

A peer-reviewed version of this preprint was published in PeerJ on 21 April 2020.

[View the peer-reviewed version](https://peerj.com/articles/8871) (peerj.com/articles/8871), which is the preferred citable publication unless you specifically need to cite this preprint.

Prieto M, Deus H, de Waard A, Schultes E, García-Jiménez B, Wilkinson MD. 2020. Data-driven classification of the certainty of scholarly assertions. PeerJ 8:e8871 <https://doi.org/10.7717/peerj.8871>

Data-driven classification of the certainty of scholarly assertions

Mario Prieto¹, Helena Deus², Anita De Waard³, Erik Schultes⁴, Beatriz García-Jiménez¹, Mark D Wilkinson^{Corresp. 1}

¹ Center for Plant Biotechnology and Genomics UPM-INIA, Pozuelo de Alarcon, Madrid, Spain

² Elsevier Inc., Cambridge, MA, United States

³ Elsevier Research Collaborations Unit, Jericho, VT, United States

⁴ GO FAIR International Support and Coordination Office, Leiden, The Netherlands

Corresponding Author: Mark D Wilkinson

Email address: markw@illuminae.com

The grammatical structures scholars use to express their assertions are intended to convey various degrees of certainty or speculation. Prior studies have suggested a variety of categorization systems for scholarly certainty. However, these have not been objectively tested for their validity, particularly with respect to representing the interpretation by the reader, rather than the intention of the author. In this study, we use a series of questionnaires to determine how researchers classify various scholarly assertions, using three distinct certainty classification systems. We find that there are three categories of certainty perceived by readers: one level of high certainty, and two levels of lower certainty that are somewhat less distinct, but nevertheless show a significant degree of inter-annotator agreement. We show that these categories can be detected in an automated manner, using a machine learning model, with a cross-validation accuracy of 89.2% relative to an author-annotated corpus, and 82.2% accuracy against a publicly-annotated corpus. This finding provides an opportunity for contextual metadata related to certainty to be captured as a part of text-mining pipelines, which currently miss these subtle linguistic cues. We provide an exemplar machine-accessible representation - a Nanopublication - where certainty category is embedded as metadata in a formal, ontology-based manner within text-mined scholarly assertions.

1 Data-driven classification of the certainty of scholarly 2 assertions

3
4

5 Mario Prieto¹, Helena Deus², Anita De Waard³, Erik Schultes⁴, Beatriz García-Jiménez¹, Mark D
6 Wilkinson¹

7

8 1. Center for Plant Biotechnology and Genomics UPM-INIA, Madrid, Spain

9 2. Elsevier Inc, Cambridge, MA, USA

10 3. Elsevier Research Collaborations Unit, Jericho, Vermont, USA

11 4. GO FAIR International Support and Coordination Office, Leiden, The Netherlands

12

13

14 Corresponding Author:

15

16 Mark D Wilkinson

17 Centre for Plant Biotechnology and Genomics UPM – INIA

18 Parque Científico y Tecnológico de la U.P.M. Campus de Montegancedo

19 Autopista M-40, Km 38, Pozuelo de Alarcón, Madrid, 28223, Spain

20 Email address: mark.wilkinson@upm.es

21

22

23 Abstract

24 The grammatical structures scholars use to express their assertions are intended to convey
25 various degrees of certainty or speculation. Prior studies have suggested a variety of
26 categorization systems for scholarly certainty. However, these have not been objectively tested
27 for their validity, particularly with respect to representing the interpretation by the reader, rather
28 than the intention of the author. In this study, we use a series of questionnaires to determine how
29 researchers classify various scholarly assertions, using three distinct certainty classification
30 systems. We find that there are three categories of certainty perceived by readers: one level of
31 high certainty, and two levels of lower certainty that are somewhat less distinct, but nevertheless
32 show a significant degree of inter-annotator agreement. We show that these categories can be
33 detected in an automated manner, using a machine learning model, with a cross-validation
34 accuracy of 89.2% relative to an author-annotated corpus, and 82.2% accuracy against a
35 publicly-annotated corpus. This finding provides an opportunity for contextual metadata related
36 to certainty to be captured as a part of text-mining pipelines, which currently miss these subtle
37 linguistic cues. We provide an exemplar machine-accessible representation - a Nanopublication -
38 where certainty category is embedded as metadata in a formal, ontology-based manner within
39 text-mined scholarly assertions.

40 Introduction

41 Narrative scholarly articles continue to be the norm for communication of scientific results.
42 While there is an increasing push from both journals and funding agencies to publish source data
43 in public repositories, the resulting article, containing the interpretation of that data and the
44 reasoning behind those conclusions, continues to be, by and large, textual. The norms of
45 scholarly writing and scholarly argumentation are learned by students as they progress through
46 their careers, with the rules of scholarly expression being enforced by journal editors and
47 reviewers. Among the unique features of scholarly writing is the tendency for authors to use
48 hedging - that is, to avoid stating an assertion with certainty, but rather to use phrases that
49 suggest that the assertion is an interpretation of experimental evidence or speculation about a
50 state of affairs, which is essential when presenting unproven propositions with appropriate
51 caution. (Hyland, 1996) For example, “*These results **suggest that** the APC is constitutively*
52 *associated with the cyclin D1/CDK4 complex and are consistent with a model in which the APC*
53 *is responsible for cyclin D1 proteolysis in response to IR...*” (Agami & Bernards, 2000); or “*With*
54 *the understanding that coexpression of genes **may imply** coregulation and participation in*
55 *similar biological processes...*” (Campbell et al., 2007). As a result, biology papers contain a
56 wide range of argumentational structures that express varying degrees of confidence or certainty.
57 These subtle linguistic structures become problematic, however, in the context of scholarly
58 citation. As discussed by De Waard & Maat (De Waard & Maat, 2012), citing papers may
59 contain reformulations of the original claims in which the degree of certainty of the original
60 claim is modified (and usually made stronger) in the absence of additional evidence (Fig. 1;
61 Latour & Woolgar, 2013). This “drift” in certainty can be very gradual over successive steps of a
62 citation chain, but the consequences may be profound, since statements with greater certainty
63 than the original author intended can be used as the basis for new knowledge. Although peer-
64 review might protect the literature from such ‘hedging erosion’, reviewers may lack the specific
65 domain knowledge required to know the legacy of a given scholarly claim. Even if they take the
66 time to follow a citation, subtle differences in expressed certainty over a single step in a citation
67 chain may not be detectable. This problem is worsened in the context of text mining algorithms
68 that currently do not richly capture the nuances of a scholarly assertion when extracting the
69 entity-relationships that make up the claim.

70

71 Given that the volume of literature published grows by approximately a half-million papers per
72 year in the biomedical domain alone, text mining is becoming an increasingly important way to
73 capture this new knowledge in a searchable and machine-accessible way. Accurate, automated
74 knowledge capture will therefore require accurate capture of the certainty with which the claim
75 was expressed. Moreover, there is increasing pressure to publish knowledge, *ab initio*, explicitly
76 for machines, in particular with the adoption of the FAIR Data Principles for scholarly
77 publishing (Wilkinson et al., 2016), and several machine-accessible knowledge publication
78 formats have recently been suggested, including NanoPublications (Groth, Gibson & Velterop,
79 2010), and Micropublications (Clark, Ciccarese & Goble, 2014). In order to capture the intent of

80 the author in these machine-readable publications, it will be necessary for them to include formal
81 machine-readable annotations of certainty.

