# The FAIR Metrics of Adherence to Citation Best Practices

Adam Craig and Carl Taswell

Brain Health Alliance | Web: www.brainhealthalliance.org | Email: ctaswell@brainhealthalliance.org | Phone: +1 (949) 481-3121

## Abstract

Measuring the merits of scholarly research articles only by citation counts and how often other research articles or social media messages cite a particular publication creates a perverse incentive for some authors to refrain from citing potential rivals. This dilemma has developed despite the historical publishing standard expected in peer review for citing and discussing related prior work. To encourage and support a countervailing incentive, research organizations should also consider metrics for how well and appropriately a scholarly article cites relevant prior work in the spirit of the classic phrase and metaphor standing on the shoulders of giants. We present a proposal for a family of such article-level metrics called the FAIR metrics and described as the FAIR Attribution to Indexed Reports or the FAIR Acknowledgment of Information Records.

## Introduction

Citation metrics introduce perverse incentives to cite or ignore prior work for non-scientific reasons [1]. Furthermore, existing plagiarism detection software does not distinguish correctly attributed ideas from plagiarized ones [2]. To address these issues, we developed metrics of adherence to good citation practices.
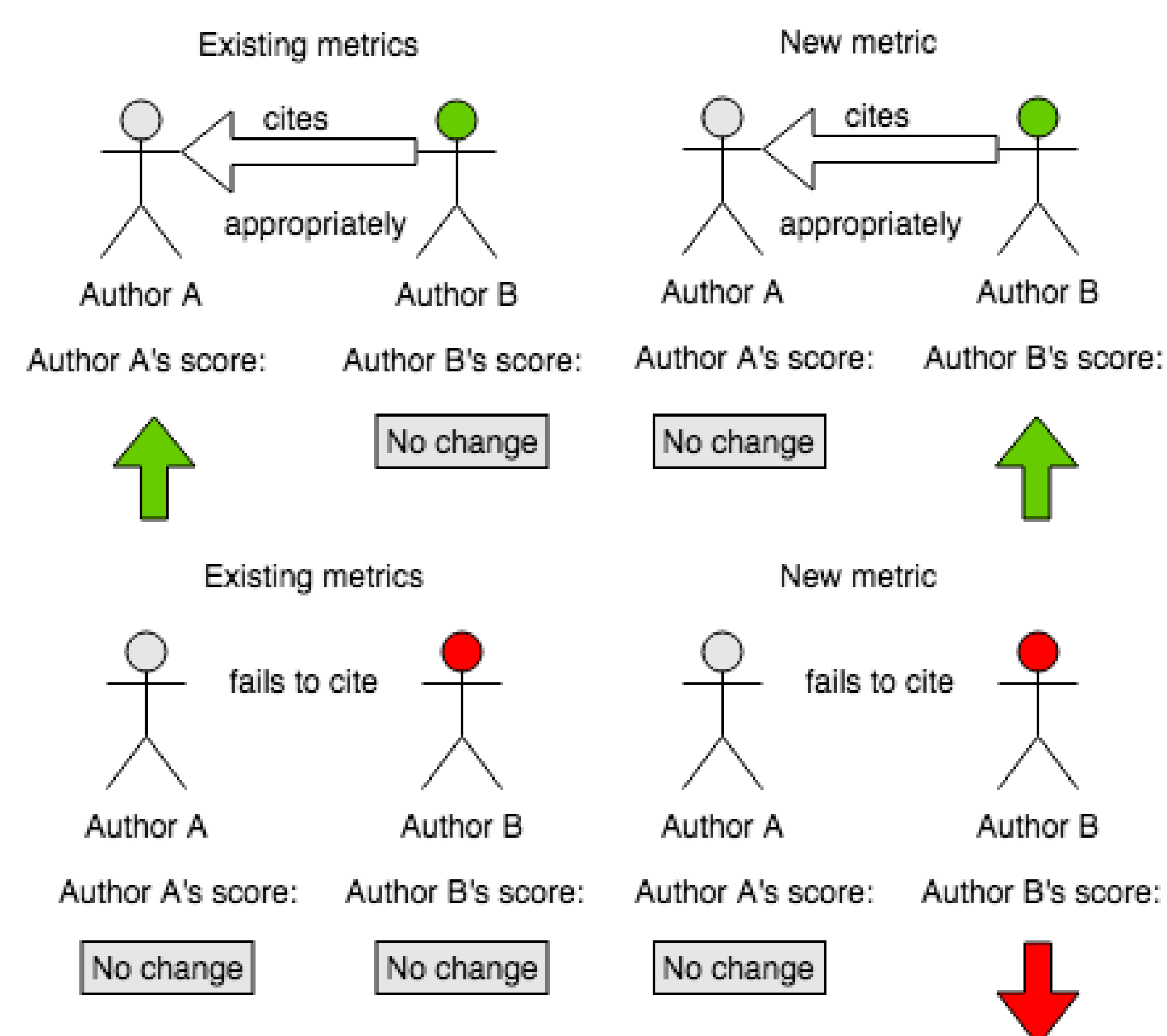


