

Please log onto Editorial Manager at <http://asem.edmgr.com> to submit your letters to the editor. If you have not already done so, you will need to register with the journal.

## Letter to the Editor re: Evaluating the Reliability of the Human Factors Analysis Classification System

Dear Editor:

The article “Evaluating the reliability of the Human Factors Analysis and Classification System,” by Cohen, Shappell, and Wiegmann,<sup>1</sup> misrepresents three of the articles<sup>3-5</sup> used in their systematic review. The authors’ listing of “Reliable” for interrater reliability (IRR) in Table III does not give a fair representation of what these three authors actually found and stated in their studies. This misrepresentation seems to be based on an overreliance on the overall IRR value. In these three articles, the authors state that although the overall IRR was acceptable, the reliability of the causal factors was not acceptable, and that “the acceptable overall reliability can be attributed to the high reliability in rejecting nanocodes that clearly did not apply to the mishaps.” Therefore, five of the six articles which specifically examined IRR found HFACS to be unreliable.

If the DoD-HFACS articles are discounted due to increased granularity in the data, three articles are left; however, as stated in Cohen et al.’s article, two of those articles claim the rating is unreliable. The remaining study, based on only two raters, found the rating reliable. Although 14 studies were included in this review, 8 did not specifically examine reliability. Furthermore, the only studies that included data from more than 10 raters were the 3 that were misrepresented.

Given that HFACS has 19 causal categories and DoD-HFACS has 149 nanocodes, it seems reasonable that IRR should be analyzed separately. To that end, all three studies (plus Hughes’ unpublished study<sup>2</sup>) that have tested DoD-HFACS specifically for IRR have found it unreliable. Two of the three studies that have tested HFACS specifically for IRR have found it unreliable. Based upon the studies to date there are insufficient data to support the reliability of HFACS.

O’Connor<sup>5</sup> recommends that coding systems be evaluated for reliability and validity prior to widespread implementation. In fact, as of 1 October 2015, the U.S. Air Force no longer codes Ground Class C and D mishaps because of the unreliability of DoD-HFACS.

**Bruce R. Burnham, CSP**  
*Air Force Safety Center, Kirtland AFB, NM*

### ACKNOWLEDGMENTS

The views expressed herein are the view of the author and do not reflect the official policy of the Department of the Air Force, the Department of Defense, or of the U.S Government.

### REFERENCES

1. Cohen TN, Wiegmann DA, Shappell SA. Evaluating the reliability of the Human Factors Analysis and Classification System. *Aerosp Med Hum Perform.* 2015; 86(8):728–735.
2. Hughes TG, Heupel KA, Musselman BT, Hendrickson E. Preliminary investigation of the interrater reliability of the Department of Defense Human Factors Accident and Classification System in USAF mishaps [Abstract]. *Aviat Space Environ Med.* 2007; 78(3):255.
3. O’Connor P. HFACS with an additional layer of granularity: validity and utility in accident analysis. *Aviat Space Environ Med.* 2008; 79(6):599–606.
4. O’Connor P, Walker P. Evaluation of a human factors analysis and classification system as used by simulated mishap boards. *Aviat Space Environ Med.* 2011; 82(1):44–48.
5. O’Connor P, Walliser J, Philips E. Evaluation of a human factors analysis and classification system used by trained raters. *Aviat Space Environ Med.* 2010; 81(10):957–960.

### In Response:

We reviewed published studies that examined the reliability of HFACS when used as a tool for classifying human factors data associated with accidents.<sup>1</sup> HFACS was not typically used during the original investigation; rather, it was used post hoc to group (code) existing causal factors into various categories. The coding process generally involved two or more coders independently performing the classification task. The agreement levels among coders (inter-rater reliability) reported in these studies was the topic highlighted in our review.

Our analysis revealed a range of reported reliability levels, some acceptable, others not. We also identified several methodological issues that could account for such discrepancies,

Reprint & Copyright © by the Aerospace Medical Association, Alexandria, VA.  
DOI: <https://doi.org/10.3357/AMHP:4555.2017>