82

83 A number of prior studies have attempted to categorize and capture the expression of scholarly
84 certainty. These, and other certainty categorization studies, are summarized, compared and
85 contrasted in Table 1, where the columns represent relevant study features that distinguish these
86 various investigations, and affect the interpretation of their outcomes. For example, the use of
87 linguistic experts, versus biomedical domain experts, will likely affect the quality of the
88 annotations, while using explicit rule-matching/guidelines will result in strict, predetermined
89 categorizations. Similarly, the use of abstracts consisting of concise reporting language, versus
90 full text which contains more exploratory narratives, will affect the kinds of statements in the
91 corpus (Lorés, 2004), and their degree of certainty.

92

93 According to Wilbur et al., “each [statement] fragment conveys a degree of certainty about the
94 validity of the assertion it makes” (Wilbur, Rzhetsky & Shatkay, 2006). While intuitively
95 correct, it is not clear if certainty can be measured/quantified, if these quantities can be
96 categorized or if they are more continuous, and moreover, if the perception of the degree of
97 certainty is shared between readers. Most studies in this domain assume that certainty can be
98 measured and categorized, though they differ in the number of degrees or categories that are
99 believed to exist, and thus there is no generally-accepted standard for certainty/confidence levels
100 in biomedical text (Rubinstein et al., 2013). Wilbur et al suggested a four category classification:
101 complete uncertainty, low certainty, high likelihood, and complete certainty/proven fact.
102 Similarly, Friedman et al. (Friedman et al., 1994) suggest that there are four categories of
103 certainty: no certainty, low, moderate, and high certainty, with an additional “cannot evaluate”
104 category. Aligning with both of these previous studies, De Waard and Schneider (De Waard &
105 Schneider, 2012) encoded four categories of certainty into their Ontology of Reasoning,
106 Certainty, and Attribution (ORCA) ontology as follows: Lack of knowledge, Hypothetical (low
107 certainty), Dubitative (higher, but short of full certainty), Doxastic (complete certainty, accepted
108 knowledge or fact). Other studies have suggested fewer or more certainty categories, and differ
109 in the manner in which these categories are applied to statements.

110

111 BioScope (Vincze et al., 2008) is a manually-curated corpus, containing 20,924 speculative and
112 negative statements from three sources (clinical free-texts, 5 articles from FlyBase and 4 articles
113 from BMC Bioinformatics) and 3 different types of text (Clinical reports, Full text articles and
114 abstracts). Two independent annotators and a chief linguistic annotator classified text spans as
115 being ‘speculative’ or ‘negative’; other kinds of assertions were disregarded. Thus, the study
116 splits certainty into two categories - speculative, or not.

117

118 Thompson et al. (Thompson et al., 2011) apply 5 meta-knowledge features - manner, source,
119 polarity, certainty, and knowledge type - to the GENIA event corpus (“GENIA Event Extraction

120 - BioNLP Shared Task”). This corpus is composed of Medline abstracts split into individual
121 sentences. With respect to certainty annotations, the corpus utilizes a classification system of
122 three certainty levels - certain, probable (some degree of speculation), and doubtful (currently
123 under investigation). Annotation was carried out by two linguistic specialists specifically trained
124 in the meta-knowledge scheme.

125
126 Light et al. (Light, Qiu & Srinivasan, 2004) investigate speculative language in biomedical
127 abstracts. Using Medline abstracts they attempt to distinguish high and low degrees of
128 speculation. Four annotators used rule-matching to classify statements. Using this annotated
129 corpus, they trained a model based on Support Vector Machines (SVM) to generate an automatic
130 classifier. This automatic classifier, therefore, is specifically tasked for speculative statements,
131 and categorizes them in a manner resembling their predefined rule-sets.

132
133 Malhotra et al. (Malhotra et al., 2013) classify hypotheses (speculative statements) in scholarly
134 text. Three annotators classified speculative statements in Medline abstracts related to
135 Alzheimer's disease using a four-class categorization, with predefined pattern-matching rules for
136 sorting statements into three speculative patterns (strong, moderate, weak) and a fourth category
137 representing definitive statements. Additionally, they explore several automated methods to
138 distinguish speculative from non-speculative statements.

139
140 Zerva et al. (Zerva et al., 2017) use a combination of the BioNLP-ST and GENIA-MK corpora -
141 both of which consist of statements manually-annotated with respect to their certain/uncertain
142 classification (degrees of uncertainty, when available, were merged resulting in a two-category
143 corpus). They applied rule induction combined with a Random Forest Classifier to create an
144 automated binary classification model. This model was run on 260 novel statements, and the
145 output classification was provided to seven annotators who were asked for simple agree/disagree
146 validation of each automated classification. The degree of disagreement between annotators was
147 in some cases surprisingly high, leading the authors to note that “the perception of (un)certainty
148 can vary among users”. In a separate experiment, two annotators ranked the certainty of 100
149 statements on a scale of 1-5. They noted low absolute annotator agreement (only 43% at the
150 statement-level), but high relative agreement (only 8% of statements were separated by more
151 than 1 point on the 5 point scale). Comparing again to the automated annotations, they found
152 high correlation at the extremes (i.e., scored by the annotators as 1 or 5) but much less
153 correlation for statements rated at an intermediate level, leading them to conclude “...looking into
154 finer-grained quantification of (un)certainty would be a worthwhile goal for future work”.

155
156 These previous works share important distinctions relevant to the current investigation. First, in
157 every case, the number of certainty categories were predetermined, and in many cases,
158 categorization rules were manually created. Second, in most cases, the work involved a small
159 number of annotators with a knowledge of linguistics, or specifically trained on the annotation

160 system, rather than being experts in the knowledge-domain represented by the statements, but
161 untrained as annotators. Third, in all cases where automated approaches were introduced, the
162 automated task was to distinguish “speculation” from “non-speculation”, rather than categorize
163 degrees of certainty. Notably, there was little agreement on the number of categories, nor the
164 labels for these categories, among these studies. Moreover, the categories themselves were
165 generally not validated against the interpretation of an (untrained) domain-expert reader. As
166 such, it is difficult to know which, if any, of these approaches could be generalized to annotation
167 of certainty within the broader scholarly literature, in a manner that reflects how domain experts
168 interpret these texts.

169

170 To achieve this would require several steps: 1) determine if there are clearly delimited categories
171 of certainty that are perceived by readers of scholarly assertions; 2) if so, determine how many
172 such categories exist; and 3) determine the fidelity of the transmission of certainty among
173 independent readers (i.e. agreement). If these are determined robustly, it should then be possible
174 to apply machine-learning to the problem of automatically assigning certainty annotations to
175 scholarly statements that would match the perceptions of human readers.

176

177 Here, we attempt a data-driven certainty categorization approach. We execute a series of
178 questionnaire-based studies using manually-curated scholarly assertions, in English, to attempt to
179 objectively define categories of perceived certainty. A different set of certainty categories are
180 provided in each questionnaire, and readers are asked to categorize each statement as to their
181 perception of its level of certainty. We use these results to examine the degree of consistency of
182 perceived certainty among readers, and run statistical tests to evaluate the degree to which the
183 categorization system provided in each survey reflects the perception of those asked to use those
184 categories. The categorization system with the highest score - that is, the one that provided the
185 highest level of agreement - was then used to manually create a corpus of certainty-annotated
186 statements. This, in turn, was used to generate a machine-learning model capable of
187 automatically classifying new statements into these categories with high accuracy. We propose
188 that this model could be used within existing text-mining algorithms to capture additional
189 metadata reflecting the nuanced expression of certainty in the original text. Finally, we provide
190 an example of a machine-accessible scholarly publication - a NanoPublication - within which we
191 have embedded this novel contextual certainty metadata.