Figure 1: Incentives in conventional vs FAIR metrics

## Methods

We started by identifying key requirements for the metric, identified the core concepts they concern, and defined abstractions to represent them.

| Feature |
| --- |
| distinguishes plagiarism from errors in citation |
| distinguishes well-citing works from ones with errors |
| consistent even when bad practices are common |
| stable against attempts at obfuscation |
| allows comparison across problem domains |
| allows for common knowledge |

Table 1: Prioritized (high-to-low) features of a FAIR Metric

While some of these goals, such as reliably keeping track of what statements are common knowledge in a given field, are challenging and may not be practical, recent developments in semantic text analysis support the feasibility of the central goal of distinguishing legitimate reports of research from acts of plagiarism. In particular, semantic analysis approaches can detect similarity of ideas despite strong obfuscation through paraphrasing [3]. To define our metrics, we first identified and defined six key concepts.

Concept a word or phrase equivalent to an item in a formal ontology or other controlled vocabulary. Define two Concepts to be equivalent if and only if they map to the same term in one or more controlled vocabularies.

Statement a statement that we can represent semantically as an ordered triple of subject, verb, and object Concepts. Define two Statements to be equivalent if they can be represented as the same *subject-verb-object* triple.

Report a scholarly work, such as a primary research article or a secondary review article, a description of which we represent as a set of Statements

Citation a triple indicating that a Statement includes a citation of a Report as its source

Index a set of Reports available for comparison

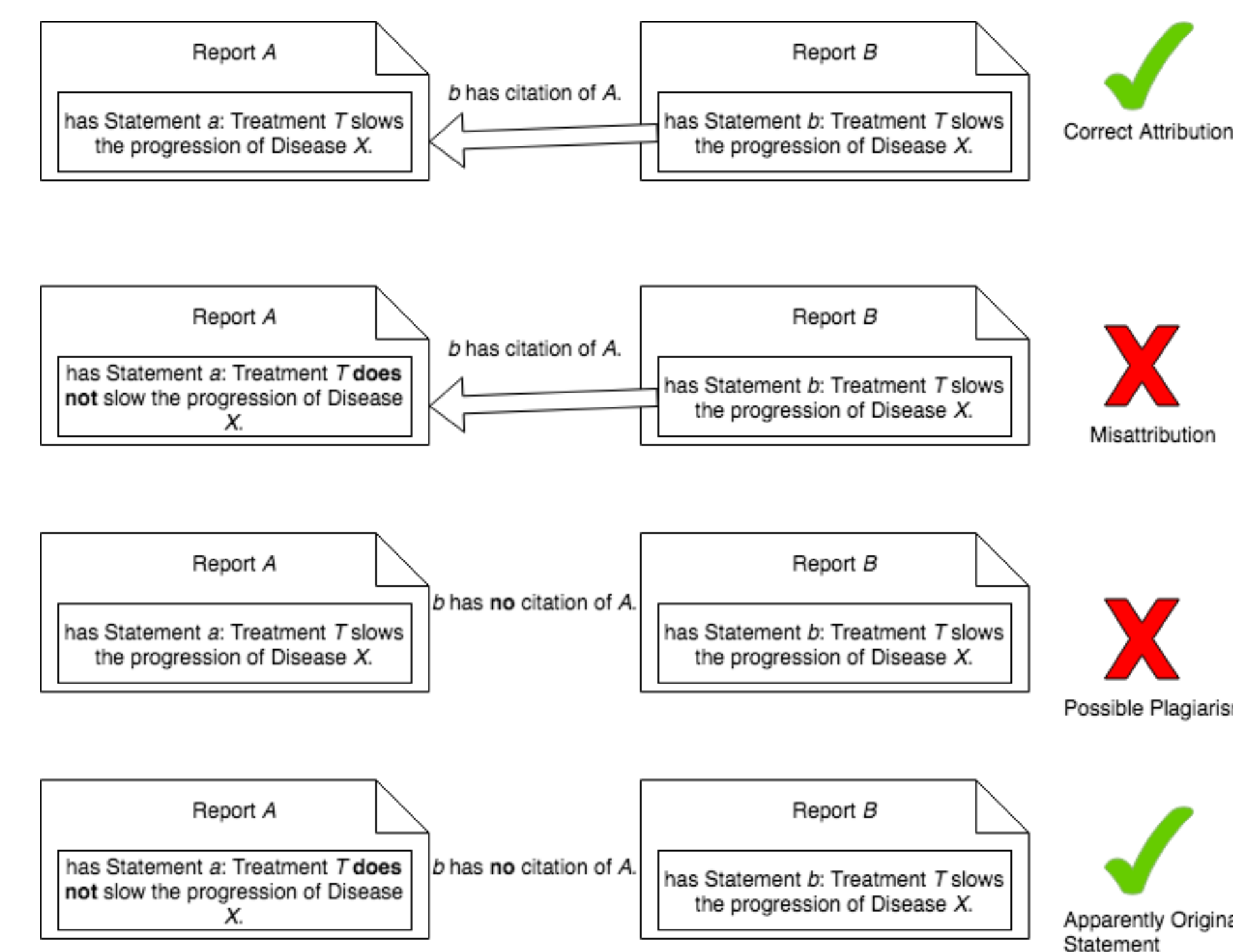Metric a function quantifying presence of a given citation practice, good or bad, in a Report

