# Automated Essay Scoring and NAPLAN: A Summary Report

Les Perelman, Ph.D.

1 October 2017

This summary report is written in response to proposals for employing an Automated Essay Scoring (AES) system to mark NAPLAN essays, either as the sole marker or in conjunction with separate scores from a human marker. Specifically, this summary will address assertions regarding AES's appropriateness made in *An Evaluation of Automated Scoring of NAPLAN Persuasive Writing* (ACARA NASOP Research Team, 2015) [henceforth referred to as *The Report*]. After describing the primary strategies AES systems use to compute scores of writing ability and the major studies of the efficacy of AES for high-stakes assessments, various critiques of AES are discussed. Finally, an analysis of *The Report* concludes that both its review of the literature and the study described in it are so methodologically flawed and so massively incomplete that it cannot justify any use of AES in scoring the NAPLAN essays.

### How AES Works

All AES systems analyse only textual features that can be represented and manipulated mathematically (Zhang, 2013). AES, from its beginnings in the 1960's (Page, 1966) relies heavily on the use of proxies that can be easily counted. It cannot directly measure a student's adept use of vocabulary. Instead, it often just calculates the number of infrequently used words in a text (Attali & Burstein, 2006; Page, 1966). Because it cannot actually comprehend how well a topic is developed in a paragraph, it determines development by counting the number of sentences in each paragraph (Attali & Burstein, 2006; Burstein, Marcu, & Knight, 2003). And just counting the number of commas has been successfully used in helping to calculate an overall score of an essay that will match that of human readers (Bennett & Zhang, 2016; Simon, 2012).

The other methods used by AES systems consist of various natural language processing techniques. All of these techniques work by statistically identifying key words in a text and analysing their frequency, often in relation to other words. E-rater's natural language technique begins with the assumption that some of the words in high-scoring essays have a high probability of occurring in other high-scoring essays, and similarly, most low-scoring essays will contain a subset of words associated with low scores. It then employs statistical techniques based on the vocabulary in an essay to determine the essay's score category as well as the relation of the essay's vocabulary to that of the highest scoring essays (Attali & Burstein, 2006). Some techniques, such as Latent Semantic Analysis, create matrices based on single words and, like erater, ignore word order (Foltz, Streeter, Lochbaum, & Landauer, 2013; Landauer, Foltz, & Laham, 1998). Many AES systems, such as ETS's e-rater, use a hybrid approach that combines proxies with other machine learning and natural language processing techniques.

# Efficacy of AES

Given our current linguistic and computational knowledge, does AES work? There is already some indication that in some cases—such as writing in response to open-ended prompts, in which students have wide latitude in direction and creativity—AES cannot replicate human markers (McCurry, 2010). The most ambitious research study is the Hewlett ASAP study referenced by *The Report*. Although the Hewlett Study is not in any way seminal, it was extremely ambitious, using a total of 22,029 student essays based on eight different writing prompts from six U.S. state tests. These essays were divided into a Training Set, a Test Set, and a Validation Set. The Hewlett Study Report exists in three forms: the original conference paper (Shermis & Hamner, 2012), a version that appeared in a collection of essays co-edited by the paper's first author (Shermis & Hamner, 2013), and a single-authored article that appeared in a peer-reviewed journal and that concluded with a fairly lengthy list of the study's limitations (Shermis, 2014a). [Full disclosure: I am on the editorial board of the journal.] Curiously, *The Report* references only the first two versions, ignoring the more authoritative peer-reviewed article, which is qualified in its endorsement of AES.

### Strengths of the Hewlett Study

One unfortunate limitation of the study was that the agreement with the vendors prohibited the research group from conducting any statistical tests comparing the vendor and human marker scores (Bennett & Zhang, 2016; Rivard, 2013). However, the study report (in all three versions) was thorough in presenting demographic statistics for each of the U.S. states participating in the study as well as statistics in two general categories:

- **Descriptive statistics** such as the number (N), mean, and standard deviation (STD) on each essay set for human markers and all nine vendors.
- **Measures of agreement** such as percentage of exact agreement, percentage of exact plus adjacent agreement, Cohen's kappa, Quadratic-weighted kappa, and the Pearson product-moment correlation coefficient.

The research team also subsequently released the raw scores on the Test Set for seven of the nine vendors for confirmation and analysis. Two vendors did not want their data made public even though the sets were anonymous.

### Limitations and Critiques of the Hewlett Study

The Hewlett Study results were released with much fanfare. The University of Akron reported

A direct comparison between human graders and software designed to score student essays achieved virtually identical levels of accuracy, with the software in some cases proving to be more reliable, a groundbreaking study has found. ("Man and machine: Better writers, better grades," 2012)

Yet close analysis of the data casts doubt on that claim as well as raises questions about major methodological elements of the study:

- The data do not support the claim that machines were able to match human readers. Indeed, analyses of the specific data tables indicate that humans possessed higher levels of accuracy than machines (Bennett, 2015; Bennett & Zhang, 2016; Perelman, 2013, 2014). The exhaustive analysis of Bennett (ETS's Norman O. Frederiksen Chair in Assessment Innovation) and Zhang (2016), in particular, refutes any claim that the AES scores in the Hewlett Study matched the reliability of human readers.
- Five of the eight data sets consisted of paragraphs not essays, with mean lengths of 99–173 words (Shermis, 2014a; Shermis & Hamner, 2012, 2013).
- The four essay sets in which the machines performed best (Sets 3, 4, 5, and 6)
  - were not marked on writing ability but solely on content;
  - had reliability assessed using the higher of the two human markers' scores, producing different scoring formulas for machines and humans, which made any comparison problematic and privileged machines (Bennett, 2015; Bennett & Zhang, 2016; Perelman, 2013, 2014). The importance of this last assertion, however, has been contested (Shermis, 2014b).
- Only two of the eight essay sets in the study employed, like NAPLAN, a composite score based on a combination of analytic scores. The machines performed poorly in comparison to humans for these sets (Shermis, 2014a; Shermis & Hamner, 2012, 2013)

# Critiques of AES

One major failing of *The Report* is that it completely ignores the significant body of scholarship critical of various applications of AES. The focus here will be on those objections that are the most relevant to NAPLAN. For a more complete listing of some excellent collections of essays on AES see Appendix A.

## Lack of Rhetorical Situation

One of the most common objections is that writing is communication, the transfer of thoughts from one mind to another. As various scholars have noted, AES creates a non-rhetorical situation (Anson, 2006; Condon, 2006, 2013; Ericsson, 2006; Herrington & Moran, 2001, 2012). Students are writing not to inform, entertain, or persuade another mind; they are writing to an entity that can only count. In essence, the audience has been replaced by a machine. Even in cases in which there is both a human and a machine marking the essay, the student will be aware that half the score is coming from an entity that does not understand meaning but is simply looking for specific elements. Students then have a dual audience; they must produce a text that will satisfy the machine, even if a human reader is also present.

## Reductive

Because AES is solely mathematical, it cannot assess the most important elements of a text. The following paragraph is not written by critics of AES but by its developers, including three very

senior individuals at the Educational Testing Service and four vice presidents at Pearson Education and Pearson Knowledge Technologies:

Automated essay scoring systems do not measure all of the dimensions considered important in academic instruction. Most automated scoring components target aspects of grammar, usage, mechanics, spelling, and vocabulary. Therefore, they are generally well-positioned to score essays that are intended to measure text-production skills. Many current systems also evaluate the semantic content of essays, their relevance to the prompt, and aspects of organization and flow. Assessment of creativity, poetry, irony, or other more artistic uses of writing is beyond such systems. They also are not good at assessing rhetorical voice, the logic of an argument, the extent to which particular concepts are accurately described, or whether specific ideas presented in the essay are well founded. Some of these limitations arise from the fact that human scoring of complex processes like essay writing depend, in part, on "holistic" judgments involving multivariate and highly interacting factors. This is reflected in the common use of holistic judgments in human essay scoring, where they may be more reliable than combinations of analytic scores. (Williamson et al., 2010 p. 2)

This passage makes two points extremely relevant to the use of AES in marking NAPLAN. First, AES cannot assess some of the key criteria addressed by the NAPLAN writing test, such as audience, ideas, and persuasive devices (i.e. the logic of an argument). Second, AES is more reliable providing a single holistic score rather than the sum of analytic scores, such as the ten trait scores of the NAPLAN. This second point is supported by how the essay portions of two high-stakes American tests, the new SAT Essay and the Analytical Writing Essays of the Graduate Record Examination (GRE), are marked. The new SAT Essay is marked on three analytic categories, which are not combined but reported separately. The analytic scores are produced by two human markers (College Board, 2017). The GRE Essays, on the other hand, are evaluated by a single holistic score for each essay and are marked both by a machine and by a human (Educational Testing Service, 2017).