192

193

194 **Materials & Methods**

195 **Broad overview:** Using TAC Biomedical Summarization Corpus (Min-Yen, 2018), we extracted
196 45 manually-curated scholarly assertions (selection process described below). Using these, a
197 total of 375 researchers in the biomedical domain, in comparable research institutes and
198 organizations, were presented with a series of assertions and asked to categorize the strength of
199 those assertions into four, three, or two certainty categories over the three independently-

200 executed questionnaires. G Index (Holley & Guilford, 1964) coefficient analysis was applied to
201 determine the degree of agreement between annotators, as a means to evaluate the power of each
202 categorization system - that is, to test the discriminatory effectiveness of the categories
203 themselves, versus the quality of the annotations or annotators. We extracted the essential
204 features of inter-rater agreement from the questionnaire data using Principal Component
205 Analysis (PCA) to guide our interpretation of the way annotators were responding to the
206 categories presented. The essential number of components identified by PCA were extracted
207 using Horn's parallel analysis, with three categories appearing to be the optimal. We then
208 clustered our collection of statements into these three categories using k-means algorithm
209 (Jolliffe, 2011; Dunham, 2006). Finally, we manually generated a corpus of statements annotated
210 using this 3-category system, and applied deep-learning techniques over this corpus to generate
211 an automated classifier model. To evaluate its accuracy, a 20-fold Cross-Validation (CV) was
212 used.

213

214 **Survey statement selection:** The 45 text blocks used in the three surveys were extracted from
215 published articles related to genetic and molecular topics, and were selected from the "Citation
216 text" and "Reference text" portions of the TAC 2014 Biomedical Summarization Track. Each
217 text block contained a sentence or sentence fragment representing a single scholarly assertion
218 that we highlighted and asked the respondents to evaluate, with the remainder of the text being
219 provided for additional context. The 45 assertions were selected using different epistemic
220 modifiers, such as modal verbs, qualifying adverbs and adjectives, references and reporting
221 verbs, which are believed to be grammatical indicators of "value of truth" statements (De Waard
222 & Maat, 2012). Given that they are intended to be used for a human survey, with the aim of
223 avoiding annotator fatigue, these were further filtered based on the length of the statement to
224 give preference to shorter ones. These were then separated into groups based on the type of
225 epistemic modifier used, and from these groups, a subset of statements were selected arbitrarily
226 to give coverage of all groups in our final statement corpus (Prieto, 2019). An example survey
227 interface presentation is shown in Fig. 2.

228

229 **Survey design:** We designed three surveys - S1, S2 and S3 - where respondents were asked to
230 assign certainty based on a number of certainty categories - 4, 2, and 3 respectively for surveys
231 S1, S2, and S3. All surveys used the same corpus of 45 scholarly assertions. To minimize the
232 bias of prior exposure to the corpus, the surveys were deployed over three comparable but
233 distinct groups of researchers, all of whom will have sufficient biomedical expertise to
234 understand the statements in the corpus.

235

236 All participants were presented a series of assertions selected randomly from the 45 in the corpus
237 - 15 assertions in S1, increased to 20 assertions in S2 and S3 in order to obtain deeper coverage
238 of the statement set. In S1, participants were asked to assess the certainty of every highlighted
239 sentence fragment based on a 4-point scale with the following response options: High, Medium

240 High, Medium Low, and Low. A 2-point scale was used for S2: Relatively High and Relatively
241 Low and 3-point numerical scale for S3: 1, 2 or 3. In addition to the assessment of certainty, for
242 each assertion, subjects were asked to indicate their impression of the basis of the assertion,
243 using a single-answer, multiple-choice question, with the options: Direct Evidence, Indirect
244 Evidence/Reasoning, Speculation, Citation or I don't know.

245

246 **Survey distribution and participant selection:** Participation in the surveys was primarily
247 achieved through personal contact with department leads/heads of 5 institutions with a focus on
248 biomedical/biotechnology research. For S1, the majority of participants came from the Centro de
249 Biotecnología y Genómica de Plantas (UPM-INIA), Spain. It was conducted between November
250 and December of 2016. S2 was executed by members from the Leiden University Medical
251 Center, Netherlands, between November and December of 2017. S3 was conducted between
252 October and November of 2018 by members of the University Medical Center Utrecht, Cell
253 Press and the Agronomical Faculty of Universidad Politécnica de Madrid. Participation was
254 anonymous and no demographic data was collected.

255

256 **Survey execution:** Participants of the surveys were engaged using the platform Survey Gizmo
257 (S1) or Qualtrics ("Qualtrics") - two online platforms dedicated to Web-based questionnaires.
258 The change in survey platform was based only on cost and availability; the two platforms have
259 largely comparable interfaces with respect to data-gathering fields such as response-selection
260 buttons and one-question-per-page presentation, with the primary differences between the
261 platforms being aesthetic (color, font, branding).

262

263 **Statistical analysis of agreement:** We evaluated each survey by quantifying the degree of
264 agreement between participants who were presented the same assertion, with respect to the level
265 of certainty they indicated was expressed by that statement given the categories provided in that
266 survey.

267

268 Agreement between participants was assessed by Holley and Guilford's G Index of agreement
269 (Holley & Guilford, 1964), which is a variant of Cohen's Weighted Kappa (K_w ; Cohen, 1968).
270 Ideally G measures the agreement between participants. It was performed based on the following
271 formula:

272

$$273 \quad G = \frac{(\text{Probability Observed}(P_o) - \text{Probability by Chance}(P_c))}{1 - P_c}$$

274

275

276 The key difference between K_w and G is in how chance agreement (P_c) is estimated. According
277 to (Xu & Lorber, 2014), "G appears to have the most balanced profile, leading us to endorse its
278 use as an index of overall interrater agreement in clinical research". G is defined *a priori*, being

279 homogeneously distributed among categories as the inverse of the number of response categories
280 (Xu & Lorber, 2014), thus making $G=0.25$ for S1; $G=0.50$ for S2; and $G= 0.33$ for S3. The
281 accepted threshold for measuring agreement and its interpretation has been suggested by Landis
282 & Koch, 1977 (Landis, Richard Landis & Koch, 1977) as follows: $0.2 = \text{Poor}$, $0.21 - 0.4 = \text{Fair}$,
283 $0.41 - 0.60 = \text{Moderate}$, $0.61 - 0.80 = \text{Substantial}$, $0.81 - 1.00 = \text{Almost Perfect}$. Anything other
284 than the 'Poor' category is considered in other studies to represent an acceptable level of
285 agreement. (Deery et al., 2000; Lix et al., 2008)

286
287 **Clustering:** We investigated the ideal number of clusters into which statements group based on
288 the profile of the annotators' responses. To estimate this, Hierarchical Clustering analysis (HCA)
289 and the Spearman correlation test were performed to determine certainty category association
290 between questionnaires (Fig. 3), using the shared classified statements in that category as the
291 metric (Narayanan et al., 2011; Campbell et al., 2010; Sauvageot et al., 2013; Narayanan et al.,
292 2014); though these represent conceptually distinct analyses, we represent them in the same chart
293 because the outputs are mutually supportive. Hierarchical clustering analysis (HCA) finds
294 clusters of similar elements, while Spearman correlation coefficient considers the weight and
295 direction of the relationship between two variables. It's worth emphasizing the importance of the
296 rank-based nature of Spearman's correlation. Spearman's formula rank the variables in order,
297 then measures and records the difference in rank for each statement/variable. Thus, "...if the data
298 are correlated, [the] sum of the square of the difference between ranks will be small" (Gauthier,
299 2001), which should be considered when interpreting the results. Interpretation of Spearman
300 correlation was as follows: Very Low ≤ 0.2 ; Low ≤ 0.5 ; Moderate ≤ 0.7 ; High ≤ 0.9 and Very
301 High > 0.9 (Dunham, 2006; Raithe, 2008). All Spearman interactions are based on hypothesis
302 testing. To determine the importance of the results, p-values were generated as an indicator of
303 the existence of correlation between certainty categories. HCA and Spearman values were
304 generated using the python libraries *seaborn* and *pandas*.

