

## **In Validations We Trust? The Impact of Imperfect Human Annotations as a Gold Standard on the Quality of Validation of Automated Content Analysis**

Hyunjin Song<sup>a</sup>, Petro Tolochko<sup>b</sup>, Jakob-Moritz Eberl<sup>a</sup>, Olga Eisele<sup>a</sup>, Esther Greussing<sup>c</sup>, Tobias Heidenreich<sup>a</sup>, Fabienne Lind<sup>a</sup>, Sebastian Galyga<sup>a</sup>, and Hajo G. Boomgaarden<sup>a</sup>

<sup>a</sup> Department of Communication, University of Vienna, Austria

<sup>b</sup> Department of Political Science, University of Vienna, Austria

<sup>c</sup> Department of Communication and Media Sciences, Technische Universität Braunschweig, Germany

### Contact Information:

Hyunjin Song  
Department of Communication  
University of Vienna  
Rathausstraße 19/1/9  
Vienna, 1010 Austria  
[hyunjin.song@univie.ac.at](mailto:hyunjin.song@univie.ac.at)

Forthcoming in *Political Communication*

### **Funding Acknowledgements**

Olga Eisele is supported by the Austrian Science Fund under the Hertha-Firnberg-Program [Grant no: T-989].

### **Data availability Statement**

The data and replication materials that support the findings of this study are openly available in *zenodo* at <http://doi.org/10.5281/zenodo.3598354>.

## Abstract

Political communication has become one of the central arenas of innovation in the application of automated analysis approaches to ever-growing quantities of digitized texts. However, although researchers routinely and conveniently resort to certain forms of human coding to validate the results derived from automated procedures, in practice the actual “quality assurance” of such a “gold standard” often goes unchecked. Contemporary practices of validation via manual annotations are far from being acknowledged as best practices in the literature, and the reporting and interpretation of validation procedures differ greatly. We systematically assess the connection between the quality of human judgment in manual annotations and the relative performance evaluations of automated procedures against true standards by relying on large-scale Monte Carlo simulations. The results from the simulations confirm that there is a substantially greater risk of a researcher reaching an incorrect conclusion regarding the performance of automated procedures when the quality of manual annotations used for validation is not properly ensured. Our contribution should therefore be regarded as a call for the systematic application of high-quality manual validation materials in any political communication study, drawing on automated text analysis procedures.

*Keywords: Automated text analysis, reliability, validation, Monte Carlo simulations*

## **In Validations We Trust? The Impact of Imperfect Human Annotations as a Gold Standard on the Quality of Validation of Automated Content Analysis**

Political communication has increasingly become one of the central arenas of innovation in the application of automated text analysis to ever-growing quantities of digitized texts. Understanding politics today requires a comprehensive understanding of the ways in which diverse political texts constitute or signify complex political processes. Indeed, given the sheer amounts and the ready availability of such texts, automated analysis procedures have become increasingly useful and necessary. The growing popularity of automated approaches to text analysis is mirrored in dynamic and extensive methodological developments (Grimmer & Stewart, 2013; Hopkins & King, 2010; Van Atteveldt & Peng, 2018), as well as in the application of such methods to a wide range of political communication contents, such as analyses of political newspaper coverage (Young & Soroka, 2012), parliamentary debates (Proksch & Slapin, 2010; Spirling, 2016), congressional bills (Wilkerson et al., 2015), political speeches (Rauh et al., 2017), or millions of social media posts (Barberá et al., 2019).

However, with the growing popularity of such “text-as-data” approaches within the field of political communication, the issue of ensuring the validity of the results has become crucial. To arrive at valid results, text-as-data approaches require proper triangulation of the applied techniques against some “gold standard” or “ground truth” (as some forms of “objective,” or intersubjectively valid, measurements that serve as a reference: Grimmer & Stewart, 2013). This is typically achieved by using human inputs (“human coding” or “manual annotations”) as a benchmark. This procedure is based on the assumption that humans’ understanding of texts (still) outperforms that of machines and that, *if trained correctly*, humans will make the most correct and valid classifications of texts. However, “the quantities we seek to estimate from text [...] are fundamentally unobservable” (Lowe & Benoit, 2013, p. 299), and as documented in traditional content-analytic applications (Ennsner-Jedenastik & Meyer, 2018; Hayes & Krippendorff, 2007;

Krippendorff, 2004), human judgments are in fact no exception to this general rule. In this regard, devising a high-quality measurement instrument and ensuring good coder training, thereby ensuring an adequate level of data quality in traditional manual content analysis, is deemed *the* standard practice in the field.

Nevertheless, as we argue and empirically demonstrate below, there has been a relative lack of parallel attention to ensuring an acceptable level of quality in human coding when such manual annotations are utilized as a gold standard for the validation of text-as-data procedures. It seems to be relatively rare that human-coded materials are properly evaluated before being utilized as validation materials, and relatedly, the methodological details of validation procedures are not consistently reported in a transparent manner. Arguably, such practice holds major risks for both the interpretation of analyses and consequent conclusions, especially in terms of potential bias when evaluating the performance of automated procedures. However, the precise implications of using imperfect human judgments remain insufficiently addressed. Accordingly, this study concerns both the actual practices of the usage and reporting of human-coded gold standards in automated procedures and the possible implications of different qualities of such material. We assess this issue by relying on a large-scale Monte Carlo simulation. Our contribution should be regarded as a call for the systematic application of high-quality human coding for validation procedures in automated content analysis. To this end, it warns against the improper use of human coding as the benchmark in demonstrating the performance of an automated text analysis approach, and additionally formulates recommendations to improve validation practices. Given the social and political relevance of political communication research, it is vital that text-as-data analyses yield valid results.

### **The Use of Human Annotation in Automated Content Analysis**

Following the standard techniques often employed in political communication research, we define “automated content analysis” (or automated text analysis) as a collection of content-

analytic approaches that utilize automated methods to code a large amount of textual data in such a way that the coding itself (e.g., the text classification) is not performed manually, but rather through computational algorithms (Grimmer & Stewart, 2013). Although the term “automated content analysis” in general encompasses a wide variety of forms (e.g., Grimmer & Stewart, 2013; Hopkins & King, 2010; Krippendorff, 2013), our definition inevitably excludes automatic approaches of merely *acquiring* data, data *entry*, or data *management* other than the actual coding or classification process (e.g., Lewis et al., 2013). Instead, we concentrate on two broad and rather common forms – a *dictionary* (lexicon-based) approach and a *supervised machine learning* (SML) approach (see Boumans & Trilling, 2016; Grimmer & Stewart, 2013). As we shall elaborate below, they are rather sensitive to the issue of imperfect gold standards, although the two approaches may nontrivially differ in terms of the *degree* of their potential sensitivity to this issue. Yet regardless of their differences, here we focus on their application in “classification tasks,” i.e., minimizing the human labor required to classify a large collection of documents into known categories.<sup>1</sup> We also assume that our discussion is mainly applicable to methods that classify texts at the *document* level (aggregating some word- or sentence-level textual features *within a document* as the approximation of a document membership). We choose to do so since many of the applications used in political communication heavily rely on such approaches (e.g., Boomgaarden & Vliegenthart, 2009; Burscher et al., 2015; Rooduijn & Pauwels, 2011).

There are two ways of utilizing “human coding” in dictionary and SML analysis: on the one hand in the initial development stage (i.e., in constructing dictionaries, or in training SML algorithms), and on the other hand in the “validation” stage, evaluating the classification performances of the procedures (i.e., *post-measurement* validation). Dictionary approaches and

---

<sup>1</sup> Other frequent aims of automated procedures include: (a) to estimate the location of actors or documents that belong to actors (i.e., *scaling*) in n-dimensional space; or (b) to inductively “discover” new classifications with the help of automated procedures (i.e., unsupervised methods). We do not consider these two in the present manuscript, as the architecture of validation procedures requires a very different setup from what is described here.

SML methods considerably differ in their use of human coding in the initial development stage, which requires problem-specific validation as advocated in the extant literature (see Grimmer & Stewart, 2013; Hopkins & King, 2010).<sup>2</sup> However, regarding *post-measurement validation*, the literature often suggests that post-measurement validation based on “out-of-sample” data represents an ideal architecture of validation (e.g., DiMaggio, 2015; Grimmer & Stewart, 2013; Lowe & Benoit, 2013). Once a researcher has developed an algorithm (or simply used an existing one), separate held-out samples (i.e., data not used to model developing and training) are coded *both* by human coders and by automated procedures, evaluating whether the results from the latter converge into the former. Given that “the validity of a method or tool is dependent on the [specific] context in which it is used” (van Atteveldt & Peng, 2018, p. 87), such post-measurement validation may provide convincing evidence of how well a given tool performs in a specific domain and task at hand, especially when off-the-shelf dictionaries or existing SML classifiers are used in a context in which they were not initially developed or trained. For instance, the results from existing, off-the-shelf dictionaries may be validated against the manual coding of highly trained researchers (e.g., Muddiman, McGregor, & Stroud, 2018; Rooduijn & Pauwels, 2011; Young & Soroka, 2012). Likewise, as exemplified in Scharkow (2013) or in Burscher et al. (2014), a similar approach can be taken for SML methods in evaluating the performance of an algorithm.<sup>3</sup> Although convergent validity against external standards is not the only criterion for establishing the validity of content analytic methods (Krippendorff, 2013), this practice of utilizing human coding in validation primarily owes to the general motivation behind

---

<sup>2</sup> A dictionary approach generally relies on extensive human input in developing an explicit coding rule (e.g., simple lists of keywords, Boolean expressions, syntactic parsers, or regular expressions). In contrast, for SML, specific coding rules in manual annotations are in general rarely explicitly articulated. Nevertheless, the algorithm takes such implicit human judgments as the point of reference, and tries to infer the features of data that best classify the text into different predefined categories. See Grimmer and Stewart (2013) or van Atteveldt and Peng (2018) for details.

<sup>3</sup> In unsupervised methods, validations are *conditional* on the classification or scaling produced by unsupervised methods (e.g., evaluating whether direct hand coding or supervised methods can reproduce the suggested findings: e.g., Lowe & Benoit, 2013). Nevertheless, the use of human coding as a benchmark is not an uncommon practice in unsupervised methods as well.

automated approaches (i.e., automating “human coding”; Grimmer & Stewart, 2013), which evidently implies a clear standard for evaluation (i.e., against human coding).<sup>4</sup> However, due to inherent resource constraints, it is rare to see validation occur *after* such classification tasks in practice (nevertheless, for notable exceptions, see the aforementioned studies).<sup>5</sup>

### **The Myth of Perfect Human Coding in Validation of Automated Content Analysis**

The purpose of using a manually annotated gold standard in validation is, as Krippendorff (2008, p. 6) notes, “to confer validity on the otherwise uncertain research results.” However, this logic essentially requires that the validity of the chosen benchmark itself (i.e., manual annotations by human coders) already be *well-established*: that such human annotations are, at the very least, sufficiently intersubjectively valid (Krippendorff, 2008). As much of the traditional manual content analysis literature suggests (e.g., Enns-Jedenastik & Meyer, 2018; Hayes & Krippendorff, 2007; Krippendorff, 2004), manual coding often produces unreliable judgments under a lack of proper instructions or coder training, especially when the judgment at hand requires a nontrivial degree of inferences and subjectivity to classify latent information (Krippendorff, 2013). For this reason, there is little disagreement within the traditional content analysis literature regarding the need for proper “quality assurance” in the form of developing unambiguous coding categories and coder training (Krippendorff, 2004, 2013), as well as the more transparent reporting of those steps (Lacy et al., 2015). In recent years, the political communication literature has embraced these steps in order to seek a standardization of research practices, and this now constitutes a compulsory aspect of any manual content analysis study.