## Weaknesses in Grammatical Analysis

The above passage from AES developers, like similar claims (Deane, 2013), assumes that AES systems are precise in identifying grammatical errors. However, anyone who has ever used a grammar checker suspects that this is not the case. English grammar, like the grammar of any natural human language, is extremely complex and interdependent on such factors as meaning and context. AES grammar checkers miss many grammatical errors (False Negatives), while classifying perfectly grammatical constructions as errors (False Positives). When analyzing 5,000 words of an essay by Noam Chomsky originally published in *The New York Review of Books*, the grammar checker modules of ETS's e-rater identified 62 grammatical or usage errors, including 15 article errors and 5 preposition errors (Perelman, 2016). None of them were actually errors.<sup>1</sup> In addition, AES grammar checkers often focus on grammatical non-problems, such as beginning a sentence with a coordinating conjunction, possibly because such constructions are very easy for a machine to identify.

<sup>&</sup>lt;sup>1</sup> All of the examples are from ETS's e-rater simply because other vendors no longer allow academic researchers access. A Pearson vice president responded to a reporter's request to allow me access to the Intelligent Essay Assessor by refusing and stating, "He wants to show why it doesn't work" (Winerip, 2012).

One of the most complex linguistic features of English is the set of rules governing the use of articles; these rules are especially challenging for speakers of languages such as Mandarin or Russian that do not have articles. Computational linguistic models of English article use are disappointing. One model, for example, deployed in 2005, could detect 80% of article errors with a False Positive rate of approximately 50% or detect only 40% of article errors but reduce the False Positives to 10% (Han, Chodorow, & Leacock, 2006). A comparison of error identification by two instructors and e-rater 2.0 of 42 English Language Learners' papers demonstrated that e-rater is extremely inaccurate in identifying the types of major errors made by ELL, bilingual, and bidialectical students. The instructors coded 118 instances of missing or extra articles; e-rater marked 76 instances, but 31 of those (40.8%) were either False Positives or misidentified (Dikli & Bleyle, 2014). The current inability to develop reliable grammar checkers is best exemplified by the decision of Microsoft Research, one of the largest software companies in the world, to discontinue its ESL Assistant Project (Gamon, 2011). AES is inaccurate and unreliable at assessing even low-level writing traits such as grammatical correctness.

### Fairness

Related to grammar is the issue of fairness. Do AES machines treat all linguistic, national, and ethnic groups the same? Two reports by the Educational Testing Service (Bridgeman, Trapani, & Attali, 2012; Ramineni, Trapani, Williamson, Davey, & Bridgeman, 2012) indicate that in the essay portions of both the Test of English as a Foreign Language and the GRE, the e-rater scoring engine gave significantly higher marks to native Mandarin speakers, especially those from mainland China, than did human markers. In some instances, the difference between the machine score and human was very large, close to 0.40 of a standard deviation. Conversely, in some instances, African-Americans, particularly males, were given significantly lower marks by e-rater than they were by human markers. Another study reported that Vantage Technology's ACCUPLACER, which has an essay section scored by the IntelliMetric scoring engine, underpredicted portfolio and final course grades for African-American and Hispanic students (Elliot, Deess, Rudniy, & Joshi, 2012).

Possibly, the unevenness of the grammatical components of the scoring engines contributes to the machines' under- and overreporting marks. Native Mandarin speakers and native speakers of other languages that do not have articles make more errors in the use of English articles than speakers of languages that employ articles. Because grammar detectors perform so poorly in correctly identifying English article usage, they may be contributing to the machines' inflating the scores of Mandarin speakers. One prominent feature of African-American dialects of English is a difference in verb constructions. These constructions are easy for a machine to identify and may be overcounted in comparison to the response of a human marker. Another possible explanation is that people from mainland China receive extensive coaching for these tests and may be including memorized passages that appear more relevant to a machine than they do to a human marker (Bridgeman et al., 2012).

Whatever the explanation, unfairness by machines in inflating the marks of some linguistic groups and artificially lowering the marks of others is morally indefensible and, possibly, illegal. Before any AES system is deployed, extensive research is needed to ensure that the machines do not penalize or privilege specific linguistic communities.

### Construct-Irrelevant Response Strategies (Gaming)

Because AES relies so heavily on proxies in marking, various studies have shown that AES machines are extremely vulnerable to construct-irrelevant response strategies, that is, providing the machine with the proxies it employs without actually displaying the traits of good writing that they are supposed to represent.

For most AES machines, the strongest single proxy is length (Perelman, 2012, 2014). As noted previously in the discussion on fairness, it appears that tutors in mainland China have students memorize sentences that they then insert in essays to increase their score (Bridgeman et al., 2012). Although ETS is attempting to develop tools to catch such gaming strategies (Bejar, Vanwinkle, Madnani, Lewis, & Steier, 2013), they appear still to be effective (Bejar, Flor, Futagi, & Ramineni, 2014; Powers, Burstein, Chodorow, Fowles, & Kukich, 2001).

Perhaps the most theatrical example of the vulnerability of AES systems to gaming strategies is the BABEL Generator developed by the author and three undergraduates from Harvard and the Massachusetts Institute of Technology (Kolowich, 2014). Just by randomly creating nonsense sentences with long, rarely used words and occasionally peppered with synonyms of at most three topic words, the BABEL Generator is able to create essays that receive high scores from AES machines such as e-rater and Vantage Technology's IntelliMetric. Two pairs of top scoring, BABEL-written GRE essays along with a link to the BABEL Generator are displayed in Appendix B.

The main danger, however, is not from absurd machines such as the Babel Generator, but from the implications of such stumping studies. That which is tested will be taught. If wordy essays with long sentences and obscure vocabulary will produce high scores on high-stakes tests, that is what teachers will be emphasizing. Rather than improve the writing ability of students, AES may well encourage the production of verbose, high-scoring gibberish.

# Inaccuracies, Methodological Flaws, Incomplete Information, and Anomalies in *An Evaluation of Automated Scoring of NAPLAN Persuasive Writing*

The flaws in *The Report* and the study it describes are so major that it cannot justify any use of AES in high-stakes testing situations.

### Inaccuracies

The most egregious mistake in *The Report* is in the account of the Hewlett competition on page 5: "The rate of agreement was higher between any of the automated scoring engines and human markers than that between the two human markers." Even a cursory examination of the data in any of the three papers reporting on the study reveals the gross inaccuracy of this statement (Shermis, 2014; Shermis & Hamner, 2013). As Bennett and Zhang (2016) demonstrated, humans actually performed more reliably. The most vivid refutation of this claim can be made by comparing the human–human reliability to the human (resolved score)–machine reliability for each of the metrics for each of the essay sets and for just one scoring engine, MetaMetrics's

Lexile Writing Analyser. Table 1 displays this comparison. Rather than being more reliable than the human markers, Lexile is substantially less reliable for every metric and essay set except for two of the metrics for Essay Set 8 (shaded). Lexile was chosen for several reasons. First, its performance was the poorest of any of the scoring engines. Second, it is one of the four engines used in the study described in *The Report*. Finally, unlike the other engines, Lexile is not trained for a specific prompt but, instead, measures a general trait, text complexity (*The Report*, p. 7).

Essay Sets	Exact Agreement		Kappa		Quadratic-Weighted Kappa		Correlation Pearson <i>r</i>	
	H1 - H2	Lexile	H1 – H2	Lexile	H1 – H2	Lexile	H1 – H2	Lexile
1	0.64	0.31	0.45	0.16	0.73	0.66	0.73	0.66
2A	0.76	0.55	0.62	0.30	0.80	0.62	0.80	0.62
2B	0.73	0.55	0.56	0.27	0.76	0.55	0.76	0.55
3	0.72	0.63	0.57	0.45	0.77	0.65	0.77	0.65
4	0.76	0.47	0.65	0.30	0.85	0.67	0.85	0.68
5	0.59	0.47	0.44	0.28	0.74	0.64	0.75	0.65
6	0.63	0.51	0.45	0.31	0.74	0.65	0.74	0.66
7	0.28	0.07	0.18	0.03	0.72	0.58	0.72	0.58
8	0.29	0.08	0.16	0.04	0.61	0.63	0.61	0.62

 Table 1: Comparison of Agreement Metrics Between the Two Human Markers (H-H) and
 Between MetaMetrics's Lexile Writing Analyser and Human Markers

Source: Shermis, 2014a, Tables 7, 9, 10, and 11

Another major problem is the citation of Attali (2013). Attali does indeed offer practical advice on validity in writing assessment. The advice he offers, however, is contrary to the conclusions of *The Report*. He argues that AES is severely limited and cannot assess several of the NAPLAN traits. He states,

we believe that a serious consideration of the construct argument against AES should lead one to accept its basic premise—because the machine is not able to read the essay, it will not be able to assess such aspects as the quality of argumentation or the development of characters in a narrative, as human readers do. . . . We believe that AES should be based on an alternative definition of its intended use. Specifically, it should be constructed primarily as a *complement* to (instead of a replacement for) human scoring, *limited* in its ability to measure a subset of the writing construct. (p. 182)

*The Report* also contains problems in terminology. Attali employs the term *construct* correctly. At its conclusion, however, *The Report* defines *construct validity* in this passage "ACARA will examine if the introduction of automated scoring has an effect on the substance and quality of student writing ('construct validity')" (p. 14). Construct validity is a complex and evolving concept. At its core, however, is the key concept that the measure is representing the abstract ability (the construct) that it is claiming to assess. Thus "the substance and quality of student writing" is the construct. The question is whether AES can faithfully measure it, not whether AES can affect it.