305
306 Normalization allows the comparison of two nominal variables on different scales. Prior to
307 Principal Component Analysis (PCA) and cluster analyses, we first normalized for the different
308 number of annotators for each statement, and centered, using the scale function from R and the
309 Python package *scikit-learn*. PCA is a widely used method for attribute extraction to help
310 interpret results. We used PCA to extract the essential features of inter-rater agreement from the
311 questionnaire data (Campbell et al., 2010; Narayanan et al., 2014). We applied PCA using *scikit-*
312 *learn* to the result-sets, and utilized K-means from the same python package to identify cluster
313 patterns within the PCA data. These cluster patterns reflect groups of similar "human behaviors"
314 in response to individual questions under all three survey conditions. The input to both statistical
315 functions was the results of the questionnaires in the form of a contingency table, where each
316 statement is represented by the profile of annotations it received from all annotators. The optimal
317 number of components was selected using Horn's parallel analysis, applied to certainty
318 categories on the 3 different surveys. Detailed output is provided in Fig. S1, S2 and S3 of the

319 supplemental information. Our decision to choose three components as the most robust number
320 to capture relevant features of our data is justified in the Results section.

321

322 To determine the optimal K (cohesion of the clusters), several indices were analyzed using the R
323 package *NbClust* (Charrad et al., 2014). *NbClust* provides 30 different indices (e.g., Gap statistic
324 or Silhouette) for determining the optimal number of clusters based on a “majority rule”
325 approach (Fig. 4; Chouikhi, Charrad & Ghazzali, 2015). Membership in these clusters was
326 evaluated via Jaccard similarity index comparing, pairwise, all three clusters from each of the
327 three surveys to determine which clusters were most alike (Table 2). This provides additional
328 information regarding the behavior of annotators between the three surveys; i.e., the
329 homogeneity of the three identified categories between the three distinct surveys. The *princomp*
330 and *paran* functions in R were utilized to execute PCA and Horn’s parallel analysis, respectively.
331 The *PCA* and *KMeans* functions from *scikit-learn* were employed to create the visualizations in
332 Fig. 5 (Pedregosa et al., 2011).

333

334 **Certainty Classification and Machine Learning Model:** We addressed the creation of a
335 machine-learning model by considering this task to be similar to a sentiment analysis problem,
336 where algorithms such as Recurrent Neural Network (RNN) with Long Short Term Memory
337 (LSTM) have been applied (Wang et al., 2016; Baziotis, Pelekis & Doulkeridis, 2017; Ma, Peng
338 & Cambria, 2018). A corpus of new statements were extracted from MedScan (Novichkova et
339 al., 2003). An initial filter was applied using the keyword 'that', since this is often indicative of
340 hedging (e.g., “*This result suggests that...*”, “*We report that...*”, “*It is thought that...*”). A total of
341 3221 statements were manually categorized using the 3 levels of certainty, based on our
342 familiarity with the classification of the 45 statements in the prior study. A 5-layer neural
343 network architecture was employed to train and validate model performance. Validation was
344 executed using 20-fold CV scheme, which is considered adequate for a corpus of this size
345 (Crestan & Pantel, 2010; Snow et al., 2008; Lewis, 2000). To design the neural network model,
346 the Python library *Keras* (Chollet and oggithers, 2015) was utilized, with *TensorFlow* (Abadi et
347 al., 2015) as the backend. Precision, recall and overall accuracy were calculated as additional
348 supporting evidence for classifier performance from a confusion matrix (Light, Qiu &
349 Srinivasan, 2004; Malhotra et al., 2013; Zerva et al., 2017), comprised of the following terms
350 and formulas: True Positive (TP); True Negative (TN); False Positive (FP); False Negative (FN);
351 Precision = $TP/(TP+FP)$; Recall = $TP/(TP+FN)$; Overall accuracy =
352 $(TP+TN)/(TP+FP+FN+TN)$. Finally, we employed Kappa as a commonly-used statistic to
353 compare automated and manual adjudication (Garg et al., 2019). Kappa was calculated using the
354 *pym* python package.

355

356 All raw data and libraries used are available in the project GitHub, together with Jupyter
357 Notebooks (both R and Python 2.7 kernels) showing the analytical code and workflows used to

358 generate the graphs presented in this manuscript and the supplemental information (Prieto,
359 2019).

360

361 **Results**

362

363 **Survey participation:** Survey 1 (S1) was answered by 101 participants of whom 75 completed
364 the survey (average of 13 responses per participant). Survey 2 (S2) had 215 participants with 150
365 completing the survey (average of 16 responses per participant). 48 of 57 participants completed
366 the entirety of Survey 3 (S3) (average of 18 responses per participant). All responses provided
367 were used in the analysis. Coverage (the number of times a statement was presented for
368 evaluation) for each of the 45 statements in the corpus was an average of: 29 for S1, 77 for S2,
369 and 23 for S3. The summary of the k-means clustering and Jaccard Similarity results over all
370 three surveys are shown in Fig. 3 and Table 2.

371

372 **Survey 1:** In S1, all statements except statement #13, scored at or above the minimum agreement
373 ($G = 0.21$; “Fair” degree of agreement on the (Landis, Richard Landis & Koch, 1977) scoring
374 system). Seven of 45 statements (15.5%) showed inter-annotator agreement achieving
375 statistically-significant scores in two certainty categories simultaneously. Table 3 shows the
376 distribution of statements among certainty categories and agreement levels. 11 of 45 statements
377 (24.4%) were classified as High certainty; 14 of 45 statements (31.1%) were Medium High;
378 Medium Low were represented by 12 of 45 statements (26.6%); and the Low certainty category
379 did not produce inter-annotator agreement for any statement

380

381 **Survey 2:** Disposition of Certainty categories and agreement levels for S2 are shown in Table 4.
382 Seven of the 45 statements (15.5%) did not achieve significant agreement for any certainty level.
383 Relatively High was selected for 19 of 45 statements (42.2%). The remaining statements (19/45;
384 42.2%) were selected as Relatively Low.

385

386 **Survey 3:** Table 5 summarizes the levels of agreement and certainty classifications observed in
387 S3. Categories were ranked numerically from 1 (the highest level of certainty) to 3 (the lowest
388 level of certainty). Minimum agreement ($G = 0.21$) or superior was observed in 41 of 45
389 statements (91.1%) with no doubly-classified statements, indicating little evidence of annotator-
390 perceived overlap between the presented categories. Four of 45 statements (8.8%) did not obtain
391 agreement for any certainty category. Category 1 was selected for 13 of 45 statements (28.8%).
392 24 of the total of 45 (53.3%) were chosen with level of Category 2. Finally, Category 3 was
393 selected for four out of 45 sentences (8.8%).