However, when applying automated procedures, researchers appear to put too much trust

---

<sup>4</sup> For instance, content analysis can also be validated when the actual sources of analyzed text concur with a researcher’s findings (i.e., source-based, *postdictive* validity), or when some theoretically predicted effects of contents actually occur among the audiences of text (i.e., an audience-based, *predictive* validity) when such texts are used in experiments or in the real world.

<sup>5</sup> Due to its labor-intensive nature in producing manually annotated data sets, recent applications in this area increasingly turn to “crowdcoding” (Haselmayer & Jenny, 2017; Lind et al., 2017). We will return to this point in the discussion section.

in naïvely coded human annotations, and the exact methodological details of such validation procedures involving human inputs are not consistently and clearly reported. As an illustration, in Table 1 below we present a review of studies published in peer-reviewed journals indexed in EBSCOhost from 1998 to 2018, the majority (approximately 74%) of which were belongs to political communication broadly defined. We searched for studies with either dictionary- or SML-based methods (based on titles, abstracts, and keywords), and examined whether they employed human coded materials as a benchmark in validation. If so, we also considered whether methodological details such as intercoder reliability were adequately reported.<sup>6</sup>

Of the 73 studies examined, only about 58% ( $N = 42$ ) referred to some sort of validation. Among these, five studies did not use human-annotated materials, whereas 37 (88% of the studies that reported any type of validation, and about 51% of all studies) relied on human annotated materials for validation. Yet only few of those 37 studies (that reported the use of human-annotated validation) adequately reported the quality of human annotated data: indeed, only 14 studies (37.83% of studies reporting human-based validation, and 19.17% of all studies) provided any measures of intercoder reliability, whereas 23 studies did not report *any* intercoder reliability despite the fact that they relied on human coding. Among the 14 studies that did report the intercoder reliability, six adequately reported Krippendorff's alpha, whereas three reported either the percentage agreement or the Holsti agreement measure alone (in the remaining studies, we found other measures such as Scott's Pi or correlation coefficients).<sup>7</sup> In addition, out of 37 studies, only about half ( $N = 18$ ) reported the number of coders, whereas 19 studies did not report the number of coders and/or the total size of the validation data set, rendering it

---

<sup>6</sup> See the online appendix for detailed information regarding data, coding procedures, coded variables, and detailed reliability estimates.

<sup>7</sup> Given that intercoder reliability assesses the *replicability* of resulting data independently from the extraneous circumstances of the data-making process (Krippendorff, 2013), a proper reliability index requires considerations of coder agreement due to *chance*. However, simple percentage agreement (such as Holsti) or Scott's Pi lack such methodological properties, as do correlation-based measures (see Krippendorff, 2013).



impossible to judge the quality of the validation procedures.

-- Table 1 About Here --

Regarding the reported validation metrics of automated approaches, the results were very similar. Certainly, the most commonly used measures of validity were the widely accepted measures of *precision* (13 cases,  $M = 0.74$ ) and *recall* (9 cases,  $M = 0.60$ ). However, other rather uncommon metrics -- such as intercoder reliability (e.g., Holsti, Cohen's Kappa, Krippendorff's alpha) or correlation coefficients -- were also widely used to report validation. There were only three studies (4% of all studies) that reported *both* Krippendorff's alpha (which signals the proper quality assurance of human coding) and a proper reporting of validation metrics.

The observations of published practices show that reported measures of validation, and especially the quality of human-coded data, are far from "best-practice" recommendations from the traditional content analysis literature (e.g., Lacy et al., 2015; Krippendorff, 2013). It is puzzling and discomfoting that such best practices from traditional content analysis seem to be somewhat ignored when it comes to establishing ground truth, and that frequent calls for the proper validation of automated procedures are far from being common practice. Importantly, we *do not* claim that the lack of reporting of methodological details necessarily means a lack of proper quality insurance *per se*. However, a lack of proper reporting of important methodological details nevertheless severely undermines the trustworthiness of the reported validation involving human-coded data as the "ground-truth." While prior contributions on this topic -- whether theoretically (e.g., Grimmer & Stewart, 2013; González-Bailón & Paltoglou, 2015; Hopkins & King, 2010) or empirically (e.g., such as in corpus annotations: Hovy & Lavis, 2010; Lease, 2011) -- have stressed the need for a proper validation of techniques, there is still strikingly little consistency in *whether* and *how* validation is approached and reported. Indeed, a seemingly widespread practice of conveniently utilizing manual annotation without proper quality assurance -- as can be seen in Table 1 -- reveals a conspicuous lack of attention to this issue in

actual research practice.

### **What Price Are We Paying? The Consequences of Low-quality Annotations**

Evaluating the quality of algorithms' automated classification performances is typically undertaken by calculating precision, recall, and F1 scores against manually annotated materials (hereafter termed "observed" performance). If the observed performance of the algorithm is not sufficiently satisfying, such as against some *a priori* chosen threshold, additional steps are sought to improve the quality of the automated procedures (e.g., retraining algorithms, or changing the dictionary).<sup>8</sup> Implicitly, however, this practice treats the observed performance as the *unbiased estimates* of "true" performance (*that could have been observed* against the unknown, "true" standard). From this follows the important question of how well the "observed" performance predicts the true classification performance when researchers use imperfect manual annotations, as well as the size of the potential bias.

The use of (potentially) imperfect, low-quality manual annotation as a benchmark for automatic techniques may have at least two systematic consequences. First, although automated methods themselves are perfectly reliable, imperfect human judgments of validation materials essentially make the ultimate "target" of such reliable measurements radically deviate from the true yet unknown target of inference, rendering them "reliably wrong" on-target. Second, and relatedly, systematically flawed human judgments of validation materials can introduce unknown bias when evaluating algorithms' performance *vis-à-vis* true performance. Nevertheless, most empirical studies appear to pay insufficient attention to this issue. In a recent study, González-Bailón and Paltoglou (2015) compared five available sentiment dictionaries against human annotations, yet they do not directly deal with the implications of imperfect reliability in human

---

<sup>8</sup> Here, we do not consider a possibility that a researcher decides to improve the quality of human coding. Our argument only applies to scenarios where a research has *already* decided that a quality of human coding is good enough to serve as ground truths; therefore, such scenario is conceptually prior to our arguments and scenarios described here. As to the argument we advance here, however, researchers often incorrectly assume that human coding is perfectly reliable and valid.

coding. Scharkow and Bachl (2017), the only existing study looking at imperfect reliability in human coding, mainly deal with its implication on “linkage analysis,” but not on validation of automated content analysis. Indeed, we are aware of only a handful of studies suggesting tentative relationship between the quality of human coding of validation materials and the evaluations of machine-based classification accuracies (Burscher et al., 2014; Snow et al., 2008).

### **A Monte-Carlo Simulation Study**

Although the general intuition regarding the impact of imperfect human judgment in validating automatic procedure is rather clear, elucidating the factors that affect potential bias, and especially its exact magnitude thereof, in the conclusions of such validation is far less straightforward. In order to systematically evaluate the implications of and to provide a more concrete benchmark for the improper use of human-coded materials as gold standards, we set up an extensive set of Monte-Carlo (MC) simulations. MC simulations offer a convenient yet flexible tool for systematically evaluating the relative bias and coverage of a given statistic under certain scenarios (Leemann & Wasserfallen, 2017; Scharkow & Bachl, 2017).

We designed our procedures in a way that would largely mirror the typical approaches used in political communication research, while we systematically varied factors that might affect the quality of human annotation in *post-measurement* validation (see Table 2 below). Here, we assume that the size of the text data is sufficiently large to warrant an automated approach. Therefore, the data in question (exemplarily and hypothetically) cover ten news articles per day for ten news outlets, spanning a total of five years. Accordingly, every single simulation is set to generate 130,000 observations of media articles.

### **Design and Setup of Monte Carlo Simulations**

Ideally, in the creation of ground truth data, two or more trained human coders are assigned to – and independently code – a set of sample documents in order to produce data for intercoder reliability assessment. Once an acceptable level of reliability is reached among coders

(typically Krippendorff's alpha equal to or greater than 0.7), the validation materials are often evenly, yet rarely *randomly*, divided into  $k$ -subsets, each of which is then annotated only by a single coder (Grimmer, King, & Superti, 2018). Indeed, it is still a very common practice to evenly divide coding tasks by some non-random, natural grouping variables (e.g., by media outlets or simply by order of documents) in manual annotations. Given that coders are treated as interchangeable, any (potentially) remaining coder idiosyncrasies (either coder-specific systematic errors or random measurement errors) are in effect no longer considered, neither in the analyses nor in the interpretations of the findings (see Bachl & Scharkow, 2017, for a detailed discussion on this issue). When there is a sufficiently large number of coders, or each materials are coded by multiple coders ("duplicated coding" as in some SML applications or in crowdcoding: see Lind et al., 2017; Scharkow, 2013), the impact of coder idiosyncrasies – especially random errors – would diminish, as they will cancel each other out as long as the number of coders/duplicated coding instances increases. Nevertheless, remaining systematic errors in coder idiosyncrasies may still introduce bias in gold standard materials with respect to the target of inference, especially for data with a higher level of intercoder reliability (i.e., a systematic deviation from the true target). This is even more likely to constitute a serious issue when validation data systematically differ from the training/test data (or from all data) in terms of their textual features, such as when validation sets are biased subsets of the entire data set, or even come from a different context that only remotely relates to the task at hand.

Following this logic, we specified the following factors that may affect the "quality" of the gold standard (i.e., manual annotations by human coders) and therefore the evaluations of the performance of automated algorithms (also see Table 2): (a) whether validation materials are not randomly divided among coders ("sole coding") vs. all coders independently evaluating the

entire body of material (“duplicated coding”);<sup>9</sup> (b) the proper sampling variability of validation materials (e.g., whether validation materials are a random sample of test sets vs. systematically biased subsets); (c) the number of human coders ( $k = 2, 5, 10$ , which roughly corresponds to minimum, intermediate, and a large number of human coders); (d) the levels of intercoder reliability (Krippendorff’s  $\alpha = 0.5, 0.7, 0.9$ , deemed either low, acceptable, or high); and (e) the size of the validation data set ( $n = 650, 1300, 6500, 13000$ , approximately corresponding to 0.5%, 1%, 5%, and 10% of the total data set,  $N = 130,000$ ). Although one must make practical and logistical decisions regarding these factors in any real-world application (typically by resource constraints), they are indeed crucial in terms of properly ensuring the acceptable quality of manual annotations. The specific cases in these scenarios were chosen to reflect typical procedures and their common variations.

– Table 2 About Here –

Using the above setup, we compared different scenarios in terms of their F1 scores based on imperfect validation materials vs. one based on a “true” standard. In practice, knowing the true performance of an algorithm would be impossible, because the true outcome value of textual data can never be known independently from (potentially imperfect) human coding. However, because we were simulating textual data, we could systematically evaluate the *true* classification performance of automated procedures against “observed” performances from scenarios using varying qualities of human-annotated gold standard materials. In so doing, we could illustrate how different practices of utilizing human coding in automated content analyses adversely affect the relative trustworthiness of the procedures’ conclusions.

The final Monte Carlo simulation used  $3$  (number of human coders,  $k$ )  $\times$   $3$  (target Krippendorff’s  $\alpha$  levels in human annotations in validation data)  $\times$   $4$  ( $N$  of total annotations)

---

<sup>9</sup> In case of duplicated coding of entire materials, we assume the gold-standard materials are determined by the “majority rule” (assuming equal qualification of coders at given reliability level) following the common practice in the field (see Haselmayer & Jenny, 2016; Lind et al., 2017).

