# Polygraph & Forensic Credibility Assessment: A Journal of Science and Field Practice

DLUME 50	2021	NUMBER
	Contents	
Different-than-Chance In Technique for Jury Select	terpretation of the Two-Alternative, Forced-Choice ion and Witness Screening	1
Emily McElfresh, I Neslihan James-K	Kevin Colwell, Cheryl Hiscock-Anisman, Laura Welc angal, Caitlin Cardenas and Brian Gavigan.	h,
Response Onset Latencie Responses During Field F Donald J. Krapohl Donnie W. Dutton	s of Electrodermal, Cardiovascular and Vasomotor Polygraph Testing I, Karen Halford, Tim Benson, Abbe Mayston and	15
Literature Review and An of Independent Target Qu Raymond Nelson,	alysis of the Multi-facet Hypothesis and the Evaluat lestions Mark Handler, and David C. Raskin	ion 28
Another Look at Electrod Donald J. Krapohl	ermal Response Ratio Minima	59

## Polygraph & Forensic Credibility Assessment: A Journal of Science and Field Practice

Editor-in-Chief: Mark Handler E-mail: <u>Editor@polygraph.org</u> Managing Editor: Nayeli Hernandez E-mail: polygraph.managing.editor@gmail.com

\*\*\*\*\*

<u>Associate Editors:</u> Réjean Belley, Ben Blalock, Kevin Colwell, Tyler Blondi, John Galianos, Don Grubin, Maria Hartwig, Charles Honts, Matt Hicks, Scott Hoffman, Don Krapohl, Thomas Kuczek, Mike Lynch, Ray Nelson, Adam Park, David Raskin, Stuart Senter, Joseph R. Stainback IV and Cholan V.

APA Officers for 2020 – 2021

President – Sabino Martinez E-mail: <u>president@polygraph.org</u>

President Elect – Roy Ortiz E-mail:<u>presidentelect@polygraph.org</u>

Chairman Darryl Starks E-mail: <u>chairman@polygraph.org</u>

Director 1 – Pamela Shaw E-mail: <u>directorshaw@polygraph.org</u>

Director 2 – Raymond Nelson E-mail: <u>directornelson@polygraph.org</u>

Director 3 – Jamie McCloughan E-mail:<u>directormccloughan@polygraph.org</u>

Director 4 – Chip Morgan E-mail: <u>directormorgan@polygraph.org</u>

Director 5 – Erika Thiel E-mail: <u>directorthiel@polygraph.org</u>

Director 6 – Donnie Dutton E-mail: <u>directordutton@polygraph.org</u> Director 7 – Lisa Ribacoff E-mail: <u>directorribacoff@polygraph.org</u>

Director 8 – Walt Goodson E-mail: <u>directorgoodson@polygraph.org</u>

Treasurer – Chad Russell E-mail: <u>treasurer@polygraph.org</u>

General Counsel – Gordon L. Vaughan E-mail: <u>generalcounsel@polygraph.org</u>

Seminar Chair – Michael Gougler E-mail: <u>seminarchair@polygraph.org</u>

Education Accreditation Committee (EAC) Manager – Barry Cushman E-mail: <u>eacmanager@polygraph.org</u>

National Officer Manager – Lisa Jacocks Phone: 800-APA-8037; (423)892-3992 E-mail: <u>manager@polygraph.org</u>

National Office Assistant - Jennifer Crawley

Subscription information: *Polygraph & Forensic Credibility Assessment: A Journal of Science and Field Practice* is published semi-annually by the American Polygraph Association. Editorial Address is <u>Editor@polygraph.org</u>. Subscription rates for 2021: One year \$150.00. Change of address: APA National Office, P.O. Box 8037 Chattanooga, TN 37414-0037. THE PUBLICATION OF AN ARTICLE IN *Polygraph & Forensic Credibility Assessment: A Journal of Science and Field Practice* DOES NOT CONSTITUTE AN OFFICIAL ENDORSEMENT BY THE AMERICAN POLYGRAPH ASSOCIATION.

## **Instructions to Authors**

#### Scope

The journal Polygraph & Forensic Credibility Assessment: A Journal of Science and Field Practice publishes articles about psychophysiological detection the of deception, and related areas. Authors are invited to submit manuscripts of original research, literature reviews, legal briefs, theoretical papers, instructional pieces, case histories, book reviews, short reports, and Special topics will be similar works. considered on an individual basis. А minimum standard for acceptance is that the paper be of general interest to practitioners, instructors and researchers of polygraphy. From time to time there will be a call for papers on specific topics.

#### **Manuscript Submission**

Manuscripts must be in English, and may be submitted, along with a cover letter, on electronic media (MS Word). The cover letter should include a telephone number, and e-mail address. All manuscripts will be subject to a formal peer-review. Authors may submit their manuscripts as an e-mail attachment with the cover letter included in the body of the e-mail to:

#### Editor@polygraph.org

As a condition of publication, authors agree that all text, figures, or other content in the submitted manuscript is correctly cited, and that the work, all or in part, is not under consideration for publication elsewhere. Authors also agree to give reasonable access to their data to APA members upon written request.

#### **Manuscript Organization and Style**

All manuscripts must be complete, balanced, and accurate. Authors should follow guidelines in the *Publications Manual* of the American Psychological Association. The manual can be found in most public and university libraries, or it can be ordered American Psychological Association from: Publications, 1200 17th Street, N.W., Washington, DC 20036, USA. Writers may exercise some freedom of style, but they will held to a standard of clarity, be organization, and accuracy. Authors are responsible for assuring their work includes correct citations. Consistent with the ethical standards of the discipline, the American Polygraph Association considers quotation of another's work without proper citation a grievous offense. The standard for nomenclature shall be the *Terminology* Reference for the Science of Psychophysiological Detection of Deception (2012) which is available from the national office of the American Polygraph Association. Legal case citations should follow the West system.

#### **Manuscript Review**

An Associate Editor will handle papers, and the author may, at the discretion of the Associate Editor, communicate directly with him or her. For all submissions, every effort will be made to provide the author a review within 4 weeks of receipt of manuscript. Articles submitted for publication are evaluated according to several criteria including significance of the contribution to the polygraph field, clarity, accuracy, and consistency.

#### Copyright

Authors submitting a paper to the American Polygraph Association (APA) do so with the understanding that the copyright for the paper will be assigned to the American Polygraph Association if the paper is accepted for publication. The APA, however, will not put any limitation on the personal freedom of the author(s) to use material contained in the paper in other works, and request for republication will be granted if the senior author approves.

## Different than Chance Interpretation of the Two-Alternative, Forced-Choice Technique for Jury Selection and Witness Screening

Emily McElfresh<sup>1</sup>, Kevin Colwell<sup>1</sup>, Cheryl Hiscock-Anisman<sup>2</sup>, Laura Welch<sup>1</sup>, Neslihan James-Kangal<sup>3</sup>, Caitlin Cardenas<sup>1</sup>, Brian Gavigan<sup>1</sup>

#### Abstract

This paper extended Forced-Choice testing to the forensically relevant settings of jury selection and screening potential witnesses to mass crimes. These new applications were made possible by changing the interpretive approach from the "Worse-than-chance" model to the "Differentthan-chance" model. The jury selection sample responded either honestly or deceptively regarding their knowledge regarding a true crime from the media. The witness sample responded honestly or deceptively regarding a mock bank robbery. FC testing made it possible to determine whether each participant was responding honestly rather than hiding or faking their knowledge of the target event. Overall, approximately 93% of participants were accurately classified, with approximately 31% improvement as a result of the Different-than-Chance model rather than the Worse-than-Chance model. The Different-than-Chance model should replace the Worse-than-Chance model for situations involving an episodic memory, and provides a more sensitive and specific mechanism for evaluating statements and informing decisions.

<sup>1</sup>Department of Psychology, Southern Connecticut State University

<sup>2</sup>Department of Psychology, National University

<sup>3</sup>Veteran's Administration, San Francisco, CA

Author Notes We have no known conflict of interest.

A portion of this research was funded by Southern Connecticut State University's Undergraduate Research Grant.

\* Address correspondence to Kevin Colwell, Ph.D., Department of Psychology, SCSU. 501 Crescent Street, New Haven, CT. 06515. Email: colwellk2@southernct.edu

1

One of the most widely used techniques in clinical and forensic assessment is the two-alternative, forced-choice test (2FC: Colwell & Sjerven, 2005; Colwell & Colwell, 2011; Hiscock & Hiscock, 1989; Pankratz, 1979). 2FC provides a sound strategy for the assessment of response style and malingering in a number of situations, including Competence to Stand Trial (Colwell & Colwell, 2011; Gottfried, Hudson, Vitacco, & Carbonell, 2017), amnesia and cognitive impairment (Colwell & Sjerven, 2005; Kapur, 1994; Schroeder, Peck, Buddin, Heinrichs, & Baade, 2012), Criminal Responsibility (Hiscock, Branham, & Hiscock, 1994), and general clinical assessment (Hiscock & Hiscock, 1989). However, relatively little attention has been given to extending 2FC to witnesses (although it has been proposed for use with witnesses for almost 20 years), and never has it been used with jury selection (Colwell, Hiscock-Anisman, Memon, Taylor, & Prewett, 2008; Colwell, Hiscock, & Memon, 2002; Jelicic, Merckelback, & van Bergen, 2004; Orthey, Vrij, Leal, & Blank, 2017).

As discussed in Colwell & Sjerven (2005), 2FC can be interpreted through either parametric or nonparametric statistics. The parametric approach gives increased sensitivity and specificity but requires a pre-tested group of relevant people and context to compare against (Colwell & Sjerven, 2005). This is not practical in situations involving the deliberate hiding of information regarding complex events. Each event is unique, and there is no way to have a pre-tested group ready upon demand. The nonparametric approach is based upon the statistics of probability, and the nature of the 2FC task. Therefore, this approach is applicable to any situation where it is possible for individual choices to be selected (Colwell & Sjerven, 2005).

#### Worse-than-Chance

Those who genuinely have no knowledge of a target event should perform at chance on a 2FC test, assuming that the test is constructed regarding information from the target event that is not self-evident or common knowledge. Those who are lying will often deliberately choose the wrong response and can therefore perform worse than chance. Those who are honest are free to choose the correct response, and therefore perform better than chance.

This model comes from the assessment of malingering, where those who know a correct answer often choose the wrong one when simulating cognitive impairment (Colwell & Colwell, 2011; C. Hiscock, Rustemier, & Hiscock, 1993). This interpretive model works well in clinical assessment, and there is a very low false-positive rate (Colwell & Sjerven, 2005).

#### **Different-than-Chance**

On a 2FC test, performing significantly worse than chance and performing significantly better than chance each convey the same information to the person who scores the test - that is, the respondent possessed information regarding the target event. To hide knowledge of a target event, it is necessary to score within the band described by chance. We propose that those who are hiding information will not know about probability or statistics (Cliffe, 1992). Instead, they are likely to approach this like more traditional deception tasks, in which one can either deny almost all knowledge of a target event, or one can deny a few critical items such that one can succeed in deception while still providing a statement that is mostly honest (Colwell, Hiscock-Anisman, Memon, Woods, & Michlik, 2006; Hines et al., 2010). The first case would likely lead to worse-thanchance performance, and the latter would likely lead to better-than-chance performance. Thus, for situations involving the attempt to hide knowledge of a complex event (such as a crime), we propose that attempts at detecting deception should be done using the Different-than-Chance model. Conceptually, the issue is how to identify those who possess information about the event. Different-than-Chance performance indicates this.

#### Deceiver's Dilemma

A 2FC test creates a dilemma for one hiding information. The only way to lie and escape detection is to miss the responses to half of the details of the event that are on the test. In missing items, the deceiver has to respond so that there are no obvious contradictions between what she or he should know and what she or he gets right or wrong on the test. The deceiver also must avoid demonstrating knowledge of any potential critical items - items that would show the deceiver to possess information that investigators want, or information that has taken on a special significance in this particular case (Colwell, Hiscock-Anisman, et al., 2008).

#### **Item Construction**

The items from existing 2FC tests can inform the process of item construction for a novel 2FC test. In the Coin-in-Hand test (Colwell & Sjerven, 2005; Kapur, 1994) the evaluator holds a coin in an open hand, and shows this to a respondent with potential memory problems. The evaluator closes her or his hand, and counts to 10. Then, the evaluator asks the respondent which hand holds the coin. This is done 10 times, switching hands. The items are relevant to the presentation of memory ability, and the items are so easy that even those with memory impairment can do them. The items are critical items, in that they directly relate to how a person chooses to present their memory for recent events.

The Test of Malingered Incompetence is administered to people who are suspected of exaggerating or completely faking mental health issues or exaggerating or faking as if they are not competent to stand trial (Colwell, Colwell, et al., 2008). The TOMI has two sets of items, created to be critical items for different portions of the Dusky Standard (Dusky v. United States, 1960). These are the General Knowledge scale (critical to the, "mental disease or defect," portion of Dusky) and the Legal Knowledge scale (critical to the, "understanding role as defendant," and, "work in own defense," prongs of Dusky). TOMI items were constructed so that they were so easy that 98% or more of honest respondents can answer each correctly.

In conclusion, items should be constructed such that: 1) Cooperative honest respondents with information will highlight themselves by responding correctly to a large proportion. 2) Those who genuinely cannot perform the task or who possess no relevant information should perform at chance. 3) The test should include items that the respondent will perceive as critical. Examples of these are, "Which hand holds the coin, Right or Left?" or, "Who had the gun, Man or Woman?"

#### **Current Project**

This project had the following purposes: 1) extend the 2FC model to jury selection in cases with significant media coverage, and to screening when there is a large number of witnesses to a single crime, and 2) demonstrate the improved effectiveness of the Different-than-Chance model when compared to the Worse-than-Chance model. The hypotheses were that using the 2FC approach with either model would outperform chance for accuracy, classifying respondents as informed or naïve, and the Different-than-Chance model would outperform the Worse-than-Chance model in detecting hidden information and overall classification accuracy.

## **Experiment 1 – Jury Selection**

## Method

#### Participants

The participants in this study were undergraduate student volunteers who received extra credit in Psychology and Criminal Justice classes at a university on the East Coast of the US. The sample comprised 44 males and 43 females (none reported trans or non-binary). Their ages ranged from 18-68 years, with a mean of 32 years (sd = 15.7 years). The self-reported Ethnic Identities were: 54.4% Caucasian/White, 17.5% African American / Black, 15.8% Latinx, and 3.5% Asian/Pacific Islander, with 8.2% "Other."

## Materials

Three different newspaper articles regarding criminal investigations/cases that were published in Los Angeles and San Diego were chosen for this study. Each of these stories was chosen because it was vivid and emotionally significant, and each involved an allegation of harm to a child. These stories were not reported on the East Coast, where the study took place, and each of the participants endorsed that they had never heard anything about the case prior to reading the article (this was one of the reasons for having 3 available, but a second choice was never required). The three articles were distributed in a balanced manner.

#### Design, Materials, and Methods

Eighty-seven participants were split into three groups of 29 and were assigned different tasks for FC testing. The goal was to highlight those respondents who would be of most interest to investigators - i.e., those who had information regarding the target event, and more especially those, who were attempting to hide this information. All participants were told to imagine that they were completing a questionnaire as part of a jury-selection process.

Participants in Group A were given an article to read and told to respond as honestly and cooperatively as possible. They were instructed that the respondent who demonstrated the most knowledge of the case would win a \$50 gift card.

Participants in Group B were given an article to read and told to imagine that they wanted to hide their knowledge of the case and respond to conceal their knowledge as a strategy to get picked for jury service. They were instructed that the respondent most able to hide information and blend in with those others without knowledge of the event, would win a \$50 gift card. Participants were not coached regarding how to blend in with others, and left to come up with their own strategy.

Participants in Group C were not given an article to read. They were told to lie in order to appear as if they had knowledge and a specific opinion of the events in order to avoid jury service. They were told the respondent who did the best at lying to appear as if they knew about the case, would win a \$50 gift card. In truth, these cards were chosen by random drawing after all FC tests were completed.

#### FC Tests

A 2FC test with 22 items was created from the information obtained in each article (3 tests, total). Each participant was given the test associated with the article that she or he read. Chance performance on each test would lead to 11 items correct. One would expect a participant with no information regarding the target event to obtain a score that falls at or near 11 items correct, and 11 items incorrect, approximately. There will be variability, and some items may appear relatively more or less likely

to some participants, which is the reason for using a confidence band rather than a specific number.

#### Manipulation Check

There was an item at the end of the demographics that asked each participant what experimental condition they were in. It read, "Which of the following best describes your experiences in this experiment: 1) I read a newspaper article and I responded honestly. 2) I read a newspaper article and I lied to act as if I had not seen it to be allowed onto the jury. 3) I did not read an article, but I tried to answer as if I had seen to avoid jury service. The response of all participants to this question aligned with their assigned experimental condition, and therefore their data was used for further analysis.

#### Results

The following statistics were computed using the Binomial Effect Size Display calculator provided by https://www.psychometrica.de/ effect\_size.html, and the Binomial Probability Calculator provided by https://stattrek.com.

According to the binomial statistic for 22 items, chance performance (alpha = .05) forms an inclusive band between 6 and 16.

#### Group A

The mean number of items correct in Group A was 20 (sd = 1.4). Thus, the items were simple enough that those who were instructed to respond honestly answered all or almost all of the items correctly.

*Worse-than-Chance.* Twenty-nine of the 29 participants in Group A obtained 16 or more correct. None obtained a score at or below chance. Because all performed better than chance, they were all were accurately classified as having information (100%). Their responding led to the same classification accuracy using the Worse-than-Chance and Different-than-Chance approaches. The accuracy rate of each of these was compared to the accuracy expected due to chance (50%) using the binomial statistic. The null hypothesis for each was that the observed rate of 100%

	Participants Below Chance (x<6)	Participants Within Chance (6≤x≤16)	Participants Above Chance (x>16)	Participants Different than Chance (x<6)+(x>16)
Group A (29) Honest information	0%	0%	100%	100%
Group B (29) Deceptive hide information	34%	21%	45%	79%
Group C (29) Honest without information	0%	97%	03%	03%

Table 1. Number of Respondents Performing Below, Within, Above, and Different From, the Chance Range.

was not a significant improvement over the rate expected chance rate of 50%. The alternate hypothesis was that the observed rate of 100% was a significant improvement over the expected rate of 50%. The results indicated that we should reject the null ( $r_{phi} = .78$ , d = 2.53, p < .01). Thus, the observed rate of 100% was a statistically-significant improvement over chance.

*Comparison to Different-than-Chance.* No comparison is presented, due to identical results.

## Group B

The mean number of items correct for this group was 9 (sd = 7.6). Ten participants scored below chance (<6), 6 scored within the chance range, and 13 scored above chance (>16).

*Worse-than-Chance.* Ten of the 29 participants who possessed information but lied to hide it were accurately classified as having information using the Worse-than-Chance model. To assess the Worse-than-Chance model, it was compared to the rate expected by chance

using the binomial statistic. The null hypothesis was that the observed rate of 34% was not a significant improvement over the expected chance rate of 50%. The results indicated that we should not reject the null ( $r_{phi} = -.16$ , d = -.33, p = .97). Thus, the performance of the Worse-than-Chance model did not lead to a significant improvement from that expected by chance.

Comparison to Different-than-Chance. Twenty-three of the 29 deceptive participants who possessed information but attempted to hide it, were accurately classified using the Different-than-Chance model (79%). The performance of the Different-than-Chance model was compared to that of the Worse-than-Chance model. The null hypothesis was that the observed rate of 79% was not a significant improvement over the previously-observed rate of 34%. The alternate hypothesis was that the observed rate of 79% was a significant improvement over the previously-observed rate of 34%. The results from the binomial effect size display calculator indicated that we should reject the null ( $r_{phi} = .34, d = .73, p < .001$ ). Thus,

the performance of the Different-than-Chance model led to a statistically-significant improvement over the performance of the Worse-than-Chance model.

Strategy Check. The mean score from those participants who scored above chance from Group B (18, sd = 1.82, 95% CI = 17.4-18.6) was compared to the mean score of Group A (20, sd = 1.4; 95% CI = 19.0-21.0). The null hypothesis was that there was no difference between groups. The alternate hypothesis was that there was a difference between groups. The results of a 1x2 ANOVA indicated a significant difference between groups (F(1,32) = 8.77, p<.01, eta squared = .22). Thus, those from Group B who chose to answer at an above-chance level were not behaving the same those in group A, who were simply telling the truth.

## **Group** C

The mean number of items correct = 12 (sd=2.8). For this group, 0 participants scored below chance (<6), 28 scored within the chance range, and 1 scored above chance (>16).

*Worse-than-Chance.* Twenty-eight of the 29 deceptive participants in this group were accurately classified as not having information (97%) using the Worse-than-Chance model.