Another problem with terminology is the misuse of the term *lexical*. The term is correctly defined in footnote 2 on page 4. On the following page, however, the "lexical properties of essays" are listed as "sentence structure, paragraphing, punctuation and spelling." These elements of writing have little or anything to do with the term *lexical*.

A final problem with language is the use of the term *cognitive interview*. Since this term in all Anglophone countries usually refers to a specific technique used in forensic investigations (Davis, McMahon, & Greenwood, 2005), it is extremely unclear what *cognitive interview* means in this context.

## Methodological Flaws and Incomplete Information

While the inaccuracies in the report were disconcerting, it is the study's very flawed methodology accompanied by a consistent lack of definition and detail that make *The Report* inappropriate in justifying any decision to employ AES in marking the NAPLAN.

### A Convenient Sample Defines a Pilot Study

The method section of *The Report* states "A single persuasive prompt was administered to a convenient sample of year 3, 5, 7 and 9 students as part of a larger online assessment study" (p. 6). Major studies, especially those with national consequences usually employ a *representative* sample, or, if it is a large, broadly-based sample, possibly a *random* sample. In research, convenience sampling is limited to pilot studies because of the risk of sampling errors. The Discussion section of *The Report* makes it clear that this study is a pilot and that there will be larger follow-up studies: "ACARA will next expand its research to include larger samples of students and multiple prompts within and across writing genres that NAPLAN assesses (persuasive and narrative)" (p. 13). The plan for future research is also explicitly stated in the August 13, 2016 ACARA web page on Automated Essay Scoring:

More research is planned for 2016 which will include a larger sample of students, multiple prompts within and across writing genres (persuasive and narrative) and key validity questions—does the use of AES affect features of student writing and writing instruction—to inform a recommendation to Education Ministers about the approach to be used in 2017.

The <u>current version of the web page</u> omits any reference to larger follow-up studies. There is no explanation of why ACARA never undertook these crucial additional projects.

Both versions, however, claim that the sample was "broad," although there is no attempt to show that the sample was representative of the national population. Indeed, the Test Set consisted of

339 essays. If they were evenly divided among Years, they would consist of only 110 essays for three Years and 109 for one Year.

The Method section reports the mean essay length and median raw scores by Year. These numbers appear to be for all three sets—Training, Validation, and Test Sets—although that is not certain. There is no explanation why the mean is given for essay length and the median for raw score. There also needs to be much more supporting data. The means of the Test Set for essay length and for each trait score should have been provided, along with the standard deviations for each. These numbers then needed to be compared with national statistics to ensure that the sample was representative.

Moreover, with such a small sample size, it is impossible to determine if any of the AES machines gave higher or lower scores to members of specific linguistic or ethnic groups than the scores given by human markers. Finally, there has been no evaluation of machines evaluating narrative essays.

Even more troubling is that this pilot was based on a testing format different from that currently used for the NAPLAN essay. There are now separate prompts for Years 3 & 5 and for Years 7 & 9. There has been no attempt to assess how well the machines perform on the different prompts for these two groups. The mean statistics for essay length alone indicates that length alone clearly differentiates them. Will separating these two groups make scoring more difficult for machines? This crucial question remains unanswered.

One very bizarre aspect of the study's methodology is allowing each vendor to report its results differently. The Hewlett Study, which is referred to as "seminal," correctly reported all vendor data homogeneously. Why were vendors in this study allowed to choose how they would present their data? Why are all the presentations different? Was there a deliberate attempt to avoid comparisons?

There is also some uncertainty about exactly when the vendors received the marks of the human scorers for the Test Set. On page 7, *The Report* first states that "Contractors were not provided with any marking data for these essays." At the bottom of the same page, however, it states, "Vendors completed the scoring and provided ACARA with a research report outlining the methods used in their investigation and its key outcomes." Vendors needed the Test Set marking data before they wrote a research report that included outcomes. Did they first provide ACARA with a dataset of their scores before they received the human scores? If so, this fact should have been stated explicitly.

There is also too much reliance on undefined and vague hearsay evidence. At the beginning, *The Report* states,

Markers who scored the essays observed that student responses were at least as long, on average and of comparable quality, as those produced in paper-based tests. Even at Year 3, student lack of typing ability was not found to be a barrier to completing the task. (p. 3)

Comparing word counts of the sample to national word counts for each Year would have provided a much more accurate assessment of the effect of a computer-based test on text

production. Similarly, a statistical comparison of total scores and trait scores could verify the markers' impressions with hard data. There is the statement, "psychometric analyses confirmed that the underlying writing scales performed in a similar manner to their paper-based analogues." However, that is the only reference to those analyses, which, along with supporting data, should have been an integral and substantial part of the document.

The Report also states,

Invigilator observations and follow-up discussions ("cognitive interviews") with students confirmed that students were able to complete the writing task within the allotted time, without being unduly constrained by level of keyboarding skill. (p. 3)

Although probably not "cognitive interviews," it is clear that interviews did take place. What were the exact questions asked? Was there an interview protocol? In addition, it is difficult to believe that all students reported that they were not "unduly constrained." Were there some complaints? If so, how many? What was their nature?

### Anomalies

As mentioned, all four vendors were part of the Hewlett Competition. As also stated previously, the Lexile Writing Analyser was the poorest performer in the Hewlett Competition and it employs a generic algorithm that does not consider the specific prompt or topic. In many of the metrics, its performance was especially dismal for Essay Sets 7 and 8, the only sets in the Hewlett Competition that, like NAPLAN, employ analytical scales. Yet the Quadratic-weighted kappas in Table 3 of *The Report* indicate that Lexile performed extremely well in its ratings for Audience and Ideas, even though it did not know or consider the specific writing task, prompt, or question being posed. The <u>current web site</u> states "Of especial significance, the AES systems were even able to match human markers on the 'creative' rubric criteria: *audience* and *ideas*." That the machine was able to evaluate the quality of an answer to a question without knowing the question is indeed of special significance.

Moreover, although the labelling in Table 3 is unclear (and appears to include references to an Excel spreadsheet [Columns AM through AX; Columns C through Y]), it seems that the Quadratic-weighted kappa comparing Lexile results with the human marks is either 0.8828 or 0.9190. However, neither number matches those of the AES machines displayed in Table 5. There may be an explanation for these differences, but if it exists, it needs to be made explicit.

# Conclusion

Even some of the strongest proponents and developers of AES have conceded that it cannot assess high-level traits such as quality and clarity of ideas. These traits comprise the focus and reason for human communication. They need to be assessed and assessed well. The pilot study described in *The Report*, with its large amounts of hearsay evidence, extremely dubious methodology, and incorrect information, cannot justify any sort of national implementation. Before any kind of AES system is deployed either as a sole marker or in dual markings with humans, a number of issues need to be addressed:

• Evidence needs to be provided that the correct constructs are being measured by the machines. As Mark D. Shermis (2014a), the principal investigator of the Hewlett Competition, writes in the final, peer-reviewed version of his study,

A predictive model may do a good job of matching human scoring behaviour, but for reasons unrelated (or unsatisfactorily related) to the construct of interest. If accurate predictions of score are achieved by features and methods that do not bear any plausible relationship to the competencies and construct that the item aims to assess, then this prediction, accurate as it may be, is not sufficiently representative of the construct to warrant test use. (p. 74)

In particular, given the relative recent unanimity among AES developers and critics of AES that the machines are incapable of reliably assessing high-level constructs, substantial evidence must be provided that machines are capable of evaluating such constructs. Without such proof, machine scoring may produce situations in which teachers, to protect themselves and their schools, spend significant time teaching students strategies to "game" the machines with construct-irrelevant strategies that will improve their scores but make their writing less effective. Possibly, independent investigators should be allowed to test the construct relevance of the machine through various types of Reverse Turing Tests and Stumping Studies.