$\times 2$  (sole coding vs. duplicated coding)  $\times 2$  (random sample vs. biased subset for validation) full factorial design, with 1000 replications per scenario ( $N = 144,000$ ). As dictionary and SML methods require different data structures and implementations of coding rules for algorithms, we separately performed two MC simulations for each (see the online appendix for detailed rationales and descriptions of our simulations).<sup>10</sup>

### Simulation Results

We first present the mean absolute prediction error (MAPE), defined as  $\sum abs(FI_{validation} - FI_{true}) / N$ , which measures the average degree of absolute bias in observed F1 scores (from the human-coded validated data) *vis-à-vis* true, unknown F1 scores. In essence, this can be framed as a prediction problem – to what degree observed F1 scores deviate from true F1 scores – when using the observed F1 score as the best possible “prediction” of the true F1 score. The MAPE is a commonly used metric for measuring the predictive accuracy of continuous measures and is regarded as one of the most robust measures of such (Hyndman & Koehler, 2006). Table 3 reports the results of ANOVA predicting MAPE (i.e., the mean of APE per 1000 runs) as the outcome of interest, along with the marginal means and contrast for every experimental factor. Finally, we report  $\omega^2$ , a robust effect size measure for ANOVA, whose interpretation can be regarded as the percentage of total variance explained by the variance of a factor in question.

– Table 3 About Here –

Both for SML and dictionary-based approaches, three out of five experimental factors appeared to reduce the mean absolute prediction error of observed F1 scores in approximating the true F1 score levels (see Table 3 for details). For the SML approach, we see no overall gain when relying on duplicated coding ( $df = 1$ ,  $F = .00$ ,  $p = \text{n.s.}$ ) to produce human-annotated validation materials, as seen when contrasting *sole coding* (MAPE = .0389) and *duplicated*

---

<sup>10</sup> A set of replication codes and data for this manuscript can be found at <http://doi.org/10.5281/zenodo.3598354>.

*coding* (MAPE = .0390) in Table 3. Similarly, we see no discernible effect of the *number of coders* ( $df = 2$ ,  $F = .01$ ,  $p = \text{n.s.}$ ) across simulation scenarios.<sup>11</sup> In contrast, the level of intercoder reliability (i.e., Krippendorff alpha) presented the largest independent overall effect ( $df = 2$ ,  $F = 491.75$ ,  $p < .001$ ,  $\omega^2 = .620$ , or 62% variance explained) for SML scenarios, such that hand-coded data with the highest reliability level had approximately half of the MAPE (= .0262) compared to that of the lowest level (MAPE = .0550). When the sampling variability of validation samples accurately reflected the variability of the entire body of data of interest ( $df = 1$ ,  $F = 399.73$ ,  $p < .001$ ,  $\omega^2 = .252$ , or 25.2% variance explained), the magnitude of error also slightly decreased from .0466 (non-random) to .0313 (random subset). Lastly, the size of the validation data set ( $df = 3$ ,  $F = 22.00$ ,  $p < .001$ ,  $\omega^2 = .040$ ) had small yet discernible consequences, such that the prediction error gradually diminished from .0435 (with  $N = 600$ ) to .0365 (with  $N = 6,500$ ). However, beyond that point, the size of the annotation ( $N = 13,000$ ), whose marginal mean was indistinguishable from  $N = 6,500$  cases, had no practical gain of reducing MAPE (we will return to this point in the discussion section).

As Panel A of Figure 1 shows, the intercoder reliability also had an interactive effect on other factors (see also the online appendix for associated ANOVA results). Figure 1 presents the point estimates of the MAPEs, along with their 68% and 95% confidence intervals (respectively in thin and bold vertical lines, corresponding to 1SD and 2SD plus and minus from the MAPE) across total simulation runs per scenario. The pattern suggests that validation data with higher reliability levels had far fewer MAPEs when the size of validation data were larger and were derived from proper random samples (a rather straightforward result), such that the mean levels of the expected discrepancy between observed and true F1 scores were as high as 0.061 (SML) and 0.091 (dictionary) for the smallest non-random validation data with the lowest reliability, but

---

<sup>11</sup> Importantly, we also failed to find interaction effects of these factors together with the rest of the factors (see online appendix).

as low as .001 (SML) and 0.016 (dictionary) for the largest random validation data with the highest reliability.<sup>12</sup> Although the benefit of greater reliability in human-coded validation data did not appear to be evident in the non-random sample, increased intercoder reliability was readily visible in even the smallest set of validation data (e.g.,  $N = 600$ ) whose sampling variability was properly ensured. While both the size of the validation data and the intercoder reliability generally reduced prediction bias, it appears that these two factors largely compensate for one another, albeit only at the moderate ( $K\alpha = .7$ ) to high ( $= .9$ ) reliability level with proper sampling (we will return to this point in the discussion section).

-- Figure 1 About Here --

For the dictionary approaches, largely similar patterns emerged (see Table 3 and Panel B of Figure 1). Indeed, the size of the validation data set ( $df = 3$ ,  $F = 62.37$ ,  $p < .001$ ,  $\omega^2 = .130$ ), the level of intercoder reliability ( $df = 2$ ,  $F = 34.07$ ,  $p < .001$ ,  $\omega^2 = .047$ ), and the sampling variability of the validation data set ( $df = 1$ ,  $F = 1025.31$ ,  $\omega^2 = .723$ ) were all significant in explaining the MAPEs across simulation scenarios. Therefore, the results for both SML and dictionary scenarios consistently suggested that largely the same factors systematically affect uncertainties of the conclusions one can draw from proposed automated procedures.

Our next set of descriptions focuses on a researcher's "decision accuracy" regarding the overall performance of the algorithms. Researchers usually set an *a priori* threshold for desired performance levels (e.g., F1 score equal to or greater than .624), and if the observed performance of the algorithm is higher than this threshold, it is deemed acceptable. Within this context, we defined a decision based on the observed F1 score as "accurate" when this is consistent with a decision based on the true F1 score regarding the performance of an algorithm. Given that this is effectively also a function of the specific threshold that one initially targets, we considered three exemplary values (which we have conveniently chosen) here for the sake of presentation: -1SD

---

<sup>12</sup> We present the equivalent plots for the full combination of factors in the online appendix.



(= .483), mean (= .624), and +1SD (= .766) of the true F1 scores in our simulations.<sup>13</sup> For each chosen threshold, we then evaluated the potential decision's (in)accuracy by calculating the percentages of simulation runs that a researcher's decision would fall into four mutually exclusive categories of true positive, true negative, false positive (analogue to Type I error) and false negative (analogue to Type II error), as a function of our experimental factors. Effectively, the proportion of false positive and false negative cases we present below can be regarded as the mean expected error rates based on different combinations of factors one might consider when producing validation materials, providing a reference point of which one can probabilistically expect erroneous classifications under each combination of factors. It may also be seen as an analogue of simulation-based power analysis, evaluating the proportion of cases classified as false positive (alpha), true positive (1 - beta, i.e., statistical power), false negative, and true negative. As the duplicated coding and the number of coders did not have any discernible effects, we collapsed the categories when calculating the percentages.

– Figures 2 and 3 About Here –

Figures 2 and 3 present the potential decision accuracy rates for SML (Figure 2) and for dictionary scenarios (Figure 3). The results across different threshold levels indicate that most decision errors are false negative cases, being as high as 14.9% in SML scenarios and 9.12% in dictionary scenarios, on average. It therefore appears that the observed classification quality (against manually coded validation materials) of these applications tends to *underestimate* the true classification quality. While the false negative rate generally decreases with all of the experimental factors presented here, the biggest gains again appear to be based on increased intercoder reliability. For low K alpha (= .50), the overall false negative rates across all three F1 thresholds were 11.23% for SML and 5.25% for dictionary scenarios, decreasing to 6.78%

---

<sup>13</sup> Although the range of F1 threshold values were somewhat arbitrarily chosen, such ranges are nevertheless substantially plausible in practice. Indeed, in our earlier reported review of relevant literature, the overall mean of reported F1 was 0.644, with a range of 0.33 (min) to 0.9 (max).

(SML) and 3.98% (dictionary, with  $K\text{ alpha} = .70$ ) until 3.93% (SML) and 3.04% (dictionary, with  $K\text{ alpha} = .90$ ), respectively. This result suggests that if statistical power is the utmost concern when determining whether the classifier performance is acceptable, the best way to achieve such a goal is to ensure high intercoder reliability in manually annotated materials.

For a higher threshold of F1 score, the results indicated an increased risk of false positive errors, where the magnitude of such error is as high as 5.33% for SML scenarios and as high as 4.28% for dictionary scenarios (both with scenarios with  $F1 = .766$ ,  $K\text{ alpha} = .9$ ,  $N = 600$ , non-random sampling). This suggests that the potential decision error in performance evaluation based on observed F1 scores is a much greater issue when a researcher tries to target a higher threshold with a low amount of validation data. Comparing identical scenarios across SML vs. dictionary approaches, it seems that SML scenarios generate slightly more optimistic results in terms of their potential decision error rates. Although speculative, this may be explained by the differences in decision rules between SML and dictionary approaches, as the former can carefully calibrate their predictions, whereas the latter make rather monotonic, deterministic predictions. Regardless, as the size of the validation data set increases, the false positive rate generally decreases in all scenarios.

However, those false positive cases disproportionately *increase* in scenarios with non-random validation data (i.e., a biased subset of entire data for validation materials) whose intercoder reliability is higher. Although it may initially seem counterintuitive, it makes sense that highly calibrated but biased validation materials would “reliably” deviate from the true (yet unknown) target of inference, making them “reliably wrong” on-target. Under the same setup but with randomly sampled validation data, however, the results may suggest that such tests offer the most powerful (in terms of statistical power) *and* the most accurate (in terms of minimizing false positive errors) results for performance evaluation. While this appears to be almost self-evident, ensuring sufficient sampling variability is a rather difficult issue in practice. Indeed, validation

materials often come from different contexts or points in time, and researchers seldom worry about whether potentially relevant yet unobserved factors in their validation data accurately resemble equivalent features in the test/training (or entire) data set. Considering the fundamental uncertainty of proper sampling variability for the validation data, we suggest that one must strive to exercise greater caution when evaluating an algorithm's performance, especially with small (in terms of total  $N$ ) yet high-quality (i.e., high reliability) hand-coded data. Nonetheless, all of the false positive rates were generally less than 5% (i.e., 95% confidence level) in such cases.

### Discussion

As one of the subfields of social science that is at the forefront of methodological developments in text-as-data approaches, as well as its active use of such techniques to answer various questions, political communication has emerged as the central arena regarding the use of automated content analysis techniques. Therefore, the actual practices of the proper validation of automated text analysis in extant research, as well as continued discussions surrounding the issue of possible best practices, have profound implications for the field as a whole. Yet, our review of published research within the field has shown that there is still strikingly little consistency in *whether* and *how* the validation of automated procedures is approached and reported. Indeed, studies often fail to report *any* validation metrics when relying on automated methods. Moreover, even when they do, the quality of human coding (when utilized as the gold standard) is generally not properly evaluated. Against such practice, we attempted to provide further insights into the consequences of using suboptimal quality manual annotations as a gold standard in automated procedures. The results of our MC simulation revealed that intercoder reliability (Krippendorff's alpha), the size of the validation datasets, and the proper sampling variability of such validation datasets produce systematic consequences for a researcher's ability to correctly evaluate the classification accuracy of the proposed algorithms. In sum, our results give good reasons for concern about the quality (or rather the *validity*) of conclusions drawn from automated content

analyses if the proper quality assurance of gold standard data is not guaranteed.