The performance of the Worse-than-Chance model was compared to that expected by chance (50%). The null hypothesis was that the accuracy rate of 97% generated by the Worse-than-Chance model does not represent a significant improvement over the 50% expected by chance. The alternate hypothesis was that the accuracy rate of 97% generated by the Worse-than-Chance model represents a significant improvement over the 50% expected by chance. The results indicated that we should reject the null ( $r_{phi} = .69$ , d = 1.90, p < .01). The Worse-than-Chance model represented a statistically-significant improvement over chance. The same is true for the Different-than-Chance model ( $r_{phi}$  = .78, d = 2.6, p < .01).

*Comparison to Different-than-Chance.* No comparison is provided, due to identical results.

## **Discussion Experiment 1**

The 2FC approach provided an effective method to evaluate the extent of a potential juror's knowledge of a case. Those who possessed information and were responding honestly were easily separated from those who did not possess information but attempted to act as if they did. Importantly, it was a simple task to find those who were attempting to avoid jury service by faking as if they knew about the case and therefore already had an opinion. Perhaps more importantly, it was relatively easy to identify 4 out of 5 of those who had read about the case and were attempting to hide this knowledge or otherwise blend with those who did not posses knowledge.

Overall, the 2FC led to an accuracy rate of 92%. This is an extremely high rate, the technique requires relatively little time, and no special equipment- only than the ability to present questions to people and record their scores. 2FC provides a way for an attorney to have quick and accurate insight into what a potential juror knows and can give insight into the thought process of each potential juror regarding the incident in question. Also, this could protect decisions to strike certain potential jurors from challenges based upon alleged discrimination or bias.

One limitation of this study was the length of the test used. However, the binomial statistic gives good sensitivity and specificity with much fewer observations (as in Experiment 1), and therefore fewer questions should be used in future applications.

#### Experiment 2 – Witnesses to Mock Crime

## Method

## Participants

The participants in this study were undergraduate student volunteers who received extra credit in Psychology and Criminal Justice classes at a university on the East Coast of the US. The sample comprised of males and females (none reported trans or non-binary). Their ages ranged from 18-56, with a mean of 29 years (sd = 14.6 years). The self-reported Ethnic Identities were: 57.4% Caucasian/ White, 16.8% African American / Black, 13.8% Latinx, 8.3% Asian/Pacific Islander, with 6.4% "Other." One important aspect of the sample

	n	Correct Worse-than-Chance	Correct Different - than-Chance	Improvement ⊿%			
Group A Honest with information	29	100%1	100%1	Na			
Group B Deceptive to hide information	29	34%1	79% <sup>2</sup>	45%			
Group C Honest without information	29	97% <sup>1</sup>	<b>97%</b> <sup>1</sup>	Na			
<ol> <li><sup>1</sup> - statistically sig</li> <li><sup>2</sup> - statistically sig</li> <li>34% from the Wo</li> </ol>	<ul> <li><sup>1</sup> - statistically significant at the .05 level, representing improvement over the chance rate of 50%</li> <li><sup>2</sup> - statistically significant at the .05 level, representing improvement over the observed rate of 34% from the Worse-than-Chance model</li> </ul>						

### Table 2. Comparison of Models.

is that well over 2/3 of the group had been trained in the theory of 2FC testing through both ACID and the TOMI. Debriefing after the study indicated that none of the participants thought to apply their knowledge from ACID or TOMI.

#### Design, Materials and Methods.

In the criminal investigation study, participants witnessed a mock bank robbery video. The first witness was questioned using a Cognitive Interview, and the information in her statement was used to create 42 FC questions. A total of 246 participants were split into three groups and assigned different response strategies for the FC test.

All participants were told to imagine that they were completing the questionnaire as part of police screening of a large number of witnesses to a bank robbery.

Participants in Group A (n=85) were shown a video of the bank robbery and told to respond as honestly and cooperatively as possible. They were instructed that the respondent who demonstrated the most knowledge of the robbery would win a \$50 gift card.

Participants in Group B (n=81) were shown the video. They were to imagine that they knew and were afraid of the thief, and to lie to hide their knowledge of the robbery. They were to respond in a manner that would keep them from standing out or be called for further investigation. They were instructed that the respondent most able to hide their knowledge and blend with those others who genuinely had no knowledge of the event would win a \$50 gift card.

Participants in Group C (n=80) were not shown the video. They were told to do their best to convince the investigators that they had witnessed the event when answering the FC items despite their lack of knowledge.

## Manipulation Check

There was an item at the end of the demographics that asked each participant what experimental condition they were in. It read, "Which of the following best describes your experiences in this experiment?" 1) I was shown a video and I responded honestly. 2) I was shown a video and I lied to act as if I had not seen it. 3) I was not shown the video, but I tried to answer as if I had seen it. The response of all participants to this question aligned with their assigned experimental condition, and therefore each participant's data was used for further analysis.

#### Results

According to the binomial statistic for 42 items, chance performance (alpha = .05) forms an inclusive band between 19 and 26.

### Group A

The mean number of items correct for this group was 41.56 (sd = 3.06). Thus, the items chosen were simple enough that those who were instructed to respond honestly answered all or almost all of the items correctly.

Worse-than-Chance. The resultant pattern of responding led to the same classification accuracy using the Worse-than-Chance and Different-than-Chance models, 100% or 85 of 85 participants accurately classified. The accuracy rate of each of these was compared to the accuracy expected due to chance (50%). The null hypothesis was that the observed rate of 100% was not a significant improvement over the expected chance rate of 50%. The alternate hypothesis was that the observed rate of 100% was a significant improvement over the expected chance rate of 50%. The results indicated that we should reject the null ( $r_{phi}$  = .81 , d = 2.83, p < .01). Thus, the observed rate of 100% was a statistically-significant improvement over chance for both models.

Comparison to Different-than-Chance. No comparison is necessary, due to identical results ( $r_{phi} = .81, d = 2.83, p < .01$ ).

#### Group B

The mean number of items correct for this group was 16.77 (sd = 11.23). Fifty-seven scored below chance (<16), 11 scored within the chance range, and 13 scored above chance (>26).

*Worse-than-Chance.* Fifty-seven of 81 participants who possessed information but lied to hide it were accurately classified as having information using the Worse-than-Chance model. This accuracy rate was compared to the accuracy expected due to chance. The null hypothesis was that the observed rate of 70% was not a significant improvement over the chance rate of 50%. The alternate hypothesis was that the observed rate of 70% was a significant improvement over the chance rate of 50%. The results indicated that we should reject the null ( $r_{phi} = .23$ , d = .46, p < .01). Thus, the observed rate of 70% was a statistical-ly-significant improvement over chance.

Comparison to Different-than-Chance. Seventy of the 81 participants who possessed informa-

	Participants Below Chance (x<16)	Participants Within Chance (16≤x≤26)	Participants Above Chance (x>26)	Participants Different than Chance (x<16)+(x>26)
Group A (85) Honest information	0%	0%	100%	100%
Group B (81) Deceptive hide information	70%	13%	16%	86%
Group C (80) Honest without information	0%	95%	05%	05%

Table 3. Number of Respondents performing Below, Within, Above, and Different From, the Chance Range.

tion but attempted to hide it were accurately classified as having information using the Different-than-Chance model. The performance of the Different-than-Chance model was compared to that of the Worse-than-Chance model. For these calculations, the null hypothesis was that the observed rate of 86% from the Different-than-Chance model was not a significant improvement over the observed rate of 70% from the Worse-than-Chance model. The alternate hypothesis was that the rate of 86% was a significant improvement over the rate of 70%. The results indicated that we should reject the null ( $r_{phi}$  = .21, d = .43, p < .01). Thus, the 86% performance of the Different-than-Chance model led to a statistically-significant improvement over the performance of the Worse-than-Chance model.

Strategy Check. The mean score for those participants who scored above chance from Group B (38.0, sd = 3.6, 95% CI = 36.4-39.6) was compared to the mean score of Group A (41.56, sd = 3.0, 95% CI = 41.1-42.0). The null hypothesis was that there was no difference in total score between Group A and those in Group B who had chosen to score above chance. The alternate hypothesis was that there was a difference between the two groups. A 1x2 ANOVA indicated a significant difference between groups (F(1,95) = 8.8, p<.01), eta squared = .084). Thus, those from Group B who chose to answer at an above-chance level were not behaving the same as those from Group A.

## Group C

The mean number of items correct = 21.41 (sd=4.23). For this group, 0 participants scored below chance (<16), 76 scored within the chance range, and 4 scored above chance (>26).

*Worse-than-Chance.* Seventy-six of the 80 deceptive participants who did not possess information but presented as if they did, were classified as not having information (95.0%). The accuracy rate was compared to that expected by chance. The null hypothesis was that the observed rate of 95% was not a significant improvement over the rate expected by chance (50%). The alternate hypothesis was that the observed rate of 95% was a significant improvement over the rate of 95%. The results indicated that we should reject the null ( $r_{phi} = .67, d = 1.80, p < .001$ ). The Worse-

than-Chance model was significantly better than chance.

Comparison to Different-than-Chance. No comparison was necessary, due to identical results (rphi = .67, d = 1.80, p < .001).

## Discussion

The 2FC approach provided an effective method for rapidly screening a large number of witnesses. The interview of the first witness required 15 minutes. It then took another 20 minutes to create the questions based upon her recorded statement. Photocopying required another 15 minutes. Finally, the participants were able to complete the questionnaire. 166 people were screened in 90 minutes. The 2FC allowed for highly effective and efficient sorting of potential witnesses. Those who possessed information and were willing to cooperate were easily separated from those who did not possess information but were attempting to cooperate. This represents a major savings of effort and time compared to taking statements from each witness, so that investigators would be able to focus their efforts.

The 2FC also allowed for detecting 3 of 4 attempts at hiding information. This is probably the most important finding of all, because these witnesses are the ones who would be most likely to be hiding within the larger group. If there was a suspect or accomplice hidden within the crowd, they would be most likely to fall into this group – i.e., they would want to hide their knowledge of the event to avoid being called as a witness.

Overall, the 2FC led to an accuracy rate of 93%. This is an extremely high rate, and the technique requires relatively little time, and no special equipment- only the ability to present questions to people and record their scores. It appears that mass crimes are not going to go away, and therefore investigators need to develop tools to address this situation.

One limitation of this study was the length of the test used. However, the binomial statistic gives good sensitivity and specificity with much fewer observations (as in Experiment 1), and therefore fewer questions should be used in future applications.

Table 4. Comparison of Models.

		Correct	Correct	Improvement				
	n	Worse-than-Chance	Different - than-Chance	⊿%				
Group A Honest with information	85	100%1	100%1	Na				
Group B Deceptive to hide information	81	70%1	86% <sup>2</sup>	16%				
Group C Deceptive without information	80	<b>95%</b> <sup>1</sup>	<b>95%</b> <sup>1</sup>	Na				
<ul> <li><sup>1</sup> - statistically sig</li> <li><sup>2</sup> - statistically sig</li> <li>70% from the We</li> </ul>	<ul> <li><sup>1</sup> - statistically significant at the .05 level, representing improvement over the chance rate of 50%</li> <li><sup>2</sup> - statistically significant at the .05 level, representing improvement over the observed rate of 70% from the Worse-than-Chance model</li> </ul>							

#### **General Discussion**

This project successfully extended 2FC testing to the forensically relevant situations of juror selection and screening large groups of witnesses. This technique is relatively simple, quick, and provides a mechanism akin to other credibility assessment techniques, in that it helps focus effort for follow-up questioning or investigation. Just as with malingering applications, these uses of 2FC do not demonstrate honesty or deception without consideration of context and motivation. However, the combination of general presentation style (claiming knowledge) plus test results (chance performance) is suggestive of deception and indicates the need for follow-up by an investigator, such as more in-depth interviewing. Similarly, claiming lack of knowledge plus different-than-chance performance is suggestive of deception, and indicates a need for follow-up. It is an efficient mechanism for decision-making and allocation of resources for more in-depth investigation/assessment.

#### **Comparison of Models**

#### Hidden Information and Chance

This project drove home the fact that those who possess information find it difficult to hide this knowledge on 2FC. Respondents describing memory for complex events find it difficult to perform in the range expected by chance. This happens even when they have been trained in the theory of 2FC testing (Orthey, Vrij, Meijer, Leal, & Blank, 2018). This happened in the current studies with those who had taken Psychology and Law.

Two basic strategies emerged, parallel to what has been seen in deception during investigative interviews. Many of the respondents simply miss a few questions related to the target information, denying knowledge of critical items but still obtaining a score that is better than chance. Those respondents missed critical items such as, "What weapon did the perpetrator use to harm the child, gun or knife?", or "Who threatened the bank teller, the man with curly hair or the woman beside him?"

Another subset missed a large proportion of items and are therefore in the worse-thanchance range. These respondents are using a strategy that is analogous to denying knowledge of the target event (Colwell et al., 2013; Derosa et al., 2019). One group is attempting to tell as much of the truth as possible, while still lying. The other is attempting to tell as little of the truth as possible. Both are effective strategies to avoid releasing sensitive information and avoiding mistakes and contradictions (Colwell, et al, 2013).

#### Parametric versus Nonparametric testing

There has been at least one paper that examined the Different-than-Chance approach using z-scores (Orthey et al., 2017). This has the potential for higher sensitivity and specificity, at the expense of practicality and efficiency. In order to use z-scores, every person (or at least a group of about 30) needed to be interviewed so that estimates for the mean and standard deviation could be created regarding the number of items answered incorrectly. Those who scored beyond a set difference score were highlighted and would be targeted for additional investigation. With a nonparametric test, the cut scores are known a priori. There is no need to test other people, or to do any on-site analysis. The number of items automatically determines the cutoffs. Thus, nonparametric interpretation is simpler, more generalizable, and requires less statistical knowledge to apply (Colwell & Sjerven, 2005)

#### **Item Creation**

This is an issue that those new to 2FC testing sometimes struggle with. It is not difficult, so much as something that causes anxiety in those who are new. The optimal item is one that: A) is so easy that just about everyone who is honest will be able to answer it correctly, B) both is relevant, and appears relevant, to the topic at hand, and C) there should be some items that appear to be critical items, so that the deceptive respondent may worry about admitting knowledge to them.

Other applications of 2FC. These studies showed that 2FC can be effective in applications regarding episodic memory. This technique has long been used as a mechanism to induce increased cognitive effort for deceivers while facilitating recall for honest respondents, and thereby triggering Differential Recall Enhancement (Colwell, Hiscock-Anisman, Memon, Taylor, & Prewett, 2008; Colwell, Hiscock-Anisman, Memon, Rachel, & Colwell, 2007; De Rosa et al., 2019). 2FC has always been proposed for use as an additional indicator of credibility when sufficient evidence is available as part of the ACID system. The current evidence further suggests 2FC as part of a multi-pronged approach (Morgan, Rabinowitz, Leidy, & Coric, 2014).

#### Weaknesses of Studies

This project combined real newspaper articles regarding real cases with a mock juror selection, and a video for a mock bank robbery with a mock investigation. Therefore, as in many credibility assessment studies, there is the potential for a difference in motivation and/ or concentration on behalf of participants. Because this happens routinely in such studies, scientists have come to recognize that the behavior of participants tends to be generally the same under lab and field settings. This is, in part, because student participants tend to be highly motivated. More importantly, in cases of genuine evaluations, increased motivation and fear can enhance the very behaviors that are being targeted. Increased motivation leads to increased impression management, and this makes attempted deception easier to detect (Colwell, Hiscock, & Memon, 2002; Colwell et al., 2008; Colwell et al., 2009; Colwell et al., 2012). This has been observed with investigative interviewing and with the assessment of malingering, and investigative interviewing and credibility assessment (Colwell, Colwell, Perry, Wasieleski, & Billings, 2008; Colwell et al., 2002; Colwell & Colwell, 2011). However, this is an empirical issue, and motivation should be manipulated directly.

## Conclusions

The 2FC approach is one of the most widely-used techniques in clinical and forensic psychology. It is easily adapted to a number of situations beyond the original assessment of malingering. These include any situation that should generate a complex episodic memory. When combined with nonparametric statistics, and a Different-than-Chance interpretive model, 2FC does not require having a pretested group, and 2FC provides a powerful tool to help focus follow-up efforts such as an investigative interview. In this way, it improves efficiency for processing large groups, and allows investigators to corroborate honest respondents and highlight those who are attempting to hide their knowledge.

Polygraph & Forensic Credibility Assessment , 2021, 50 (1)

#### References

- Binder, L. M., Larrabee, G. J., & Millis, S. R. (2014). Intent to fail: Significance testing of forced choice test results. *Clinical Neuropsychologist*, 28(8), 1366–1375. https://doi.org/10.1080/ 13854046.2014.978383
- Cliffe, M. J. (1992). Symptom-validity testing of feigned sensory or memory deficits: A further elaboration for subjects who understand the rationale. *British Journal of Clinical Psychology*, 31(2), 207–209. https://doi.org/10.1111/j.2044-8260.1992.tb00985.x
- Colwell, K., Colwell, L. H., Perry, A. T., Wasieleski, D., & Billings, T. (2008). The test of malingered incompetence (TOMI): A forced-choice instrument for assessing cognitive malingering in competence to stand trial evaluations. *American Journal of Forensic Psychology*, *26*(3).
- Colwell, K., Hiscock, C. K., & Memon, A. (2002). Interviewing techniques and the assessment of statement credibility. *Applied Cognitive Psychology*, 16(3). https://doi.org/10.1002/acp.788
- Colwell, K., & Sjerven, E. R. (2005). The "coin-in-hand" stratagem for the forensic assessment of malingering. *American Journal of Forensic Psychology*, 23(1).
- Colwell, K, Hiscock-Anisman, C., Memon, A., Taylor, L., & Prewett, J. (2008). Assessment Criteria Indicative of Deception (ACID): An integrated system of investigative interviewing and detecting deception. *Journal of Investigative Psychology and Offender Profiling*, 4(3), 167–180.
- Colwell, K, Hiscock-Anisman, C., Memon, A., Woods, D., & Michlik, P. (2006). Strategies of impression management among deceivers and truth-tellers: How liars attempt to convince. American Journal of Forensic Psychology, 24(2), 1–9.
- Colwell, K, Hiscock, C., & Memon, A. (2002). Interviewing techniques and the assessment of statement credibility. *Applied Cognitive Psychology*, *16*(3), 287–300.
- Colwell, Kevin, Hiscock-Anisman, C. K., Memon, A., Rachel, A., & Colwell, L. (2007). Vividness and spontaneity of statement detail characteristics as predictors of witness credibility. *American Journal of Forensic Psychology*, 25(1), 5–30.
- Colwell, Kevin, Hiscock-Anisman, C., Memon, A., Colwell, L. H., Taylor, L., & Woods, D. (2009). Training in Assessment Criteria Indicative of Deception to Improve Credibility Judgments. *Journal of Forensic Psychology Practice*, 9(3), 199–207. https://doi.org/10.1080/15228930902810078
- Colwell, Kevin, & Sjerven, E. R. (2005). The "coin-in-hand" stratagem for the forensic assessment of malingering. *American Journal of Forensic Psychology*.
- Colwell, L., Colwell, K., Hiscock-Anisman, C. K., Hartwig, M., Cole, L., Werdin, K., & Youschak, K. (2012). Teaching Professionals to Detect Deception: The Efficacy of a Brief Training Workshop. *Journal of Forensic Psychology Practice*, 12(1), 68–80.
- Colwell, L.H., & Colwell, K. (2011). Assessing feigned cognitive impairment in defendants hospitalized for competency restoration: Further validation of the TOMI. *Journal of Forensic Psychology Practice*, 11(4). https://doi.org/10.1080/15228932.2011.562804
- Colwell, Lori H., & Colwell, K. (2011). Assessing Feigned Cognitive Impairment in Defendants Hospitalized for Competency Restoration: Further Validation of the TOMI. Journal of Forensic Psychology Practice, 11(4), 293–310. Retrieved from http://www.tandfonline.com/doi/abs/1 0.1080/15228932.2011.562804



De Rosa, J., Hiscock-Anisman, C., Blythe, A., Bogaard, G., Hally, A., & Colwell, K. (2019). A comparison of different investigative interviewing techniques in generating differential recall enhancement and detecting deception. *Journal of Investigative Psychology and Offender Profiling*, *16*(1). https://doi.org/10.1002/jip.1519

Dusky v. United States (1960). U.S.C.