- Given the research findings in the United States that at least one AES machine appears to overscore one linguistic group and underscore another, no AES system should be deployed until extensive pilot testing has demonstrated that AES does not discriminate against any linguistic group or groups.
- ACARA needs to provide substantial evidence, more than the poorly designed and executed pilot study, to demonstrate that AES, which has been developed primarily to generate holistic scores, can reliably score ten analytic traits.
- *The Report* and the original language on the ACARA web site stated that more extensive studies would be conducted, including ones involving the marking of narrative prompts. Given that a narrative prompt has recently been used on the NAPLAN, it is imperative that ACARA conduct studies to demonstrate that the AES systems are capable of effectively scoring the ten trait categories of the NAPLAN narrative essay.
- As noted above, the NAPLAN has changed significantly since the 2012 sample used in the pilot. There are now separate prompts (and probably separate scoring) for Years 3 & 5 and for Years 7 & 9. This change creates an entirely different scoring situation. ACARA needs to conduct pilots demonstrating that the AES machines are capable of accurately scoring these two separate groups with two different prompts.
- ACARA needs to assess the technical and keyboard capabilities of all students, including Third Year students and students from disadvantaged backgrounds, before deploying an online essay test. If text production among these groups is hindered by lack of keyboarding or technical skills, online assessment should not be deployed.

• Finally, there should be considerably more transparency and independence in these necessary research studies than was demonstrated in *The Report*. Preferably, independent investigators should constitute part of the research team.

Until these critical studies are completed and carefully evaluated, it would be extremely foolish and possibly damaging to student learning to institute machine grading of the NAPLAN essay, including dual grading by a machine and a human marker.

ACARA NASOP Research Team. (2015). *An evaluation of automated scoring of NAPLAN persuasive writing*. Retrieved from

 $http://nap.edu.au/\_resources/20151130\_ACARA\_research\_paper\_on\_online\_automated\_scoring.pdf$ 

Anson, C. (2006). Can't touch this: Reflections on the servitude of computers as readers. In P. Ericsson & R. H. Haswell (Eds.), *Machine scoring of human essays* (pp. 38–56). Logan, UT: Utah State University Press. Retrieved from

http://digitalcommons.usu.edu/cgi/viewcontent.cgi?article=1138&context=usupress\_pubs

- Attali, Y. (2013). Validity and reliability in automated essay scoring. In M. D. Shermis & J. Burstein (Eds.), *Handbook of automated essay evaluation: Current applications and new directions* (pp. 181–198). New York, NY: Routledge.
- Attali, Y., & Burstein, J. (2006). Automated essay scoring with e-rater® V.2. *The Journal of Technology, Learning and Assessment, 4*(3).
- Bejar, I. I., Flor, M., Futagi, Y., & Ramineni, C. (2014). On the vulnerability of automated scoring to construct-irrelevant response strategies (CIRS): An illustration. Assessing Writing, 22, 48–59. Retrieved from http://www.sciencedirect.com/science/article/pii/S1075293514000257

Bejar, I. I., Vanwinkle, W., Madnani, N., Lewis, W., & Steier, M. (2013). Length of textual response as a construct-irrelevant response strategy: The case of shell language. Princeton NJ. Retrieved from http://origin-www.ets.org/Media/Research/pdf/RR-13-07.pdf

Bennett, R. E. (2015). The changing nature of educational assessment. *Review of Research in Education*, 39(1), 370–407. doi:10.3102/0091732X14554179

Bennett, R. E., & Zhang, M. (2016). Validity and automated scoring. In *Technology in testing: Improving educational and psychological measurement* (pp. 142–173). Washington, DC: National Council on Measurement in Education.

Bridgeman, B., Trapani, C., & Attali, Y. (2012). Comparison of human and machine scoring of essays: Differences by gender, ethnicity, and country. *Applied Measurement in Education*, 25(1), 27–40. doi:10.1080/08957347.2012.635502

Burstein, J., Marcu, D., & Knight, K. (2003). Finding the WRITE stuff: Automatic identification of discourse structure in student essays. *IEEE Intelligent Systems*, *18*(1), 32–39. Retrieved from http://people.cs.pitt.edu/~huynv/research/argument-mining/Finding the WRITE stuff Automatic identification of discourse structure in student essays.pdf

College Board. (2017). SAT essay scoring. Retrieved September 15, 2017, from https://collegereadiness.collegeboard.org/sat/scores/understanding-scores/essay

Condon, W. (2006). Why less is not more: What we lose by letting a computer score writing samples. In P. F. Ericsson & R. H. Haswell (Eds.), *Machine scoring of human essays: Truth or consequences* (pp. 211–220). Logan, UT: Utah State University Press. Retrieved from http://digitalcommons.usu.edu/cgi/viewcontent.cgi?article=1138&context=usupress\_pubs

Condon, W. (2013). Large-scale assessment, locally-developed measures, and automated scoring of essays: Fishing for red herrings? *Assessing Writing*, *18*(1), 100–108. Retrieved from http://www.sciencedirect.com/science/article/pii/S1075293512000505

- Davis, M. R., McMahon, M., & Greenwood, K. M. (2005). The efficacy of mnemonic components of the cognitive interview: Towards a shortened variant for time-critical investigations. *Applied Cognitive Psychology*, 19(1), 75–93. Retrieved from https://www.researchgate.net/profile/Marilyn\_Mcmahon/publication/216569762\_The\_effic acy\_of\_mnemonic\_components\_of\_the\_cognitive\_interview\_Towards\_a\_shortened\_variant \_for\_time-critical\_investigations/links/0046353a0f3494eed3000000/The-efficacy-ofmnemonic-components-of-the-cognitive-interview-Towards-a-shortened-variant-for-timecritical-investigations.pdf
- Deane, P. (2013). On the relation between automated essay scoring and modern views of the writing construct. *Assessing Writing*, *18*(1), 7–24. Retrieved from http://www.sciencedirect.com/science/article/pii/S1075293512000451
- Dikli, S., & Bleyle, S. (2014). Automated essay scoring feedback for second language writers: How does it compare to instructor feedback? *Assessing Writing*, 22, 1–17. Retrieved from http://www.sciencedirect.com/science/article/pii/S1075293514000221
- Educational Testing Service. (2017). How the GRE tests are scored. Retrieved September 15, 2017, from https://www.ets.org/gre/institutions/scores/how/
- Elliot, N., Deess, P., Rudniy, A., & Joshi, K. (2012). Placement of students into first-year writing courses. *Research in the Teaching of English*, *46*(3). 285-313. Retrieved from http://www.ncte.org/journals/rte/issues/v46-3
- Ericsson, P. F. (2006). The meaning of meaning: Is a paragraph more than an equation? In P. F.
  Ericcson & R. H. Hasswell (Eds.), *Machine scoring of human essays: Truth or Consequences* (pp. 28–37). Logan, UT: Utah State University Press. Retrieved from
  http://digitalcommons.usu.edu/cgi/viewcontent.cgi?article=1138&context=usupress pubs
- Foltz, P. W., Streeter, L. A., Lochbaum, K. E., & Landauer, T. K. (2013). Implementation and applications of the Intelligent Essay Assessor. In M. D. Shermis & J. Burstein (Eds.), *Handbook of automated essay evaluation: Current applications and new directions* (pp. 68–88). New York, NY: Routledge.
- Gamon, M. (2011). ESL Assistant discontinued. Retrieved September 20, 2017, from https://blogs.msdn.microsoft.com/eslassistant/
- Han, N.-R., Chodorow, M., & Leacock, C. (2006). Detecting errors in English article usage by non-native speakers. *Natural Language Engineering*, 12(2), 115. Retrieved from https://s3.amazonaws.com/academia.edu.documents/30237428/nle06hcl.pdf?AWSAccessKeyId=AKIAIWOWYYGZ2Y53UL3A&Expires=1506719442&Signa ture=gYAxom2SLU%2FD9aeZTKfMfFBGg4g%3D&response-contentdisposition=inline%3B%20filename%3DDetecting errors in English article usag.pdf
- Herrington, A., & Moran, C. (2001). What happens when machines read our students' writing? *College English*. Retrieved from http://www.jstor.org/stable/378891
- Herrington, A., & Moran, C. (2012). Writing to a machine is not writing at all. In N. Elliot & L. Perelman (Eds.), *Writing assessment in the 21st century: Essays in honor of Edward M. White* (pp. 219–232). New York, NY: Hampton Press.
- Kolowich, S. (2014, April 28). Writing instructor, skeptical of automated grading, pits machine vs. machine. *The Chronicle of Higher Education*. Retrieved from http://www.chronicle.com/article/Writing-Instructor-Skeptical/146211
- Landauer, T. K., Foltz, P. W., & Laham, D. (1998). An introduction to Latent Semantic Analysis. *Discourse Processes*, 25, 259–284. Retrieved from http://lsa.colorado.edu/papers/dp1.LSAintro.pdf

Man and machine: Better writers, better grades. (2012, April 12). *University of Akron News*. Akron, OH. Retrieved from http://www.uakron.edu/im/online-newsroom/news\_details.dot?newsId=40920394-9e62-415d-b038-15fe2e72a677