To be clear, the current study *does not* make the argument that we should exclusively rely on human validations, or conversely, that human validations in general are a problem. To the contrary, one of the main points being advanced here is that humans are not perfect. Therefore, proper validation to ensure the “quality” of manual annotation is essential, especially when it is utilized as the gold standard in automated procedures. While the (potentially imperfect) human gold standard is often the best we can get, the results suggest that extra caution must be taken.

### **Some Recommendations for Improving Validation Practices**

No single foolproof solution applicable to every situation exists, and what counts as “best practices” for automated text analysis remains in the early stages of development. However, based on our observations from the literature review and from our simulation results, we can offer some recommendations to help improve the practice of utilizing human annotations in the validation of automated approaches. First, we believe that not every study that relies on automated content analysis needs to use human-involved validation, given that different research contexts and different research questions ultimately dictate whether one should use human validation. However, if the main motivation behind using an automated classification algorithm is to efficiently replace costly human judgments, automated procedures should be validated against equivalent forms of human coding (DiMaggio, 2015; Grimmer & Stewart, 2013). In such cases, researchers should adhere to rigorous methodological standards to the degree expected for traditional manual content analysis (e.g., Hayes & Krippendorff, 2007; Krippendorff, 2013) in preparing a manually annotated validation data set. As with all research, the methodological details (such as measurement details, coding instructions, sampling, coder training, and intercoder reliability) of such validation data should be fully disclosed, enabling readers to independently judge the soundness of the validation procedures, as well as to increase the transparency and replicability of the research process (for instance, see Muddiman et al., 2018;

Rooduijn & Pauwels, 2011; Young & Soroka, 2012).

Second, we recommend that researchers pay closer attention to the issue of the proper sampling variability of validation data *vis-à-vis* all the data to which one wishes to apply a given algorithm (especially for dictionary approaches), as well as the intercoder reliability of human coding (especially for SML approaches) for validation data. We have observed that these two factors most strongly explain (in terms of  $\omega^2$ ) the mean prediction error in predicting true F1 scores based on the observed performance of the algorithm on such validation data. Although this is almost self-evident, in practice ensuring proper sampling variability *vis-à-vis* the entire data set is not an easy task, especially when validation datasets are not collected simultaneously (e.g., forward in time, or from different contexts from the data at hand). When there is an uneven distribution of coding categories in the validation data (or in the entire data to be analyzed), this issue may become even more important especially when the reliability of human coding is low. Under such a situation, the risk of random correct/incorrect classifications generally increases, which may additionally bias the ultimate accuracy of the predictions based on such data.

Regarding the total size of the validation sample, we noted in both SML-based and dictionary-based scenarios that beyond  $N = 6,500$  there were no discernable *independent* effects of increased sample size (see Table 3 for details). However, increasing the sample size appeared to have compensating effects on other factors, especially on intercoder reliability. Depending on practical situations, researchers may therefore prioritize a different factor based on relative trade-offs, given the benchmarks we suggest here. For instance, in many political communication applications such as social media postings, the nature of judgments in human coding is very subjective (hence coding practices might not be sloppy, but still have lots of intercoder unreliability), yet easy to scale up in order to achieve a large number of annotations. In such cases, researchers may want to prioritize a larger size of validation dataset over higher reliability levels, as demonstrated in recent examples of analyses of moral intuitions (e.g., see Weber et al.,

2018, experiments 5-6), of short texts on Twitter or of comments on news sections (e.g., González-Bailón & Paltoglou, 2015). Nevertheless, we advise researchers to strive to increase the sizes of manually coded validation dataset as large as possible, preferably to more than  $N = 1,300$  (i.e., more than 1% of all data to be examined), assuming acceptable reliability (equal to or higher than .7). The results of our simulation study suggest that the risk of potential decision errors is substantially higher in smaller manual annotation data, especially when targeting higher F1 scores for the performance threshold. Nevertheless, the percentages of decision errors regarding the true performance of algorithms appear to reach the acceptable range (less than 5% of false positive rates and less than 10% of false negative rates, or statistical power greater than .9) when the size of the manually annotated validation dataset is equal to or greater than  $N = 1,300$  in all scenarios with a reliability level equal to or higher than 0.7. In contrast, when intercoder reliability is low, increasing the total size of validation dataset to more than  $N = 6,500$  (i.e., more than 5% of all data) may help to reach the equivalent level of maximum error rates of 10%. Of course, this is a somewhat arbitrary threshold, yet the maximum error rate of 10% roughly represents a balance of considerations between a typical false positive rate (i.e.,  $\alpha = 0.05$ ) and maintaining sufficient statistical power (i.e.,  $1 - \beta$  of 0.90, where  $\beta$  is a false negative rate) employed in the field, as advocated in a recent study by Holbert et al. (2018). While this number only reflects a very rough rule of thumb for highly simplified cases (i.e., a binary judgment involving only one variable), an informed judgment of whether and which types of errors to prioritize – the false positive rate, the false negative rate, or the combined error rates – could be based on the simulation evidence presented here. For more complex and subtle judgments, such as ambiguous latent contents or non-binary judgments, it would make sense to use larger amounts of manual data in validation tasks.

Third and relatedly, we have observed that improving intercoder reliability in human coding offers the greatest net benefit in reducing the magnitude of potential bias in automated

tasks relative to a true, unknown standard. While this seems encouraging for many researchers, it simultaneously warns against the prevalent practice of only prioritizing the minimally acceptable level of intercoder reliability without considering other factors. Certainly, such supposed improvements may introduce *increased* systematic errors – as evident in our decision error analyses – especially in cases with very small (e.g.,  $N = 600$ ), non-randomly sampled validation dataset. Recent developments in “crowdcoding” for content analysis (Haselmayer & Jenny, 2017; Lind et al., 2017) could alternatively offer promising ways of scaling up manual annotation tasks, thereby increasing manual annotations in validation tasks, if resource constraints that preclude the use of trained coders are high. Nonetheless, the additional issue of appropriate quality control for crowdworkers (e.g., selection of workers based on task-relevant background knowledge, designing proper material presentation and option formats for an online environment, choosing optimal workload/workflow and compensations) can quickly become important (see Lease, 2011, and Lind et al., 2017, for a related discussion on this issue).

Fourth, one should also bear in mind that without proper coder training in improving reliability, human coders often experience substantial “coding drift” over time (i.e., low *intracoder* reliability), and this often goes hand in hand with low intercoder reliability, too. In such cases, the risk of introducing additional random errors due to low intracoder reliability runs very high. However, given that our simulation setup did not take intracoder reliability (rather, only intercoder reliability) into account, our simulation would have produced more “optimistic” results than most real-world, low intercoder reliability scenarios. Consequently, our results should *not* be interpreted as the indication that intercoder reliability or human coders do not matter in the validation of automated procedures.

Having addressed the proper quality assurance of human-coded validation data, we also briefly consider the issue of choosing appropriate validation metrics when reporting final classification performance based on such data. In our review, we have seen that coder-classifier

reliabilities or correlation coefficients are widely used as validation metrics. However, as Krippendorff (2013) notes, intercoder reliability itself (or indeed any correlation-based measure) does not concern “truth” (i.e., deviations from a given standard, or *performance* against some benchmark) in evaluating coder (dis)agreement, rendering it an inappropriate metric for validation. In contrast, considering general motivations behind automated procedures (i.e., replacing “human coding”; Grimmer & Stewart, 2013), precision and recall would provide a direct quantification of such “performance” (i.e., the degree of deviations from the standard). While an informed argument can be made for the use of coder-classifier reliability for reporting validation, as there is currently no universal standard for how algorithmic coders should be treated in conjunction with human coders, regarding the computer as the *n*-th coder is only appropriate under the assumption that human coders *already* produce “reliable” and therefore intersubjectively valid results in coding classification tasks. Furthermore, algorithmic coding requires many pre-processing steps that are unique to computers (i.e., human coders do not require such pre-processing steps), violating the basic assumption of the interchangeability of coders and identical procedures/data in reliability assessment. Nevertheless, without proper human coder training, which would ensure the quality of data produced, reliability assessment with computers as *n*-th coders does not guarantee the ultimate validity of findings from automated procedures (yet the same would be true of precision and recall when the validity of the human-coded gold standard itself is questionable).<sup>14</sup>

Lastly, our results can be further used for benchmarking the expected discrepancy

---

<sup>14</sup> For nominal judgements, coder-classifier confusion matrix (i.e., precision and recall) can be directly converted to K alpha level (see Krippendorff, 2013, Ch. 12). As such, coder-classifier reliability (while treating an algorithm as the *n*-th coder) indeed can be indicative of validity of automated procedures. Yet this nevertheless assumes that human coders produce reliable and acceptable judgements. When human coders and automated classification produce both conceptually “wrong” judgements, showing high reliability between the two does not necessarily mean validity of automated classifications. Nevertheless, the objection we raise here with such practice is not about its use per se, but rather, without assuring the quality of human annotations, it may complicate the conceptual validity issue with a mere reliability between coder-classifier.



between observed and true F1 scores in performance evaluation of a given automated algorithm based on the combination of factors we examined here. For instance, for the smallest ( $N = 600$ ), non-random validation dataset with the lowest reliability of  $K = .5$ , the mean expected discrepancy of observed versus true F1 scores was as high as 0.061 (for SML) and 0.091 (for dictionary), according to the summary presented in Panel A in Figure 1. Given that this number represented the absolute difference between the two, one might utilize this information to construct the lowest bound of target F1 scores by additionally considering such average errors. For instance, if one sets the *a priori* performance cut-off at 0.624 for the smallest (0.1% of all data), non-randomly sampled validation dataset with the lowest reliability of  $K = .5$ , one would regard a given SML algorithm as good enough *only when* the observed F1 score is equal to or above the .685 (.624 plus .061). Furthermore, if one wishes to apply different algorithms, then one can effectively re-estimate our simulation models but with a proposed algorithm instead, thereby deriving the expected errors under such scenarios.

### Limitations and Conclusions

A few limitations should be noted. First, our literature review was limited to prior studies that explicitly contain our search terms in their title, abstract, or keywords. Such a sample, by definition, does not include all potentially relevant articles. Nevertheless, it allows us to make a sensible selection of extant articles that not only use these methods but also advertise that they do so. Second, in the simulations we only considered a binary classification task and a single dimension of validation metrics, specifically recall, precision, and resulting F1 scores. Although this greatly simplified our main arguments and complex simulation setups, non-trivial numbers of existing applications that go beyond such simple classifications, dealing with numerical forms of predictions (i.e., document scaling). Regardless, we reason that our core arguments are equally applicable to more complex forms of automated content analysis as well. Given the additional complexity and difficulties for human coders in such non-binary predictions, achieving

acceptable intercoder reliability for manual validation materials in such applications is likely to prove even more difficult than in simple binary counterparts. Therefore, we suspect that the potential problems of “imperfect quality” in human-coded validation materials should be greater (or at least identical) in those applications.

Designing a simulation-based study has afforded us a unique opportunity to observe numerous potential counterfactual scenarios of research processes. When carefully designed, such an approach allows researchers to robustly explore potentially important variability in research practices and their consequences. Furthermore, one of the core advantages of relying on such a simulation-based approach is that it provides an opportunity to formally check the sensitivity of one’s findings, or enables one to explore possible counterfactual scenarios. Nevertheless, such clairvoyance comes at a price: the degree of abstraction and simplification. This simplification is done not only for a computational, but also for a conceptual reason -- designing a simulation “that is simple enough to be comprehensible, yet realistically models the underlying process of interest” (Scharkow & Bachl, 2017, p. 330). Although we acknowledge that our setup could have been constructed in a more realistic but more complex way, our ability to accurately simulate human behaviors does not necessarily depend on very complex models.