- Gottfried, E. D., Hudson, B. L., Vitacco, M. J., & Carbonell, J. L. (2017). Improving the Detection of Feigned Knowledge Deficits in Defendants Adjudicated Incompetent to Stand Trial. Assessment, 24(2), 232–243. https://doi.org/10.1177/1073191115605631
- Hines, A., Colwell, K., Hiscock-Anisman, C., Garrett, E., Ansarra, R., & Montalvo, L. (2010). Impression management strategies of deceivers and honest reporters in an investigative interview. European Journal of Psychology Applied to Legal Context, 2(1).
- Hiscock, C., Branham, J., & Hiscock, M. (1994). Detection of feigned cognitive impairment: The two-alternative forced-choice method compared with selected conventional tests. *Journal* of Psychopathology and Behavioral Assessment, 16, 95-110. Retrieved from http://www. springerlink.com/index/L1440H515L126261.pdf
- Hiscock, C., Rustemier, P., & Hiscock, M. (1993). Determination of Criminal Responsibility Application of the Two-Alternative Forced-Choice Stratagem. *Criminal Justice and Behavior*, *20*(4), 391– 405.
- Hiscock, M., & Hiscock, C. (1989). Refining the forced-choice method for the detection of malingering. *Journal of Clinical and Experimental Neuropsychology*, *11*(6), 967–974. Retrieved from http:// www.tandfonline.com/doi/abs/10.1080/01688638908400949
- Jelicic, M., Merckelback, H., & van Bergen, S. (2004). Symptom validity testing of feigned amnesia for a mock crime. *Archives of Clinical Neuropsychology*, 19, 525–531.
- Kapur, N. (1994). The coin-in-the-hand test: A new "bedside" test for the detection of malingering in patients with suspected memory disorder [8]. Journal of Neurology Neurosurgery and Psychiatry. BMJ Publishing Group. https://doi.org/10.1136/jnnp.57.3.385
- McElfresh, E.N., Hiscock-Anisman, C.K., James-Kangal, N., Maciel, V., Gavigan, B., Colwell, K. Forced-Choice Testing for Deception Detection in Jury Selection and Investigation of Mass Crimes (2020). A poster presented at the American Psychology and Law Society. New Orleans.
- Morgan, C., Rabinowitz, Y., Leidy, R., & Coric, V. (2014). Efficacy of Combining Interview Techniques in Detecting Deception Related to Bio-threat Issues. *Behavioral Sciences and the Law, 32*, 269–285.
- Orthey, R., Vrij, A., Leal, S., & Blank, H. (2017). Strategy and Misdirection in Forced Choice Memory Performance Testing in Deception Detection. *Applied Cognitive Psychology*, 31(2), 139–145. https://doi.org/10.1002/acp.3310
- Orthey, R., Vrij, A., Meijer, E., Leal, S., & Blank, H. (2018). Resistance to coaching in forced-choice testing. *Applied Cognitive Psychology*, 32(6), 693–700. https://doi.org/10.1002/acp.3443
- Pankratz, L. (1979). Symptom validity testing and symptom retraining: procedures for the assessment and treatment of functional sensory deficits. *Journal of Consulting and Clinical Psychology;* .... Retrieved from http://psycnet.apa.org/journals/ccp/47/2/409/



Schroeder, R. W., Peck, C. P., Buddin, W. H., Heinrichs, R. J., & Baade, L. E. (2012). The coin-inthe-hand test and dementia: More evidence for a screening test for neurocognitive symptom exaggeration. *Cognitive and Behavioral Neurology*, 25(3), 139–143. https://doi.org/10.1097/ WNN.0b013e31826b71c1



## Response Onset Latencies of Electrodermal, Cardiovascular and Vasomotor Responses During Field Polygraph Testing.

Donald J. Krapohl<sup>1</sup> Karen Halford<sup>2</sup> Tim Benson<sup>3</sup> Abbe Mayston<sup>4</sup> Donnie W. Dutton<sup>5</sup>

#### Abstract

One of the important characteristics for determining whether a physiological response is associated with a stimulus is whether it is timely and can arguably be attributed to the presented stimulus. While there are scoring policies taught in most polygraph schools that address onset latencies, we could find very little published normative data on which to support those policies. To address this open question, field polygraph examiners manually measured and recorded onset latencies associated with relevant and comparison questions from 154 field cases for the electrodermal, cardio-vascular and vasomotor channels. From those measurements we summarized the data to suggest recommendations for polygraph Response Onset Windows in manual scoring of 1.2 - 8.0 seconds for electrodermal responses, 1.0 - 9.0 seconds for cardiovascular responses, and 2.0 - 9.0 for vasomotor responses. The electrodermal and vasomotor onsets were consistent with previously reported results in research on the Orienting Response. All windows used question onsets as reference for measurements of response onset latency for reasons discussed in the article.

<sup>4</sup> Ms. Abbe Mayston is a Detective Constable and Polygraph Officer, Essex Police, UK.

<sup>5</sup> Mr. Dutton is an APA Past President and the Vice President of the Capital Center for Credibility Assessment.

The views expressed in this article are solely those of the authors and do not necessarily represent those of the APA or the authors' agencies and organizations.

15

 $<sup>^{1}</sup>$  Mr. Don Krapohl is an APA Past President, Education Director for the Capital Center for Credibility Assessment (C<sup>3</sup>A) and School Director for the Behavioural Measures United Kingdom (BMUK) Polygraph Training Centre. Comments and questions can be directed to him at APAkrapohl@gmail.com. Reasonable requests for the original data from this project will be honored.

<sup>&</sup>lt;sup>2</sup> Ms. Karen Halford is a Detective Constable and Polygraph Officer, Cochester Police Station, Essex Police, UK.

<sup>&</sup>lt;sup>3</sup> Mr. Tim Benson was amongst the first police polygraph examiners to be trained for the UK police, in 2014. He works as a full-time polygraph examiner for Hertfordshire Constabulary and sits on the UK police National Polygraph Working Group. Although primarily working with post-conviction sex offenders, he pioneered the expansion of polygraph into other areas of policing including domestic abuse and initiatives aimed at prolific burglars.

We are grateful to the UK police examiners who the contributed cases for this study. In addition to the three UK co-authors there were Alex Tills, Kerry Clarke, Leanne Powell, and Ned Kelly. This project would not have been possible without them. We also thank the reviewer for helpful comments and suggestions.

In polygraph testing, inferences of veracity are derived through a systematic process in which physiological responses elicited by different kinds of test questions are compared to one another. The types of physiological responses found predictive of veracity are well established and described elsewhere (Bell et al., 1999; Handler et al., 2010; Kircher & Raskin, 1988; Kircher et al., 2005; Nelson et al., 2008). One characteristic of those responses that has less agreement in the polygraph literature is response onset latency, that is, the normal delay between the presentation of the test question and the initiation of a concomitant physiological response. It is the window of time for helping distinguish those reactions that are associated with the test question from those that are spontaneous or unrelated.

To our knowledge there has not been published research in the polygraph context which has led to recommendations for minimum and maximum latency periods, commonly referred to as the Response Onset Window (ROW) in polygraph testing. As such there remain differing opinions among polygraph examiners, schools and writers on what those ROWs should be.

In a different, non-polygraph context there have been numerous investigations of minimum response onset latency in the psychophysiological literature for two of the common polygraph data channels. Researchers investigating the orienting response (OR) have reported some onset latency statistics for the electrodermal (e.g., Sjouwerman & Lonsdorf, 2018) and vasomotor responses (Biferno & Dawson, 1978; Furedy, 1968). We were unable to locate similar research in the OR literature for the pneumograph and cardiograph channels.

The available psychophysiological research has not directly assessed ROWs for polygraph questions. Rather, researchers have to date used tones, lights, images, words and other abrupt stimuli to measure the period between stimulus and response. With these kinds of stimuli, it is possible to determine the absolute shortest onset latency because they require minimal cognitive processing. Indeed, the responses are almost entirely reflexive. When OR research uses startle prompts, a subsequent response may simply signal the examinee has noticed the stimulus. The OR onset may appear even before the test subject knows what the stimulus is.

Polygraph testing, in contrast, relies on the examinee assessing the salience of the test questions, the stimuli. The appraisal process to determine salience requires the examinee have some threshold amount of information about the question as it is being presented during testing. An understanding of the test question may not take place before a significant portion of it has been presented. Because the assumptions about polygraph stimuli are markedly different from those in research on onset latency for the OR, it seems likely there will be some differences in response onset windows between these two paradigms.

A physiological response may either be an indicant of a psychological event or it may be a purely random reaction. To be of value in polygraph testing there must be confidence the response is linked to the test question. As a response onset latency departs from the normative range there can be less certainty the test question is responsible for the physiological response. It became our interest to seek the normative range for response onsets for three of the most easily measured polygraph data channels: electrodermal, cardiovascular and vasomotor. To obtain data we took advantage of our access to field cases conducted by police polygraph examiners in the United Kingdom (UK).

#### Method

#### Cases

Measurements were taken from polygraph cases performed by UK police polygraph examiners conducting voluntary polygraph examinations of sex offenders living in the community. Seven police examiners volunteered to participate in the study and were requested to provide measurements taken from the first 25 polygraph examinations they conducted in the calendar year 2020. Due to the shutdown from the pandemic that year only four of the seven examiners completed 25 examinations in the specified time period. The final tally was 154 cases, all conducted on LX5000 Lafayette computerized instruments.

The average age of the examinees was 46.2 years (sd = 14.8), with a range of 21 to 78 years. The digital voice was used in 72 of the cases during the test phase and in the remaining 82 cases the examiners presented the questions using their own voice.

The cases were all screening examinations using a version of the Air Force Modified General Question Technique with either two or three relevant questions and two, three or four comparison questions. Probable-lie comparison questions were employed in 133 cases and directed lies in 21.

#### **Exclusionary** Criteria

Cases that resulted in fewer than three charts were excluded as were those affected by detected countermeasures. Likewise, tracings that were distorted by movements or reactions that could have been caused by sources unrelated to the test questions were not used.

## Measurements

Only onset latencies of relevant and comparison questions were recorded. The polygraph examiners conducting the cases made all the measurements and entered their data into standardized Microsoft Excel spreadsheets. All data were anonymized in the spreadsheet using alphanumeric codes to represent the cases. The measurements used the timestamp at the beginning of the question presentation to the point where the physiological response began. The examiners were instructed on the use of the caliper feature in the Lafayette software for more precise and reliable measurements. For the electrodermal channel the measurement period began at 1.2 seconds after question onset to the beginning of the inflection point signaling the start of the reaction, but not later than 10.0 seconds from where the question began. The minimum onset of 1.2 seconds was selected because of the evidence in the OR literature that EDR onsets have not been shown to be shorter than this period. The use of the 10-second maximum was arbitrarily set, based on the experience of the researchers. In the cardiograph the measurement period was from 0.2 seconds to 10.0 seconds from question onset. The cardiograph measurement period was broad owing to a lack of OR research to guide the choice of window. The measurement period for the vasomotor response was from 2.0 seconds to 9.0 seconds after question onset, again, based on the previously cited reports from the OR literature. Instructions given to the examiners can be found in Appendix A.

## Procedure

The sample consisted of onset measurements for 2,011 electrodermal, 1,720 cardiovascular and 1,068 vasomotor responses. We performed standard descriptive statistical treatments to the onset latency data using Microsoft Excel statistical packages.

## Results

The median onset latency of electrodermal responses was 3.3 seconds. The standard measure of variability, the standard deviation, was not applied to these data because the frequency distribution of latencies, as can be seen in Figure 1, had a significant skew  $(1.57)^6$ . As an alternative, a simple count found 99% of the latencies were 8.0 seconds or less. Figure 1 shows the frequency of all 2,011 EDR onset latencies in 0.2-second increments.



<sup>&</sup>lt;sup>6</sup> The Excel software formula for calculations of skewness is  $\frac{n}{(n-1)(n-2)}\sum_{s} \left(\frac{x_{s}-\bar{x}}{s}\right)^{s}$ . Whether the degree of skewness is a cause for

concern can be estimated by the formula  $\sqrt{6/N}$ . If the measure of skewness is more than twice this value, the distribution is considered non-symmetrical. From https://www.stattutorials.com.



Figure 1. Frequency of Response Onset Latencies for 2,011 Electrodermal Responses.

The cardiograph data showed a median value for onset latency of 3.2 seconds and a skew of 0.88. A count showed 99% of the latencies were 8.4 seconds or less. Figure 2 is a frequency distribution of all 1,720 cardiograph response latencies partitioned into 0.2 second increments.

Figure 2. Frequency of Response Onset Latencies for 1,720 Cardiovascular Responses.



There was a relative delay in vasomotor responses compared to the other two data channels. The vasomotor data had a median onset latency of 4.3 seconds and a positive skew at

0.61. A count showed 99% of the latencies were 8.8 seconds or less. Figure 3 displays the frequency of vasomotor onset latencies for the 1,068 measurements in this sample.



Figure 3. Frequency of Response Onset Latencies for 1,068 Vasomotor Responses.

#### Discussion

The present findings for the electrodermal channel suggest onset latencies in polygraph testing are not identical to those found in the OR literature. The polygraph electrodermal onset latencies appear to be longer and the distribution more skewed. In their recent article on the OR for the electrodermal response, Sjouwerman and Lonsdorf (2019) reported a mean onset latency to startle prompts of 1.9 seconds. By way of comparison, our polygraph data produced a median latency of 3.3 seconds. The longer latency period with polygraph data might be expected when considering that polygraph responses are induced by the personal significance of the test question, which requires the presentation of enough of the test question for the examinee to appreciate the degree of personal significance the question contains. The threshold for recognition of the intent of a polygraph test question can be a function of the examiner's speaking rate or whether the earliest words in the polygraph question signal the meaning of the entire question. Typically, these are not factors in OR research that use stimuli such as tones, lights, images, or single words. Longer electrodermal response onset latencies are therefore expected in polygraph testing.

The polygraph electrodermal data, combined with that from the OR research may provide guidance on polygraph scoring rules as they apply to onset latency. The OR research, as well as investigations of nerve conduction transmission speeds (Lim et al., 2003), do not support the beginning of electrodermal responses recorded at the extremities in less than 1.2 seconds from stimulus onset. At the other end, the present research found that 99% of all electrodermal responses to polygraph questions began within 8.0 seconds of the question onset. Therefore, we submit that an electrodermal ROW from 1.2 seconds to 8.0 seconds is consistent with what is currently known about electrodermal response onset latencies.

To our knowledge there is no previously published normative data on onset latency for the manner the cardiovascular data are recorded in polygraph testing. Consequently, there is no benchmark against which to compare the present findings. We can report that virtually all 1,720 cardiovascular responses in our field

19

data took place within 8.4 seconds of question onset. If this statistic generalizes, 8.4-second onset latency or the next rounded value (9.0 seconds) could be used as the outer bound of the cardiovascular ROW in polygraph testing when conducting manual scoring. Estimating minimum onset latency, however, proved more challenging.

The response onset of the cardiograph is recognized by the upward rise of the pulse. Pulses are intermittent events unlike the continuous data in the electrodermal channel. Determining a precise response onset time is constrained by the rate in which these events take place. For example, consider a sitting pulse rate of 72 beats per minute. At this rate the individual pulses are about 0.8 seconds apart. Therefore, the onset of a phasic cardiovascular response cannot be visually resolved for any period less than 0.8 seconds. Even the relatively high pulse rate of about 90 beats per minute found in criminal polygraph tests of confirmed deceptive examinees (Ansley & Krapohl, 2000), response onset cannot be resolved more precisely than about 0.7 seconds. Returning to the present data, 99% of the cardiovascular responses took place at or after 0.6 seconds from question onset, a value lower than the temporal resolution of the data. We concede our visually derived data cannot precisely identify a minimum onset latency for this channel, at least for human interpreters alone. It can be agreed that examinees need sufficient information about a test question to assess its personal significance before producing a cardiovascular response. It would be optimistic to suppose the minimum required information can be presented under field conditions in the period between two heartbeats. For this reason, we tentatively propose a minimum of at least 1 second (roughly the period between heartbeats) from question onset is necessary to initiate a cardiovascular response prompted by a test question. In the present data set 97% of all reported latencies were 1 second or longer.

The lower limit of the vasomotor response onset period is not firmly established in the existing literature. Psychophysiologists who have investigated the OR with the vasomotor response have reported using minima of 2.0 seconds (Biferno & Dawson, 1978; Furedy, 1968) or four heartbeats (Stern & Anschel, 1968). Researchers on the Concealed Information Test place the minimum vasomotor onset latency at 3.0 - 4.0 seconds (Matsuda & Nittano, 2018). What has not been reported, however, is the research basis for using these onset minima.

The present data found 99% of the vasomotor response onsets required 1.8 seconds or more. Like the cardiovascular data, however, the vasomotor responses rely on individual pulses to demark response onset, and thereby inherit the same problem of temporal resolution. For this reason the use of the same minimum proffered by psychophysiologists, 2.0 seconds, is not contested by the polygraph data and seems reasonable for use in polygraph testing.

As the maximum onset latency, 99% of the vasomotor responses in the polygraph data began within 8.8 seconds of question onset. This outer bound is not meaningfully different from that reported by Biferno et al. (1978) research, who used 9.0 seconds. Based on the present data and the practices of psychophysiologists, a vasomotor ROW of 2.0 - 9.0 seconds can be defended.

Within these ROWs may be another consideration, that of response stereotype. It is generally recognized that individuals are relatively consistent within themselves in the manner in which they react physiologically. Therefore, the value of long or short response onset latencies in polygraph testing can be further judged according to whether they occur in a window of time that is normal for the person being tested. Normative within-person variability in response latencies for polygraph reactions is the focus of a separate investigation.

As a general observation we would offer that the preferred reference point for all ROWs would be the question onset as opposed to other events, such as the point of the examinee's answer. This suggestion is based on the undisputed fact the polygraph is not a lie detector – it is not the act of lying but the personal significance of the test question that evokes the response. Polygraph testing is a Stimulus – Response paradigm where the stimulus is the test question. As such, anchoring the timing of the physiological response to the onset of the test question is the most reasonable approach.

#### Limitations

All measurements in this project were based on visual assessments by the testing polygraph examiners. Though the examiners were instructed how to use the caliper tool in the Lafayette polygraph software for making the measurements there was no obligation for them to use this tool. Consequently, individual differences may reside in the interpretations of onset latency.

A second limitation is that we used simple counts to establish the normal range of onset latencies. The skewness of the distributions limited the choices for statistical characterizations of data central tendency and distribution. We suspect advanced statistical treatments of the relatively course data are unlikely to have produced substantially different ROWs for use by human examiners relying on visual discernment. For greater confidence in these findings an independent replication is recommended.

#### Summary

With large sample sizes of physiological responses we recorded response onset latencies for the electrodermal, cardiovascular and vasomotor responses in field cases. From the present findings, which are compatible with the findings from related research, the following ROWs can be considered normal for human scorers of polygraph data, starting at question onset: 1.2 to 8.0 seconds for electrodermal response onsets, 1.0 to 9.0 seconds for cardiovascular response onsets, and 2.0 to 9.0 seconds for vasomotor response onsets. We also argue that response onset latency should only be measured against the onset of the test question.

#### References

- Ansley, N. & Krapohl, D.J. (2000). The frequency of appearance of evaluative criteria in field polygraph charts. *Polygraph*, 29(2), 169-176.
- Bell, B.G., Raskin, D.C., Honts, C.R. and Kircher, J.C. (1999). The Utah Numerical Scoring System. *Polygraph, 28*(1), 1 9.
- Biferno, M.A., and Dawson, M.E. (1978). Elicitation of subjective uncertainty during vasomotor and electrodermal discrimination classical conditioning. *Psychophysiology*, *15*(1), 1 8.
- Furedy, J.J. (1968). Human orienting reaction as a function of electrodermal versus plethysmographic response modes and single versus alternating stimulus series. *Journal of Experimental Psychology*, 77, 70 – 80.
- Handler, M., Nelson, R., Goodson, W., and Hicks, M. (2010). Empirical Scoring System: A crosscultural replication and extension study of manual scoring and decision policies. *Polygraph*, 39(4), 200 – 215.
- Kircher, J. C., Kristjansson, S. D., Gardner, M. K. & Webb, A. (2005). *Human and computer decision*making in the psychophysiological detection of deception. University of Utah. Final Report.
- Kircher, J. C. & Raskin, D. C. (1988). Human versus computerized evaluations of polygraph data in a laboratory setting. *Journal of Applied Psychology*, 73, 291-302.
- Lim, C.L, Seto-Poon, M., Clouston, P.D., and Morris, J.G.L. (2003). Sudomotor nerve conduction velocity and central processing time of the skin conductance response. *Clinical Neurophysiology*, 114, 2172 – 2180.
- Matsuda, I., and Nittono, H., (2018). Physiological responses in the Concealed Information Test. In (J.P. Rosenfeld, Ed) Detecting Concealed Information and Deception: Recent Developments. Academic Press: New York.
- Nelson, R., Krapohl, D.J., and Handler, M. (2008). Brute force comparison: A Monte Carlo study of the Objective Scoring System version 3 (OSS-3) and human polygraph scorers. *Polygraph*, 37(3), 185 – 215.
- Sjouwerman, R., and Lonsdorf, T.B. (2019). Latency of skin conductance responses across stimulus modalities. *Psychophysiology*, 56(4), 1-11.
- Stern, R.M. and Anschel, C. (1968). Deep inspirations as stimuli for responses of the autonomic nervous system. *Psychophysiology*, 5(2), 132 141.



