286$, $F(1,58) = 5.174$, $p \leq .027$, $r^2 = .082$).

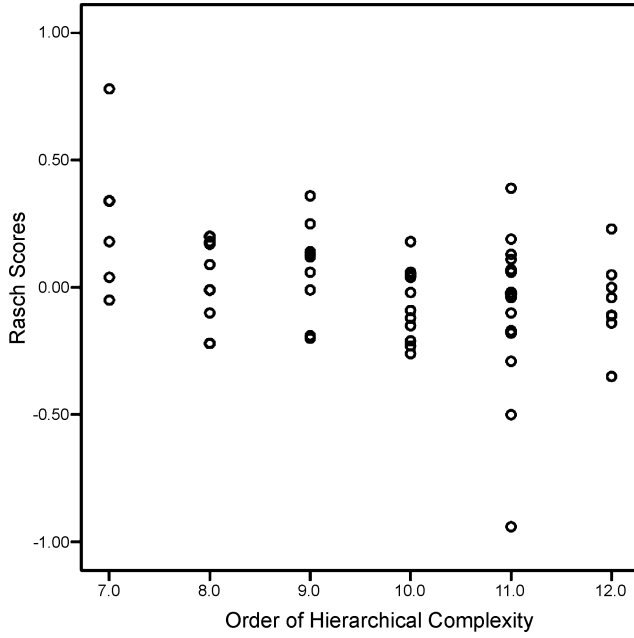


Figure 2. Regression of order of hierarchical complexity versus Defining Issues Rasch scores ($r(58) = .372$, $F(1,58) = 9.333$, $p \leq .003$, $r^2 = .139$).

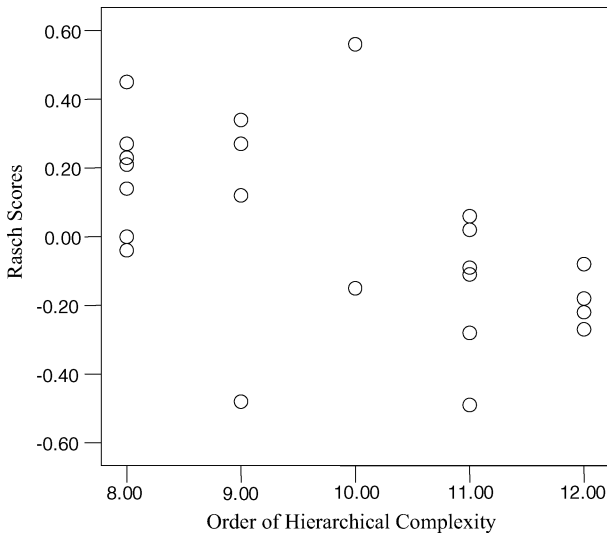


Figure 3. Regression of Rasch Scores on order of hierarchical complexity for the Moral Judgment items ($r(22) = .557$, $F(1,22) = 9.912$, $p \leq .005$, $r^2 = .311$).

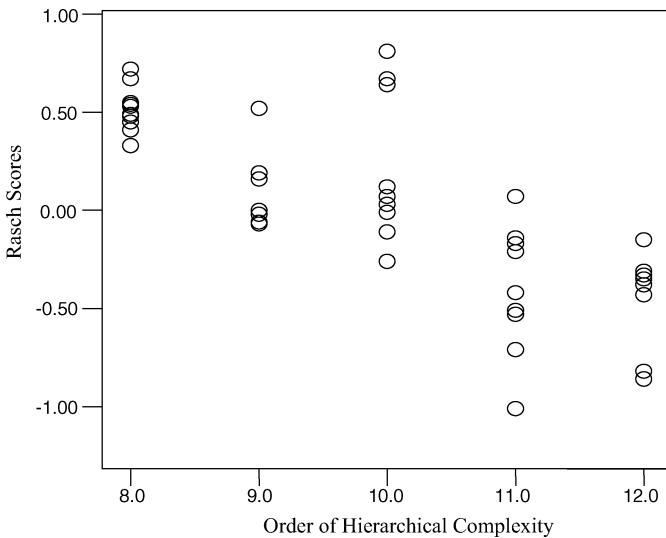


Figure 4. Regression of Rasch Scores on order of hierarchical complexity for the Politician-Voter items ($r(43) = .767$, $F(1,43) = 61.331$, $p \leq .000$, $r^2 = .588$).

Figure 4 depicts scattergram results for the politician–voter items, $r(43) = .767$. The items of the Politician–Voter Problem vignettes rated a .74 person reliability score, and a .97 item reliability score, indicating an extremely good overall test reliability.

The higher the r , the more valid the instrument. That is because these instruments are supposed to reflect hierarchical complexity of the items.