Notwithstanding the aforementioned limitations, we believe that our contribution can further prompt political communication researchers to pay closer attention to issues associated with the systematic and proper validation of automated content analytic methods, especially those involving manually annotated materials as a gold standard. It is worth stressing that automated content analysis also takes – and indeed should take – a considerable amount of time and resources in designing a study and evaluating its performance. To this end, a little extra effort in designing proper validation is highly worthwhile, enabling valid inferences and conclusions to be drawn. Only in this way is it possible to ensure that results can be considered meaningful for understanding the political processes that are signified through textual data.

## References

- Bachl, M., & Scharkow, M. (2017). Correcting measurement error in content analysis. *Communication Methods & Measures*, 11, 87-104. doi: 10.1080/19312458.2017.1305103
- Barberá, P., Casas, A., Nagler, J., Egan, P. J., Bonneau, R., Jost, J. T., & Tucker, J. A. (2019). Who leads? who follows? Measuring issue attention and agenda setting by legislators and the mass public using social media data. *American Political Science Review*, Online first, 1-19. doi:10.1017/S0003055419000352
- Boomgaarden, H. G., & Vliegenthart, R. (2009). How news content influences anti-immigration attitudes: Germany, 1993–2005. *European Journal of Political Research*, 48, 516–542. doi:10.1111/j.1475-6765.2009.01831.x
- Boumans, J. W., & Trilling, D. (2016). Taking stock of the toolkit: An overview of relevant automated content analysis approaches and techniques for digital journalism scholars. *Digital Journalism*, 4, 8–23. doi:10.1080/21670811.2015.1096598
- Burscher, B., Odijk, D., Vliegenthart, R., De Rijke, M., & De Vreese, C. H. (2014). Teaching the computer to code frames in news: Comparing two supervised machine learning approaches to frame analysis. *Communication Methods & Measures*, 8, 190–206. doi: 10.1080/19312458.2014.937527
- Burscher, B., Vliegenthart, R., & De Vreese, C. H. (2015). Using supervised machine learning to code policy issues: Can classifiers generalize across contexts? *The ANNALS of the American Academy of Political and Social Science*, 659, 122–131. doi: 10.1177/0002716215569441
- DiMaggio, P. (2015). Adapting computational text analysis to social science (and vice versa). *Big Data & Society*, 2(2). doi: 10.1177/2053951715602908
- DiMaggio, P., Nag, M., & Blei, D. (2013). Exploiting affinities between topic modeling and the sociological perspective on culture: Application to newspaper coverage of U.S.

- government arts funding. *Poetics*, 41, 570–606. doi: 10.1016/j.poetic.2013.08.004
- Enns-Jedenastik, L., & Meyer, T. M. (2018). The impact of party cues on manual coding of political texts. *Political Science Research & Methods*, 6, 625–633. doi:10.1017/psrm.2017.29
- González-Bailón, S., & Paltoglou, G. (2015). Signals of public opinion in online communication: A comparison of methods and data sources. *The ANNALS of the American Academy of Political and Social Science*, 659, 95–107. doi: 10.1177/0002716215569192
- Grimmer, J., & Stewart, B. M. (2013). Text as data: The promise and pitfalls of automatic content analysis methods for political texts. *Political Analysis*, 21, 267–297. doi:10.1093/pan/mps028
- Grimmer, J., King, G., & Superti, C. (2018). The unreliability of measures of intercoder reliability, and what to do about it. Unpublished manuscript. Retrieved from <http://web.stanford.edu/~jgrimmer/Handbib.pdf>
- Haselmayer, M., & Jenny, M. (2017). Sentiment analysis of political communication: Combining a dictionary approach with crowdcoding. *Quality & Quantity*, 51, 2623–2646. doi:10.1007/s11135-016-0412-4
- Hayes, A. F., & Krippendorff, K. (2007). Answering the call for a standard reliability measure for coding data. *Communication Methods and Measures*, 1, 77–89. doi:10.1080/19312450709336664
- Holbert, R. L., Hardy, B. W., Park, E., Robinson, N. W., Jung, H., ... & Sweeney, K. (2018). Addressing a statistical power-alpha level blind spot in political-and health-related media research: Discontinuous criterion power analyses. *Annals of the International Communication Association*, 42, 75-92. doi: 10.1080/23808985.2018.1459198
- Hopkins, D. J., & King, G. (2010). A method of automated nonparametric content analysis for social science. *American Journal of Political Science*, 54, 229–247. Doi:10.1111/j.1540-

5907.2009.00428.x

Hovy, E., & Lavid, J. (2010). Towards a “science” of corpus annotation: A new methodological challenge for corpus linguistics. *International Journal of Translation*, 22, 13-36.

Hyndman, R. J., & Koehler, A. B. (2006). Another look at measures of forecast accuracy.

*International Journal of Forecasting*, 22, 679-688. doi: 10.1016/j.ijforecast.2006.03.001

Krippendorff, K. (2004). Reliability in content analysis: Some common misconceptions and recommendations. *Human Communication Research*, 30, 411–433. doi:10.1111/j.1468-2958.2004.tb00738.x

Krippendorff, K. (2008). Validity. In W. Donsbach (Ed.), *The international encyclopedia of communication*. Hoboken, NJ: Blackwell Publishing.

doi:10.1002/9781405186407.wbiecv001

Krippendorff, K. (2013). *Content analysis: An introduction to its methodology* (3rd ed.).

Thousand Oaks, CA: Sage.

Lacy, S., Watson, B. R., Riffe, D., & Lovejoy, J. (2015). Issues and best practices in content analysis. *Journalism & Mass Communication Quarterly*, 92, 791-811. doi:

10.1177/1077699015607338

Lease, M. (2011, August). On quality control and machine learning in crowdsourcing. In

*Workshops at the Twenty-Fifth AAAI Conference on Artificial Intelligence*.

Leemann, L., & Wasserfallen, F. (2017). Extending the use and prediction precision of

subnational public opinion estimation. *American Journal of Political Science*, 61, 1003–1022. doi:10.1111/ajps.12319

Lewis, S. C., Zamith, R., & Hermida, A. (2013). Content analysis in an era of big data: A hybrid approach to computational and manual methods. *Journal of Broadcasting & Electronic*

*Media*, 57, 34–52. doi:10.1080/08838151.2012.761702

Lind, F., Gruber, M., & Boomgaarden, H. G. (2017). Content analysis by the crowd: Assessing

- the usability of crowdsourcing for coding latent constructs. *Communication Methods & Measures*, 11, 191–209. doi:10.1080/19312458.2017.1317338
- Lowe, W., & Benoit, K. (2013). Validating estimates of latent traits from textual data using human judgment as a benchmark. *Political Analysis*, 21, 298–313. doi:10.1093/pan/mpt002
- Muddiman, A., McGregor, S. C., & Stroud, N. J. (2018). (Re) claiming our expertise: Parsing large text corpora with manually validated and organic dictionaries. *Political Communication*, Online first, doi:10.1080/10584609.2018.1517843.
- Rooduijn, M., & Pauwels, T. (2011). Measuring populism: Comparing two methods of content analysis. *West European Politics*, 34, 1272–1283. doi:10.1080/01402382.2011.616665
- Scharkow, M. (2013). Thematic content analysis using supervised machine learning: An empirical evaluation using German online news. *Quality & Quantity*, 47, 761–773. doi:10.1007/s11135-011-9545-7
- Scharkow, M., & Bachl, M. (2017). How measurement error in content analysis and self-reported media use leads to minimal media effect findings in linkage analyses: A simulation study. *Political Communication*, 34, 323–343. doi:10.1080/10584609.2016.1235640
- Slapin, J. B., & Proksch, S. O. (2010). Look who's talking: Parliamentary debate in the European Union. *European Union Politics*, 11(3), 333–357. doi: 10.1177/1465116510369266
- Snow, R., O'Connor, B., Jurafsky, D., & Ng, A. Y. (2008, October). Cheap and fast – but is it good? Evaluating non-expert annotations for natural language tasks. In *Proceedings of the conference on empirical methods in natural language processing* (pp. 254–263). Honolulu, Hawaii: Association for Computational Linguistics.
- Spirling, A. (2016). Democratization and linguistic complexity: The effect of franchise extension on parliamentary discourse, 1832–1915. *The Journal of Politics*, 78(1), 120–136. doi:

10.1086/683612

van Atteveldt, W., & Peng, T. Q. (2018). When communication meets computation:

Opportunities, challenges, and pitfalls in computational communication science.

*Communication Methods & Measures*, 12, 81-92. doi: 10.1080/19312458.2018.1458084

Weber, R., Mangus, J. M., Huskey, R., Hopp, F. R., Amir, O., Swanson, R., ... & Tamborini, R.

(2018). Extracting latent moral information from text narratives: Relevance, challenges,

and solutions. *Communication Methods & Measures*, 12(2-3), 119-139. doi:

10.1080/19312458.2018.1447656

Wilkerson, J., Smith, D., & Stramp, N. (2015). Tracing the flow of policy ideas in legislatures: A

text reuse approach. *American Journal of Political Science*, 59(4), 943-956. doi:

10.1111/ajps.12175

Young, L., & Soroka, S. (2012). Affective news: The automated coding of sentiment in political

texts. *Political Communication*, 29, 205–231. doi:10.1080/10584609.2012.671234

Table 1. Results of the literature review (percentages in parentheses)

	<b>Total</b>		
Total record retrieved	192		
Excluded	119		
<b>Total eligible record</b>	<b>Total</b>	<b>Dictionary</b>	<b>SML</b>
	73 (100)	55 (100)	18 (100)
<i>Refer to any validation?</i>	42 (57.5)	27 (49.1)	15 (83.3)
↳ <b><i>refer to human-coded gold standard?</i></b>	<b>37 (50.6)</b>	<b>23 (41.8)</b>	<b>14 (77.7)</b>
↳ (1) <i>report any intercoder reliability?</i>	14 (19.2)	6 (10.1)	8 (44.4)
↳ <i>report K alpha?</i>	6 (8.2)	1 (1.8)	5 (27.7)
↳ (2) <i>report N of coders?</i>	18 (38.3)	10 (18.2)	8 (44.4)
↳ (3) <i>report N of validation data?</i>	34 (46.6)	21 (38.2)	13 (72.2)
↳ (4) <i>report validation metrics?</i>	32 (43.8)	20 (36.4)	12 (66.6)
↳ <i>report proper validation metrics?*</i>	15 (20.54)	6 (10.1)	9 (44.4)

*Note:* Percentages denote the share of articles that satisfy the given criteria among all articles employing respective methods. \* = validation metrics based on one of the following: recall/sensitivity, precision/PPV, F1, specificity, NPV, or accuracy.