394

395 **Clustering:** As shown in Fig. 3, HCA and Spearman correlation rank revealed three primary
396 clusters when executed over the three surveys combined. The first branch of the HCA, left-top
397 side of Fig. 3 clustered S3-1, S1-High and S2-Relatively High. Numbers inside the squares of

398 this cluster show significant Spearman correlation (**S1-High/S2-Relatively High**: $r = 0.81$, p-
399 value = $2.3e-11$; **S1-High/S3-1**: $r = 0.72$, p-value = $3.2e-8$; **S2-Relatively High/S3-1**: $r = 0.79$,
400 p-value = $1.3e-10$). The other branch of the hierarchical tree is split again into 2 main sub-trees,
401 including the center and right sides of the figure. The central cluster in the figure, differentiated
402 by excellent Spearman correlation, contains S1-Medium Low, S2-Relatively Low and S3-3
403 categories (**S1-Medium Low/S3-3**: $r = 0.78$, p-value = $2.5e-10$; **S2-Relatively Low/S3-3**: $r =$
404 0.81 , p-value = $1.3e-11$; **S1-Medium Low/S2-Relatively Low**: $r = 0.83$, p-value = $1.5e-12$).
405 Finally, the smaller cluster identified by HCA on the right side of Fig. 3 comprises S1-Medium
406 High and S3-2 with moderate Spearman correlation ($r = 0.55$, p-value = $7.8e-5$), confirming that
407 a third certainty category has sufficiently strong support.

408

409 Supporting the previous cluster tests, using the majority rule approach, *NbClust* results (Fig. 4)
410 indicate that:

- 411 • 16 indices proposed 3 as the optimal number of clusters for the results of S1 (Fig. 4A)
- 412 • 11 indices proposed 2 as the optimal number of clusters for the results of S2 (Fig. 4B)
- 413 • 6 indices proposed 3 as the optimal number of clusters for the results of S2 (Fig. 4B)
- 414 • 14 indices proposed 3 as the optimal number of clusters for the results of S3 (Fig. 4C)

415

416 Note that, surprisingly, the second-most optimal number of clusters for Survey 2 was three (Fig.
417 4B), despite S2 having only two possible responses. This will be discussed further in the
418 Discussion section.

419

420 Standard Deviation (row 1), Proportion of Variance (row 2) and Cumulative Proportion (row 3)
421 are summarized in Table 6 for S1, for each principal component. Table 6 additionally supplies
422 the information to explain each component and its relative weighting, requisite to understanding
423 all components. Horn's parallel analysis on S1 retained optimally 2 factors, though 3 factors was
424 also within acceptable boundaries. S2 also retained 2 factors, and for S3 3 factors were retained.
425 Given the S3 parallel analysis results, and given the more robust separation of, and cohesion
426 within categories in this third survey, we believed that the optimal number of components to
427 retain was 3. Detailed output is provided in Fig. S1, S2 and S3 of supplemental information. The
428 first 3 components explain 97% of the variance of the data. Figure 5A shows the graph resulting
429 from a principal components analysis (PCA) of responses to statements from S1, clustered by K-
430 Means (colored dots). Red lines represent the eigenvectors of each variable (here the certainty
431 categories) for PC1 against PC2. A coefficient close to 1 or -1 indicates that variable strongly
432 influences that component. Thus, the High category has a strong influence on PC1 (0.62),
433 Medium High primarily influences PC2 (0.79), and Medium Low and Low have a notably strong
434 negative relationship with PC1 (-0.62 and -0.44, respectively). Additionally, Low influences PC2
435 in a negative relationship (-0.54). The same approach was followed for S2 and S3, with the
436 results shown in Fig 5B, C, and D. For survey S2, we show the K-Means clustering results for
437 both a three-cluster solution (Fig. 5B), and a two-cluster solution (Fig. 5C).

438

439 **Machine learning:** The corpus of 3221 author-categorized statements was used to build a neural
440 network model. This was validated using a 20-fold CV due to the size of the dataset, with the
441 result indicating that it achieved 89.26% +/- 2.14% accuracy. A test of its performance relative to
442 the highest-scoring dataset (Survey 3 majority rule classification of the publicly-annotated 45
443 statements) showed 82.2% accuracy (see Table 7, right). A further test was done to validate the
444 author-categorized corpus compared to the publicly annotated dataset (see Table 7, left).
445 Majority rule vs. the author's classification gave a kappa value of 0.649 (substantial), while
446 comparison with the model's classification gave a kappa of 0.512 (moderate)

447

448 Discussion

449

450 Evidence to support three levels of certainty in scholarly statements

451

452 In S1, we began with a four-category classification system, since this is the highest number
453 presumed in earlier studies (Zerva et al. used a 5 point numerical scale, but we do not believe
454 they were proposing this as a categorization system). In the absence of any agreed-upon set of
455 labels between these prior studies, and for the purposes of asking untrained annotators to
456 categorize scholarly statements, we labelled these categories Low, Medium Low, Medium High
457 and High. The results of this survey revealed statistically significant categorization agreement for
458 37 of the 45 statements (82.2% of total), with seven statements being doubly-classified and one
459 statement showing poor inter-annotator agreement, for a total of eight 'ambiguous'
460 classifications. The G index (Holley & Guilford, 1964) with only four categories is small, and
461 the statistical probability of chance-agreement in the case of ambiguity is therefore high, which
462 may account for the high proportion of doubly-classified statements. Interestingly, the category
463 Low was almost never selected by the readers. We will discuss that observation in isolation later
464 in this discussion; nevertheless, for the remainder of this discussion we will assume that this
465 category does not exist in our corpus of statements, and will justify this in these later, more
466 detailed arguments.

467

468 With respect to the categories themselves, the category of High had robust support using the G
469 index statistic, indicating that it represents a valid category of certainty based on agreement
470 between the annotators on the use of that labelled category. Support for the other two, medium-
471 level, categories was less robust. This could be interpreted in two ways - one possibility is that
472 these two categories are not distinct from one another, and that readers are selecting one or the
473 other "arbitrarily" with statistical significance, because there were only two choices. This would
474 suggest that there are only two certainty categories used in scholarly writing. The other option is
475 that the labels assigned to these two non-high categories do not accurately reflect the perception
476 of the reader, and thus that the categorizations themselves are flawed, leading to annotator
477 confusion.

478

479 In Survey 2, with only two categories (Relatively High and Relatively Low), statistical support
480 for these two categories was evident, but deeper examination of the results suggests that these
481 categories may still not accurately reflect the reader's perception. For example, seven of the 45
482 statements (15.5%) showed no inter-annotator agreement. Of the remainder, Table 4 shows a
483 clear pattern of association between the strength of certainty perceived by the reader, and the
484 degree to which the readers agreed with one another. Effectively, there was greater agreement on
485 the categorization of high-certainty statements, than low-certainty statements. This mirrors the
486 observations from Survey 1, where the category High generated the highest levels of agreement
487 among annotators. Since this binary categorization system lacks an intermediate category, the G
488 index in this survey is 0.5, meaning that agreement by chance is high. It appears that statements
489 that would have been categorized into a middle class from Survey 1 became distributed between
490 the two Survey 2 categories, rather than being categorized uniformly into the lower category.
491 This would indicate that the two-category explanation for Survey 1 is not well-supported, and
492 possibly, that the labelling of the categories themselves in both Survey 1 and Survey 2 confounds
493 the analysis and does not reflect the perception of the reader. In other words, the category
494 High/Relatively High seems to match a perception that exists in the minds of the readers, but the
495 categories Medium High (S1), Medium Low (S1) and Relatively Low (S2) might not correspond
496 to the perception of the readers for the lower certainty statements, which is why they are less
497 consistent in the selection of these categories.