## Appendix A

Thank you for volunteering to participate in the onset latency project. With your help we can develop evidence-based onset windows for the cardio, EDA and PPG signals. This is the first project of its type, so whatever we discover will be new and also useful to our fellow examiners.

Keep this page somewhere handy in case you have any questions. In it I will try to describe all of the steps in sufficient detail that you will feel comfortable with what you are doing. Let's start first with the minimum criteria for the cases we will use.

#### Exclusionary Criteria for Cases

- 1. Fewer than 3 interpretable charts
- 2. Countermeasure cases

#### **Exclusionary** Criteria for **EDA Tracings**

- 1. Distorted by movements
- 2. EDRs induced by DBs
- 3. EDRs that began before the question was presented
- 4. For EDRs, reactions that begin in less than **1.2 seconds or after 10 seconds** from question onset taken at the inflection point when recorded in manual mode. See below.



- 5. Exclusionary Criteria for Cardio
- a. Distorted by movements
- b. Rises that begin in less than 0.2 seconds or more than 10 seconds from question onset
- c. Rises that immediately follow a PVC. See below.



23

d. Rises that appear to be part of a repeating cycle and not necessarily associated with the test question. See below:



- 6. Exclusionary Criteria for PPG
  - a. Distorted by a DB or movements

b. Constrictions that begin in less than 2 seconds or more than 9 seconds from question onset.

c. Constrictions after PVCs or that are part of a recurring pattern. See below.



Displaying the data

- 1. For the EDA, turn up the gain until the response onset inflection is clearly seen.
- 2. For cardio and PPG, adjust the gain such that the pulse amplitude averages about 2-3 chart divisions.

Use of the calipers

- 1. Right click on the chart.
- 2. On the dropdown menu click on Show Calipers
- 3. Our interest is only in the timing of the reactions and so only the two vertical caliper bars will be used.
- 4. Move the calipers to the question with the responses you want to measure. This can be done by clicking on either the <Previous or Next> button in the Caliper Statistics dropdown menu.
- 5. Move the right vertical bar to where the reaction began. The left bar will remain at the question onset point. You can expand the chart if needed to help find response onsets.



- 6. Record the time, which can be found on the dropdown menu under Dimensions, the Width. See below for examples.
  - a. EDA: First inflection between 1.2 and 10 seconds after question onset.



b. Cardio: The first risen diastolic point for a rise of three or more pulses.



c. PPG: First constricted pulse of at least three consecutive constricted pulse taking place between 2 and 9 seconds from question onset.

Ca	liper Stati	istics	_	$\times$			
	Dimensions		Line Lengt	h			
$\leq$	Width: Height:	4.6 sec 1.2 div	P1:	0 0			
E	DA (Ohms	;) (1	CP (mmHG	)			
· 1	Avg: Min: Max:	NA NA NA	Avg: Min: Max:	NA NA NA			
	leart Rate	(bpm)	Resp. Rate	e (cpm)	cert	1-1-1-1-	
	Avg:	NA	Avg:	NA			
0	Cardio Bas	eline Chara	cteristics		A-A-A-		<b>   </b>
M f	Area unde Time to reti Rate of arc	r curve: urn: usal:	NA NA	blks sec blk/s	AAA	VVVVV	WWW
	< Previ	ous	Nex	d >		R4	

#### Judgment

- 1. There will be instances where the data are somewhat noisy, and you may experience difficulties pinpointing the reaction onset. If you are not confident in making a decision, or the data are too contaminated, record an "N" in the data sheet. Note, the PPG will likely have the most Ns and EDA the least.
- Similarly, if there is no discernible reaction, or the reaction began before or after the on set window (EDA=1.2 to 10 sec, cardio=0.2 sec to 10 sec, PPG=2.0 - 9.0 sec), record an "N" in the data sheet.

#### **Data Sheet**

- 1. Use only the first three scorable charts from your case.
- 2. The data sheet is set up to record latencies for the EDA, cardio and PPG for up to 3 RQs and 4 CQs.
- 3. There is no need to fill in the blocks where there is no data: e.g., third block for an RQ when there was only 2 RQs asked on the chart.
- 4. Use the numbers you see in the Caliper Statistics, which will always be time recorded in 10ths of a second.
- 5. An RQ latency is an RQ latency and a CQ latency is a CQ latency. It does not matter in what order you record the latency just so long as you correctly place it in an RQ block or a CQ block. For example, if for the first chart you record the latency information for R4 in the first RQ block on the data sheet you can record R6 in the first RQ block for the next chart. Our plan is to average the latencies of all RQs separately from the latencies for all CQs, but we don't care about whether it is R4 or R6. An example of a filled-in data sheet is below.

kotoš le	inve 👓 🗄 🤊 Home insert	- C - D B G = Page Layout Formulas	Data Review V	lew Help Acrobat Яse	arch	Data sheet - Exc	:6				Donald Krapoh
ermal	Page Break Page Cur Preview Layout W Workbook Views	dam Gridines Headings	Bar Q Doom 100% Zoom Stele	The New Arrange Freeze Divide	CO View Side by 5 D) Synchronous 5 ie ân Reset Window Window	ide Scroling Position Windows *	Macros				
	• I ×	√ ≴ R	C	D	F	E	G	Ц	1		K
1	Case	Case ID	Chart	Channel	C	R	C	R	C	R	C
2	1	DK246	1	EDA	1.8	N	6.4	2.4	2.7		
3				Cardio	4.0	3.0	2.1	2.8	4.1		
4				PPG	Ν	Ν	2.7	3.8	Ν		
5			2	EDA	5.0	2.4	2.2	1.9	Ν		
6				Cardio	3.9	1.9	3.0	4.4	5.5		
7				PPG	4.5	Ν	4.6	2.8	N		
8			3	EDA	7.0	2.8	1.6	Ν	2.8		
9				Cardio	4.9	5.2	5.2	N	4.1		
10				PPG	3.1	Ν	N	Ν	4.3		



- 6. The data sheet should be printable in landscape layout in case you want to record the numbers manually and enter them into the Excel file later.
- 7. I am asking for 25 cases from each volunteer. They should be 25 consecutive cases beginning on January 1, 2020. If a case is not useable because of cms or the session ended before three charts are collected, it still counts toward your 25. Simply enter "N" in every block in the data sheet for that case.

Questions:

1. If you have any questions about a measurement, just make a comment on the data sheet.

## Literature Review and Analysis of the Multi-facet Hypothesis and the Evaluation of Independent Target Questions

Raymond Nelson,

Mark Handler,

and David C. Raskin

#### Abstract

Published literature was reviewed for studies investigating the Multi-Facet Hypothesis (MFH). It postulates responses to PDD target questions vary independently when the questions address different levels of involvement in known or alleged incident. Implicit within the MFH is responses to individual questions serve to discriminate deception and truth-telling at rates greater than chance. The MFH suggests PDD test questions will not only discriminate guilty from innocent persons, but may also discriminate the behavioral role or level of involvement of a guilty person. The MFH is essentially the hypothesis of a multiple issue diagnostic exam, and is based on decision rules using question subtotal scores. Overall accuracy with the MFH was (.78) and lower than accuracy using grand total scores (.89). These results provided limited support for high test sensitivity with guilty persons when classifications were made under the MFH, but with test specificity for innocent persons not exceeding the range of chance. In contrast, decisions using grand total scores provided test sensitivity and specificity that were significantly greater than chance. The hypothesis of effective role discrimination using multiple issue diagnostic exams is inconsistent with the clear trend in the published scientific literature and is therefore not supported. Furthermore, these results do not support polygraph field practices that attempt to determine that a person has been deceptive to one or more questions and also truthful to one or more questions within a single examination. This literature review and associated analysis supports the evidence that multi-facet questions are nonindependent.

The authors are indebted and grateful to Dr. Joe Stainback and Dr. Stewart Senter for their reviews, edits, comments and critique on earlier versions of this manuscript.



The Multi-Facet Hypothesis (MFH) posits that responses to psychophysiological detection of deception (PDD) test questions vary independently when the test questions address different behavioral roles or different levels of involvement in a known or alleged incident. This hypothesis suggests that the effectiveness of discrimination of guilt and innocence can be optimized through the use of decision rules that treat the questions independently. Multi-facet exams are a form of diagnostic exam<sup>1</sup> conducted in the context of known or alleged events. They differ from other event-specific diagnostic examinations in that they attempt to assess both guilt vs. innocence and the level of involvement of guilty examinees. The MFH holds that asking about different roles or actions is sufficient for responses to vary independently<sup>2</sup>. For example, a thief, accomplice, and lookout may all be involved in a crime though they have different roles. Multi-facet test questions may also be formulated to describe either primary/direct involvement or secondary/indirect involvement, such as knowledge of the identity of the guilty criminal, plans or details regarding a known or alleged crime, or the details of a known or alleged crime. Similarly, use of words such as helping, planning and participating in a crime

may be thought of as three different roles (or different levels of involvement) that may be distinct from, or secondary to, primary or direct involvement in a behavioral action such as a theft, robbery, assault, or other crime.

Field practices have traditionally emphasized the scoring and interpretation of multi-facet criminal investigation exams in the same manner as multiple-issue screening exams for which both criterion variance and response variance<sup>3</sup> of the target questions are assumed to be independent (Department of Defense, 2006; Department of Defense, 2006b). In practical terms, the MFH is the hypothesis of a multiple issue diagnostic exam .<sup>4</sup> The null-hypothesis to the MFH will maintain that the use of different action verbs is insufficient to achieve independent response variance to PDD test stimuli, and that interpreting responses to test questions with an assumption of independent variance does not optimize the discrimination between guilt and innocence, compared to interpretation with no assumption of independence, nor does it identify the behavioral role or level of involvement among guilty persons at rates greater than chance (Department of Defense, 2006; Department of Defense, 2006b).

<sup>&</sup>lt;sup>1</sup>Similar to other forms of testing, psychophysiological detection of deception (PDD) examinations can be thought of as belonging to the category of either diagnostic or screening exams (American Polygraph Association, 2011; Krapohl & Stern, 2003). PDD diagnostic tests are limited in scope to a single known or alleged behavioral incident, and interpreted at the level of the test as a whole. In contrast, PDD screening examinations are those tests conducted in the absence of a known or alleged problem. Screening exams are often formulated to simultaneously investigate multiple issues of concern for which the external criterion state is assumed to vary independently.

<sup>&</sup>lt;sup>2</sup>An assumption of independent variance, in the scientific context, holds that responses to individual questions are influenced by no factors in common with each other – that responses to different stimuli have no shared sources of variance. This assumption is inherently compromised by the fact that responses to all individual questions have in common the examinee. The practical result of this is that it is not possible to conclude, with any scientific confidence, that an examinee has been deceptive to one or more questions while being truthful to other questions within the same examination. All comparison question test formats represent a form of single-subject, repeat-measures design, in which differential responses to different types of test stimuli permit the analysis of response variance with each examinee serving as a control set for oneself.

<sup>&</sup>lt;sup>3</sup>Response variance is distinct from criterion variance, and criterion variance and criterion state refer to the actual disease state of the issue we are attempting to diagnose, (i.e., guilt or innocence). Variance can be thought of as either independent (i.e., having no causal factors in common with, and therefore not influenced by, the other test items) or non-independent (i.e., subject to potential influence from factors in common with other test items) when the requirements for independence are not satisfied.

<sup>&</sup>lt;sup>4</sup>The MFH in field practice has been a source of ambiguity and discussion. And example of this is the Federal ZCT format, for which the third relevant question is formulated as an evidence-connecting or indirect involvement question. Despite the fact that results for the Federal ZCT format are given at the level of the test as a whole, there has been discussion among field practitioners about whether this format is an example of a multi-facet exam, with the implication that the third relevant question may vary independently. As with the MFH in general, validity of this idea would require evidence of increased accuracy as a result of interpretation at the question subtotal compared to results at the level of the test as a whole.

Deception and truth-telling are amorphous phenomena that are without physical substance and which therefore cannot be physically measured. Response features used when scoring and interpreting PDD examinations are criteria that are correlated with deception and truth-telling, and can be recorded, aggregated, normed, and interpreted categorically based on the probabilistic strength of the information. Numerous published studies have shown that responses to relevant and comparison stimuli vary significantly as a function of the criterion states of deception and truth-telling (APA, 2011; Honts, Thurber & Handler, 2020; Kircher & Raskin, 1988; National Research Council, 2003; Offe & Offe, 2007; Podlesny & Truslow, 1993; Raskin, Kircher Honts & Horowitz, 1988). PDD results are both categorical and probabilistic statements for which response features serve as proxy data that are correlated with the criterion states of the issues being tested - similar to the way that hormones or antibodies can serve as proxies for the presence of a pregnancy or disease state. All correlations are imperfect, and all physiological responses are multi-purposed (i.e., all physiological activities serve multiple functions and therefore have multiple external correlates. An inevitable feature of probabilistic information is that test data is always a combination of diagnostic variance, also referred to as controlled variance, explained variance or signal information, and error variance, also referred to as random variance, uncontrolled variance or noise.

In contrast to probabilistic tests and probabilistic models, are those that are deterministic for which randomness, uncertainty, and uncontrolled variance play no role. Deterministic models are based on physical substance and thus offer the potential for mechanical measurement – subject only to measurement error. Probabilistic test results are based on proxy data that cannot be interpreted with deterministic expectations of perfection, and for this reason scientific tests are not expected to be infallible. Understanding the accuracy of probabilistic models is a matter of quantifying the degree of uncertainty associated with a model or test result. All probabilistic test results based on proxy data will include some margin of potential error. Scientific tests achieve their effectiveness by quantifying whether the margin of error conforms to stated requirements for accuracy, generally expressed in the form of the tolerance for error or alpha boundary. Practical effect sizes can be described as the observed or expected proportions of correct or incorrect test results, and can also be described as the degree of improvement over chance outcomes. The MFH suggests that error variance can be minimized through the assumption that test question response variance varies independently, commensurate with the criterion state. Implicit within the multi-facet independence hypothesis is the expectation that responses to individual relevant stimuli will also vary significantly from responses to comparison stimuli as a function of guilt or innocence<sup>5</sup>.

In field practice, assumptions about the independence of reactions to relevant questions are operationalized through the use of decision rules that use subtotal scores. Ideally, these subtotal scores would be considered with cut-scores that satisfy requirements for statistical significance. Using the subtotal scores for multi-facet or multiple issue exams, the examinee is considered deceptive if the subtotal score for *any* individual relevant question equals or exceeds the normative cutscore, while truthful classifications are made when *all* subtotal scores equal or exceed the normative cut-score for a truthful result. If

<sup>&</sup>lt;sup>5</sup> PDD test effectiveness, as in other forms of diagnosis and screening, is achieved through development and validation of structural models composed of statistically optimal combinations of physiological response features that have been shown individually to have significant correlation with the criterion. Implicit in this correlation is an understanding that no relationship between test data and criterion can be uniform or perfect, and that test results are probabilistic and not deterministic. Test data analysis, test accuracy, and test validation are a process of developing predictive classification models by quantify differential responses in physiology that occur in response to test stimuli. Measured physiological responses are aggregated mathematically and used to make probabilistic calculations of how well a test result fits our expectations according to a statistical reference model that may represent either the null-hypothesis or the hypothesis. The principles of scientific hypothesis testing and decision theory can then be used to make categorical and probabilistic classifications of the test results regarding deception and truth-telling.



neither of these two conditions is satisfied, the result is considered inconclusive. [Refer to Nelson (2018a) for a discussion of various decision rules used in polygraph field practice.] If the subtotal score for any question meets the criterion for deception, then the numerical results of other question subtotal scores that do not meet the criterion for deception are uninterpretable.<sup>6</sup> This contrasts with decisions based on the grand total score, in which deceptive and truthful classifications are made by comparing the grand total score to normative cut-scores and statistical reference models without regard to subtotal scores.<sup>7</sup>

The MFH suggests several potential benefits, including: 1) increased discrimination between deception and truth-telling, 2) discrimination at the level of the individual target questions regarding the behavioral role or level of involvement among guilty examinees, 3) increased confidence that a reasonable and complete array of stimuli was presented to the examinee during testing, and 4) improved context setting for post-polygraph discussion, interview, and investigation. The third and fourth of these hypothesized effects can be thought of as staging effects intended to clarify for examinees and referring professionals the semantic and logical meaning of the test stimuli. The literature review performed herein is limited to the first two, which can be thought of as criterion effects because they are related to criterion accuracy, (i.e., whether decisions at the level of the individual questions can increase the criterion accuracy of the test, and whether the test can determine involvement or non-involvement in the distinct behavioral targets of the investigation). Criterion accuracy can be measured along several dimensions, including decision accuracy with guilty and innocent cases, sensitivity and specificity rates, false-positive and false-negative error rates, inconclusive rates with guilty and innocent

persons, and Positive and Negative Predictive Values (outcome confidence in the test result). Consideration of these dimensions of test accuracy is important because different risk-assessment and risk-management contexts may prioritize different aspects of the risk-benefit ratios that may be informed by the test result.

The following research questions were investigated in this literature review: 1) does the interpretation of test data at the level of the individual target question provide any increase or advantage to criterion accuracy, 2) do subtotal scores for individual target questions discriminate deception and truth-telling, and 3) do subtotal scores for individual target questions discriminate the behavioral role or level of involvement of guilty examinees. Investigation of these effects will require the identification of published studies that reported results using both grand total and subtotal scores.

Validity of the MFH will be supported if published evidence reveals increased test effectiveness as a result of the interpretation of physiological responses at the level of the individual target questions when compared to the interpretation of responses at the level of the test as a whole. The null-hypothesis (responses to individual questions do not vary independently and decisions based on subtotal scores do not optimize the discrimination of deception and truth-telling) can be rejected if subtotal scores discriminate deception vs. truth-telling at rates significantly greater than chance, and if increases in test accuracy are observed as a result of the scoring and interpretation of subtotal scores in comparison to results based on grand total scores. The related hypothesis regarding discrimination of the behavioral role or level of involvement of guilty persons can be supported and the corresponding null hypothesis (responses to individual questions do not effectively discriminate the behavioral role or



<sup>&</sup>lt;sup>6</sup>The decision rule is described here in the context of a norm-referenced statistical hypothesis test. However, traditional cut-scores appear to have sometimes been determined through a combination of empirical and heuristic analysis without regard for statistical reference distributions or the level of statistical significance. It is possible that test accuracy might be further optimized through the use of statistically determined norm-referenced cut-scores.

<sup>&</sup>lt;sup>7</sup>Combinations of the SSR and GTR have also been described, including the simultaneous use of the SSR and GTR (Department of Defense, 2006), and the sequential combination of the SSR and GTR (Bell, Raskin, Honts & Kircher, 1999; Handler & Nelson, 2008; Senter, 2003; Senter & Dollins, 2003;). All assumptions about the criterion, including assumptions about independence and non-independence, are embedded in pragmatic and procedural test operations.

level of involvement of guilty persons) can be rejected if the observed combinatoric decision accuracy rate for individual questions is significantly greater than chance for both deceptive and truthful classifications.

#### History and evolution of independent target questions

Keeler (1930) provided one of the earliest examples of the use of a multi-facet target selection and question formulation approach.<sup>9</sup> He described a procedure in which four or five irrelevant <sup>10</sup> questions are presented, e.g., Do you live in Chicago? Following the irrelevant questions, several relevant questions would be asked that described multiple facets of a crime incident. Keeler gave the following as examples of relevant questions: Did you dine with Jones Tuesday night?, Did you return to Jones' apartment that night?, Did you owe Jones some money?, Did you discuss this indebtedness?, Have you ever been in California?, Did you shoot Jones?, and Do you own a Savage forty-five?

Keeler (1933) described a burglary in Los Angeles during 1923, in which the owner of a second story apartment came home to surprise a burglar attempting to open a safe. Reportedly, the burglar attempted to open a window to the fire escape, became entangled in the curtains before shooting and killing the owner, and then left through the door. Keeler described the presentation of the following questions that had no bearing on the case: Is your name Jones?, Do you live in Los Angeles?, Do you own an automobile? The following relevant questions were also presented: Do you live on Maple Street? (the street of the burglarized apartment.) Do you live in a second story apartment?, Have you heavy draperies in your apartment?, Have you a safe in your apartment?<sup>11</sup> Keeler explained that innocent persons showed no difference in reactions to the relevant questions compared to the other questions, while the guilty person reacted stronger to the questions about the burglarized apartment.

Summers (1939) described the use a series of significant questions<sup>12</sup> and provided these examples from a mock crime laboratory experiment: Do you know who took the money?, Did you take the money?, Have you the money on your person? These questions were included in sequence with these matter of fact questions:<sup>13</sup> Are you wearing a black coat?, *Did you eat breakfast this morning?* Summers also described the use of relevant questions regarding whether an examinee had helped to kill X, including: Were you in the home of X on the day of the murder?, Did you kill X?, and Do you know who killed X? Importantly, Summers included questions called emotional standards,<sup>14</sup> for which he gave the following examples: Were you ever arrested, Are you living with your wife?, "Do you own a revolver? <sup>15</sup> Emotional standard questions preceded each of the significant questions in the test question sequence, and the sequence of questions was presented three times.<sup>16</sup> Summers explained that if an examinee showed generally greater reactions to the significant questions than to the emotional standards, the examinee was attempting to deceive the examiner; if the examinee showed reactions to the significant questions that were generally not greater than reactions to the emotional standards, the examinee was answering truthfully. This

<sup>&</sup>lt;sup>8</sup> In the combinatoric approach to the hypothesis of role discrimination, a test result is correct if the results are correct for all individual questions.