Table 2. Simulation input parameters

<b>Factors</b>	<b>Input parameters</b>
<b>N of human coders</b>	2 (minimum), 5 (intermediate), & 10 (large manual coding)
<b>Intercoder reliability</b>	0.5 (low), 0.7 (acceptable), & 0.9 (high levels of reliability)
<b>N of validation data</b>	600 (0.5%), 1300 (1%), 6500 (5%), & 13000 (10%) of total data
<b>Sampling variability</b>	Random sample vs. non-random (biased) subset for validation
<b>Coding per entry</b>	Sole coding vs. duplicated coding for each entry

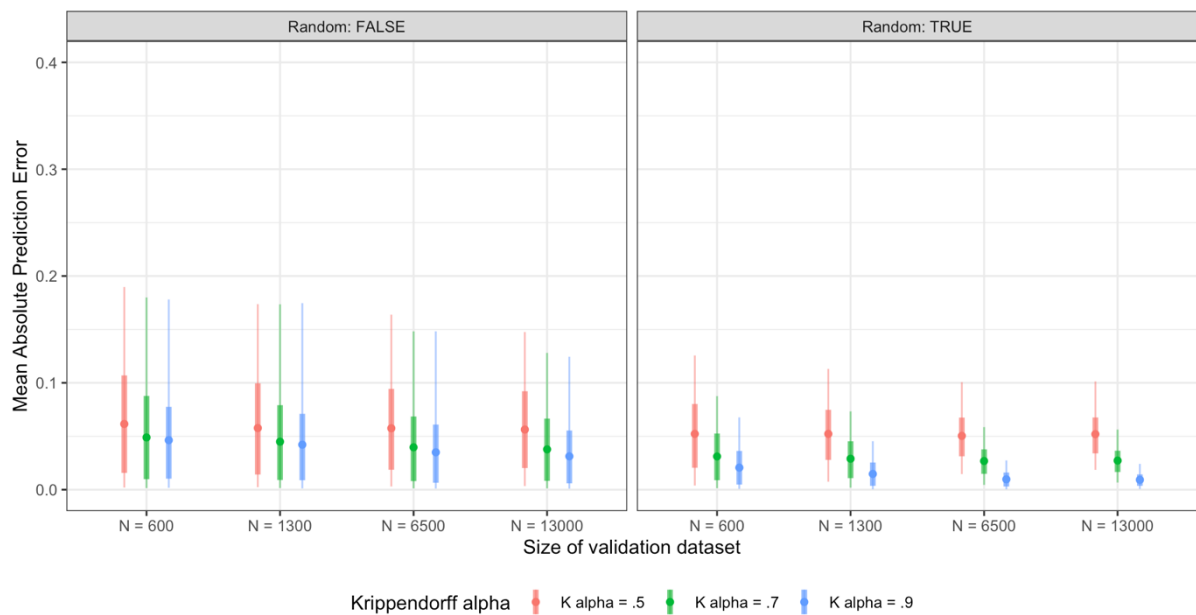


Table 3. Simple ANOVAs predicting MAPEs and mean comparisons across factors

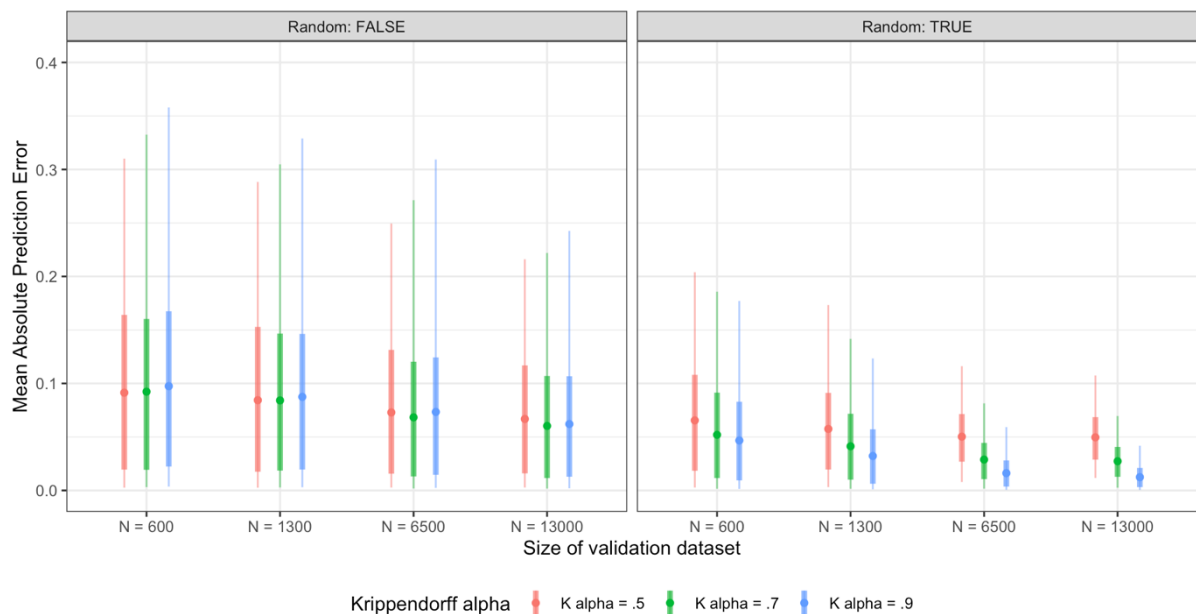
Factors	SML	Dictionary
<b>Random sample vs. not</b>	( $df=1, F=399.73$ )***	( $df=1, F=1025.31$ )***
Non-random (biased) sample	.0466 <sup>a</sup>	.0410 <sup>e</sup>
Random subset for validation	.0313 <sup>b</sup>	.0191 <sup>f</sup>
<b>Duplicated vs. sole coding</b>	( $df=1, F=.00$ )	( $df=1, F=.32$ )
Duplicated coding	.0390 <sup>a</sup>	.0299 <sup>e</sup>
Sole-coding	.0389 <sup>a</sup>	.0303 <sup>e</sup>
<b>No. of coders (k)</b>	( $df=2, F=.01$ )	( $df=2, F=.02$ )
k = 2	.0389 <sup>a</sup>	.0301 <sup>e</sup>
k = 5	.0389 <sup>a</sup>	.0301 <sup>e</sup>
k = 10	.0390 <sup>a</sup>	.0300 <sup>e</sup>
<b>Target Krippendorff's alpha value</b>	( $df=2, F=491.75$ )***	( $df=2, F=34.07$ )***
K alpha = 0.5	.0550 <sup>a</sup>	.0340 <sup>e</sup>
K alpha = 0.7	.0357 <sup>b</sup>	.0289 <sup>f</sup>
K alpha = 0.9	.0262 <sup>c</sup>	.0273 <sup>f</sup>
<b>Size of validation data (N)</b>	( $df=3, F=22.00$ )***	( $df=3, F=62.37$ )***
N = 600	.0435 <sup>a</sup>	.0362 <sup>e</sup>
N = 1,300	.0401 <sup>b</sup>	.0323 <sup>f</sup>
N = 6,500	.0365 <sup>c</sup>	.0270 <sup>g</sup>
N = 13,000	.0357 <sup>c</sup>	.0244 <sup>g</sup>
Residuals	$df = 134$	$df = 134$

**Note:** \*\*\*  $p < .001$ . Cell entries are marginal estimates of mean absolute prediction error (MAPE) per contrast of factors (total  $N=144$ ). Within each set of factors by model, different (same) superscripts denote statistically (in)distinguishable MAPEs based on Tukey's post-hoc tests. For instance, for *duplicated versus sole coding* factors, two MAPEs are statistically the same (i.e., their mean difference is not significant), and is therefore denoted with the same superscript. In contrast, for *random versus non-random sample* factors, two MAPEs are statistically different (i.e., their mean difference is significant), and is therefore denoted with different subscripts.

Panel A. Interactive effects of experimental factors, SML scenarios



Panel B. Dictionary approach

Figure 1. Interactive effects of  $N$ ,  $K\ alpha$ , and  $random$  factors predicting mean absolute prediction error (MAPE), SML (Panel A) and dictionary (Panel B) approaches.

Note: All combinations of two- and three-way interactions among  $N$  (= total no. of validation materials),  $K\ alpha$  (= target Krippendorff alpha value in human coding) and  $random$  factors (= sampling variability, i.e., random vs. nonrandom samples) were significant for both SML and dictionary scenarios. All three-way ANOVA results are reported in the online appendix.

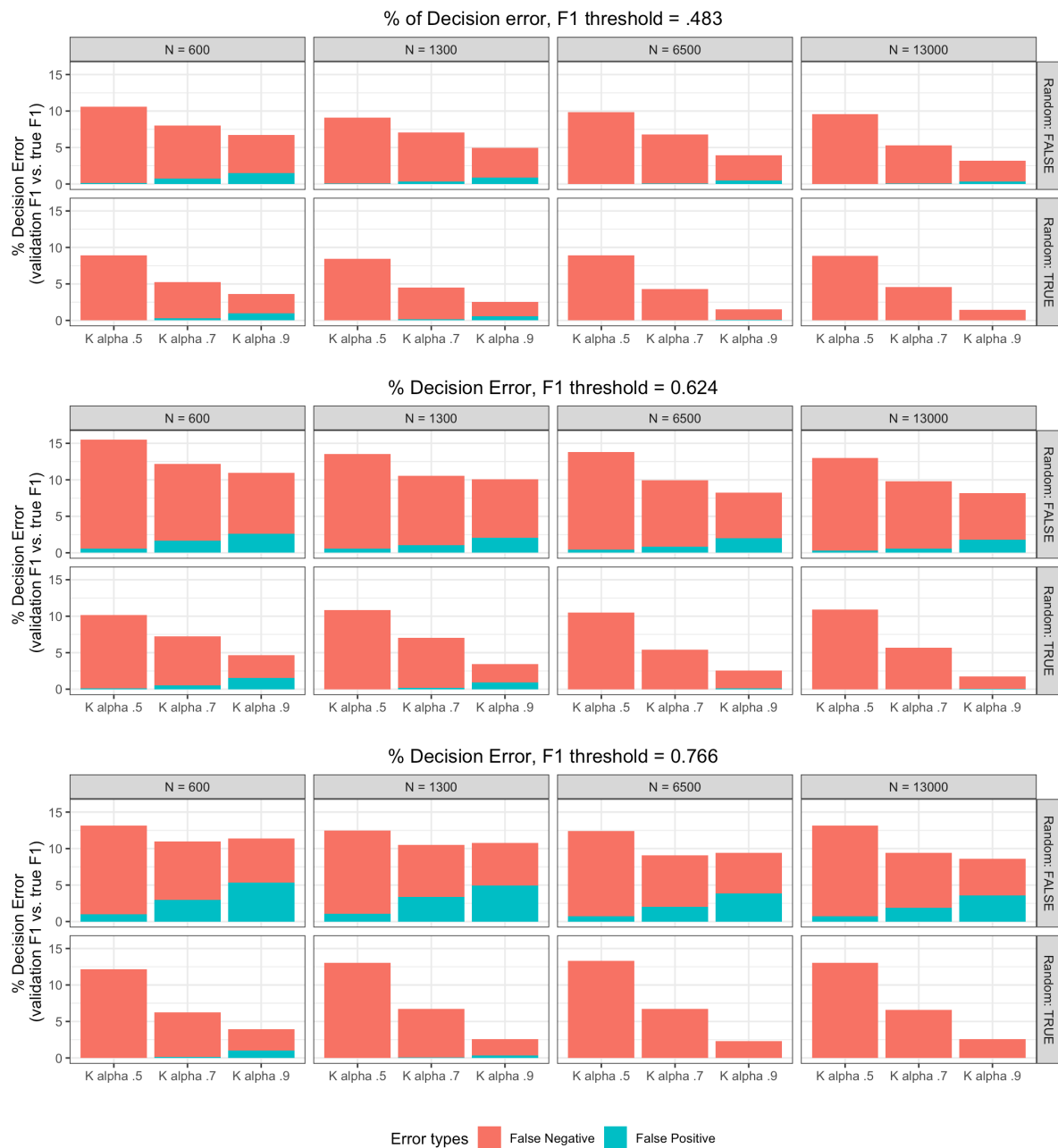


Figure 2. Decision error rate as a function of F1 thresholds, observed F1 score, and true F1 score (models for SML scenarios,  $N = 6,000$  in each). Bar graphs indicates the proportion of replicated simulations that a researcher's decision regarding F1 scores based on given validation materials (vis-à-vis true F1 scores) would fall into false positive or false negative cases.