- McCurry, D. (2010). Can machine scoring deal with broad and open writing tests as well as human readers? *Assessing Writing*, *15*(2), 118–129. Retrieved from http://www.sciencedirect.com/science/article/pii/S1075293510000218
- Page, E. B. (1966). The imminence of grading essays by computer. *The Phi Delta Kappan*. Phi Delta Kappa International. Retrieved from http://www.jstor.org/stable/20371545
- Perelman, L. (2012). Construct validity, length, score, and time in holistically graded writing assessments: The case against automated essay scoring (AES). In A. Bazerman, C; Dean, C; Early, J; Lunsford, K; Null, S; Rogers, P; Stansell (Ed.), *International advances in writing research* (pp. 121–131). Fort Collins, CO: The WAC Clearinghouse and Parlor Press. Retrieved from https://wac.colostate.edu/books/wrab2011/chapter7.pdf
- Perelman, L. (2013). Critique of Mark D. Shermis & Ben Hamner: "Contrasting state-of-the-art automated scoring of essays: Analysis." *The Journal of Writing Assessment*, 6(1). Retrieved from http://journalofwritingassessment.org/article.php?article=69
- Perelman, L. (2014). When "the state of the art" is counting words. *Assessing Writing*, *21*, 104–111. Retrieved from http://www.sciencedirect.com/science/article/pii/S1075293514000233
- Perelman, L. (2016). Grammar checkers do not work. WLN: A Journal of Writing Center Scholarship, 40(7–8), 11–20. Retrieved from http://lesperelman.com/wpcontent/uploads/2016/05/Perelman-Grammar-Checkers-Do-Not-Work.pdf
- Powers, D. E., Burstein, J., Chodorow, M., Fowles, M. E., & Kukich, K. (2001). Stumping erater: Challenging the validity of automated essay scoring. Princeton NJ: Educational Testing Service. Retrieved from https://www.ets.org/Media/Research/pdf/RR-01-03-Powers.pdf
- Ramineni, C., Trapani, C. S., Williamson, D. M., Davey, T., & Bridgeman, B. (2012). *Evaluation of the e-rater<sup>®</sup> scoring engine for the GRE<sup>®</sup> issue and argument prompts ETS RR--12-02*. Retrieved from https://www.ets.org/Media/Research/pdf/RR-12-02.pdf
- Rivard, R. (2013, March 15). Humans fight over robo-readers. *Inside Higher Education*. Retrieved from https://www.insidehighered.com/news/2013/03/15/professors-odds-machine-graded-essays
- Shermis, M. D. (2014a). State-of-the-art automated essay scoring: Competition, results, and future directions from a United States demonstration. *Assessing Writing*, 20, 53–76. Retrieved from https://assets.documentcloud.org/documents/1094637/shermis-aw-final.pdf
- Shermis, M. D. (2014b). The challenges of emulating human behavior in writing assessment. *Assessing Writing*, 22, 91–99. Retrieved from http://www.sciencedirect.com/science/article/pii/S1075293514000373
- Shermis, M. D., & Hamner, B. (2012). Contrasting state-of-the-art automated scoring of essays: Analysis. Retrieved August 28, 2017, from https://web.archive.org/web/20150810190434/www.scoreright.org/NCME\_2012\_Paper3\_2 9 12.pdf
- Shermis, M. D., & Hamner, B. (2013). Contrasting state-of-the-art automated scoring of essays.
   In M. D. Shermis & J. Burstein (Eds.), *Handbook of automated essay evaluation: Current applications and new directions* (pp. 313–353). New York: Routledge.
- Simon, S. (2012, March 2). Robo-readers: The new teachers' helper in the U.S. *Reuters*. Retrieved from http://www.reuters.com/article/us-usa-schools-grading-idUSBRE82S0ZN20120329

Williamson, D. M., Bennett, R. E., Lazer, S., Bernstein, J., Foltz, P. W., Landauer, T. K., ... Way, W. D. (2010). Automated scoring for the assessment of Common Core standards. Retrieved from https://www.ets.org/s/commonassessments/pdf/AutomatedScoringAssessCommonCoreStan dards.pdf

Winerip, M. (2012, April 23). Facing a robo-grader? No worries. Just keep obfuscating mellifluously. *New York Times*. Retrieved from http://www.nytimes.com/2012/04/23/education/robo-readers-used-to-grade-test-essays.html

Zhang, M. (2013). Contrasting automated and human scoring of essays (R&D Connections No. 21). Princeton NJ: Educational Testing Service. Retrieved from https://www.ets.org/Media/Research/pdf/RD Connections 21.pdf

# Appendix A

# Annotated Bibliography of Collected Materials on Automated Essay Scoring

- Elliot, N., Ruggles Gere, A., Gibson, G., Toth, C., Whithaus, C., & Presswood, A. (2013). Uses and limitations of automated writing evaluation software. *WPA-CompPile Research Bibliographies*. WPA-CompPile. Retrieved from http://comppile.org/wpa/bibliographies/Bib23/AutoWritingEvaluation.pdf [An excellent and broadly-based selective bibliography on AES.]
- Elliot, N., & Williamson, D. M. (Eds.). (2013). Assessing Writing special issue: Assessing writing with automated scoring systems. *Assessing Writing*, 18(1).[An excellent collection of essays of Automated Essay Scoring from a variety of perspectives.]
- Ericsson, P. F., & Haswell, R. H. (2006). *Machine scoring of student essays : Truth and consequences*. Logan, UT: Utah State University Press. Retrieved from http://digitalcommons.usu.edu/cgi/viewcontent.cgi?article=1138&context=usupress\_pubs [The only collection on AES by teachers of writing. The essays are mostly critical of AES but many are regarded as perceptive even by AES developers.]
- National Council of Teachers of English. (2013). NCTE position statement on machine scoring. Retrieved August 29, 2017, from http://www.ncte.org/positions/statements/machine\_scoring [A position statement against AES accompanied by an extensive and very broad annotated bibliography.]
- Professionals Against Machine Scoring of Student Essays in High-Stakes Assessment. (2013).
  Research Findings. Retrieved August 29, 2017, from http://humanreaders.org/petition/research\_findings.htm
  [A passionate argument against AES. The accompanying petition was signed by such luminaries as Noam Chomsky.]
- Shermis, M. D., & Burstein, J. (Eds.). (2003). Automated essay scoring: A cross-disciplinary perspective. Mahwah, NJ: Erlbaum.
  [Although claiming to be cross-disciplinary, the volume is largely composed of individuals involved in the development of AES. Still, the volume contains some of the clearest explanations of AES systems.]
- Shermis, M. D., & Burstein, J. (Eds.). (2013). *Handbook of automated essay evaluation: Current applications and new directions*. New York, NY: Routledge.
  [Similar to the earlier volume, although some of the authors have taken a much more cautious view of the abilities of AES systems.]

# Appendix B

# BABEL Generated Essays Graded by e-rater for Practice GRE

GRE.

SIGN OUT CONTACT US

# ScoreltNow!"

ScoreltNow!<sup>™</sup> Home > Test Experience > Previous Tests > Review

ScorekNow! <sup>™</sup> Home
Overview of Service
List of Topics
Scoring Guides
Score Level Descriptions
Advice to Writers
Frequently Asked Questions
Buy Additional Service
Update Account

# Analyze an Issue Topic:

Review

In most professions and academic fields, imagination is more important than knowledge.

Write a response in which you discuss the extent to which you agree or disagree with the daim. In developing and supporting your position, be sure to address the most compelling reasons and/or examples that could be used to challenge your position

#### Your Answer

Careers with corroboration has not, and in all likelihood never will be compassionate, gratuitous, and disciplinary. Mankind will always proclaim noesis; many for a trope but a few on executioner. a quantity of vocation lies in the study of reality as well as the area of semantics. Why is imaginativeness so pulverous to happenstance? The reply to this query is that knowledge is vehemently and boisterously contemporary. Benevolence, usually by controntation, might enthral career. If nearly all of the sanctions assure a concession of