498

499 To reveal patterns that may clarify what defines these lower categories, we utilized a variety of
500 clustering approaches (Figures 3 and 4). That there are three, rather than two, categories is
501 supported by the hierarchical clustering of all three surveys, shown in Fig. 3 (see clusters along
502 the top edge) which reveals three primary clusters in the data, where high is strongly
503 differentiated from non-high categories. The output from NbClust's "majority rule" approach to
504 selecting the optimal number of clusters was executed on individual surveys. The results for S1
505 and S2 are shown in Fig. 4A and Fig. 4B. The majority rule indicates that there were three
506 discernable clusters in S1. Survey 2 was assessed by the 30 NbClust indices (Charrad et al.,
507 2014) (Fig. 4B). Surprisingly, we found that, while 11 indices recommended only two clusters,
508 six indices suggested that there were three clusters. Since a cluster represents a pattern of
509 categorization-behavior among all evaluators, we take these results as further indication that
510 there are three discernable annotator responses when faced with a certainty categorization task.

511

512 To further explore the meaning of these clusters, we executed a feature reduction analysis using
513 Principal Components. The PCA of Survey 1 revealed three primary components accounting for
514 ~97% of the variability. The main component, accounting for more than half (~56%) of the
515 variation, is characterized by a strong positive influence from the category labelled High, and a
516 strong negative influence from the categories labelled Medium Low and Low. This lends support
517 to our earlier interpretation that there is little ambiguity among annotators about what statements

518 are classified as highly certain, and moreover, when faced with a high-certainty statement
519 annotators will almost never select one of the low categories. The second and third components
520 (accounting for ~32% and ~10% respectively) are more difficult to interpret. Component 2 is
521 characterized by a strong positive influence from the category Medium High, and a strong
522 negative influence from the category Low; Component 3's "signature" is distinguished through a
523 positive influence from the category Low, though as stated earlier, this category was rarely
524 selected, and showed no significant agreement among annotators, making this difficult to
525 interpret. The lack of clarity regarding the interpretation of these second and third components
526 may reflect ambiguity arising from the labelling of the non-high certainty categories in the
527 questionnaire; effectively, the words used for the labels may be confusing the readers, and/or not
528 aligning with their impressions of the statements.

529

530 In an attempt to gain additional evidence for a three-category classification system, we undertook
531 a third survey (S3) in which the reader was offered three categories, ordered from higher to
532 lower, but with numerical labels (1, 2, or 3). The rationale for this was twofold. First, we could
533 not think of three suitable labels that would not inherently bias the results (for example, 'high',
534 'medium', and 'low' would not be suitable because we have already determined that the category
535 'low' is almost never selected). In addition, we wished to know if category labels were a
536 potential source of bias, and therefore more semantically neutral labels might lead to a stronger
537 correspondence between the annotators. Indeed, Survey 3 generated the most consistent
538 agreement of the 3 questionnaires, where only four of the 45 statements did not meet the cutoff
539 level for annotator agreement, and none were doubly-classified. It is not possible to disambiguate
540 if this enhanced agreement is due to the annotators being presented with a "correct" number of
541 categories, or if it supports the suggestion that the presentation of meaningful (but non-
542 representative) category labels caused annotators to behave inconsistently in S1 and S2, or
543 perhaps a combination of both. As with S1, NbClust's "majority rule" proposes 3 clusters for S3
544 (Fig. 4C).

545

546 In Fig. 3 we present the correlation matrix to show how the categories relate to one another
547 between the three surveys, using a Spearman Correlation. High (S1) is clearly correlated with
548 Relatively High (S2) and Category 1 (S3). Medium Low (S1), Relatively Low (S2) and Category
549 3 (S3), are also highly correlated. Low (S1) only has moderate correlation with Relatively Low
550 (S2) and Category 3 (S3). The intermediate values Medium High (S1) and Category 2 (S3), are
551 found on the negative side of Principal Component 1 (Fig. 5A & 5D), which supports the
552 interpretation that a High certainty category is strongly supported, and strongly distinct from
553 other categories. The non-high categories appear as distinct blocks within the correlation matrix,
554 but with more ambiguity or inconsistency, though the Jaccard similarity index was sufficient to
555 support the existence of these two lower-certainty categories. Additionally, the clusters identified
556 by the Spearman analysis (3 clusters) are supported by the results of the HCA analysis (3
557 branches).

558

559 One general source of inconsistency we noted in the data could be described as a “tendency
560 towards the middle”. When a category is removed, statements from that category tend to
561 distribute to adjacent categories. We presume this reflects some form of “central tendency bias”,
562 a behavioral phenomenon earmarked as a preference for selecting a middle option.
563 (Hollingworth, 1910; Huttenlocher, Hedges & Vevea, 2000; Duffy et al., 2010). Nevertheless,
564 this did not appear to be sufficiently strong in this investigation to mask the detection of distinct
565 clusters of categorization behavior.

566

567 In summary, the results suggest that there are three categories of certainty in the minds of the
568 readers of scholarly assertions. One category is clearly distinguished as representing high-
569 certainty statements. The other two categories, representing non-high certainty statements, are
570 distinct from one another in the minds of the annotators, however, seem to not be reflected well
571 by the labels “moderately/relatively + high/low”. Nevertheless, they do appear to represent a
572 higher-to-lower spectrum, since the replacement of textual labels with a numerical range resulted
573 in stronger annotator agreement about these two lower categories.

574

575 **The absence of a Low certainty category**

576

577 Several studies that preceded this one (Friedman et al., 1994; Wilbur, Rzhetsky & Shatkay, 2006;
578 De Waard & Schneider, 2012) suggested four categories of certainty, with one of those being a
579 category that would represent the lowest certainty. In this study, we identify only three. The
580 category that seems to be absent from our data is this lowest category - generally described as
581 “no knowledge” in these three precedent studies. We examined our corpus and, given the
582 grammatical cues suggested by De Waard (De Waard & Maat, 2012) we identified two
583 statements in our corpus that, by those metrics, should have scored in the Low category. Those
584 are Statement 3, “*However, this was not sufficient for full blown transformation of primary*
585 *human cells, which also required the collaborative inhibition of pRb, together with the*
586 *expression of hTERT, RASV12.*”, and Statement 4, “*Hence, the extent to which miRNAs were*
587 *capable of specifically regulating metastasis has remained unresolved.*” Looking at the results in
588 Table 3, 4 and 5, these two statements were annotated with considerable agreement as high-
589 certainty statements - the opposite of what would have been predicted. One explanation for this
590 is that the statements are making a negative claim, with high certainty, and thus are being
591 categorized as high-certainty assertions by our annotators. If that is the case, then the category of
592 “no knowledge” may not be a category that lies anywhere on the spectrum of certainty, and may
593 reflect a distinct feature of scholarly communication discourse, or (more likely) a combination of
594 the meta-knowledge facets of certainty and polarity.