 $<sup>^{9}</sup>$  Keeler's technique, though he does not appear to have used the term at the time, forms the basis of what is referred to today as the relevant-irrelevant technique (Department of Defense, 2006).

<sup>&</sup>lt;sup>10</sup> Referred to today as neutral questions.

<sup>&</sup>lt;sup>11</sup>Although Keeler's method of target selection is unlikely to be used today, the use of seemingly neutral relevant questions is noteworthy as an innovative early solution to challenges of test data analysis.

<sup>&</sup>lt;sup>12</sup> Referred to today as relevant questions. Use of the term significant generally implies statistical significance.

<sup>&</sup>lt;sup>13</sup> Similar to the irrelevant questions described by Keeler and neutral questions of today.

is one of the earliest examples of the current conception of the polygraph technique that strength of reactions to one group of questions or the other (relevant or comparison) is a function of deception or truth-telling to the relevant questions (Kircher & Raskin, 1988; Nelson, 2015a; Nelson, 2015b).

Reid (1947) introduced a revised questioning technique that included relevant questions, irrelevant questions, and two types of comparison questions.<sup>17</sup> He suggested the use of the following irrelevant questions: Have you ever been called 'Red'? (the examinee had used one or more alias names including the name Red). Did you have something to eat today?, and Did you ever smoke? (the examiner had seen the examinee smoking). Affirmative truthful answers were solicited from the examinee regarding these irrelevant questions. Reid provided these examples of comparative response questions to which the examiner was reasonably sure the examinee would lie: Have you ever been arrested before?, Since you got out of the penitentiary, have you committed any burglaries?, Have you ever lied on the witness stand?, Have you ever stolen anything?, Have you ever cheated on your income tax returns?, Have you ever committed adultery?, and Beside that five dollars you told me about, have you ever stolen any money? Reid stated, The examiner must feel reasonably sure, as a result of his preliminary interrogation, that the subject will answer "No" to any of the above suggested questions used for comparative response pur-

poses. Reid gave the following as examples of relevant questions: Do you know who shot John Jones?, About two months ago did you kill a man during a burglary at 112 State Street?, Did you stay in Chicago last night?, Do you know who shot John Jones?, Did you kill John Jones last Saturday night?, Did you steal a diamond ring from John Jones' room last Saturday night?, and Were you present when John Jones was shot? These questions represent a retention of earlier attempts to stimulate the guilty examine with a variety of target questions that describe distinct aspects of a known or alleged incident under investigation.<sup>18</sup> However, field practice with the Reid technique involved one decision about the test as a whole, and published studies on this technique involved uniform decisions at the level of the question and therefore did not attempt to interpret the individual questions with an assumption that they may vary independently.

Kubis (1962) reported on a series of three experiments on various aspects of lie detection, including whether response variance could be scored and interpreted independently for individual test questions. One experiment suggested the potential for role discrimination in which blind scorers were able to differentiate between thief, lookout, and innocent suspects in a simulated theft. Relevant questions were: *Were you an accomplice to the thief?*, *Did you take the money from the coin box?*, *Do you have the coin box money with you?* A separate series asked these questions: *Do you know how much* 

<sup>14</sup>Later referred to as emotional controls, then control questions and more recently comparison questions.

<sup>15</sup>These questions are unlikely to be viewed today as satisfactory comparison questions.

<sup>16</sup>Summers' question formulation method, and the technique as a whole - with three relevant target questions, three comparison questions, and three iterations of the test question sequence - bears early resemblance to contemporary versions of the family of zone comparison techniques.

<sup>17</sup>One type of comparison question was a guilt complex question based on a fictitious crime of the same type being investigated, and suggested that any reaction to the fictitious crime that is greater than or about the same as the actual crime question would indicate the examinee is innocent. Reid further suggested that reaction to the crime questions without reaction to the fictitious crime would indicate guilty knowledge or responsibility rather than nervousness or other factors. Raskin, Barland, and Podlesny (1978) reported that guilt complex questions were less effective than comparison questions at determining truthfulness or deception. These questions are not used in current PDD techniques.

<sup>18</sup>Backster (1963) recognized that the use of several behaviorally distinct questions introduces complex and troublesome implications regarding attention, salience, interpretation of test data and test accuracy. Backster began to advocate for the use of a series of single-issue test questions that forgo attempts to differentiate role involvement and instead describe the most important aspect of a known or alleged event. Barland Honts and Barger (1989) later showed that a series of single issue exams offered no advantages – or may suffer from the same psychological and statistical multiplicity complications – compared to the use of multiple issue exams.

money was in the coin box?, Did you act as a lookout for the thief?, and Do you have the coin box money with you? However, blind scorers in two later experiments, a crime scenario and a classified-information scenario, were informed that the examinees were deceptive to two of four target questions, but were unable to discriminate deceptive from truthful questions.<sup>19</sup>

The Department of Defense (2006a; 2006b) described the use of a variety of types of relevant test questions in the context of the Comparison Question Test (CQT) format, a generic category of the test that encompasses variations of the Modified General Question Technique (MGQT), including versions developed by the United States Army and United States Air Force (US Customs and Border Protection, 2010). Stimulus questions in this format are formulated to serve different purposes. These include both primary relevant questions that test possible direct involvement and secondary relevant questions that test possible involvement in a crime such as helping, planning, or any participation. Secondary relevant questions can also describe issues of guilty knowledge such as seeing, hearing, or knowing specific crime details. Consistent with the trend toward quantitative decision models, these examination formats employ a structured numerical scoring method with standardized protocols for feature identification and numerical transformations. Interpretation of CQT/MGQT test question stimuli is accomplished through structured decision rules that assume reactions vary independently for each relevant question. The test results are reported deceptive if reactions to any relevant question are determined to be indicative of deception, and are reported as truthful when reactions to all question are determined to be indicative of truth-telling.

## Method of Literature Review and Evaluation

Ten published studies were found in the literature to provide quantitative information regarding effect sizes for decisions using both grand total and subtotal scores. One of these studies described the result of two experiments using two different samples, and the sample data from one study was used in another larger analysis. In all, eleven separate experiments were identified as providing results using both grand total and subtotal scores. Included studies are listed below:

- 1 Horvath and Reid (1971)
- 2. Hunter and Ash (1973)
- 3. Slowik and Buckley (1975)
- 4. Wicklander and Hunter (1975)
- 5. Raskin, Kircher, Honts, and Horoitz (1988)
- 6. Barland, Honts, and Barger (1989)
- 7. Podlesny and Truslow (1993)
- 8. Krapohl and Norris (2000)
- 9. Senter (2003)
- 10.Senter and Dollins (2003)

Two-way unbalanced ANOVAs were calculated using the method described by Cohen (2002) to test the level of significance of differences among results using grand total and subtotal scores for criterion guilty and criterion innocent groups, including decision accuracy, sensitivity, specificity, false-positive and false-negative errors, and inconclusive results. Unbalanced ANOVAs were used due to difference in study sample sizes. Standard

<sup>&</sup>lt;sup>19</sup>Kubis (1962) is also interesting for a number of other practical and historical reasons, including the use discriminate analysis to show that the structural correlation of electrodermal responses was greater than respiratory and cardiovascular responses. Kubis also provided one of the earliest published references to numerical scoring, including the use of a Likert type numerical transformation to assign integer scores 0, 1, 2, and 3 to relevant and comparison stimuli using a rubric of "non-significant," "doubtfully significant," "significant," and "very significant."

deviation estimates were calculated using binomial approximation to the normal distribution, using the arithmetic mean of the sample sizes.<sup>20</sup> One-way, post-hoc ANOVAs were used as needed to investigate significant interactions. Random chance probability of a correct or incorrect classification was assumed to be .5 for decisions at the levels of the test as a whole and the individual questions. Statistical tests were conducted with statistical significance set at  $\alpha$ = .05.

## Chronology of research pertaining to the independence hypothesis

Horvath and Reid (1971) reported the criterion accuracy of blind evaluation of confirmed field examinations conducted by Horvath using the Reid Technique (Reid, 1947; Reid & Inbau, 1977).<sup>21</sup> This study involved 10 examiners who provided blind subjective judgments of a sample of confirmed criminal investigation exams after removing exams that were regarded as very easy or very difficult to interpret. Half of the examinees (n = 20) were reportedly confirmed guilty<sup>22</sup> and the other half (n = 20) were reportedly confirmed innocent.<sup>23</sup> Unweighted decision accuracy was reported as 87.5% for the test as a whole. In addition, blind evaluators provided categorical judgments for 164 individual questions, for with the criterion state was coded uniformly for each case. The

external criterion was coded innocent for half of the questions and the other half was coded as guilty. Decision accuracy for individual test questions was reported as 88.2% .<sup>24</sup>

Hunter and Ash (1973) reported the results of a study in which seven examiners from the staff of John E. Reid and Associates (Chicago, Illinois) provided blind judgments of deceptive or truthful, regarding a sample of N= 20 confirmed field cases that were conducted using the Reid Technique. Field cases included: theft, official misconducted, sexual assault, and homicide. Half of the sample cases were reportedly verified innocent, and the other half were reportedly verified guilty. Judgments were made by evaluating whether the consistency of physiological response was greater to the relevant or comparison questions. Unweighted decision accuracy was reported to be 87.1% for judgments made at the level of the test as a whole. Decision accuracy was reported to have been 92.1% for judgments regarding the individual relevant questions with criterion states coded uniformly for all questions within each case.<sup>25</sup>

Slowik and Buckley (1975) reported the results of a study in which seven examiners from the staff of John E. Reid and Associates provided blind categorical judgments of deception and truth-telling regarding a sample of

35

<sup>&</sup>lt;sup>20</sup>Mean sample size was used to calculate the binomial approximation of the sample standard deviation for both grand total and subtotal scores. This method was preferred because the assumptions of frequentist hypothesis testing assert that sampling error is a function of sample size when the sample is randomly selected from independent members of the population. Because subtotal scores within each examination violate the assumption of independence, attempts to use the N of subtotals were expected to overestimate the precision of the sampling statistic.

<sup>&</sup>lt;sup>21</sup>This technique is not presently being taught at any accredited polygraph school that we are aware of, but forms part of the basis of the family of MGQT formats.

<sup>&</sup>lt;sup>22</sup>The terms guilty and innocent are used to refer to the criterion state to avoid potential confusion when using the terms deceptive and truthful to discuss categorical test results. These terms are not intended to convey assumptions or implications about legal or criminal culpability in either laboratory or field settings.

<sup>&</sup>lt;sup>23</sup>Among the important differences between this technique and contemporary polygraph techniques is that the Reid technique employed a clinical approach in which judgments were made via consideration of a combination of whether physiological responses to test stimuli were interpreted as loaded onto either relevant or comparison stimuli (without fixed numerical cutscores), together with evaluation of the case background information and information from behavioral observation of the examinee during testing.

<sup>&</sup>lt;sup>24</sup>Uniform coding of the criterion state in this manner, together with interpretation at the level of the individual question, is not an expression of the MFH or the assumption that the both the criterion state and responses of individual questions will vary independently, and cannot address the hypothesis of role discrimination among the guilty cases.

<sup>&</sup>lt;sup>25</sup>Uniform coding does not address MFH.

N = 30 field cases that were selected from the reportedly verified case files at John E. Reid and Associates. Half of the sample cases were reported as confirmed as innocent, and half were reported as confirmed guilty. Cases involved theft, industrial sabotage, drug abuse, and sexual assault. The sample consisted of 141 individual relevant questions for which 71 questions were coded as confirmed innocent, and 70 were coded as confirmed guilty. Decision accuracy was reported as 86.7% for judgments made at the level of the individual test questions, with criterion states coded uniformly for all questions within each case.<sup>26</sup> Accuracy was reported to have been 88.9% for judgments made at the level of the overall test.

Wicklander and Hunter (1975) reported the results of a study in which six examiners from the staff at John E. Reid and Associates provided categorical judgments regarding a sample of 20 reportedly confirmed field cases that were conducted by the authors. Half of the sample cases were reported as confirmed guilty, and half were reported as confirmed innocent. Sample examinations consisted of 89 individual relevant questions, for which 43 were reported as verified innocent and 46 were reported as verified guilty. Criterion accuracy for categorical judgments of individual relevant questions was reported to have been 92.8% when the criterion state for individual questions was coded uniformly for each case. <sup>27</sup> Criterion accuracy for judgments made at the level of the overall test was reported as 93.9%.

Raskin et al. (1988) analyzed 76 examinations that had been conducted by the United States Secret Service during FY83, FY84, and FY85. The criterion states were confirmed (through an exhaustive records search, described in the study) by a combination of admissions and reliable evidence independent of the polygraph examinations. The examinations were independently scored by six experienced examiners employed by the United States Secret Service and one psychophysiologist-examiner. The cases were divided into Pure Verification and Mixed Verification cases of which 37 Pure Verification examinees were confirmed deceptive to all questions and 26 were confirmed truthful to all questions. The Mixed Verification cases consisted of 13 cases for which the examinees were confirmed as deceptive to at least one question and also confirmed as truthful to at least one question. There were 23 questions for which the examinees were confirmed deceptive and 19 questions for which the examinee were confirmed truthful. Accuracy of blind numerical scores for Pure Verification cases in which questions were not independent was 90.1%. Unweighted average accuracy for the Mixed Verification cases was 75.6%<sup>28</sup> for decisions at the level of the individual test questions. Raskin et al. concluded that test accuracy is maximized by formulating test questions that can be interpreted as non-independent.

Barland et al., (1989) investigated the ability to detect deception and truth-telling at the level of the individual questions. Federal polygraph instructors tested 100 participants who completed none, one, two, or three mock espionage and sabotage behaviors. Participants were randomly assigned to guilty or innocent groups. Fifty examinees were tested using a multiple-issue examination consisting of criterion independent target questions describing possible involvement in different mock espionage and sabotage activities. The other 50 examinees were tested with a series of three single-issue exams, each of which described

 $<sup>^{28}</sup>$ Also noteworthy in this study was that there was no effect for experience or type of training, with both highest and lowest criterion accuracy rates resulting from blind scores provided by experienced field examiners. The main effect for decisions was significant (F [1, 212] = 1340.26, p < .001), but the main effect for confirmation was not significant (F [1, 212] = 1.57, ns) suggesting that decisions did not vary as a function of confirmation. The authors also reported that the computer algorithm model generally outperformed blind numerical scores and most human scorers. They recommended the use of algorithms as a form of quality control.



<sup>&</sup>lt;sup>26</sup>Uniform coding does not address the potential for role discrimination among guilty cases.

<sup>&</sup>lt;sup>27</sup>Uniform coding does not address MFH and cannot be used to investigate the potential for role discrimination among guilty cases.

a single crime with two relevant questions in a sequence repeated three times. Unweighted accuracy was 87.3% for the participants tested using the series of single-issue exams and 83.9% for those participants who were evaluated using a multiple-issue examination format.<sup>29</sup> The difference in accuracy was not significant for tests using target questions that varied independently vs. decisions made using a series of examinations for which decisions were made at the level of the test as a whole. However, unweighted accuracy was only 47.8% when decisions were made at the level of the individual target questions, and only 33% of the results on individual crimes were correct when attempting to interpret the results of individual target questions. Although overall accuracy was high for guilty subjects, they were unable to determine which behavior had been committed by the guilty examinees.

Podlesny and Truslow (1993) reported the results of a laboratory study designed to investigate the ability of the polygraph test to identify guilty and innocent examinees and to determine the role or level of involvement of each examinee . Participants were recruited from the local community through a temporary employment agency and were paid for their participation, and were reported as not suffering from lack of sleep or illness. Ninety-six participants were randomly assigned equally to innocent, perpetrator, accomplice, or confidant and were tested using a Modified General Question Technique. Relevant questions included items about direct involvement, secondary involvement, knowledge, and innocence or truth-telling in general. Data were evaluated manually by 10 polygraph instructors from the Department of Defense.<sup>31</sup> Innocent participants produced a significantly positive mean grand total

and significantly positive mean subtotal scores to all relevant questions except knowledge. Each of the guilty groups differed significantly from the innocent group, but there was no significant difference among the guilty groups. Guilty perpetrators produced significantly negative mean scores to questions about perpetration, general truth, and knowledge, but not to participation/involvement. Guilty accomplices also had significant negative mean scores for questions regarding general truth, perpetration, and knowledge, but not participation/involvement. Guilty confidants produced significant negative scores to knowledge questions but failed to produce significant negative scores to questions about general truth. Guilty confidants also produced a higher rate of false negative errors and failed to produce significant negative mean total scores

Podlesny and Truslow reported that questions about participation/involvement did not discriminate among any of the guilty roles.<sup>32</sup> Unweighted accuracy was 89.7% using the grand total scores and 79.1% for subtotal (i.e., question) scores. The authors concluded that the evidence provided general support for the hypothesis that guilty participants would produce responses to individual relevant questions that varied significantly from responses to comparison stimuli, but did not generally support the hypothesis of role discrimination. The authors cautioned that although some discrimination among guilty groups might be possible, attempts to sub-categorize deception or to differentiate perpetration from participation/involvement and guilty knowledge may place overambitious demands on the testing context and examinee.



<sup>&</sup>lt;sup>29</sup>Barland, Honts and Barger (1989) also reported the results of a two-way ANOVA showing a significant main effect for criterion state (F [1, 96] = 30.4, p < .001) but no significant interaction between criterion state and type of testing (i.e., multiple issue exam or series of single issue exams) for decisions made at the level of the test as a whole.

<sup>&</sup>lt;sup>30</sup>Podlesny and Truslow (1993), in this study, also replicated the feature extraction and structural model development of Kircher and Raskin (1988) and Raskin, Kircher, Honts and Horowitz (1988), and provided replication and extension of the significant loading effect and numerical score differences for members guilty and innocent groups in a CQT test paradigm.

<sup>&</sup>lt;sup>31</sup>Data were also scored using a computer algorithm similar to the one described by Kircher and Raskin (1988), and hypothesis testing was completed using linear discriminate analysis.

<sup>&</sup>lt;sup>32</sup>Podlesny and Truslow (1993) also reported the results of a multivariate analysis that showed a significant interaction for role x relevant question (F [9, 219), = 1.57, p < .05), and a significant main effect for group (F [3, 92] = 30.2, p < .01).

Krapohl and Norris (2000) reported the results with 16 reportedly confirmed innocent and 16 reportedly confirmed guilty criminal investigation tests that were conducted using the Army MGQT format (Department of Defense, 2006b).<sup>33</sup> The series of test questions for the Army MGQT exam is structurally similar to that of the Reid Technique.<sup>34</sup> These examinations consisted of several different types of relevant questions including: 1) primary relevant questions that address direct involvement in the crime or issue under investigation, 2) secondary relevant questions that describe helping, planning or participating indirectly in the crime or issue, or secondary involvement such as seeing hearing, or knowing details about the crime or issue, 3) evidence-connecting relevant questions designed to determine if the examinee was involved with any of the evidence of the crime, or is aware of the nature or location of various items of evidence, and 4) guilty-knowledge relevant questions that are used to determine if the examinee has any knowledge of who committed the incident under investigation (Department of Defense, 2006a). Data were scored by three experienced examiners using a seven-position scoring model. Decisions made at the level of the test as a whole resulted in an accuracy rate of 75.5% but only 56.9% when decisions were made at the level of the individual questions.

Senter (2003) published the results of a study of decision rules with the Army MGQT format (Department of Defense, 2006b). Data were 205 verified criminal investigation tests and included 161 reportedly confirmed guilty cases and 44 reportedly confirmed innocent cases that were used in four previous research samples. Twenty-six of the exams were conducted by the United States Army Criminal Investigative Division (USACID). Forty-seven exams were conducted by the Bureau of Alcohol Tobacco and Firearms (BATF). Thirty-two exams were used in a sample used by Krapohl and Norris (2000). One-hundred MGQT cases were from a stratified matched sample used by Blackwell (1998). Cases from the BATF and USACID were scored by the original examiners. Results from the three examiners in the Krapohl and Norris (2000) study and the three examiners from the Blackwell (1998) study were averaged to achieve a single set of scores for each of those samples. There were 644 questions for which the criterion state was coded as guilty and 176 questions coded as innocent, with the criterion state coded uniformly for all questions within each case. Using traditional cutscores and decision rules, in which categorical results are determined at the level of the individual target questions, the unweighted mean for criterion accuracy of guilty and innocent groups was 66.1%. Unweighted accuracy using only the grand total scores was 83.8%.