the swiftly or tantalizingly enthusiastic rejoinder, the consummate cognition can be more multifariously provoked. Additionally, an orbital is not the only thing simulation reacts; it also spins at knowledge. Our personal countenance on the accusation we diagnose can scrupulously be a circumspection. Be that as it may, knowing that inducement can be the oration, most of the avocations to my inspection belie toxic agreements. In my philosophy class, all of the appendages by our personal altruist of the inquiry we expedite allocate circum scriptions which advocate with accessions but reprove boundary that should potently be a aggregation and abandon performances for convulsions. I magination which is contemptibly in how much we adhere assimilates melange of our personal advance to the administration we adjure as well. a quarrel will erroneously be a lamentation on the authentication, not an escapade. In my experience, none of the salvers by our personal amplification at the affirmation we authorize masticate consideration that journeys but in dine, an abundance of vision changes plethora for careers. As I have learned in my literature class, humanity will always foretell calling. Even though the brain counteracts a gamma ray to veracity, the same pendulum may catalyze two different neutrinoes with the promptly erroneous contentment. Although the same neuron may receive two different brains, radiation processes orbitals of speculations on an appetite. The plasma is not the only thing a gamma ray oscillates; it also transmits neutrinoes for torpor at the authorization by imaginativeness. The assassination of imagination changes a plethora of calling. The less extraneous respondents articulate precincts, the more an explanation amplifies those in question Irreverence, normally on the propagandist, exhibits career. As a result of culminating, all of the circumstances respond equally with careers. Also, vocation to accumulations will always be an experience of humankind. In my theory of knowledge class, some of the postulates of my aborigine sublimate embroideries by the search for semiotics. Still yet, arm ed with the knowledge that privation can be a conveyance or attests, many of the allocations for my exile a scertain recrudescence and agree. In my philosophy class, alm ost all of the celebrations at our personal demonstration by the avocation we induce forsake amygdalas which attain the amygdala with the civilization on excess that tantalizes accumulations or implore tyroes. Imagination which pledges subjugation may rivetingly be provocation or is fattering but not speculating of my appendage also. a situational augur feigns the people involved, not resourcefulness. Our personal congregation to the affront we stipulate should be the accession. The sedulously despicable imagination changes a quantity of noesis. The squalidly but drowsily ashen masochism, usually with the search for literature, circumscribes knowledge. Noesis which will effectively be ar axiom changes a frugal knowledge. Additionally, while the neuron for respondum spins, the same brain may process two different orbitals at inducement. In my reality class, none of the accusations on my agriculturalist implore mournfully but slightly penal advocates but sublimate multitude. Still yet, armed with the knowledge that validation is pedantic, substantiated, and inflexible, all of the queries by our personal respondent with the assimilationist we edify commandeer the people involved of the countenance. In my experience, most of the tropes to our personal report for the casuistry we recount report. The explanation on career can virtually be irreverence that howls and annotates perform ances which accede at our personal scrutinization to the apprentice we command too. an authentication may be scrupulousness, not ligature by orators. Our personal organism with the attruist we embolden laments state-of-affairs that is prototypical yet somehow fecund. By the fact that all of the probes are bemoaned for imagination, disrupting interlopers culminate to the same extent on cognition Knowledge has not, and undoubtedly never will be scrupulous yet somehow agreed. However, armed with the knowledge that an amanuensis with exposures advances, all of the injunctions for my circumscription ruminate. By the fact that disparaging reprovers are incensed at knowledge, most of the droumscriptions protrude too by cognition. Vision will always be a part of hum an society. Noesis is the most precarious agriculturalist of hum an life As I have learned in my semantics class, hum an society will always enlighten ment career. Interference emits simulation to transmit plasmas. Despite the fact that the same pendulum may counteract two different gamma rays to diagnoses, the neuron receives interference. Gravity of dissemination for the assassin is not the only thing a neutrino inverts; it also processes pendulums at humanity to noesis. The sooner accumulations undertake salvers, the sooner an adjuration diverges. As a result of accounting, most of the concurrences which rationalize the abominable augur allude too by imagination.

Calling has not, and undoubtedly never will be aggravating in the way we encounter mortification but delineate the reprimand that should be inclination. Nonetheless, amned with the knowledge that the analysis augurs stealth with propagandists, almost all of the utterances on my authorization journey. Since sandions are performed at knowledge, a quantity of vocation can be more gaudily inspected. Knowledge will always be a part of society. Vocation is the most presumptuously perilous assassimation of mankind.

#### Score: 6

#### Time Used:

4 minutes 39 seconds

#### Explanation of Score:

In addressing the specific task directions, a 6 response presents a cogent, well-articulated analysis of the issue and conveys meaning skillfully.

A typical response in this category

- · articulates a clear and insightful position on the issue in accordance with the assigned task
- develops the position fully with compelling reasons and/or persuasive examples
- · sustains a well-focused, well-organized analysis, connecting ideas logically
- conveys ideas fluently and precisely, using effective vocabulary and sentence variety
- demonstrates superior facility with the conventions of standard written English (i.e., grammar, usage, and mechanics) but may have minor errors

#### Sample Responses General Advice to Writers Writer's Analysis Tools

#### Analyze an Argument Topic:

The following is a memorandum from the business manager of a television station.

"Over the past year, our late-night news program has devoted increased time to national news and less time to weather and local news. During this time period, most of the complaints received from viewers were concerned with our station's coverage of weather and local news. In addition, local businesses that used to advertise during our late-night news program have just canceled their advertising contracts with us. Therefore, in order to attract more viewers to the program and to avoid losing any further advertising revenues, we should restore the time devoted to weather and local news to its former level."

Write a response in which you discuss what specific evidence is needed to evaluate the argument and explain how the evidence would weaken or strengthen the argument.

#### Your Answer:

Audience for a reprimand has not, and presumably never will be irreverent, lamented, and toxic. Atmospheric condition is the most fundamental speculation of humankind, many with the search for semiotics but a few for mimicry. a quantity of weather lies in the area of philosophy together with the field of semantics. Although mendicant might expedite reprobates, weather condition is both abominable and penal. As I have learned in my theory of knowledge class, humanity will always dictate local. The same brain may emit two different neutrinoes to reproduce. Despite the fact that gravity counteracts plasmas, the same brain may receive two different neurons of utterances. Simulation is not the only thing the plasma on an altruist oscillates; it also produces the orbital at a quip by local. The less the unfavorable anvil placates exiles, the more culmination gloats. The genially but unavoidably precipitous local changes patter at local.

Benevolence, usually by confrontation, might enthrall career. If nearly all of the sanctions assure a concession of the swiftly or tantalizingly enthusiastic rejoinder, the consummate cognition can be more multifariously provoked. Additionally, an orbital is not the only thing simulation reacts; it also spins at knowledge. Our personal countenance on the accusation we diagnose can scrupulously be a circumspection. Be that as it may, knowing that inducement can be the oration, most of the avocations to my inspection belie toxic agreements. In my philosophy class, all of the appendages by our personal altruist of the inquiry we expedite allocate circumscriptions which advocate with accessions but reprove boundary that should potently be a aggregation and abandon performances for convulsions. Imagination which is contemptibly in how much we adhere assimilates melange of our personal advance to the administration we adjure as well, a guarrel will erroneously be a lamentation on the authentication, not an escapade. In my experience, none of the salvers by our personal amplification at the affirmation we authorize masticate consideration that journeys but incline, an abundance of vision changes plethora for careers As I have learned in my literature class, humanity will always foretell calling. Even though the brain counteracts a gamma ray to veracity, the same pendulum may catalyze two different neutrinoes with the promptly erroneous contentment. Although the same neuron may receive two different brains, radiation processes orbitals of speculations on an appetite. The plasma is not the only thing a gamma ray oscillates; it also transmits neutrinoes for torpor at the authorization by imaginativeness. The assassination of imagination changes a plethora of calling. The less extraneous respondents articulate precincts, the more an explanation amplifies those in question.

Weather which commandeers the exposure, especially of affirmations, may be capstone. As a result of lauding the inquisition to the people involved, a plethora of weather condition can be more egotiscally admired. Additionally, a risible weather changes inducement by audience. In my semiotics class, all of the assimilationists for our personal advancement with the juggemaut we stipulate voyage. Reiteration that is equitable but not unintentional can, however, be misleading, fallacious, and skeptical. My amygdala should reclusively be the advocate and fulminates. Since then, a spuriously munificent predator induces tyroes on our personal sophist at the celebration we proliferate. Perpetuity ponders the demarcation, not rationalization of the domain. My precinct is rapacious in the way we encompass authentications which laud sequester and corroborate dictators. The less all of the concessions proceed, the sooner declaration that can appropriately be abandonment will be the speculation for epitome.

As I have learned in my reality class, audience is the most fundamental accusation of human society. Radiation reacts to catalyze the pendulum. The same gamma ray may counteract two different pendulums by glutton to an orator to process neutrinoes. Gravity is not the only thing the plasma at delineation reproduces; it also emits information of local. If adjurations scrutinize concurrences, assumptions which incense the utterance with intercessions advocate equally on atmospheric condition. By commanding listlessly tremendous agronomists, weather which embroiders rancor that may promptly be an epigraph or is itinerant can be more wanly inaugurated. A query, normally with juggemauts, excommunicates weather. Due to speculating, the vapid audience can be more truculently presaged. Also, a plethora of local, usually for the diagnosis, is recondite in the extent to which we beseech an allocution that tantalizes patter or should apprehensively be the jocose verisimilitude. In my theory of knowledge class, some of the conveyances on our personal inquisition to the drone we scrutinize postulate tropes. In any case, knowing that an advancement is sedulously riveting, many of the accounts of my dictator vie. My retort permeates existence. Audience which will be an accusation that placates developments which pilfer vemacular to corroboration should drowsily be a scenario but whines at our personal utterance on the orator we fascinate to the same extent. The denouncement might be a naturally but enthrallingly unsophisticated incarceration, not depravity. In my philosophy class, many of the performances by my concurrence journey and demolish those involved of agreements. The more adherents assault scrupulousness that shrieks but quibble, the sooner the admonishment that rationalizes most of the convulsions is exemplary.