595

596 **Application of this categorization system**

597

598 As indicated in the Introduction, a primary motivation for this study is its application to the
599 automated capture of metadata related to the certainty being expressed in text-mined scholarly
600 assertions, or to identify or monitor ‘hedging erosion’. To demonstrate how the outcomes of this
601 study can be applied, we have used the data described here to generate, by machine-learning, an
602 automated certainty classifier capable of assigning new scholarly statements into one of the three
603 certainty categories. Two exemplar outputs from this classification system are shown in Figs. 6
604 and 7. Figure 6 shows three sets of statements, color-coded by the category of certainty detected
605 by our classifier - green (Category A, associated with High certainty), orange (Category B, non-
606 high/moderate), and red, (Category C non-high/low). Two citation chains relate to the
607 accumulation of beta-APP in muscle fibers of Alzheimer’s Disease patients (Fig. 6A & 6B),
608 while Fig. 6C shows a longer citation chain identified by Greenberg as being problematic with
609 respect to ‘citation-distortion’ (Greenberg, 2009). The panels reveal that the degree of certainty
610 can change through citation, becoming higher (Fig.6A & 6B). Fig. 6C reveals a similar trend
611 toward increasing certainty, with the exception of one author who used a clearly high-certainty
612 assertion four years before others in the community expressed the same idea with certainty.

613

614 Figure 7 demonstrates how this certainty classification could be used to enhance the quality of
615 machine-extracted information. The figure shows a block of machine-readable information
616 following the NanoPublication model for scholarly publishing. The sentence which has been
617 extracted in this exemplar is from the article with DOI ‘10.1371/journal.pone.0073940’, and the
618 specific sentence “Consequently miRNAs have been demonstrated to act either as oncogenes
619 (e.g., miR-155,miR-17–5p and miR-21) or tumor suppressors (e.g., miR-34,miR-15a,miR-16–1
620 and let-7)” Following the rules of NanoPublications, a single scholarly assertion is captured - in
621 this case, that “miR-34 has the function of tumor suppressor” (red text). The provenance block
622 contains information showing the degree of certainty being expressed (Category A, which maps
623 to the highest certainty category in our classifier; blue text). Finally, there is a block of citation
624 information regarding the NanoPublication itself, such that the author of the certainty
625 classification can be properly cited (green text).

626

627 **Tools for researchers, authors, reviewers, and data miners**

628

629 As discussed in the introduction, researchers may lack the knowledge required to assess the
630 legitimacy of claims that are not directly in their domain, or may be unaware of the history of a
631 claim if they have not followed a citation chain to its roots. Similarly, when acting as peer
632 reviewers, there is little tooling to assist them in evaluating the validity of assertions in the
633 submitted manuscript or funding proposal. In parallel with research into automated identification
634 of reference-spans (Saggion, Ronzano & Others, 2016), the availability of a certainty classifier
635 would make it possible to automate the creation of annotated citation chains such as shown in
636 Fig. 6. Reviewers could then use these to determine if a claim was being made with unusually
637 high (or low) certainty - like the Magstalia statement from 2003, shown in Fig. 6C - and thus

638 enhance the confidence of their reviews. Similarly, such tools could become an important part of
639 the scholarly planning process. During the preparation of a paper or proposal, researchers could
640 be made aware of dubious assertions, and avoid relying on these as the bases for their hypothesis.
641 In the context of automated data mining, assuming that incremental steps towards certainty
642 should be associated with the existence of supporting data, the automated detection of “certainty
643 inflection points” could be used by data mining algorithms to identify the specific dataset
644 containing data supporting (or refuting) a given claim. Together with the use of certainty
645 classification in the context of text-mining discussed above, the use of such a classification
646 system may become an important part of the scholarly publishing lifecycle.

647

648 **Future investigations to elucidate perceptions of certainty**

649

650 A variety of future studies could provide additional insight into how researchers communicate
651 and perceive certainty. The results presented here seem to suggest that words like “medium” and
652 “low” do not align well with the perception held by researchers as they read statements that fall
653 into non-high certainty categories. Future studies could extract additional information in the
654 questionnaire, such as questions related to the basis upon which an assertion was made (e.g.
655 speculation, direct or indirect observation, etc.), as it may be that the distinction between the two
656 lower certainty categories is being made based on other kinds of implicit information, rather than
657 being specifically “medium” or “low” expressions of certainty. It would also be interesting to
658 capture demographic information, to determine if perception of certainty changes as a researcher
659 becomes more experienced, if it differs between different linguistic groups, or if it is associated
660 with other demographic variables

661

662

663 **Conclusions**

664

665 This study attempted to derive a data-driven certainty classification system, using statements
666 from scholarly literature in the biological sciences. We found support for three categories of
667 certainty within the dataset of 45 scholarly statements we selected. These consisted of one very
668 distinct High Certainty category, and two lower-certainty categories that were seemingly not
669 well-described using textual labels, but were clearly distinguishable from one another using
670 statistical algorithms. We suggest that a fourth category described in previous studies - best
671 described as “lack of information” - likely does not belong in the same categorization system,
672 and is likely a measure of a different discourse feature than “certainty”. Finally, we show how
673 this categorization system could be used to capture key contextual information within text-
674 mining pipelines, to improve the quality of automated information capture. Work on the
675 machine-learning models leading to such an automated classifier are well underway, and are
676 already showing a high degree of accuracy, indicating that machines may be capable of detecting
677 and distinguishing the subtle linguistic cues of certainty that we have observed in this study.

678 While this study was limited to biomedical statements, and thus may be applicable only in this
679 domain, it nevertheless seems likely that the results will be more generalizable, at least within
680 the sciences where these kinds of grammatical structures are commonly used.

681

682 **Acknowledgements**

683

684 We wish to thank all of the anonymous volunteers who donated their time to answering these
685 questions to the best of their ability. We wish to thank Leiden University for hosting us during a
686 student exchange, and for providing free access to their Qualtrics questionnaire platform. The
687 authors would like to acknowledge Dr. Ron Daniel Jr. for useful discussions around survey
688 design and the staff at Cell Press for allowing use of their Boston offices and providing feedback
689 on the surveys. We would also like to thank the Foundation DTL DP (Data Projects) for their
690 support of this initiative.

691

692 **References**

- 693 Abadi M, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, Corrado GS, Davis A, Dean J,
694 Devin M, Ghemawat S, Goodfellow IJ, Harp A, Irving G, Isard M, Jia Y, Józefowicz R,
695 Kaiser L, Kudlur M, Levenberg J, Mané D, Monga R, Moore S, Murray DG, Olah C,
696 Schuster M, Shlens J, Steiner B, Sutskever I, Talwar K, Tucker PA, Vanhoucke V,
697 Vasudevan V, Viégas FB, Vinyals O, Warden P, Wattenberg M, Wicke M, Yu Y, Zheng X.
698 2016. TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems.
699 *CoRR* abs/1603.04467.
- 700 Agami R, Bernards R. 2000. Distinct initiation and maintenance mechanisms cooperate to induce
701 G1 cell cycle arrest in response to DNA damage. *Cell* 102:55–66.
- 702 Baziotis C, Pelekis N, Doulkeridis C. 2017. Datastories at semeval-2017 task 4: Deep lstm with
703 attention for message-level and topic-based sentiment analysis. In: *Proceedings of the 11th*
704 *international workshop on semantic evaluation (SemEval-2017)*. 747–754.
- 705 Campbell J, Narayanan A, Burford B, Greco M. 2010. Validation of a multi-source feedback tool
706 for use in general practice. *Education for primary care: an official publication of the*