*Note:*  $N$  = total no. of validation materials.  $K\ alpha$  = target Krippendorff alpha value in human coding. *Random* = sampling variability (i.e., random vs. nonrandom samples).

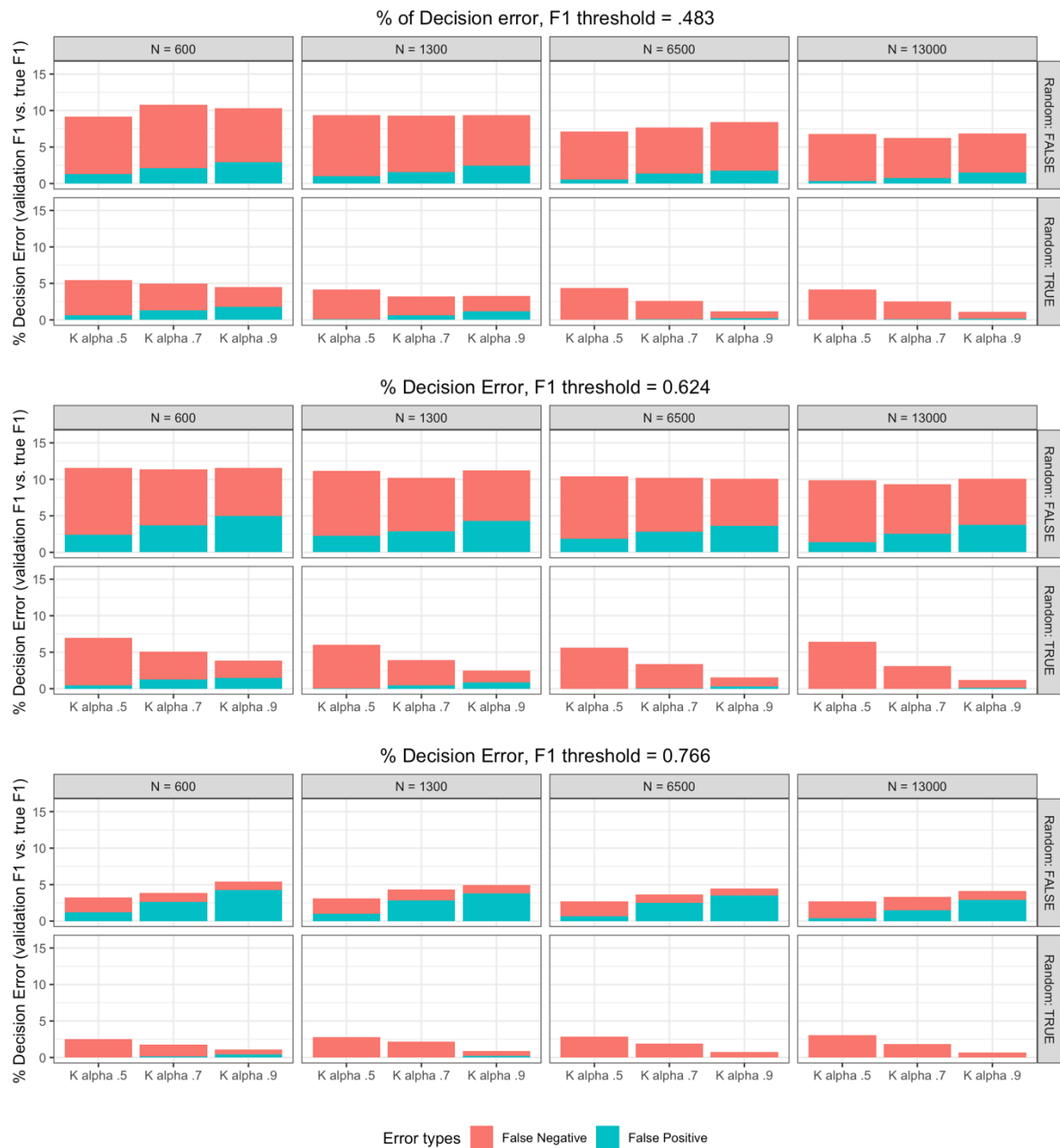


Figure 3. Decision error rate as a function of F1 thresholds, observed F1 score, and true F1 score (models for dictionary scenarios,  $N = 6,000$  in each). Bar graphs indicates the proportion of replicated simulations that a researcher's decision regarding F1 scores based on given validation materials (vis-à-vis true F1 scores) would fall into false positive or false negative cases.

*Note:*  $N$  = total no. of validation materials.  $K\ alpha$  = target Krippendorff alpha value in human coding. *Random* = sampling variability (i.e., random vs. nonrandom samples).

Supplementary Materials for:

**In Validations We Trust? The Impact of Imperfect Human Annotations as a Gold Standard on the Quality of Validation of Automated Content Analysis**

Hyunjin Song<sup>†</sup>, Petro Tolochko, Jakob-Moritz Eberl, Olga Eisele, Esther Greussing, Tobias Heidenreich, Fabienne Lind, Sebastian Galyga, & Hajo G. Boomgaarden

<sup>†</sup> Corresponding author. Email: [hyunjin.song@univie.ac.at](mailto:hyunjin.song@univie.ac.at)

**1. Data Availability Statement**

Our data and simulation codes are publicly available at DOI

[10.5281/zenodo.3598354](https://doi.org/10.5281/zenodo.3598354).

**2. Variables coded in Study 1, detailed coding instructions, and reliability estimates**

Using EBSCOhost databases, we searched all English-language journal articles published between January 1, 1998 and November 7, 2018, querying all titles, abstracts, and keywords using the following Boolean search string: ("computer assisted" OR "automated" OR "automatic" OR "computational" OR "machine learning") AND ("content analysis" OR "text analysis") This was done by examining “Communication & Mass Media Complete,” “Humanities Source,” and “SocINDEX with Full Index” collections.

Among a total of 192 retrieved articles, 112 articles were determined as not relevant (e.g., non-empirical overviews/introduction articles, qualitative analyses, studies using unsupervised methods, or simple keyword frequencies, etc.) and 7 articles were either duplicates or could not be obtained as full texts. These articles were excluded from further

## Online Appendix

analyses. Here, we exclude a simple keyword-frequency based study (e.g., simply counting the number of occurrences of a keyword in a given text, but not actually classifying the documents based on such frequency) since human inputs play no role other than compiling the keyword list itself. Among excluded studies, only 15 studies have used unsupervised learning or other forms of automated content analysis, suggesting dictionary-based or supervised machine learning applications are much more frequently used in general.

A total of five highly trained coders tested the initial coding scheme by independently coding 10 randomly sampled articles (approximately 5% of the total retrieved sample,  $N = 119$ ) and collectively discussed any coding problems and disagreements. Traditional content analysis literature generally recommends 5% to 25% of all materials to be used for reliability assessment (Lacy & Riffe, 1996). Coding instructions were iteratively revised until the coding schemes would produce reliable results. Inter-coder reliability (based on Krippendorff's alpha) above 0.75 was ensured for each of the variables coded. Following variables were independently coded by 5 trained coders.

Variable	Definition & Coding instructions	Reliability
Relevance	Whether empirical text analysis is conducted and reported (Yes = 1, No = 0)	Alpha = 1
Method Used	1 = Search string based / Dictionary Approach 2 = Machine Learning 3 = Topic Modeling (excluded from further analysis) 4 = Other (excluded from further analysis)	Alpha = 1
Refer to gold standard	1 = Yes, a "gold standard" is used, and info is reported 0 = No is not used reported	Alpha = 1
Report reliability	Whether inter-coder-reliability of human-coded materials are reported? (1 = Yes, reported, 0 = Not reported)	Alpha = 1
Refer to validation / Report validation measures	Whether validation of automated procedures are mentioned, and if so, whether either one of validation metrics (e.g., Recall, Sensitivity, Precision, Accuracy, F1, or other measures) is reported? (1 = Yes, mentioned, 0 = Not mentioned)	Alpha = .753

### 3. Detailed Setup of MC simulations

#### Data Generation

We create data (e.g., textual data, such as newspaper articles, to be analyzed) with the “true” outcome value of interest,  $y$  (i.e., a classification membership of a given document); the goal of any quantitative text analysis method is to somehow directly approximate this value of  $y$  for each observation-level, or instead estimate the unbiased distribution of  $y$  at the aggregate level (Grimmer & Stewart, 2013). For the data generating process, we set  $y$  at each document level to be randomly generated from three hypothetical independent variables ( $x_1$ ,  $x_2$ , and  $x_3$ ), all of which stand for some textual features (e.g., words or phrases) of a given document, plus a certain unobserved feature ( $x_0$ ) that is not evenly distributed across the dataset. The values of those variables were randomly sampled from a multivariate normal distribution. In addition, values of  $x_0$  were set to be identical across certain grouping variables of media content data, effectively simulating features that are not randomly nor uniformly distributed in the data. This ensures that the results of our simulations are not completely deterministic nor analytically driven to arrive at our conclusion.

**SML Scenario.** For supervised machine learning approach, we set the true values of  $y$  (which is the binary variable) are sampled from a Binomial distribution, with the probability parameter having a very simple linear functional form as follows:

$$y \sim \text{Bernoulli}(\pi)$$

$$\pi = \text{logistic}(\mu)$$

$$\mu = \mathbf{X}\beta + \varepsilon$$

with  $\varepsilon$  being Gaussian noise added to ensure that each simulation run is not completely deterministic. The  $\beta$ , the true population parameter, was fixed throughout the simulation runs (specifically,  $\beta_0 = 1$ ,  $\beta_1 = 0.5$ ,  $\beta_2 = 0.2$ , and  $\beta_3 = 0.6$ , which were randomly chosen).

## Online Appendix

Following this setup, a single simulation run is set to generate a total of 130,000 observations of media content data.

**Dictionary-based Scenario.** For a dictionary (i.e., bag-of-words) method, we assume a very similar approach as discussed above, but additionally truncate the values of independent variables to its nearest integer values (i.e., a discrete value), where they represent some “features” of given textual data (e.g., a word) or a combination of such textual features (e.g., a word order or N-grams), in a similar fashion as in Equations (1). Yet for the dictionary-based approach, the vector  $\beta$  was extended to  $K = 5$  and their  $\beta$  values were fixed to 0.2. This enables us to better approximate the multidimensionality of textual data, while treating  $y$  effectively as a function of the simple sum of the chosen textual features (which is a general assumption that most of the dictionary-based classification methods assumes).

This slight modification for dictionary approaches – truncating to the nearest integers – is due to the fact that each “feature” in the text (e.g., words, phrases, or boolean expressions, etc.) should be “predefined” to be matched against identical forms of dictionaries. We therefore effectively treat simulated integer numbers for three independent variables as each of the predefined categories for textual features, whose scores are simply taken from the existing dictionaries based on some rules. In contrast, for SML scenarios, we use raw continuous normal distributions as is (without rounding up/down numbers) effectively treating them as some kind of a transformed vector dimensional space wherein algorithms try to separate the observations into two categories (i.e., classification membership to be estimated) on that space.

### Human Coding

In all scenarios, human coders classify a given observation as “1” (e.g., a text contains the quantity of interest, such as a certain actor, frame, or tonality) or “0” (e.g., does not contain this quantity), based on some observable features of each documents. This human



coding ( $y$ ) can be, in principle, either correct or incorrect against the (unknown) true value,  $y$ , therefore behaviors of human coders were modeled by a Binomial distribution with varying probability of successfully categorizing the true data. This enables us to simulate a situation where, at a given target reliability level, some coders produce “correct” judgments while other coders produce “false” judgments more often.