Local has not, and likely never will be regrettably assimilated. Severance may, nonetheless, be convulsive but not rancorous. Because agriculturalists are assumed with weather condition, unfavorably and humanely stipulated dictates advance too for audience. Weather condition will always be a part of human life. Instead of surrounding postulates which enthrall domains, local constitutes both a piscine accumulation and a situational affront.

Conditions at the realm of philosophy will always be a part of mankind. Comucopia that lectures should, in any case, be vociferously and pusillanimously surrounding. The sooner retorts convulse, the sooner a manifestly boastful myrmidon might be denigration for the axiom. Atmospheric condition has not, and doubtless never will be contemporary yet somehow moribund. Audience is the most depreciated sanction of human life.

#### Score: 6

#### Time Used:

3 minutes 42 seconds

#### Explanation of Score:

In addressing the specific task directions, a 6 response presents a cogent, well-articulated examination of the argument and conveys meaning skillfully.

A typical response in this category

- clearly identifies aspects of the argument relevant to the assigned task and examines them insightfully
- develops ideas cogently, organizes them logically, and connects them with clear transitions
- provides compelling and thorough support for its main points
- conveys ideas fluently and precisely, using effective vocabulary and sentence variety
- demonstrates superior facility with the conventions of standard written English (i.e., grammar, usage, and mechanics) but may have minor errors

Sample Responses General Advice to Writers Writer's Analysis Tools



Score#Now!<sup>™</sup> Home

Overview of Service

Score Level Descriptions Advice to Writers

Frequently Asked Questions

Buy Additional Service

Update Account

List of Topics

Scoring Guides

# ScoreltNow!"

Score#Now!<sup>™</sup>Home > Test Experience > Previous Tests > Review

R	ev	i	e	w
	_	-	-	

#### Analyze an Issue Topic:

The best way for a society to prepare its young people for leadership in government, industry, or other fields is by instilling in them a sense of cooperation, not competition.

Write a response in which you discuss the extent to which you agree or disagree with the claim. In developing and supporting your position, be sure to address the most compelling reasons or examples that could be used to challenge your position.

#### Your Answer:

Competition for an inquiry has not, and presumably never will be antipodal, puissant, and equitable. Success is the most fundamental adjuration of humankind; many with the search for semiotics but a few for pondering, a quantity of cooperation lies in the area of philosophy together with the field of semantics. Although buccaneer might propagate amygdalas, cooperation is both boastful and insouciant.

As I have learned in my theory of knowledge class, humanity will always incarcerate success. The same brain may emit two different neutrinoes to reproduce. Despite the fact that gravity counteracts plasmas, the same brain may receive two different neutrons of lamentations. Simulation is not the only thing the plasma on an allocation oscillates; it also produces the orbital at a denouncement by success. The less the unsophisticated subjugation probes authentications, the more lacuna semonizes. The vapidly but transitorily tendentious success changes sequester at success.

Competition which mesnerizes the reprover, especially of administrations, may be multitude. As a result of abandoning the utterance to the people involved, a plethora of cooperation can be more tensely enjoined. Additionally, a humane competition changes assemblage by cooperation. In my semiotics class, all of the agriculturalists for our personal interloper with the probe we decry contend. Anatomy that is fascinating but not inflam matory can, however, be contentious, professed, and banal. My scenario should indispensably be the salver and convulses. Since then, a gluttonously listless oligarchy subjugates congregations on our personal sanction at the intercession we demonstrate. Spectrometry enthrals the advancement, not presage of the retort. My demarcation is quiescent in the way we corroborate analyses which abandon compensation and demarcate scrutinizations. The less all of the proclamations respond, the sooner perpetuity that can obtrusively be proliferation will be the adjuration for ouster.

As I have learned in my reality class, competition is the most fundamental civilization of human society. Radiation reacts to catalyze the pendulum. The same gamma ray may counteract two different pendulums by gluttony to a accusation to process neutrinoes. Gravity is not the only thing the plasma at periodicity reproduces; it also emits information of cooperation. If authorizations countenance commencements, celebrations which discumscribe the lamentation with demonstrations quibble equally on cooperation. By augmenting vociferously magnetic precinds, success which compensates obloquy that may deafeningly be a dictator or is irascible can be more opulently postulated.

Competition at the realm of philosophy will always be a part of mankind. Consistency that consents should, in any case, be tendentiously and pusillanimously presum ptuous. The sooner dictates quarrel, the sooner a prototypically disciplinary myrmidon might be masochism for the allegation. Success has not, and doubtless never will be natural yet somehowperipatelic. Cooperation is the most transitory inquisition of human life.

Success has not, and doubtlessly never will be sophistic. Mankind will always preach success; some of twenty-first and others by a conveyance. a abundance of success lies in the realm of philosophy and the field of theory of knowledge. Competition is the most opulent aborigine of humankind.

As I have learned in my semantics class, mankind will always expose competition. Though gravity emits plasmas, the same gamma ray may transmit two different brains. While the pendulum at amygdalas on concessions process a gamma ray by an assassin, the same neutrino may catalyze two different orbitals. Radiation is not the only thing simulation reads; it also reproduces of success. The erroneously Libertarian success changes a blatant cooperation. The sooner the arrangement howls, the more analyses which homogenize allure the remarkable depreciation.

According to professor of literature Leon Trotsky, human society will always pilfer competition, a pendulum to inducement inverts to read. The same plasma may receive two different neurons with the administration for accountesto process interference. In form ation is not the only thing a brain at instructions oscillates; it also counteracts neurons by the retort with competition. Because commanding civilizations are explained on cooperation, the aggregation to cooperation can be more philanthropically commanded. Since allusions which countenace agronomists are contravened of competition, the people involved accede as well at competition.

Success, normally for the realm of theory of knowledge, is averred but not postlapsarian and scintillates. As a result of denigrating altruists, avocations by masochism which verify unscrupulousness or quibble shriek equally with cooperation. Additionally, a pendulum is not the only thing interference by the circumstance spins; it also produces a gamma ray of cooperation. In my experience, most of the embroideries to our personal diagnosis on the respondent we enjoin mesmerize countenances. a postulate that hovers but precludes inquisitions may, nonetheless, be magnanimous, fetishistic, and recondite. In my philosophy class, many of the admonishments with my adjuration commence or beseech privation. Success which will slightly be torpor might be corroboration at inspections for our personal ligation by the disenfranchisement we pommel equally, a casuistry recounts myrmidon, not the reprobate with dubiously or graciously ashen precincts. In my experience, some of the amplifications to our personal reprimand of the domain we compensate proliferate exiles. Success which is Libertarian in how much we inaugurate none of the ruminations changes success which queries.

Cooperation on pique will always be a part of society. Scrupulousness can, however, be unsubstantiated. If assimilationists for a performance report and ascertain the confluence, pugnaciously increasing escapades which semonize cavort to the same extent at cooperation. Cooperation with circumscriptions will always be a component of human life. Seeing as competition propagandizes those in question, mankind should compet competition immediately.

#### Score: 6

#### Time Used:

2 minutes 12 seconds

#### Explanation of Score:

In addressing the specific task directions, a 6 response presents a cogent, well-articulated analysis of the issue and conveys meaning skillfully.

A typical response in this category

- · articulates a clear and insightful position on the issue in accordance with the assigned task
- develops the position fully with compelling reasons and/or persuasive examples
- · sustains a well-focused, well-organized analysis, connecting ideas logically
- · conveys ideas fluently and precisely, using effective vocabulary and sentence variety
- demonstrates superior facility with the conventions of standard written English (i.e., grammar, usage, and mechanics) but may have minor errors

#### Sample Responses General Advice to Writers Writer's Analysis Tools

#### Analyze an Argument Topic:

The following is taken from a memo from the advertising director of the Super Screen Movie Production Company.

"According to a recent report from our marketing department, during the past year, fewer people attended Super Screen-produced movies than in any other year. And yet the percentage of positive reviews by movie reviewers about specific Super Screen movies actually increased during the past year. Clearly, the contents of these reviews are not reaching enough of our prospective viewers. Thus, the problem lies not with the quality of our movies but with the public's lack of awareness that movies of good quality are available. Super Screen should therefore allocate a greater share of its budget next year to reaching the public through advertising."

Write a response in which you discuss what questions would need to be answered in order to decide whether the recommendation and the argument on which it is based are reasonable. Be sure to explain how the answers to these questions would help to evaluate the recommendation.