Senter and Dollins (2003) studied the criterion accuracy of decisions based on grand total and subtotal scores using the Federal ZCT test format (Department of Defense, 2006a; 2006b) which makes use of primary relevant questions, that describe the possible direct involvement of an examinee in a known incident or allegation, along with secondary relevant questions that describe indirect involvement or guilty knowledge regarding an issue under investigation. Laboratory sample data were aggregated from previous studies, including data from Senter and Dollins (2004) <sup>35</sup> involving a sample of 50 criterion guilty cases and 50 criterion innocent cases that were scored by five examiners, along with data from two

<sup>&</sup>lt;sup>33</sup>This technique is no longer used in field settings, and was not included in the American Polygraph Association (2011) meta-analytic survey of criterion accuracy. Results are included here because the study design addressed the issues surrounding assumptions about independent or non-independent variance, and decision rules.

<sup>&</sup>lt;sup>34</sup>One important difference between the Reid format and the Army MGQT is that decisions are made by comparison of subtotal scores for individual target stimuli to fixed cut-scores when using the Army MGQT. This is in contrast to clinical judgments made at the level of the test as a whole when using the Reid technique. Both formats may use a combination of target questions developed around both primary and secondary involvement in addition to evidence-connection or guilty-knowledge questions.

<sup>&</sup>lt;sup>35</sup>Note that the publication dates are out of sequence due to apparent differences in project and publication timelines.

experiments described by the Department of Defense Polygraph Institute Research Division Staff (2001) involving a sample of 32 Federal ZCT exams (16 guilty cases and 16 innocent cases) that were scored by three examiners, and another sample of 32 Federal ZCT exams (16 guilty cases 16 innocent cases) that were scored by seven examiners. Laboratory studies included a total of 82 guilty examinees and 82 innocent examinees, whose data were scored by 15 different examiners. Data consisted of 492 individual questions, including 246 questions for which the criterion state was guilty and 246 questions for which the criterion was innocent. Decisions made with the grand total resulted in an unweighted average accuracy rate of 91.8%. Unweighted accuracy was 86.8% for decisions based subtotal scores, with the criterion state coded uniformly for all questions within each case.

Senter and Dollins (2003) also studied decisions based on grand total and subtotal scores using data from two Department of Defense archival samples of reportedly confirmed field exams that were conducted using the Federal ZCT format. Field cases consisted of examinations scored by six different examiners, with 115 guilty examinees and 85 innocent examinees. One hundred of the Federal ZCT cases from the sample were previously used by Blackwell (1998) and included 65 criterion deceptive cases and 35 criterion innocent cases that were scored by three examiners. An additional 100 cases were from the sample used by Krapohl, Dutton and Ryan (2001). There were a total of 600 individual questions, including 345 questions for which the criterion state was guilty and 255 questions for which the criterion state was innocent. Unweighted accuracy for the Federal ZCT exams was 91.5% for decisions based on grand total scores, but unweighted accuracy was only 79.8% for decisions based on subtotal scores.<sup>36</sup>

## Results

## Discrimination of deception and truth-telling.

Data from the 11 studies was scored by 85 examiners who evaluated N = 888 individual cases consisting of n = 549 criterion guilty cases and n = 339 criterion innocent cases. There were a total of N = 2870 individual target questions, including n = 1748 criterion guilty questions and n = 1122 criterion innocent questions.<sup>37</sup> Sensitivity and specificity rates for these studies are shown in Appendix A, while false-positive and false-negative error rates are shown in Appendix B. Information from the sensitivity, specificity, and error rates, together with the sample size, provide the information for this analysis.

The weighted mean of unweighted accuracies for grand total scores was 88.9% (SEM = .011, 95% CI = .869 to .910) and the weighted mean of unweighted accuracy rates for decisions based on subtotal scores was 79.0% (SEM = .014, 95% CI = .763 to .816).<sup>38</sup> Decisions based on an assumption of non-independent response variance, using the grand total score, were significantly more accurate than decisions based on an assumption that physiological responses to individual questions, and resulting subtotal scores, will vary independently, F (1,1775) = 16.737, (p < .001).

Table 1 shows the weighted means, standard errors, and confidence intervals for guilty and innocent cases. Figure 1 shows the mean plot for guilty and innocent cases. Results from a two-way unbalanced<sup>39</sup> ANOVA showed a significant interaction between criterion state and decision rule, F (1,1772) = 40990.342, (p < .001). Post-hoc one-way unbalanced ANO-VAs showed a significant one way effect for criterion state for decisions based on grand



<sup>&</sup>lt;sup>36</sup>Senter and Dollins (2003) also reported results using a sequential combination of the grand total score and subtotal scores and that would provide similarly high decision accuracy (87.7% for laboratory cases and 86.4% for field cases) with increased test sensitivity and reduced inconclusive results.

 $<sup>^{37}</sup>$ Data from Krapohl and Norris (2000) are excluded from these totals because the sample was included in the data reported by Senter (2003).

<sup>&</sup>lt;sup>38</sup>The weighted mean is preferred in this situation, giving more importance to the statistical estimates from larger studies, due to the premise of frequentist inference that the precision and error with which a sampling statistic describes the population is a function of the sample size, with the assumption that samples are representative of the population if the sample cases are independent of each other and are randomly selected (i.e., all members of the population have an equal chance of being selected).

total scores, F (1,837) = 7.032, (p = .008) and for subtotal scores, F (1,837) = 39.724, (p < .001). The difference in accuracy with guilty and innocent cases was significant for decisions based on grand total scores, F (1,1097) = 9.923, (p = .002) and subtotal scores, F (1,677) = 31.154, (p < .001). Decisions based on grand total scores were more accurate with innocent cases, while decisions using subtotal scores were more accurate with guilty cases. However, unweighted decision accuracy was significantly greater than chance for both grand total and subtotal scores with both guilty and innocent cases.

		Criterion guilty	Criterion innocent
Results based on grand	Weighted mean	.841	.924
total scores	SEM	.016	.014
	95% confidence	.811 to .872	.895 to .952
	interval		
		Criterion guilty	Criterion innocent
Results based on	Weighted mean	.926	.694
subtotal scores	Standard deviation	.011	.025
	95% confidence	.904 to .948	.647 to .745
	interval		

## Figure 1. Mean plot for accuracy of guilty and innocent cases based on grand total and subtotal scores.



<sup>&</sup>lt;sup>39</sup>Use of unbalanced ANOVAs, using the harmonic mean of sample sizes, was required due to inequality in the cell sizes. Unbalanced ANOVAs provide a slight reduction in statistical power, and are thought to be a more cautious test of significance under these circumstances.



Table 2 shows the sensitivity and specificity rates of decisions based on grand total and subtotal scores. Figure 2 shows the mean plot for test sensitivity and specificity. Results from a two-way unbalanced ANOVA showed a significant interaction between decisions and decision rule, F (1,1772) = 41642.631, (p < .001). Post-hoc one-way ANOVAs showed that the difference in test sensitivity and specificity was not significant for decisions based on grand total scores, F (1,837) = 1.05, (p =

0.306), but was statistically significant for results based on subtotal scores, F (1,837) = 72.713, (p < .001). The difference in test sensitivity to deception for grand total and subtotal scores was statistically significant, F (1,1097) = 32.047, (p < .001), and the difference in test specificity for grand total and subtotal scores was also significant, F (1,677) = 15.469, (p < .001). Use of subtotal scores was more accurate with guilty cases, while grand total scores were more accurate with innocent cases.

#### Table 2. Sensitivity and specificity cases using grand total and subtotal scores.

		Sensitivity	Specificity
Results based on grand	Weighted mean	.660	.706
total scores	Standard deviation	.020	.025
	95% confidence	.620 to .699	.658 to .755
	interval		
		Sensitivity	Specificity
Results based on subtotal scores	Weighted mean	.834	.465
	Standard deviation	.015	.027
	95% confidence	.831 to .889	.449 to .555
	interval		

#### Figure 2. Mean plot for sensitivity and specificity.



False-negative and false-positive errors are shown in Table 3, and Figure 3. A two-way unbalanced ANOVA showed a statistically significant interaction between errors and decision rule, F (1,1772) = 7979.41, (p < .001). Posthoc analysis showed that false-negative and false-positive errors did not differ significantly at the .05 level for results using grand total scores, F (1,837) = 2.709, (p = .100). However, the difference in error rates was significant for results based on subtotal scores, F (1,837) = 25.567, (p < .001). False-negative error rates did not differ significantly for results based on grand total and subtotal scores, F (1,1097) = 1.374, (p = 0.241), but the difference in false-positive errors was significant for decisions based on grand total and subtotal scores, F (1,677) = 5.034, (p = 0.025). An assumption of independent variance contributed to a statistically significant increase in FP errors.

## Table 3. False-negative and False-positive errors.

		False-negative errors	False-positive errors
Results based on grand	Weighted mean	.072	.119
total scores	Standard deviation	.011	.018
	90% confidence	.051 to .094	.085 to .154
	interval		
		False-negative errors	False-positive errors
Results based on	Weighted mean	.048	.209
subtotal scores	Standard deviation	.009	.022
	90% confidence interval	.031 to .066	.166 to .252

Figure 3. Mean plot for false-negative and false-positive errors.





Inconclusive rates for results based on grand total and subtotal scores are shown in Table 4 and Figure 4. A two-way unbalanced ANO-VA showed a significant interaction between the criterion state and decision rule for inconclusive results, F (1,1772) = 11102.575, (p < .001). Post-hoc one-way unbalanced ANO-VAS showed that the difference in inconclusive rates for guilty and innocent cases was not significant for results based on grand total scores, F (1,837) = 0.218, (p = .641), but was statistically significant when results were based on subtotal scores, F (1,837) = 29.029, (p < .001). Inconclusive rates the innocent cases did not differ significantly for results based on grand total and subtotal scores, F (1,677) = 1.518, (p = .218). However, the difference in inconclusive rates was statistically significant for guilty cases when comparing results using grand total and subtotal scores, F (1,1097) = 16.39, (p < .001). Subtotal scores produced fewer inconclusive results with guilty cases.

#### Table 4. Inconclusive rates for results based on grand total and subtotal scores.

		Criterion guilty	Criterion truthful
Results based on grand	Weighted mean	.218	.237
total scores	Standard deviation	.018	.023
	90% confidence	.184 to .253	.192 to .283
	interval		
		Criterion guilty	Criterion truthful
Results based on	Weighted mean	.120	.323
subtotal scores	Standard deviation	.024	.044
	90% confidence	.070 to .119	.248 to .345
	interval		

Figure 4. Inconclusive rates for results based on grand total and subtotal scores.



43

Unweighted decision accuracy shown in Table 1 was significantly greater than chance with both guilty and innocent cases for grand total scores and for subtotal scores. For decisions using grand total scores, both test sensitivity to deception and test specificity to truth-telling were significantly greater than chance (0.5), with no significant difference in sensitivity and specificity. However, decisions using subtotal scores showed an increased in test sensitivity along with a reduction of mean test specificity to chance.

False-positive and false-negative error rates did not differ significantly for results based on grand total scores, and the mean false-negative error rate was not significantly different for results based on subtotal scores when compared with results of grand total scores. However, decisions based on subtotal scores produced significantly more false-positive errors than decisions using grand total scores. Inconclusive results did not differ significantly for guilty and innocent cases when results were based on grand total scores but were significantly lower for guilty cases when results were based on subtotal scores.

Results of these analyses provide general support for the hypothesis that PDD examination data can discriminate deception and truth-telling at rates significantly greater than chance for the test as a whole and at the level of individual question. The results do not support the validity of the MFH that interpretation of independent response variance for individual target questions optimizes or increases the effectiveness of event-specific polygraph examinations. Results using grand-total scores, based on an assumption of non-independent response variance produced greater mean accuracy and more balanced error rates compared to results using subtotal scores. However, observed differences were not uniform for the criterion guilty and criterion innocent cases. Based on an assumption of independent response variance, subtotal scores showed increased effectiveness with guilty cases but also showed a disproportionately larger decrease in effectiveness with innocent cases.

## Discrimination of behavioral role or level of involvement.

Analysis of effect sizes for role discrimination or level of involvement of guilty persons requires that the guilt vs. innocence criterion state is known or set independently for each individual relevant question, in addition to a requirement that decisions are made using subtotal scores. Although results were reported for the study samples using both grand total and subtotal scores, a majority of these samples consisted of target questions for which the criterion state was coded uniformly (i.e., non-independently) for all questions within each case. In other words, the criterion state was coded as guilty to all questions or innocent to all questions within these sample cases. Uniform coding in this manner makes no assumption of independent criterion variance, and this imposes a substantial barrier to the interpretation of independent responses and the effect of the MFH for these sample cases.

Samples from three of the included studies (Barland et al., 1989; Podlesny & Truslow, 1993; Raskin et al., 1988) did include target questions for which the criterion states of the individual target questions were coded independently. An ancillary analysis was completed using data from these three studies in comparison to data from the seven studies (Horvath & Reid, 1971; Hunter & Ash, 1973; Krapohl & Norris, 2000; Senter, 2003; Senter & Dollins, 2003; Slowik & Buckley, 1975; Wicklander & Hunter 1975) for which the criterion states of the individual questions was coded uniformly.

Mean unweighted accuracy of studies for which the criterion state of the individual target questions was coded independently was .892, SEM = .021 (95% CI = .851 to .934) for decisions using grand total scores, and .792 SEM = .028 (95% CI = .737 to .847) for decisions using subtotal scores. Mean accuracy of study samples for which the criterion state of individual target questions was coded uniformly for within each case was .889, SEM .012 (95% CI = .865 to .912) for decisions using grand totals, and .789, SEM = .016 (95%) CI = .758 to .820) for decision using subtotal scores. Differences were not statistically significant for grand total decisions of questions coded independently or uniformly, nor was the

difference statistically significant for subtotal decisions of questions coded independently or uniformly.

To test for differences in criterion accuracy for independent and uniform coding of guilt and innocence within each case, a series of two-way ANOVA contrasts was conducted for decision accuracy, sensitivity and specificity, false-negative and false-positive errors, and inconclusive results for guilty and innocent cases. Figure 5 shows the mean plot for unweighted decision accuracy, Figure 6 shows the mean plot for sensitivity and specificity, Figure 7 shows the mean plot for errors, and Figure 8 shows the mean plot for inconclusive results.

#### Figure 5.







## Figure 7.



#### Figure 8.





The factor of interest for the planned contrasts was the method of coding of the criterion state for the individual target questions within each case, i.e., whether the criterion variance was independently coded or uniformly coded. The effect for differences in decision accuracy for independent versus uniform question coding, shown in Figure 5, was not significant, F (1,1768) = 0.132, (p = .717). Decision accuracy was similarly effective with cases for which the criterion state was coded independently for each question compared to cases coded uniformly. The one-way effect for test sensitivity and specificity, shown in Figure 6, was not significant at the .05 level, F(1,1768) = 3.435, (p = .064), The difference in errors, shown in Figure 7, was also not statistically significant for independent vs uniform coding of question criterion states, F (1,1769) = 0.348, (p = .555). Mean decision accuracy was lower for innocent cases when decisions were based on subtotal scores regardless of the method through which the criterion state was coded.

The one-way difference in inconclusive rates, shown in Figure 8, was statistically significant, F(1,1768) = 4.077, (p = .044) at the .05 level. Cases for which the criterion state of individual questions was coded independently produced more inconclusive results compared to cases for which questions were coded uniformly, and this difference was loaded on innocent cases. This is not surprising when considering that standard field practice when making decisions using subtotal scores is that the examinee is considered deceptive if the subtotal score for any individual relevant question equals or exceeds the normative cut-score, while truthful classifications are made when all subtotal scores equal or exceed the normative cutscore for a truthful result. Inconclusive results occur when neither of these two conditions is satisfied. A corollary to this procedure is that field polygraph practices prohibit examiners from rendering both deceptive and truthful classifications within a single exam. [Refer to Nelson, Blalock & Handler (2019) for a discussion of how test results are parsed for individual questions and for the test as a whole for single-issue and multiple-issue polygraphs.] As a consequence, when any single question within a test is classified as deceptive, results are deemed inconclusive for all other questions that are not statistically significant for deception. In other words, when the criterion

state was coded independently for the relevant questions within each exam, and when the criterion states were mixed for the individual questions, with one or more questions coded guilty and one or more questions coded innocent, questions were classified as inconclusive if they were coded innocent in an exam for which any single questions was classified as deceptive.

As shown in Figure 5, decision accuracy was slightly higher for guilty cases when decisions were made using subtotal scores, but with a disproportionately larger decrease in accuracy with innocent cases. In addition, an interaction effect for decision accuracy can be observed in the results with grand total scores, shows in Figures 5, 6 and 8, between criterion state and whether the criterion states of individual questions were coded independently or uniformly when decisions were based on grand total scores. The increase in accuracy that resulted from the use of grand total scores with truthful cases was not observed when the criterion state was set independently for individual target questions but was observed when target questions were coded uniformly. When the criterion states of the individual questions were coded independently the use of grand total decisions resulted in less effective discrimination of deception and truth-telling compared to decisions using the grand total score with cases for which the criterion state of individual questions was coded uniformly.

## Discussion

Important limitations of the present study include the form of analysis and the use of a series of unbalanced multivariate ANOVAs to evaluate several dimensions of test accuracy by using proportional statistics that describe categorical results. Calculation of variance estimates, statistical confidence intervals, and sums of squares was accomplished with binomial approximation of the normal distribution. A more precise analysis might be obtained through bootstrapping or Monte-Carlo simulation or other method. It is also possible that meta-analytic methods might lead to more interesting insights than were observed in this analysis.

Fundamentally, any analytic method that re-

lies on null-hypothesis significance testing is based on the assumptions that samples are collected in a random manner and are representative of the population. It is also assumed that sampling bias is a function of sample size, i.e., a sufficiently large sample will approximate the population with reasonable accuracy. However, large random samples are difficult to obtain. Furthermore, other sources of bias or error that may be introduced to the sample data and study results include sampling methodology, test administration (including test question construction and inter-rater reliability)and test data analysis. This analysis is premised on and assumes that the published sampling information is in some way representative of the larger population.

Some of the samples included field investigation cases for which case selection was contingent on non-random criteria involving the availability of confirmation data (i.e., confession or investigative evidence) that may not be independent of the PDD examination results. In other words, the PDD result may have contributed to the resolution of the field cases, and this can be expected to inflate the correlation between the test result and criterion state. A further limitation of field sampling procedures is the potential for the systematic exclusion of both false-positive and false-negative errors for which confirmatory confession or evidence may be difficult or impossible to obtain. Non-random sampling methodology may inflate observed test accuracy. Despite these known sampling confounds, field investigation samples offer the advantage of external and ecological validity.

Although ecological validity is not required to achieve the more important goal of external validity, laboratory samples have been subject to criticism because they may incompletely represent field testing contexts. Laboratory studies offer the important advantage of scientific control over the research questions, which cannot be achieved in field research. Results of laboratory studies have been shown to be more conservative but not significantly different from field studies (Honts, Thurber & Handler, 2020). Correspondence between the result of field and laboratory sampling data provides additional insight as to the stability or reproducibility of observed results. An equally important limitation of the present analyses is that the results of the included studies were universally derived without regard for the normative distributions of the scores of guilty and innocent persons. This means that decision cut-scores were statistically under-informed and potentially sub-optimal as a result of unknown alpha levels for some of the sampling data. Also, statistical corrections were not applied to decision alpha nor and some of the included studies were completed without fixed numerical cutscores.

Individual field examinations represent a form of single-systems experiment, and examinations involving the simultaneous interpretation of multiple, potentially independent, sources of response variance represent the execution of multiple statistical comparisons within each test experiment. There are known statistical multiplicity effects that can occur when making multiple statistical comparisons, including the potential for inflation of alpha boundaries corresponding increase in Type I (i.e., false-positive) errors when making deceptive classifications, the potential for deflation of alpha boundaries, and corresponding increase in Type II (i.e., false-negative) errors when requiring multiple statistically significant truthful scores in order to make truthful classifications. For this reason, there is an unquantified possibility that sampling data from experiments with norm-referenced cut-scores and statistically informed alpha boundaries can produce different results.

Another limitation for this literature review is the small number of published studies that have evaluated the MFH. A number of the included studies made use of subtotal scores with uniform criterion coding that could not fully evaluate the MFH. Those that did attempt to code the criterion states of individual target issues point to a general conclusion against the notion of effective role discrimination, and there is a clear trend in the literature that scores based on uniform coding and grand total scores can maximize the criterion accuracy of PDD test data.

It was not possible to investigate whether the use of multi-facet questions contributes to any semantic increase in test sensitivity – which might result from describing a broader range of crime-event related stimuli to which guilty persons may respond. Similarly, it was not possible to evaluate whether the staging and presentation of multi-facet test questions leads to increased effectiveness and resolution of post-test discussion and investigation activities or increased satisfaction among referring professionals. This review is limited to quantitative effects only. These hypothesized effects are beyond the scope of the present study but are nonetheless important and should be investigated in future studies.