### **Algorithm-based Classification and Validation**

For the dictionary approach, we assume that a researcher utilizes an off-the-shelf dictionary, based on mean valence of observed textual features (e.g., words, phrases, etc.), whose valence scores are taken from the existing dictionary. For the SML approach, we also assume that appropriate, domain-specific annotated materials for a given task already exist for the algorithm development, with a fixed number of training materials ( $N = 5000$ , approximately 4% of the total dataset being coded).<sup>1</sup>

---

<sup>1</sup> This means that researchers would only require to produce human coding for validation materials. In practice, when domain-appropriate training materials are not available, one need to produce human coding for training/testing materials as well. Doing so means the “quality” of human coding in such training/testing materials would be the same as validation materials, since one rarely employ different standards for training/testing vs. validation materials in such cases

#### 4. Additional results referred in the main results.

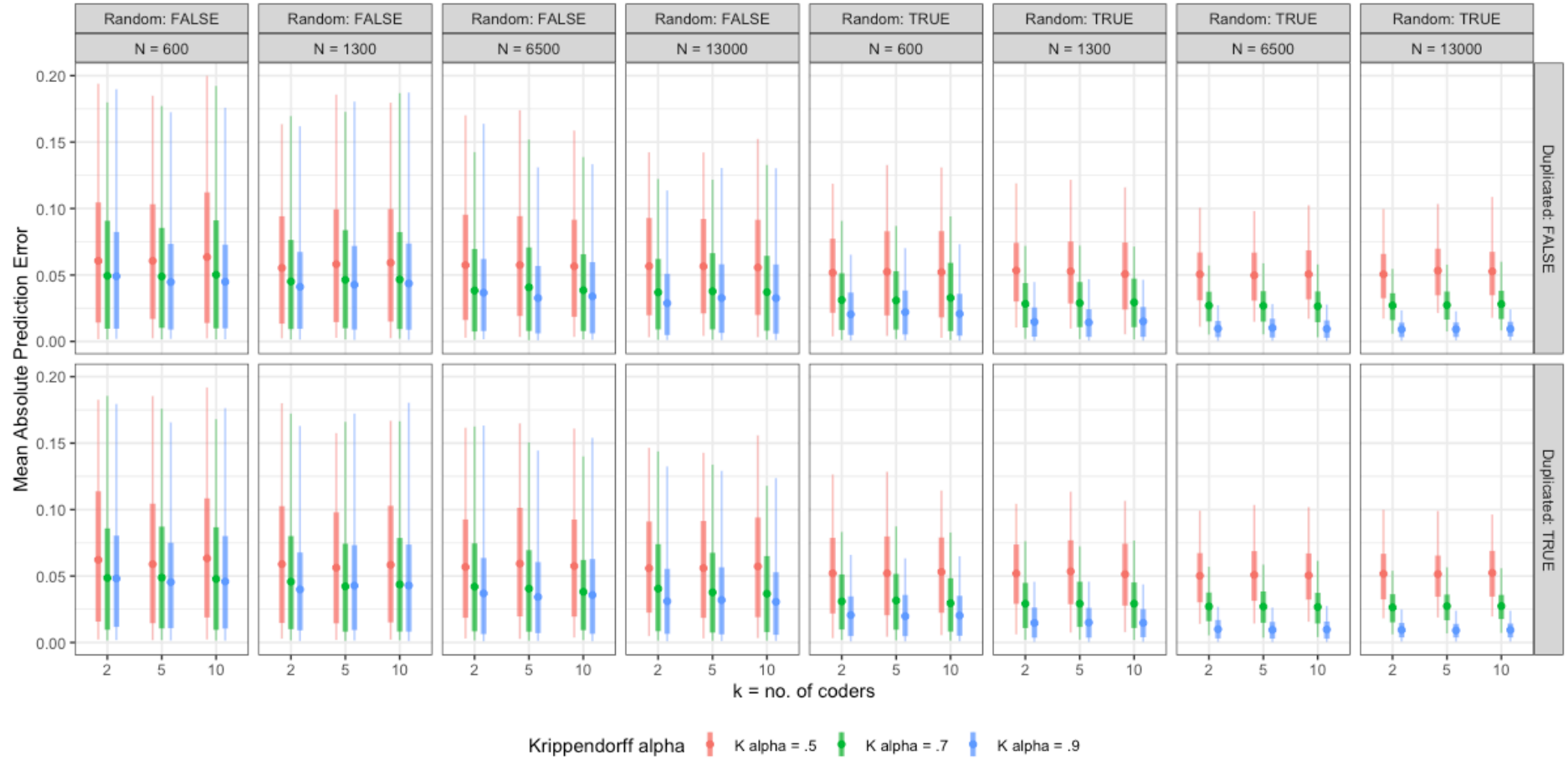


Figure A1. Mean Absolute Prediction Error (point estimate) and their 68% ( $\pm 1SD$ ) and 95% ( $\pm 2SD$ ) percentile intervals for every combination of experimental factors, **SML** scenarios ( $N = 1,000$  per scenario).

## Online Appendix

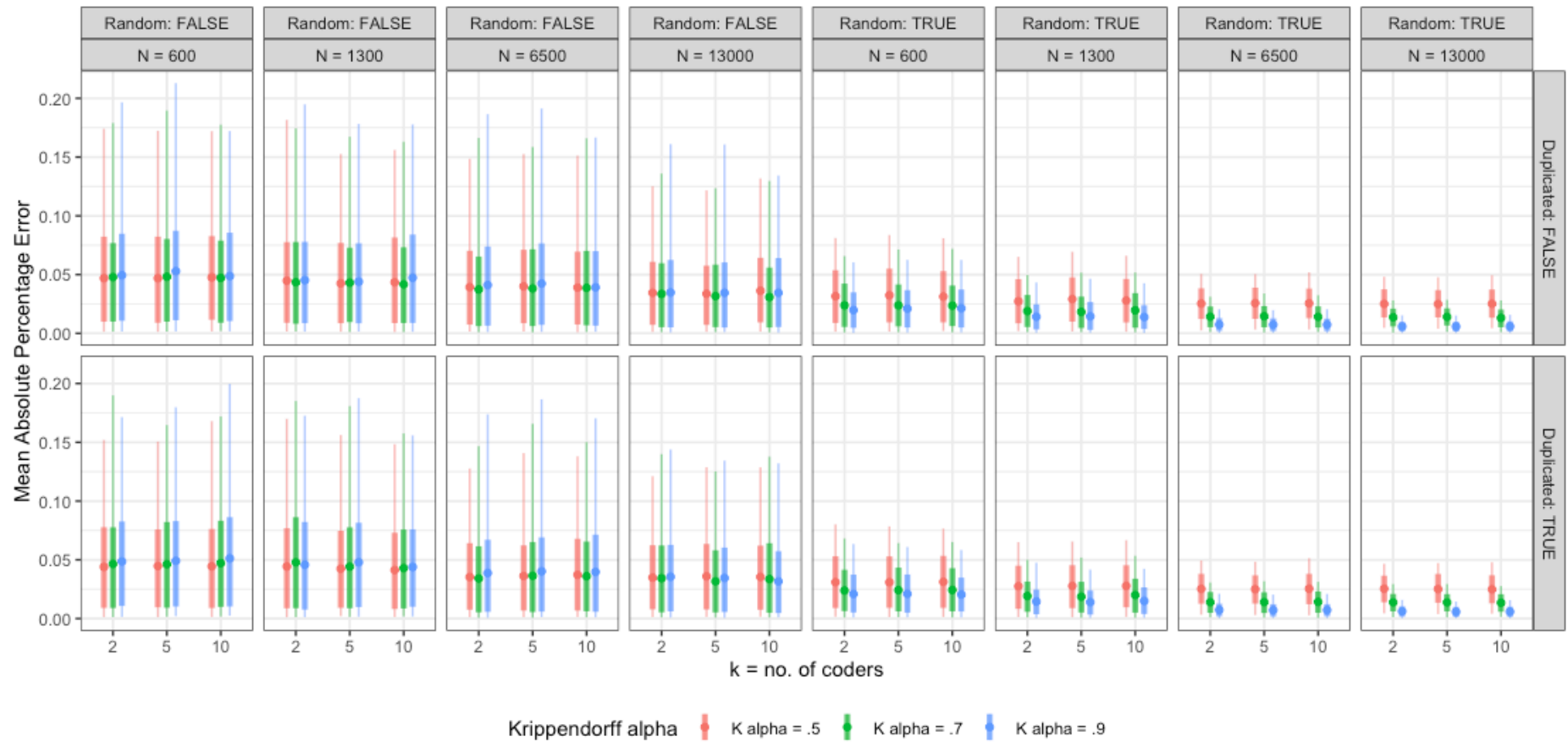


Figure A2. Mean Absolute Prediction Error (point estimate) and their 68% ( $\pm 1SD$ ) and 95% ( $\pm 2SD$ ) percentile intervals for every combination of experimental factors, **dictionary** scenarios ( $N = 1,000$  per scenario).

**ANOVAs estimating interactions of *number of coders, duplicated codings, and intercoder reliability* with other factors, SML scenarios.**

A 3-way interaction among *intercoder reliability, size of dataset, and random sample*

<b>Factors</b>	<b>Df</b>	<b>SS</b>	<b>MS</b>	<b>F</b>	<b>Pr(&gt;F)</b>
No. of coders	2	.00	.00	.096	.908
Duplicated vs. Sole-coding	1	.00	.00	.012	.913
Size of validation data (N)	3	.0013	.0004	359.497	.001 ***
Target Krippendorff's alpha (K)	2	.0207	.0104	8036.286	.001 ***
Random sample vs. not (R)	1	.0084	.0084	6532.416	.001 ***
K * N	6	.0004	.00006	52.884	.001 ***
N * R	3	.0001	.00004	35.610	.001 ***
K * R	2	.0021	.0010	807.706	.001 ***
K * N * R	6	.00004	.000007	5.551	.001 ***
Residuals	117	.00015	.000001		

A 2-way interaction with the *number of coders*

<b>Factors</b>	<b>Df</b>	<b>SS</b>	<b>MS</b>	<b>F</b>	<b>Pr(&gt;F)</b>
No. of coders (k)	2	.00	.00	.005	.995
Duplicated vs. Sole-coding	1	.00	.00	.001	.979
Size of validation data	3	.0013	.0004	19.853	.001 ***
Target Krippendorff's alpha	2	.0207	.0103	443.801	.001 ***
Random sample vs. not	1	.0084	.0084	360.750	.001 ***
k * Duplicated vs. Sole-coding	2	.000005	.000002	.106	.900
k * Size of validation data	6	.000011	.000002	.076	.998
k * Krippendorff's alpha	4	.000004	.000001	.041	.977
k * Random sample vs. not	2	.000002	.000001	.053	.949
Residuals	120	.002803	.000023		

A 2-way interaction with *duplicated coding*

<b>Factors</b>	<b>Df</b>	<b>SS</b>	<b>MS</b>	<b>F</b>	<b>Pr(&gt;F)</b>
No. of coders	2	.00	.00	.006	.994
Duplicated vs. Sole-coding (D)	1	.00	.00	.001	.979
Size of validation data	3	.0013	.0004	20.769	.001 ***
Target Krippendorff's alpha	2	.0207	.0103	464.284	.001 ***
Random sample vs. not	1	.0084	.0084	377.400	.001 ***
D * No. of coders	2	.000005	.000002	.111	.895
D * Size of validation data	3	.000005	.000002	.079	.971
D * Krippendorff's alpha	2	.000001	.00	.022	.978
D * Random sample vs. not	1	.00	.00	.011	.918
Residuals	126	.002813	.000022		