Correlational Analysis

A correlational analysis was performed on the Rasch scores of content area test items. The results are presented in Table 1. Significant correlations were found among all of the grouped items except the defining issues items and the politician–voter items.

Correlations were also analyzed from the different content areas of the moral reasoning measures. The results are shown in Table 2.

No significant correlations were found between the defining issues content areas and the politician–voter items. However, there were significant correlations between each of the five defining issues content areas. This may indicate that the defining issues items and politician–voter items are measuring different constructs. Interestingly, none of the defining issues content areas correlated with the anti-arms items, the pro-worker, or the pro–right to die items. These content areas would typically be considered politically liberal. There was, however, a significant positive correlation between the defining issues and the anti-worker, and anti–right to die content areas. Both of these content areas (anti-worker, and anti–right to die) would typically be categorized as conservative political positions. This means

Table 1
Content Area Rasch Score Correlations

	School Board	Reporter Duty	Famine to Die	Right to Die	Demonstration	Pro-Arms	Anti-Arms	Pro-Worker	Anti-Worker	Pro-Right to Die	Anti-Right to Die	Pol.-Voter
School Board	1	.570**	.526**	.451**	.490**	.109	.023	.069	.270**	.053	.074	.100
Reporter Duty		1	.492**	.502**	.483**	.206*	.110	.151	.145	-.034	.321**	.040
Famine to Die			1	.371**	.452**	.142	.147	.103	.248**	.023	.106	.085
Right to Die				1	.470**	.107	.056	.002	.281**	-.125	.331**	.155
Demonstration					1	.063	.203	.015	.189*	-.127	.230*	.071
Pro-Arms						1	-.187	-.003	.062	.195	.093	.325**
Anti-Arms							1	.207	.047	-.104	-.029	.067
Pro-Worker								1	-.217*	.338**	.121	.207*
Anti-Worker									1	-.003	.265**	.267**
Pro-Right to Die										1	-.340**	.389**
Anti-Right to Die											1	-.012
Pol.-Voter												1

Table 2
Correlations among Grouped Items of the Measures of Moral Reasoning

	Politician-Voter	Moral Judgment	Defining Issues	Right to Bear Arms
Politician-Voter	1	.451**	.108	.393**
Moral Judgment		1	.268**	.499**
Defining Issues			1	.279**
Right to Bear Arms				1

* $p \leq .05$, ** $p \leq .01$.

that when participants rated an anti-worker, or anti-right to die argument as low on the Likert scale (1 = no) their Rasch score increased (toward the negative end of the scale), and vice versa. This could be evidence that some defining issues content areas are biased toward liberal responses.

The pro-right to die, and the anti-right to die items were negatively correlated ($r = -.340$, $p \leq .01$). The pro-worker and anti-worker items were moderately negatively correlated ($r = -.217$, $p \leq .05$). This suggests that bias plays a role as suggested by LaLave (2006). The pro-arms and anti-arms were not correlated. Due to the pro and con relationship of these content areas a significant negative correlation would only be predicted by bias. If content did not matter, as was the case in the right to bare arms, there would not have been a negative correlation but a positive one.

Factor Analysis

A factor analysis was performed on all four grouped moral reasoning items to examine the interrelationships among the variables and to explain these variables in terms of their common underlying dimensions. As expected, all of the grouped items loaded significantly on the first component, or the moral reasoning factor. A Kaiser, Meyer, and Olkin (KMO) test was conducted to determine if the items were measuring a common factor as suggested. The KMO value was .693 indicating that moral reasoning accounts for a fair amount of variance, but not a substantial amount. A factor analysis was also performed on the specific content areas of each of the moral reasoning measures (Table 3).

The defining issues content areas, the anti-right to die, and anti-worker content areas loaded significantly on the primary moral reasoning factor. The politician-voter, and the pro-right to die items loaded significantly on component two, and the pro- and anti-right to bear arms, the pro-worker, and the anti-worker all loaded significantly on component three. Although the pro-worker and pro-right to arms loaded on component three they also were significantly correlated with component two. The anti-worker content area was not only correlated with component three, but also with the moral reasoning factor. The KMO value was .636, which indicates a mediocre degree of common variance. KMO assess the degree of multicollinearity. There is a KMO statistic for each individual variable, and their sum is the KMO overall statistic. KMO varies from 0 to 1.0 and KMO

Table 3
Factor Analysis of Moral Measure Content Areas

Content Area	Component		
	1	2	3
Defining Issues Reporter	.784	-.021	.120
Defining Issues School	.751	.071	.047
Defining Issues Right to Die	.734	-.123	-.102
Defining Issues Demonstration	.729	-.151	.152
Defining Issues Famine	.711	.067	.140
Anti-Right to Die	.420	-.340	-.182
Pro-Right to Die	-.060	.851	.042
Politician-Voter	.235	.690	-.212
Anti-Right to Arms	.177	-.057	.646
Pro-Worker	.120	.484	.598
Pro-Right to Arms	.248	.445	-.486
Anti-Worker	.441	-.025	-.460

Note: Principal Component Analysis. 3 Components extracted.

overall should be .60 or higher to proceed with factor analysis. Certainly these results make it clear that some of the specific content areas from these instruments are measuring very different constructs.