Generalizability of these results to multiple-issue screening exams is neither completely certain nor completely unknown. Multiple-issue screening exams are conducted in a context where there is a strong assumption of independence of the criterion state for the individual questions. Greater independence of response variance in multi-issue testing may lead to differences in the ability of individual questions to discriminate deception and truth-telling. However, the degree to which the target questions of multiple-issue screening exams may include an increase in the independence of response variance is unknown. Practical experience in field polygraph settings has revealed potential problems with over-reliance on an assumption of independent response variance. As a result, field polygraph examiners are generally prohibited from making both deceptive and truthful classifications within a single exam. Also unknown are the effects of variable prior probabilities of guilt for multi-issue screening polygraphs compared to the priors associated with event-specific diagnostic exams. Few practicing field polygraph examiners are prepared to consider the effect that base-rates may have on Positive or Negative Predictive Values of the individual test question results.

## Conclusion

The MFH holds that response variance is independent when the relevant questions employ different behavioral action verbs that will provide valid and reliable information about a person's behavioral role or level of involvement in a known allegation or incident. The MFH suggests that, in addition to the ability to partition response variance into the larger dimensions of relevant and comparison stimuli, response variance can also be partitioned into sub-dimensions for individual relevant questions. This is similar to that of repeated-measures or within-subjects experimental designs. Within-subjects designs are useful for their efficiency and statistical power but have the potential disadvantage of complex analysis.<sup>40</sup> An important difference is that omnibus analytic methods are not routinely expected to provide granular conclusions, while the MFH in field polygraphy includes an implicit expectation of granular precision. In practical terms, the MFH is the hypothesis of a multiple-issue diagnostic exam for which results are classified as deceptive or truthful using subtotal scores for individual relevant questions. The MFH assumes that test outcomes are optimized by making decisions using the question subtotals instead of the overall test total. Consistent with the known limitations of omnibus analysis methods, field practitioners have adopted standards and policies that preclude conclusions that a person has been deceptive to one or more questions and truthful to one or more questions within a single exam. Nevertheless, the expectation that multi-facet questions will isolate the response variance associated with deception has been difficult to counter even though the published literature contains a number of studies that directly and indirectly address the MFH.

Published evidence at this time does not support the hypothesis that test accuracy with event-specific exams can be optimized through the interpretation of independent response variance at the level of the individual target questions. There is support for the more general hypothesis that both grand total and subtotal scores may discriminate deception and truth-telling at rates better than chance but with considerable imbalance with guilty and innocent persons. Results from these analyses indicate that overall decision

<sup>&</sup>lt;sup>40</sup>Which is largely mitigated through the use of computerized statistical analysis.

accuracy of event-specific polygraph exams may be optimized by using grand total scores, thus forgoing the assumption that responses will vary independently for individual relevant questions – including when those questions employ different action verbs or describe different behavioral roles in a known or alleged incident. <sup>41</sup>

Accuracy with subtotal scores was found to be generally lower than for grand total scores, and imbalanced with guilty and innocent persons. Published information was insufficient to investigate the degree to which this may be improved through the use of statistically referenced cut-scores and mathematical corrections for multiplicity effects. However, evidence does not preclude the use of multi-facet questions without an assumption of independence. Effect sizes for grand total scores of examinations consisting of a combination or primary, secondary and behavioral role-descriptive questions are equivalent to published information for other single issue exams.

These results are consistent with previous reports (e.g., Senter & Dollins, 2003) that decisions using grand total and subtotal scores may provide different advantages. Overall decision accuracy was greater with grand total scores, with better balanced accuracy for guilty and innocent cases and better specificity to truth-telling compared to decisions using subtotal scores. However, decisions using subtotal scores provided increased test sensitivity to deception and reduced inconclusive rates for guilty examinees, though with disproportionately weaker accuracy with innocent persons. Considering the imbalance of effects with guilty and innocent cases and the fact that criterion accuracy for innocent cases was reduced to less than .5, these results do not support the hypothesis that test questions can be used to determine behavioral role or level of involvement of guilty examinees. Furthermore, these results do not support polygraph field practices that attempt to determine that a person has been deceptive to one or more questions and also truthful to one or

more questions within a single examination.

Prior to the advent of numerical evaluation and the use of normative data, the use of test questions describing several different behavioral aspects of a known or alleged crime would have been viewed as an important strategic innovation. It is easy to understand the strategy of presenting stimulus questions with an array of factual information designed to focus the attention and responses of the guilty person - while the innocent person is expected to have no factual or behavioral association. While this strategy seems to have potential face validity, evidence-based professional practices require replicable empirical support for the continued use of a method. Published evidence suggests that the MFH contributes to a small increase in test sensitivity is associated with decisions using subtotal scores accompanied by a larger decrease in test specificity and lower overall accuracy.

Use of the term multi-facet cannot be easily found in testing contexts outside the polygraph profession. One of the earliest references to the term within the polygraph profession was in the terminology reference by Krapohl and Sturm (1997). The concept of a multiple-facet test appears to have emerged within the polygraph profession as a useful way of distinguishing between multiple-issue screening polygraphs and event-specific diagnostic polygraphs that make use of both primary and secondary relevant questions. Both multiple-issue screening polygraphs and multi-facet diagnostic polygraphs have been traditionally scored and interpreted with the assumption that responses to test questions vary independently. Although somewhat useful for its ability to remind us of the difference between investigative polygraphs of known or alleged incidents and screening tests, the term multi-facet has the potential for confounding our knowledge and understanding of the decision-theoretic and statistical concerns that determine test effectiveness - specifically, whether we assume independent or non-independent variance of the relevant questions.

<sup>&</sup>lt;sup>41</sup>Neither Senter, Dollins, and Krapohl (2008), nor recent studies on the AFMGQT (Nelson & Blalock, 2012; Nelson, Blalock & Handler, 2011; Handler & Nelson, 2012; Nelson, Handler, Morgan, & O'Burke, 2012) attempted to interpret independent response variance or the criterion accuracy of the individual target questions.

Traditional practice has been to interpret these specific-incident exams using the same decision rules as with multi-issue screening exams, premised on an assumption of independent variance. Support for the validity of the multi-facet hypothesis requires that the interpretation of independent response variance leads to increased criterion accuracy. Despite a number of attempts to investigate its contribution to polygraph test effectiveness, there remains insufficient evidence to support the validity of the MFH. Evidence has consistently indicated that multi-facet questions are non-independent.

Studies dating back to the early 1970s have consistently informed us that the polygraph test is effective at discriminating between deception and truth-telling, but less effective at determining individual questions to which a person is lying or truthful. In practical terms this should be taken to mean that examinees can be said to pass or fail the test as a whole, but not the individual questions.<sup>42</sup>

Despite the absence of quantitative empirical support for the multi-facet hypothesis, use of multi-facet test questions may provide qualitative advantages that are not captured by a numerical or statistical scoring model. These unquantified advantages may pertain to the staging of pretest discussion and to confidence that complete and adequate test stimuli were presented to the examinee during testing. Also, field examiners may find multi-facet questions useful for developing a post-test interview strategy. Referring professionals should be cautioned against naïve expectations that multiple roles or behavior can be successfully targeted and parsed with granular precision during a single examination.

In response to anticipated arguments that it is the *utility* of the test that is optimized and optimization of test accuracy is not an objective of the MFH, we caution that utility is not optimized when the test results are more likely to produce an error that misleads an investigator or consumer. In response to numerous anecdotal stories regarding individual examinations during which an examinee exhibited strong response to a particular question with subsequent confirmation, anecdotal stories comprise a miniscule fraction of examinations that are actually conducted Evidence-based polygraph testing necessarily focuses on what happens most of the time.

Field examiners who find it useful to use multi-facet questions to stage the discussion or investigation before, during or after the polygraph test, may wish to do so using the event-specific formats developed at the University of Utah (Handler & Nelson, 2008; Kircher & Raskin, 1988; Nelson 2018b; Raskin, Honts, Nelson & Handler, 2015). These formats consist of three or four relevant stimulus questions regarding a single known or alleged incident. Relevant questions can be expressed as a single issue or can also be used to describe multiple facets of a single known or alleged incident. An example of the use single issue questions is as follows: Did you rob that bank located at \_\_\_\_ in Austin?, Did you rob that bank located at \_\_\_\_\_ in Austin last Thursday?, and, Did you rob that bank at \_\_\_\_ on (date)? (Nelson & Handler, 2008). An example of multiple facet test questions, involving both primary and secondary relevant questions is as follows: Did you rob that bank at on (date)?, Did you plan with anyone to rob that bank at ?, and Did you participate in that robbery of that bank? Regardless of the approach to target selection and question formulation, testing protocols must conform to the published literature if the test data are interpreted with the assumption that response variance is non-independent.

We caution against any expectation of perfect test accuracy, and remind that scientific tests and scientific test results are fundamentally probabilistic. The polygraph, like other tests, is a useful though imperfect tool. Recall that Podlesny and Truslow (1993) wrote:

Users of results obtained with the present technique should be cautioned



<sup>&</sup>lt;sup>42</sup>This practical approach may differ for multi-issue screening tests for which the strength of the assumption of independent criterion variance is thought to be greater. More research is needed with multiple issue exams regarding the independence of response variance. [Refer to Nelson, Blalock and Handler (2019) for additional discuss about how categorical results are parsed and reported in different ways for event-specific diagnostic exams and multiple-issue screening polygraphs.]

that errors in classifying guilty and innocent subjects are not unlikely. The results further suggest that attempts to subcategorize deception using present MGQTs are not advisable. Where MGQT examinations are used to detect deception versus no deception in distributed-crime-role contexts, decision methods based on overall question totals might reduce inconclusives and improve accuracy with innocent subjects in comparison with methods based on individual question scores. (p.795).

A similar caution was provided previously by Raskin et al., (1988), when they wrote:

> A related problem is raised by the finding of higher false positive rates for questions answered truthfully by suspects who were also deceptive to at least one relevant question in the same test. It appears that answering deceptively to at least one relevant question in the test tends to weaken the reactions to the control questions, thereby making it difficult for them to produce reactions that are larger than those to relevant questions that are answered truthfully. Therefore, field polygraph examiners should attempt to devise sets of relevant questions that the suspect can be expected to answer all truthfully or all deceptively. The case information and the importance of each relevant question should be carefully considered in formulating the set of relevant question to be asked, and separate question series should be used whenever it seems likely that the suspect might answer some of the relevant questions truthfully and some of them deceptively.

Despite the pragmatic recommendations of Raskin *et* al. (1988) regarding the formulation of test questions, and of Podlesny and Truslow in interpreting the results, the practical appeal of multiple-issue tests has resulted in the continued use of polygraph tests for which responses test questions are assumed to vary independently. And, this practice continues to occur in both screening and diagnostic testing contexts. While there is obvious convergence of evidence regarding the MFH and the independence of questions withing multi-facet polygraph exams, conducted in the context of a specific incident or allegation, less is known about the variance of multi-issue screening polygraphs conducted in the absence of any known allegation or incident.

Barland et al. (1989) wrote:

One interesting finding of Experiment 2 was that the examinations did not detect deception at the level of the individual crimes. This result has important implications for examiners who must test on multiple relevant issues, as it suggests that the numerical scores associated with individual relevant issues may be a poor guide in choosing issues for interrogation. This result suggests that when deception is inferred, the interrogator may need to address all of the relevant issues on the examination with the interrogation.

The available evidence at this time indicates that the MFH is false. However, this does not negate field practices that make use of combinations of primary and secondary relevant questions when it is understood that responses to multi-facet test stimuli regarding a single known or alleged incident do not vary independently. Test accuracy is optimized by interpreting responses to multi-facet questions in a manner similar to that of other event-specific examinations (i.e., assuming that the variance of responses to relevant stimuli is non-independent). In practical terms, this may require reevaluating decision policies used to achieve the final categorical result that is interpreted from the numerical data. Accurate PDD test results contribute to guide decisions about the need for subsequent discussion and investigation, but inaccurate PDD results may contribute to inaccurate risk-assessment and risk management decisions.

## References

- American Polygraph Association (2011). Meta-analytic survey of validated polygraph techniques. Polygraph, 40(40), 203-305. [Electronic version] Retrieved January 6, 2012, from http:// www.polygraph.org.
- Backster, C. (1963). Standardized polygraph notepack and technique guide: Backster zone comparison technique. Cleve Backster: New York.
- Barland, G. H., Honts, C. R., & Barger, S. D. (1989). Studies of the accuracy of security screening polygraph examinations. Department of Defense Polygraph Institute.
- Barland, G. H. & Raskin, D.C. (1975). Psychopathy and detection of deception in criminal suspects. *Psychophysiology*,(12), 224.
- Bell, B. G., Raskin, D. C., Honts, C. R. & Kircher, J.C. (1999). The Utah numerical scoring system. *Polygraph*, 28(1), 1-9.
- Blackwell, J. N. (1998). PolyScore 33 and psychophysiological detection of deception examiner rates of accuracy when scoring examination from actual criminal investigations. Available at the Defense Technical Information Center. DTIC AD Number A355504/PAA. Reprinted in Polygraph, 28, (2) 149-175.
- Cohen, B. H. (2002). Calculating a factorial ANOVA from means and standard deviations. Understanding Statistics., 1(3), 191-203.
- Department of Defense (2006). *Federal Psychophysiological Detection of Deception Examiner Handbook.* (Retrieved from https://www.antipolygraph.org/documents/federal-polygraphhandbook-02-10-2006.pdf on 2-11-2010). Reprinted in Polygraph, 40(1), 2-66.
- Department of Defense (2006). Psychophysiological Detection of Deception Analysis II -- Course #503. Test data analysis: DoDPI numerical evaluation scoring system. Available from the author. (Retrieved from https://www.antipolygraph.org/documents/dodpi-numerical-scoring-08-2006.pdf on 6-28-2007).
- Department of Defense Polygraph Institute Research Division Staff. (2001). Test of a mock theft scenario for use in the psychophysiological detection of deception: IV (DoDPI00-R-0002). Fort Jackson, SC: Department of Defense Polygraph Institute.
- Handler, M. (2006). The Utah PLC. Polygraph, (35), 139-148.
- Handler, M. & Nelson, R. (2008). Utah approach to comparison question polygraph testing. *European Polygraph,(2),* 83-110. Reprinted in Polygraph 38(1) 15-33.
- Handler, M. & Nelson, R. (2012) Criterion Validity of the United States Air Force Modified General Question Technique and Three Position Scoring, *Polygraph*, 41 (1).
- Handler, M., Nelson, R., Goodson, W. & Hicks, M. (2010). Empirical Scoring System: A cross-cultural replication and extension study of manual scoring and decision policies. *Polygraph,(39)*, 200-215.
- Honts, C. R., Thurber, S. & Handler, M. (2020). A comprehensive meta-analysis of the comparison question polygraph test. *Applied Cognitive Psychology*, 2021, 1-17.

53

- Horvath F. S. (1988). The utility of control questions and the effects of two control question types in field polygraph techniques. *Journal of Police Science and Administration*, 16(3), 198-209. Reprinted in Polygraph, 20, 7-25.
- Horvath, F. & Palmatier, J. (2008). Effect of two types of control questions and two question formats on the outcomes of polygraph examinations. *Journal of Forensic Sciences*, *53(4)*, 1-11.
- Horvath, F. S. & Reid, J.E. (1971). The reliability of polygraph examiner diagnosis of truth and deception. *Journal of Criminal Law, Criminology and Police Science*,(62), 276-281.
- Hunter, F. L. & Ash, P. (1973). The accuracy and consistency of polygraph examiners' diagnosis. Journal of Police Science and Administration,(1), 370-375.
- Keeler, L. (1930). A method for detecting deception. American Journal of Police Science, 1, 38-52. Reprint in Polygraph, 23, (2), 134-144.
- Keeler, L. (1933). Scientific methods of crime detection with a demonstration of the polygraph. Kansas Bar Association Journal, (12), 22-31.
- Kircher, J. C. & Raskin, D.C. (1988). Human versus computerized evaluations of polygraph data in a laboratory setting. *Journal of Applied Psychology*, 73, pp. 291-302.
- Krapohl, D. J., Dutton, D. W. & Ryan, A.H. (2001). The rank order scoring system: Replication and extension with field data. *Polygraph,(30)*, 172-181.
- Krapohl, D., Gordon, N. & Lombardi, C. (2008). Accuracy demonstration of the Horizontal Scoring System using field cases conducted with the Federal Zone Comparison Technique. *Polygraph,(37)*, 263-268.
- Krapohl, D., Handler, M. & Sturm, S. (2012). PDD *Terminology Reference*. American Polygraph Association: Nashville, TN.
- Krapohl, D. & McManus, B. (1999). An objective method for manually scoring polygraph data. *Polygraph, 28, 209-222.*
- Krapohl, D. J. & Norris, W. F. (2000). An exploratory study of traditional and objective scoring systems with MGQT field cases. *Polygraph*, 29, 185-194.
- Krapohl, D. J. & Stern, B.A. (2003). Principles of Multiple-Issue Polygraph Screening: A model for applicant, post-conviction offender, and counterintelligence testing. *Polygraph*, 32, pp. 201-210.
- Krapohl, D., & Sturm, S. (1997). Terminology Reference for the Science of Psychophysiological Detection of Deception. Chattanooga, TE: American Polygraph Association.
- Kubis, J. F. (1962). Studies in Lie Detection: Computer Feasibility Considerations. RADC-TR 62-205, Contract AF 30(602)-2270. Air Force Systems Command, U.S. Air Force, Griffiss Air Force Base. New York: Rome Air Development Center.
- Light, G. D. (1999). Numerical evaluation of the Army zone comparison test. Polygraph, 28, 37-45.
- Nelson, R. (2015). Scientific basis for polygraph testing. Polygraph 41(1), 21-61.
- Nelson, R. (2015). Scientific (analytic) theory of polygraph testing. APA Magazine, 49(5), 69-82.



- Nelson, R. (2018). Practical polygraph: a survey and description of decision rules. *APA Magazine*, *51(2)*, 127-133.
- Nelson, R. (2018). Credibility assessment using Bayesian credible intervals: a replication study of criterion accuracy using the ESS-M and event-specific polygraphs with four relevant questions. *Polygraph & Forensic Credibility Assessment 47 (1)*, 85-90.
- Nelson, R. & Blalock, B. (2012). Extended analysis of Senter, Waller and Krapohl's USAF MGQT examination data with the Empirical Scoring System and the Objective Scoring System, version 3. *Polygraph, 42,*.
- Nelson, R., Blalock, B. & Handler, M. (2011). Criterion validity of the Empirical Scoring System and the Objective Scoring System, version 3 with the USAF Modified General Question Technique. *Polygraph*, 40(11).
- Nelson, Blalock, B. & Handler, M. (2019). Practical polygraph: how to parse categorical results for test questions of diagnostic and screening polygraphs. *APA Magazine*, *52(3)*, 60-65.
- Nelson, R., Handler, M., Morgan, C., & O'Burke, P., (2012). Short Report: Criterion validity of the United States Air Force Modified General Question Technique and Iraqi scorers. *Polygraph*, 41(1).
- National Research Council (2003). The Polygraph and Lie Detection. National Academy of Sciences.
- Offe, H. & Offe, S. (2007). The comparison question test: does it work and if so how? *Law and human behavior*, *31*, 291-303.
- Podlesny, J. A. & Truslow, C. M. (1993). Validity of an expanded-issue (modified general question) polygraph technique in a simulated distributed-crime-roles context. *Journal of Applied Psychology*, 78, 788-797.
- Raskin, D. C., Barland, G. H., & Podlesny, J. A. (1978). *Validity and reliability of detection of deception*. Washington, D.C.: U.S. Government Printing Office.
- Raskin, D. C., Kircher, J. C., Honts, C. R. and Horowitz, S. W. (1988) A Study of Validity of Polygraph Examinations in Criminal Investigation, Grant number 85-IJ-CX-0040. Salt Lake City: Department of Psychology, University of Utah.
- Raskin, D. C., Honts, C. R., Nelson, R. & Handler, M (2015). Monte Carlo estimates of the validity of four relevant question polygraph examinations. *Polygraph*, 44(1), 1-278.
- Reid, J. E. (1947). A revised questioning technique in lie detection tests. *Journal of Criminal Law and Criminology*, (37), 542-547. Reprinted in Polygraph 11, 17-21.
- Reid, J. E. & Inbau, F.E. (1977). *Truth and deception: The polygraph ('lie detector') technique (2nd ed).* Williams & Wilkins.
- Senter, S. M. (2003). Modified general question test decision rule exploration. *Polygraph*, 32(4), 251-263.
- Senter, S. M., & Dollins, A. B. (2003) New Decision Rule Development: Exploration of a Two-Stage Approach. Department of Defense Polygraph Institute, Report No. DoDPI01-R-0007. Reprinted in Polygraph, 37(2), 149-164.