## Results



Figure 2: Scenarios a FAIR metric must differentiate

Let $I$ be an Index and $b$ be a Report in $I$.

Correctly Attributed Statements $\kappa(b, I) \equiv$ the number of Statements in $b$ that it attributes to any other Report that contains an equivalent Statement.

Misattributed Statements $\mu(b, I) \equiv$ the number of Statements in $b$ that it attributes to any Report that does not contain an equivalent Statement.

Potentially Plagiarized Statements $\rho(b, I) \equiv$ the number of Statements in $b$ that it does not attribute to a prior Report but that have equivalent Statements in at least one prior Report in $I$.

Apparently Original Statements $\alpha(b, I) \equiv$ the number of Statements in $b$ that it does not attribute to a prior Report and that have no equivalent Statement in any prior Report in $I$.

For example, manually parsing the abstract of retracted report $b$ ([4]) into *subject-verb-object* triples revealed $\alpha(b, I) = 7$ apparently original statements and $\rho(b, I) = 22$ unattributed statements equivalent to ones derived from the abstract of Report $a$ ([5]) using the same approach. Intuitively, these values suggest significant plagiarism occurred, reflecting identical experimental procedure and numerical data, despite differences in vegetable used, model fitted to the data, and time and place of the experiment.

## Future Work

To evaluate the utility of the metrics we have described here, we will need to evaluate them on a set of scholarly articles, including ones free of citation errors, ones with correctable mistakes, and ones that have been retracted for plagiarism. Once we have a suitable sampling of results, we will use appropriate statistical tests to judge which metrics differ significantly among these three categories. As part of this effort, we are developing tools to evaluate the metrics over large sets of articles, using the Nexus-PORTAL-DOORS System to manage article metadata and semantic descriptions of statements [6].

## References

[1] M Ryan Haley.
On the inauspicious incentives of the scholar-level h-index: an economist's take on collusive and coercive citation. *Applied Economics Letters*, 24(2):85–89, 2017.

[2] Taiseer Abdalla Elfadil Eisa, Naomie Salim, and Salha Alzahrani.
Existing plagiarism detection techniques: A systematic mapping of the scholarly literature. *Online Information Review*, 39(3):383–400, 2015.

[3] Deepa Gupta et al.
Study on extrinsic text plagiarism detection techniques and tools. *Journal of Engineering Science & Technology Review*, 9(5), 2016.

[4] Fahim Ullah, Min Kang, Mansoor Khan Khattak, and Said Wahab.
Retracted: Experimentally investigated the asparagus (asparagus officinalis l.) drying with flat-plate collector under the natural convection indirect solar dryer. *Food Science & Nutrition*, 6(6):1357–1357.

[5] SK Sansaniwal and M Kumar.
Analysis of ginger drying inside a natural convection indirect solar dryer: An experimental study. *Journal of Mechanical Engineering and Sciences*, 9(unknown):1671–1685, 2015.

[6] Carl Taswell.
A distributed infrastructure for metadata about metadata: the HDMM architectural style and PORTAL-DOORS system. *Future Internet*, 2(2):156–189, 2010.