CONCLUSIONS

Overall results indicated that with a few specific exceptions, moral stage did not predict political affiliation. Findings did support that the test items were measuring moral reasoning stages. Education-stage and household income were found to be significant predictors of political affiliation, supporting the findings of Emler, Renwick, and Malone (1983). Stage of religiosity was correlated with and found to be a significant predictor of one's identification as a liberal or a conservative. These results are in keeping with Kohlberg's notion that moral stage is not about content of judgment but the structure of the cognitive process. The results also indicated that the test items used to measure moral reasoning in this study were significant. They loaded on the first factor of moral stage. The correlations may have been somewhat low because the restricted range of the participant stages.

Based on these results it appears that when the moral reasoning items are grouped as separate measures, they are all measuring the same underlying variable—that is, moral reasoning stage.

It is not clear what the second component is measuring. The politician-voter, and the moral judgment pro-right to die items loaded significantly on it, but the items do not appear to have much in common. Even when rotated, both the politician-voter and pro-right to die load significantly on component two. The pro-right to die items typically represent a liberal attitude and therefore might

be expected to load on component three but did not do so. The politician–voter problem is a hierarchically complex design concerning how well politicians inform their constituency about a proposed method for resolving a community issue. It does not appear to have either liberal or conservative attitudinal content. It is possible that component three measures some other psychological or social value that is not readily apparent.

Major reasons for the finding of no role of moral stage in predicting political behavior in results between this study and previous studies are twofold. First, this study did not use the traditional moral reasoning tests, or scoring procedures that were used in previous studies. The participants completed the standard DIT-2. But the raw data was not sent to Minneapolis for scoring. Instead, specific story items were extracted from the DIT-2, and assigned an order of hierarchical complexity from Concrete stage 8 to Metasystematic stage 12. Then the raw data were transformed into a linear Rasch scale. Hierarchical complexity was used to base decisions regarding participant performance, making the scoring procedure more objective. Previous studies based results of a participant's performance using subjective scoring criteria, whereas hierarchical complexity is based on mathematical principles. The downside to using a non-standard scoring system was the issue of reliability and validity, which could not simply be assumed, or provided by the test authors. However, the end result is hopefully a more objective approach to the study of moral reasoning.

Secondly, by using the Rasch Model for data analysis, results were not determined by analyzing raw data, or counts, but by constructing objective, additive scales. According to Bond and Fox (2001), the only way objectively to construct scales that are separable from the distribution and the attributes they measure is to use the Rasch Model that has become popular in the social sciences. It allows one objectively to examine the processes underlying why people and items behave in a particular way rather than simply how a person performed on a particular item. This is of primary importance in the measurement of moral reasoning because it eliminates the possibility of biased and subjective scoring of participant responses to particular test items.

Certainly other variables may have played a part in these results. Both item and participant reliability and validity were generally good. The item reliabilities were very high, .95. However, some items demonstrated unacceptable reliability. These issues are discussed later with the moral reasoning measures. It was also discovered that the participant sample was constricted in response ranges because of the community college sample. Rasch models need a wide range of responses from low to high in order to construct a proper scale (Bond and Fox, 2001). This sample did not provide enough responses at the systematic and metasystematic stages. As a result, the scales were slightly skewed toward the lower ordered Rasch scores.

These results are not entirely surprising from a Kohlbergian point of view. The core of Kohlberg's cognitive-developmental position was that cognitive stages are qualitative differences in modes of thinking that form an invariant sequence, and that each of these sequential modes of thought form a structured whole. Therefore, a given stage–response on a task does not just represent

a specific response determined by knowledge of that task, but actually represents an underlying organization of thought (Kohlberg, 1984). Previous studies have relied on the identification of particular conceptual content, via moral measures such as the DIT, rather than the direct identification of the underlying thought structures. Therefore, the relationship between stage and content may be confounded so that stage is defined in terms of that content rather than the structures that form the basis of cognitive developmental theory (Dawson, 2000).

As Gross (1996) pointed out, because Kohlberg's theory focuses on the structure rather than the content of moral reasoning, one would not expect that the reasoning of liberals and conservatives would necessarily be different. Results to the contrary had cast doubt on the structural integrity of cognitive development theory, suggesting that moral stages reflect a bias in liberal democratic norms. The idea that individual differences in adult moral reasoning were actually a reflection of politico-moral ideology was also the position of Emler et al. (1983). Their claim was that moral and political attitudes are overlapping domains, and stage differences between liberals and conservatives are merely that of ideological content rather than structural complexity.