- Senter, S. M., & Dollins, A. B. (2004). Comparison of three versus three or five question series contingency rules: A replication. Department of Defense Polygraph Institute, Report No. DoDPI01-R-0008. Reprinted in Polygraph, 33(4), 223-233.
- Slowik, S. M. & Buckley, Joseph P, III (1975). Relative accuracy of polygraph examiner diagnosis of respiration, blood pressure and GSR recordings. *Journal of Police Science and Administration,(3)*, 305-309.

Summers, W. G. (1939). Science can get the confession. Fordham Law Review, 8, pp. 334-354.

- US Customs and Border Protection (2010). U.S. Customs and Border Protection Office of Internal Affairs Credibility Assessment Division Policy and Reference Manual with Appendices from DACA and Federal Standards. Retrieved from https://antipolygraph.org/documents/cbp-polygraph-handbook-2010-01-07.pdf.
- Wicklander, D. E. & Hunter, F.L. (1975). The influence of auxiliary sources of information in polygraph diagnosis. *Journal of Police Science and Administration,(3)*, 405-409.



## Appendix A.

#### Table 3. Sensitivity and specificity of grand total and subtotal scores.

Study	Grand total Scores		Subtotal scores	Subtotal Scores
	Sensitivity (sd) {90% Cl}	Specificity (sd) {90% Cl}	Sensitivity (sd) {90% Cl}	Specificity (sd) {90% Cl}
Horvath and Reid, 1971	.850 (.059)	.905 (.048)	.814 (.064)	.827 (.062)
	{.791 to .909}	{.857 to .953}	{.75 to .878}	{.765 to .889}
Hunter and Ash, 1973	.871 (.055)	.857 (.058)	.833 (.061)	.839 (.060)
	{.816 to .926}	{.799 to .915}	{.772 to .894}	{.779 to .899}
Slowik and Buckley, 1975	.838 (.061)	.905 (.048)	.733 (.073)	.881 (.053)
	{.777 to .899}	{.857 to .953}	{.660 to .806}	{.828 to .934}
Wicklander and Hunter, 1975	.942 (.038)	.867 (.056)	.905 (.048)	.808 (.065)
	{.904 to .98}	{.811 to .923}	{.857 to .953}	{.743 to .873}
Raskin et al., 1988	.651 (.078)	.515 (.082)	.478 (.082)	.316 (.076)
	{.573 to .729}	{.433 to .597}	{.396 to .56}	{.24 to .392}
Barland et al., 1989	.816 (.064)	.417 (.081)	.667 (.078)	.545 (.082)
	{.752 to .88}	{.336 to .498}	{.589 to .745}	{.463 to .627}
Podlesny and Truslow, 1993	.694 (.076)	.750 (.071)	.847 (.059)	.375 (.080)
	{.618 to .770}	{.679 to .821}	{.788 to .906}	{.295 to .455}
Krapohl and Norris, 200043	.604 (.080)	.479 (.082)	.896 (.050)	.083 (.045)
	{.524 to .684}	{.397 to .561}	{.846 to .946}	{.038 to .128}
Senter, 200344	.530 (.082)	.746 (.072)	.944 (.038)	.212 (.067)
	{.448 to .612}	{.674 to .818}	{.906 to .982}	{.145 to .279}
Senter and Dollins, 2003	.600 (.081)	.714 (.074)	.766 (.070)	.549 (.082)
(experiment 1)	{.519 to .681}	{.64 to .788}	{.696 to .836}	{.467 to .631}
Senter and Dollins, 2003	.713 (.074)	.671 (.077)	.861 (.057)	.423 (.081)
(experiment 2)	{.639 to .787}	{.594 to .748}	{.804 to .918}	{.342 to .504}
Weighted means <sup>45</sup>	.660 (.078)	.707 (.075)	.833 (.061)	.483 (.082)
	{.582 to .738}	{.632 to .782}	{.771 to .894}	{.401 to .565}

<sup>43</sup>Data from Krapohl and Norris (2000) was included in the data and analysis from Senter (2003), and are excluded from the weighted mean calculations in Table 2 to avoid redundancy.

<sup>44</sup>The N of individual questions was imputed from the number of reported sample cases based on an assumption that this test format was commonly used with four target questions at the time the sample data were collected.

<sup>45</sup>Confidence intervals were calculated using a variance estimate derived from the binomial approximation to the normal distribution using the number of the number of cases, not the n of the individual questions. This approach is thought to be more conservative, with a larger variance estimate and wider corresponding confidence intervals.



## Appendix B.

#### Table 2. Error rates for decisions based on grand total and subtotal scores.

Study	Grand total Scores		Subtotal scores	Subtotal Scores
	False-negative rate (sd) {90% Cl}	False-positive rate (sd) {90% Cl}	False-negative rate(sd) {90% Cl}	False-positive rate (sd) {90% Cl}
Horvath and Reid, 1971	.095 (.029)	.150 (.036)	.125 (.033)	.096 (.029)
	{.047 to .143}	{.091 to .209}	{.071 to .179}	{.048 to .144}
Hunter and Ash, 1973	.100 (.030)	.143 (.035)	.115 (.032)	.144 (.035)
	{.051 to .149}	{.085 to .201}	{.063 to .167}	{.086 to .202}
Slowik and Buckley, 1975	.152 (.036)	.067 (.025)	.182 (.039)	.065 (.025)
	{.093 to .211}	{.026 to .108}	{.119 to .245}	{.024 to .106}
Wicklander and Hunter, 1975	.083 (.028)	.050 (.022)	.054 (.023)	.078 (.027)
	{.038 to .128}	{.014 to .086}	{.017 to .091}	{.034 to .122}
Raskin et al., 1988	.058 (.023)	.148 (.036)	.091 (.029)	.170 (.038)
	{.020 to .096}	{.090 to .206}	{.044 to .138}	{.108 to .232}
Barland et al., 1989	.051 (.022)	.182 (.039)	.188 (.039)	.468 (.050)
	{.015 to .087}	{.119 to .245}	{.124 to .252}	{.386 to .55}
Podlesny and Truslow, 1993	.125 (.033)	.042 (.020)	.056 (.023)	.208 (.041)
	{.071 to .179}	{.009 to .075}	{.018 to .094}	{.141 to .275}
Krapohl and Norris, 2000 <sup>48</sup>	.25 (.043)	.104 (.031)	.001 (.003)	.521 (.050)
	{.179 to .321}	{.054 to .154}	{<.001 to .006}	{.439 to .603}
Senter, 200347	.037 (.019)	.432 (.050)	.006 (.008)	.432 (.050)
	{.006 to .068}	{.351 to .513}	{<.001 to .019}	{.351 to .513}
Senter and Dollins, 2003 (experiment 1)	.088 (.028)	.027 (.016)	.051 (.022)	.139 (.035)
	{.041 to .135}	{<.001 to .054}	{.015 to .087}	{.082 to .196}
Senter and Dollins, 2003	.07 (.026)	.055 (.023)	.049 (.022)	.228 (.042)
(experiment 2)	{.028 to .112}	{.017 to .093}	{.013 to .085}	{.159 to .297}
Weighted means <sup>48</sup>	.072 (.035)	.119 (.053)	.050 (.036)	.207 (.067)
	{.030 to 0.115}	{.066 to .173}	{.014 to .086}	{.140 to .273}

<sup>46</sup>Data from Krapohl and Norris (2000) was included in the data and analysis from Senter (2003), and are excluded from the weighted mean calculations in Table 2 to avoid redundancy.

<sup>47</sup>The number of individual questions was imputed from the number of reported sample cases based on an assumption that this test format was commonly used with four target questions at the time the sample data were collected.

<sup>48</sup>All standard deviations were calculated using binomial approximation to the normal distribution using the number of the number of cases. for both grand total and subtotals. This approach is thought to be more conservative, with a larger variance estimate and wider corresponding confidence intervals.



#### Another Look at Electrodermal Response Ratio Minima

Donald J. Krapohl<sup>1</sup>

#### Abstract

In separate papers Nelson (2020) and Krapohl (2020) conducted analyses of measurements of electrodermal responses to assess whether there was an optimal minimum ratio of response amplitudes for scoring between those elicited by relevant questions and those from comparison questions. Both evaluated the same data set and though the statistical treatments were different, both showed compatible patterns across ratios. Graphical representations of the data indicated performance of the electrodermal channel peaked when there was a minimum of 20% difference between the amplitudes of the relevant question and the comparison question against which it was scored. The effect was modest. Nelson proposed the differences were not meaningful. Krapohl opined that even a limited effect might be exploited to incrementally improve polygraph decision accuracy. To have confidence the improvement is genuine, though, the trend would have to be found in other samples. Here we analyzed data from laboratory studies conducted at two independent research centers. Both suggested a small improvement in EDR score assignment using minimum differences of about 20% up to 50% between EDR amplitudes over simply basing a score on all differences greater than 0%. These and the previous data point to a robust-if-modest effect when requiring a minimum difference in electrodermal response amplitudes for score assignment. We argue that minimum differences of about 20% may be the "sweet spot" for scoring for single-issue testing. More work remains for mixed-issue testing.



<sup>&</sup>lt;sup>1</sup>APA Past President and regular contributor to this publication. Comments, questions and requests for the raw data can be directed to APAkrapohl@gmail.com

The author wishes to express appreciation to Dr. John Kircher for making available the dissertation data analyzed in this paper. I am also grateful for the critically helpful suggestions and comments of the anonymous reviewer.

#### Introduction

The present effort is a replication and extension of Krapohl's (2020) pursuit of an optimal minimum ratio of electrodermal response (EDR) amplitudes between those of relevant questions and of comparison questions against which they are scored. We had previously hypothesized that basing scoring decisions on small differences between EDR amplitudes from relevant and comparison questions may incidentally capture excessive random variability that, in the long run, would degrade the value of the resultant manual scores. Conversely, a requirement of a large minimum difference between EDR amplitudes may come at the expense of a loss of diagnostic information. Somewhere between these extremes would be a hypothetical "sweet spot," a proposed minimum ratio that might perform better than others. That spot, according to our earlier data, was where the difference between EDR amplitudes was about 20%. The effect was limited, with the 20% minimum threshold producing a Detection Efficiency Coefficient (DEC, Kircher, Horowitz & Raskin, 1988) of 0.785 against a DEC of 0.727 for scores based on absolutely any difference in EDRs irrespective of size.

In an article in the same journal Nelson (2020) undertook a sophisticated statistical analysis of the same data set, looking not only at EDRs but respiratory and cardiovascular measurements. Using a different but related statistical approach to detection efficiency, Nelson's analyses indicated a peak efficiency of 0.450 when scoring relevant and comparison question EDR amplitudes with 20% differences versus 0.411 when scores were assigned to any absolute difference. Nelson's careful study addressed automated analysis of electrodermal data, but his findings had direct implications for manual scoring, as well. For automated analysis Nelson concluded the small differencees provided no clear advantage for imposing any minimum EDR ratio.

Both the Nelson and Krapohl studies found the Bigger is Better Rule (BIBR) for the electrodermal channel is effective across a broad range of minimum ratios, even at the smallest possible ratio. The modest improvement at a 20% minimum ratio found by both analyses may be tantalizing, but unless a similar finding were to come from other data sets there would be little reason to suppose there could be a benefit for using this or any other minimum ratio in manual scoring. The findings required replication with other data sets.

Available to us were two high quality and untapped archival data sources. One data set came from three Ph.D. dissertation projects from the University of Utah which focused on polygraph testing. A second was from an unpublished polygraph screening study conducted by the National Center for Credibility Assessment. As in the earlier Krapohl (2020) analysis, we were interested in the impact of different minimum EDR ratios on the effectiveness of that data channel.

#### Method

#### Data Set 1

In 2001 the Department of Defense Polygraph Institute (now the National Center for Credibility Assessment, or NCCA) conducted an unpublished multiple-issue screening study (Dollins, Senter & Pollina, 2001). Volunteers from the community were recruited to commit various mock crimes and then undergo polygraph testing to determine whether the examiners could discern those individuals from other volunteers who had not committed those acts. There were 102 volunteer examinees, 52 of them programmed deceptive. As part of the polygraph examination process each volunteer was given two separate test series using the Air Force Modified General Question Technique (AFMGQT) with three relevant questions in each series covering one issue (Krapohl & Shaw, 2015). Each AFMGQT series had three charts. Polygraphs used to record physiological data were made by Axciton Systems (Interface Version S7.1) and Lafayette Instrument Company (Model LX2000). EDR amplitudes were measured using Extract software (Extract version 4.0, Johns Hopkins University).

To create ratios for the present project, the EDR amplitude from each relevant question was divided by the amplitude of the stronger EDR response from one of two adjacent comparison questions. However, the third relevant question was preceded but not followed by a comparison question, and the EDR amplitude of that question was divided by the EDR amplitude of the immediately preceding comparison question. Both the relevant and comparison questions were rotated across the three charts in each case. There were 3,672 possible measurements to create 1,836 ratios (3 ratios per chart X 3 charts per case X 102 examinees X 2 series per examinee = 1,836).

## Data Set 2

Electrodermal amplitudes from polygraph dissertation projects of Drs. John Podlesny (Podlesny & Raskin, 1978), John Kircher (1983) and Paul Bernhardt (2005) of the University of Utah were analyzed. All the polygraph examinations in those projects addressed the examinee's involvement in a mock crime. Skin conductance was recorded with either a Beckman Type R Dynograph or a CPS-LAB system. Collectively there was a total of 255 cases (127 programmed deceptive), each with three charts, three relevant questions and three comparison questions for a total of 4590 EDR measurements.

To create ratios, the amplitude of each relevant question was divided by the amplitude of the comparison question immediately preceding it in the testing sequence. The comparison questions were systematically rotated across the three charts in each case. There were 2,295 ratios (3 ratios per chart X 3 charts per case X 255 cases = 2,295).

#### Procedure

Score assignment for given ratios followed the procedure in Krapohl (2020). Briefly, we systematically changed the minimum ratio required for score assignment. The bottom ratio was operationally defined as any difference in EDR amplitudes greater than 0%. We then increased the minimum in 10% increments up to a maximum of 80% difference in EDR amplitudes. If the larger response was on the relevant question a -1 was assigned. A +1 was given if the comparison question EDR was the larger, and a 0 if the difference did not exceed the thresholds we tested. For missing values we assigned 0 as the score. The scores were then summed per case. Point bi-serial correlations were calculated between total case scores and ground truth for each of the minimum ratios. Ground truth was coded +1 for programmed truthful and -1 for programmed deceptive.

Correlation statistics were calculated using online tools available at www.socialstatistics. com.

## Results

## Data Set 1

Figure 1 represents the relationship between score totals across nine EDR ratio minima for the laboratory screening study. As has been shown in previous work on this issue, the BIBR functions well regardless of the minimum. For all EDR amplitudes different by >0% the total by-case scores correlated with ground truth well above chance  $[r_{pb}(203) =$ 0.633, p < .001]. The lowest correlation among those tested was for a minimum ratio of >80% which also proved to be significantly greater than chance  $[r_{pb}(203) = 0.623, p < 0.001]$ . As Figure 1 shows, there was a plateau of better performance using minimum EDR differences between 10% and 50%.

Figure 2 shows the proportion of the 204 test series in which the total score for a case was 0. The proportions ranged from 0.044 to 0.113.





Figure 2. Proportion of cases in Data Set 1 which had total scores of 0 at escalating minimum ratios between >0% and >80% in 10% increments.





#### Data Set 2

Figure 3 shows correlation coefficients between total case scores and ground truth for nine minimum ratios. Again, the BIBR works quite well for all ratios. Even the most liberal ratio at >0% produced a correlation coefficient significantly greater than chance [rpb.(254) = 0.646, p < .001]. As Figure 3 shows, the best EDR performance takes place when requiring minimum ratios between 20% and 60% for this data set.

Figure 4 shows the proportion of the 255 examinations in which the total score for a case was 0. The proportions ranged from 0.012 to 0.086.

## Figure 3. Point bi-serial correlation coefficients between ground truth and total EDR scores for Data Set 2 at escalating minimum ratios between >0% and >80% in 10% increments.





## Figure 4. Proportion of cases in Data Set 2 which had total scores of 0 at escalating minimum ratios between >0% and >80% in 10% increments.



#### Discussion

Once again, the BIBR has been shown to be a good principle in scoring EDRs. Both data sets revealed significant correlations between scores and ground truth across all tested minimum ratios.

The trend seen in Figure 1 for Data Set 1 suggests EDR effectiveness is highest when using minimum ratios between 10% and 50% over the use of a ratio lower than 10%. This finding is compatible with Krapohl's (2020) earlier data. It might be recalled that the Krapohl (2020) procedure involved scoring each relevant question against the immediately preceding comparison question in the test sequence. The procedure used in the present analysis with Data Set 1 was to score against the stronger of two adjacent comparison questions for two of the relevant questions, and against the preceding comparison question for the final relevant question. That the two approaches should both show EDR performance is better at higher ratios than it is when using any ratio greater than 0% suggests the trend reported in Nelson (2020) and Krapohl (2020) could be robust.

Figure 3 for Data Set 2 shows the greatest EDR performance takes place between minimum ratios of 20% to 60% over the use of higher or lower ratios. The ratios in Data Set 2 were based on scoring each relevant question to the immediately preceding comparison question, the same method employed in the Krapohl (2020) analysis and therefore the present findings might not be unexpected. The trend in Data Set 2 also aligns with that of Data set 1 in that some EDR differences seem to be better than others.

Considering the present laboratory data as well as the analysis of field cases in Nelson (2020) and Krapohl (2020) there appears to be a modest increase in EDR performance using minimum EDR amplitude differences of about 20% for score assignment. Trend lines found across all analyses showed the 20% minimum requirement boosted EDR effectiveness over the use of a minimum difference greater than 0%. The trend consistently appeared in three independent samples, with both lab and field data, and using different rules for choosing comparison questions against which to score relevant questions. The consistent pattern may suggest a robust effect.

#### Limitations

The present findings are relevant only to 3-position and Empirical Scoring Systems. Optimal EDR ratios for 7-position scoring were offered previously (Krapohl, 2002) and are different from the present findings. There are no similar investigations to our knowledge for scoring systems based on rank ordering.

Our findings are not expected to generalize to mixed-issue screening examinations. Data Sets 1 and 2 both used test questions in which the examinee was either truthful to all three relevant questions or deceptive to all of them. This permitted the summing of scores across all relevant questions. Mixed-issue testing, where examinees may be deceptive to only some of the relevant questions, would not allow adding together scores for all relevant questions. Decision rules for mixed-issue testing relies on the sum of scores of individual relevant questions, which are expected to show more instability, and may require higher minimum differences between EDR amplitudes.

We conducted no tests of significance between or among correlation coefficients, and given the modest differences there is no expectation any would achieve significance with the available sample sizes. Samples of sufficiently large size will always result is statistically significant differences, which is an issue of ongoing discussion in scientific publications. Our suggestion that EDR score assignment based on a 20% difference in EDR amplitudes could outperform score assignment based on ratios merely greater than 0% is founded solely on the consistency of trends in the three independent samples, each containing at least 200 polygraph examinations. The difference in practical outcomes such as correct or incorrect results may be the greatest source of information when considering a recommendation for field practice. Because only electrodermal responses were investigated, accuracy of test outcomes were beyond the scope of this paper.



#### References

- Bernhardt, Paul, C. (2005). Effects of Prior Demonstrations of Polygraph Accuracy on Outcomes of Probable-Lie and Directed-Lie Polygraph Tests. Doctoral Dissertation, Department of Educational Psychology, University of Utah.
- Dollins, A.B., Senter, S.M., and Pollina, D.A. (2001). A Test of the Counterintelligence Screening Polygraph Process. Report No. DoDPI01-R-002. Ft. Jackson, SC. Unpublished.
- Kircher, J.C. (1983). Computerized Decision-Making and Patterns of Activation in the Detection of Deception. Doctoral Dissertation, Department of Psychology, University of Utah.
- Kircher, J.C., Horowitz, S.W., and Raskin, D.C. (1988). Meta-analysis of mock crime studies of the control question polygraph technique. *Law and Human Behavior*, 12(1), 79 90.
- Krapohl, D.J. (2002). Short Report: Update for the Objective Scoring System. *Polygraph*, 31(4), 298-302.
- Krapohl, D.J. (2020). Electrodermal responses: When is bigger really better? Polygraph & Forensic Credibility Assessment. 49(2), 104 109.
- Krapohl, D.J., and Shaw, P. (2015). *Fundamentals of Polygraph Practice*. Academic Press: San Diego, CA.
- Nelson, R. (2020) Bigger is better for automated scoring: Analysis of minimum constraints for RQ/ CQ ratios. *Polygraph & Forensic Credibility Assessment.* 49(2), 110 – 120.
- Podlesny, J.A., and Raskin, D.C. (1978). Effectiveness of techniques and physiological measures in the detection of deception. *Psychophysiology*, 15(4), 344 359.