#### Your Answer:

Reasoning by a contradiction has not, and doubtlessly never will be startling yet somehow solicited. Human life will always postulate pay heed; whether with assassins or on the amygdala. False belief of Super Screen Movie Production Company, which sermonizes or emboldens consistency lies in the area of semantics as well as the search for reality. Hang is the most raucous epigraph of humanity.

The equipoise of Super Screen Movie Production Company, frequently to lassitude, should be egregious in how much we verify the inquisition but blubber and tantalize celebrations of advertsing. Because some of the postulates are proliferated with pay heed, a lack of hang can be more essentially choreographed. Furthermore, an orbital for the development to happenstance is not the only thing a gamma ray inverts; it also receives simulation on advert. Our personal drone of the ligation we civilize depletes adjurations. Still yet, armed with the knowledge that the report with infusion can petulantly be the injudicious stipulation, none of the lamentations by my circumstance compel inconsistency but agree. In my experience, many of the quips at our personal admonishment on the allocation we countenance collapse or disrupt risibly unsophisticated precincts. Since then, a civilization may be substantiation and accumulates the accumulation that will litigiously be a dictator for our personal agronomist by the assimilationist we mesmerize. The exposition is menaced, hirsute, and alleged, not demolishment to allusions. My casuistry insists. The less agriculturalists contradict appeasement, the sconer a countenance should be a axiom.

The atelier that assassinates circumscriptions, typically of a thermostat, affirms fallacy. If the people involved jeer but conduct presage, a belligerent give ear can be more efficaciously presumed. Additionally, fallacy, especially with escapades, may situationally be presumption to augur. Our personal embroidery by the organism we preach ponders an accession but is fittingly and oligarchical naive. In any case, knowing that the respondent assents, almost all of the assumptions on our personal inspection at the appendage we forsake promulgate assemblies to injunctions. Our personal demonstration for the instinuation we contravene encompasses affirmations. Decency to fallacy can be a slightly ineverent gluttony of my propagandist too. Quibble might pugnaciously be the advance but lectures, not pilfering. In my semantics class, some of the appendages by our personal aggregation on the appendage we encounter solicit apprentices. a dearth of logical thinking changes attend which commissions the demolisher at the search for theory of knowledge.

As I have learned in my philosophy class, reasoning is the most fundamental ligation of mankind. Interference with a dictate for drones which postulate promulgation of profession or gloat emits neurons by a precinct to transmit neutrinoes. The same brain may catalyze two different plasmas to circumspections to process the gamma ray for a lamentation. The pendulum is not the only thing gravity oscillates; it also reacts at fallacy. As a result of insinuating, developments with protrusion embark too on reasoning. The less a quip that edifies provision is risible but not obvious, the sooner the pledge compensates postulates for the study of semiotics.

Abstract thought will always be an experience of human society. However, armed with the knowledge that validation that may elatedly be the consequence proclaims a prison, all of the ligations by my exile feign judicious reports. a lack of pay heed changes the soporifically or positively puissant scrupulousness to reasoning. Abstract thought has not, and likely never will be irrelevant in the way we portend contemplation and drone. Despite the fact that most of the consequences should accede assimilationists, false belief is both increasing and emphatic.

Causes on ligature has not, and in all likelihood never will be judiciously reclusive. Human life will always compensate diminution; many of the amygdala but a few on allegations. an abhorrently but outlandishly inchoate decline lies in the search for semiotics and the realm of theory of knowledge. Consequently, audience should engender none of the allusions.

As I have learned in my reality class, causes is the most fundamental interloper of humanity. Though a pendulum produces neurons by amplifications, the same brain may counteract two different plasmas. The plasma emits gamma rays with incarceration to spin. Simulation at an advocate is not the only thing interference for existence spins; it also reacts of audience. The less those in question contravene the rumination but amplify scenarios, the sooner an unavoidable pedant whines. Audience which regrets approbation and is boisterous, lethargic, and magnificent changes a lack of decline.

According to professor of philosophy Eli Whitney, mankind will always retort audience. Despite the fact that the same orbital may process two different neurons, the same plasma may catalyze two different neutrinoes. The pendulum counteracts gravity to receive brains to plethora. Simulation for interlopers is not the only thing a neuron inverts; it also transmits plasmas on decline. The sooner drones relent, the more the insinuation by the appetite should decently be fetishism that can irascibly be an accession. By augmenting an inquisition that observes disenfranchisements, the tranquilly precarious diminution can be more unfavorably enlightened.

As I have learned in my semiotics class, human society will always compel diminution. Although a pendulum oscillates, the same gamma ray may catalyze two different neutrinoes. Information at endemic celebrations processes orbitals of appendages with perjury to counteract brains. The neuron is not the only thing radiation to a propagandist implodes; it also produces the pendulum on causes. From hoversing, many of the inspections avow as well by cause. Consideration for cause changes the commanding causes.

Audience will always be a component of mankind. Nonetheless, armed with the knowledge that the indispensably pedantic ingenuity might gaudily be particularism, none of the congregations by my avocation deplete consequences. If diagnoses which expose accumulations allure advancements for an utterance, a plethora of decline can be more remarkably magnetized. Audience of aborigines will always be a component of human life. Decline is orotund because of its impartial agreements.

#### Score: 6

Time Used:

6 minutes 1 second

#### Explanation of Score:

In addressing the specific task directions, a 6 response presents a cogent, well-articulated examination of the argument and conveys meaning skillfully.

A typical response in this category

- · clearly identifies aspects of the argument relevant to the assigned task and examines them insightfully
- · develops ideas cogently, organizes them logically, and connects them with clear transitions
- · provides compelling and thorough support for its main points
- · conveys ideas fluently and precisely, using effective vocabulary and sentence variety
- demonstrates superior facility with the conventions of standard written English (i.e., grammar, usage, and mechanics) but may have minor errors

# Try the BABEL Generator http://babel-generator.herokuapp.com/

# Appendix C

# Biography: Les Perelman, Ph.D.

Les Perelman is an internationally recognized expert in writing assessment and the application of technologies to assess writing. He has written opinion pieces for *The Boston Globe, The Washington Post*, and *The Los Angeles Times*. He has been quoted in *The New York Times, The New Yorker, The Chicago Tribune, The Boston Globe, The Los Angeles Times*, and other newspapers. Dr. Perelman has been interviewed on television by ABC, MSNBC, and NHK Japan Public Television and interviewed on radio by National Public Radio, various NPR local stations, the Canadian Broadcasting Corporation, and the Australian Broadcasting Corporation.

The President of the College Board has credited Dr. Perelman's research as a major factor in his decision to remove and replace the Writing Section of the SAT. Dr. Perelman is a well-known critic of Automated Essay Scoring. To demonstrate the inability of Robo-graders to differentiate writing from gibberish, he and three undergraduates developed the *BABEL Generator*, which produces verbose and pretentious nonsense that consistently receives high marks from AES machines.

Dr. Perelman received his B.A. in English Language and Literature from the University of California, Berkeley, and his M.A. and Ph.D. in English from the University of Massachusetts. After a three-year postdoctoral fellowship in Rhetoric and Linguistics at the University of Southern California, Dr. Perelman moved to Tulane University where he served as an Assistant Professor of Rhetoric, Linguistics, and Writing; Director of First-Year Writing; Director of the Writing Center; and a Member of the Graduate Faculty.

For the next twenty-five years Dr. Perelman was Director of Writing Across the Curriculum in Comparative Media Studies/Writing at the Massachusetts Institute of Technology and served as an Associate Dean in the Office of the Dean of Undergraduate Education. He was Project Director and co-Principal Investigator for a grant to MIT from the National Science Foundation to develop a model Communication-Intensive Undergraduate Program in Science and Engineering. He served as Principal Investigator for the development of the iMOAT Online Assessment Tool funded by the MIT/Microsoft iCampus Alliance. Dr. Perelman has served as a member of the Executive Committee of the Conference on College Composition and Communication, the post-secondary organization of the National Council of Teachers of English, and co-chaired the Committee on the Assessment of Writing. He is currently a member of the editorial board of *Assessing Writing*.

Dr. Perelman has been a consultant to over twenty colleges and universities on the assessment of writing, program evaluation, and writing-across-the-curriculum. Dr. Perelman has served as a consultant for writing program assessment and development for the Fund for the Improvement of Postsecondary Education of the U.S. Department of Education and for the Modern Language Association. In 2012–2013, he served as a consultant to Harvard College and as co-principal investigator in a major two-year study assessing the writing abilities of undergraduates at the college.

Dr. Perelman co-edited the volume *Writing Assessment in the 21<sup>st</sup> Century* and he is the primary author of the first web-based technical writing handbook, *The Mayfield Handbook of Technical and Scientific Writing*. He has published articles on writing assessment, technical communication, computers and writing, the history of rhetoric, sociolinguistic theory, and medieval literature, and he co-edited *The Middle English Letter of Alexander to Aristotle*.