Results of this study generally found no significant differences between an individual's moral stage and their self-reported political affiliation. These results suggest that there is no liberal bias in cognitive developmental theory, or even in the sets of grouped items that make up the moral reasoning measures, but rather there may be bias in the scoring systems that are currently used to determine moral stage, thus finding differences that do not really exist.

Further research is needed to compare the results of traditional moral reasoning scoring systems with the hierarchical complexity model used in this study. Hierarchical complexity may prove to be a valuable tool in objectively measuring individual differences in other realms of social science. Some of the instruments developed here need to be improved. Only the politician-voter instrument worked perfectly.

REFERENCES

- Alker, H. A., and Poppen, P. J. 1973. Personality and ideology in university students. *Journal of Personality* 41(4): 652-671.
- Bond, T. G., and Fox, C. M. 2001. *Applying the Rasch Model: Fundamental measurement in the human sciences*. Mahway, NJ: Lawrence Erlbaum.
- Commons, M. L., Miller, P. M., Goodheart, E. A., and Danaher-Gilpin, D. 2005. Hierarchical Complexity Scoring System (HCSS): How to score anything. Unpublished manual. Dare Institute, Cambridge, MA.
- , and Robinett, T. L. 2006. *Moral measures based on hierarchical complexity*. Unpublished manuscript, Dare Institute, Cambridge, MA.
- Crowson, M. H. 2002. *Is the defining issues test a measure of moral judgment development: A test of competing claims?* PhD dissertation, University of Alabama.
- Davidson, M. L. 1979. The internal structure and psychometric properties of the Defining Issues Test. In *Development in judging moral issues*, Ed. Rest, J., 223-245. Minneapolis: University of Minnesota Press.

- Dawson, T. 2000. Layers of structure: A comparison of two approaches to development assessment. *Genetic Epistemologist* 29(4): 2–14.
- Emler, N., Renwick, S., and Malone, B. 1983. The relationship between moral reasoning and political orientation. *Journal of Personality and Social Psychology* 45(5): 1073–1080.
- Fishkin, J., Keniston, K., and MacKinnon, C. 1973. Moral reasoning and political ideology. *Journal of Personality and Social Psychology* 27(1): 109–119.
- Gross, M. L. 1996. Moral reasoning and ideological affiliation: A cross-national study. *Political Psychology* 17(2): 317–338.
- Kohlberg, L. 1984. *The psychology of moral development: The nature and validity of moral stages*. Vol. 2. San Francisco: Harper and Row.
- LaLlave, J. 2006. Moral judgment competence and attitude as moderators of decisions concerning war, through preferences of frames and arguments on the Iraq war. Dissertation Zur Erlangung des Doktorgrades der Naturwissenschaften (Dr. rer nat) Fachbereich Psychologie, Universität Konstanz. <http://w3.ub.uni-konstanz.de/v13/volltexte/2006/1883/pdf/LaLlaveDiss.pdf>.
- Linacre, J. M., and Wright, B. D. 2000. *Winsteps: Multiple-choice, rating scale, and partial credit Rasch analysis*. Computer software. Chicago: MESA Press.
- Lind, G. 1985. The theory of moral-cognitive development: A socio-psychological assessment. Trans Wren, T. E. In *Moral development and the social environment*, Eds. Lind, G., Hartmann, H. A., and Wackenhut, R., 21–53. Chicago: Precedent.
- Lind, G. 2004. *The meaning and measurement of moral judgment competence: A dual aspect model* [Electronic Version]. Unpublished manuscript, University of Konstanz, Germany. <http://www.uni-konstanz.de/ag-moral/pdf/Lind-2005>.
- Rest, J. R. 1975. Longitudinal study of the Defining Issues Test of moral judgment: A strategy for analyzing developmental change. *Developmental Psychology* (11): 738–748.
- , Thoma, S., and Edwards, L. 1997. Designing and validating a measure of moral judgment: Stage preference and stage consistency approaches. *Journal of Educational Psychology* 89(1): 5–28.
- Thoma, S. J. 2002. An overview of the Minnesota approach to research in moral development. *Journal of Moral Education* 31(3): 225–245.
- , Narvaez, D., Rest, J., and Derryberry, P. 1999. Does moral judgment development reduce to political attitudes or verbal ability? Evidence using the Defining Issues Test. *Educational Psychology Review* 11(4): 325–341.
- Walker, L. J. 2002. The model and the measure: An appraisal of the Minnesota approach to moral development. *Journal of Moral Education* 31(3): 353–367.