- 707 *Association of Course Organisers, National Association of GP Tutors, World Organisation*
708 *of Family Doctors* 21:165–179.
- 709 Campbell PA, Perez-Iratxeta C, Andrade-Navarro MA, Rudnicki MA. 2007. Oct4 targets
710 regulatory nodes to modulate stem cell function. *PloS one* 2:e553.
- 711 Charrad M, Ghazzali N, Boiteau V, Niknafs A. 2014. NbClust: AnRPackage for Determining the
712 Relevant Number of Clusters in a Data Set. *Journal of statistical software* 61. DOI:
713 10.18637/jss.v061.i06.
- 714 Chollet F and others. 2015. Keras. Available at <https://keras.io> (accessed May 25, 2019).
- 715 Chouikhi H, Charrad M, Ghazzali N. 2015. A comparison study of clustering validity indices. In:
716 *2015 Global Summit on Computer Information Technology (GSCIT)*. 1–4.
- 717 Clark T, Ciccarese PN, Goble CA. 2014. Micropublications: a semantic model for claims,
718 evidence, arguments and annotations in biomedical communications. *Journal of biomedical*
719 *semantics* 5:28.
- 720 Cohen J. 1968. Weighted kappa: nominal scale agreement with provision for scaled
721 disagreement or partial credit. *Psychological bulletin* 70:213–220.
- 722 Crestan E, Pantel P. 2010. Web-scale Knowledge Extraction from Semi-structured Tables. In:
723 *Proceedings of the 19th International Conference on World Wide Web*. WWW '10. New
724 York, NY, USA: ACM, 1081–1082.
- 725 Deery C, Wagner ML, Longbottom C, Simon R, Nugent ZJ. 2000. The prevalence of dental
726 erosion in a United States and a United Kingdom sample of adolescents. *Pediatric dentistry*
727 22:505–510.
- 728 De Waard A, Maat HP. 2012. Epistemic modality and knowledge attribution in scientific
729 discourse: A taxonomy of types and overview of features. *Proceedings of the 50th Annual*

- 730 *Meeting of the Association for Computational Linguistics*, Jeju Island, Korea — July 08-14,
731 2012.
- 732 De Waard A, Schneider J. 2012. Formalising uncertainty: An ontology of reasoning, certainty
733 and attribution (ORCA). In: *Proceedings of the Joint Workshop on Semantic Technologies*
734 *Applied to Biomedical Informatics and Individualized Medicine (SATBI+ SWIM 2012)*.
735 Boston MA, USA, Nov. 11–15, 2012.
- 736 Duffy S, Huttenlocher J, Hedges LV, Crawford LE. 2010. Category effects on stimulus
737 estimation: shifting and skewed frequency distributions. *Psychonomic bulletin & review*
738 17:224–230.
- 739 Dunham MH. 2006. *Data Mining: Introductory And Advanced Topics*. Pearson Education India.
- 740 Friedman C, Alderson PO, Austin JH, Cimino JJ, Johnson SB. 1994. A general natural-language
741 text processor for clinical radiology. *Journal of the American Medical Informatics*
742 *Association: JAMIA* 1:161–174.
- 743 Garg R, Oh E, Naidech A, Kording K, Prabhakaran S. 2019. Automating Ischemic Stroke
744 Subtype Classification Using Machine Learning and Natural Language Processing. *Journal*
745 *of stroke and cerebrovascular diseases: the official journal of National Stroke Association*.
746 S1052-3057(19)30048-5 (epub ahead of print). DOI:
747 10.1016/j.jstrokecerebrovasdis.2019.02.004.
- 748 Gauthier TD. 2001. Detecting Trends Using Spearman’s Rank Correlation Coefficient.
749 *Environmental Forensics* 2:359–362.
- 750 GENIA Event Extraction (GENIA) - BioNLP Shared Task. Available at <http://2011.bionlp->
751 [st.org/home/genia-event-extraction-genia](http://2011.bionlp-st.org/home/genia-event-extraction-genia) (accessed May 13, 2019).
- 752 Greenberg SA. 2009. How citation distortions create unfounded authority: analysis of a citation

- 753 network. *BMJ* 339:b2680.
- 754 Groth P, Gibson A, Velterop J. 2010. The anatomy of a nanopublication. *Information services &*
755 *use* 30:51–56.
- 756 Holley JW, Guilford JP. 1964. A Note on the G Index of Agreement. *Educational and*
757 *psychological measurement* 24:749–753.
- 758 Hollingworth HL. 1910. The Central Tendency of Judgment. *The Journal of Philosophy,*
759 *Psychology and Scientific Methods* 7:461–469.
- 760 Huttenlocher J, Hedges LV, Vevea JL. 2000. Why do categories affect stimulus judgment?
761 *Journal of experimental psychology. General* 129:220–241.
- 762 Hyland K. 1996. Writing Without Conviction? Hedging in Science Research Articles. *Applied*
763 *Linguistics* 17:433–454.
- 764 Jolliffe I. 2011. Principal Component Analysis. In: *International Encyclopedia of Statistical*
765 *Science*. 1094–1096.
- 766 Landis JR, Richard Landis J, Koch GG. 1977. The Measurement of Observer Agreement for
767 Categorical Data. *Biometrics* 33:159.
- 768 Latour B, Woolgar S. 2013. *Laboratory Life: The Construction of Scientific Facts*. Princeton
769 University Press.
- 770 Lewis RJ. 2000. An introduction to classification and regression tree (CART) analysis. In:
771 *Annual meeting of the society for academic emergency medicine in San Francisco,*
772 *California*. San Francisco, California, USA. May 22-25.
- 773 Light M, Qiu XY, Srinivasan P. 2004. The language of bioscience: Facts, speculations, and
774 statements in between. In: *HLT-NAACL 2004 Workshop: Linking Biological Literature,*
775 *Ontologies and Databases*.

- 776 Lix LM, Yogendran MS, Shaw SY, Burchill C, Metge C, Bond R. 2008. Population-based data
777 sources for chronic disease surveillance. *Chronic diseases in Canada* 29:31–38.
- 778 Lorés R. 2004. On RA abstracts: from rhetorical structure to thematic organisation. *English for*
779 *Specific Purposes* 23:280–302. DOI: 10.1016/j.esp.2003.06.001.
- 780 Malhotra A, Younesi E, Gurulingappa H, Hofmann-Apitius M. 2013. “HypothesisFinder:” a
781 strategy for the detection of speculative statements in scientific text. *PLoS computational*
782 *biology* 9:e1003117.
- 783 Ma Y, Peng H, Cambria E. 2018. Targeted aspect-based sentiment analysis via embedding
784 commonsense knowledge into an attentive LSTM. In: *Thirty-Second AAAI Conference on*
785 *Artificial Intelligence*.
- 786 Min-Yen KAN. The Computational Linguistics Scientific Summarization Shared Task (CL-
787 SciSumm 2018). Available at <http://wing.comp.nus.edu.sg/~cl-scisumm2018/> (accessed
788 January 21, 2019).
- 789 Narayanan A, Greco M, Powell H, Bealing T. 2011. Measuring the quality of hospital doctors
790 through colleague and patient feedback. *Journal of Management & Marketing in*
791 *Healthcare* 4:180–195.
- 792 Narayanan A, Greco M, Reeves P, Matthews A, Bergin J. 2014. Community pharmacy
793 performance evaluation: Reliability and validity of the Pharmacy Patient Questionnaire.
794 *International Journal of Healthcare Management* 7:103–119.
- 795 Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer
796 P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M,
797 Duchesnay É. 2011. Scikit-learn: Machine Learning in Python. *Journal of machine learning*
798 *research: JMLR* 12:2825–2830.

- 799 Prieto M. 2019. Certainty Corpus. Available at
800 https://github.com/Guindillator/Certainty/blob/master/Corpus/Complete_statements.txt
801 (accessed March 5, 2019).
- 802 Prieto M. 2019. Guindillator/Certainty. Available at <https://github.com/Guindillator/Certainty>
803 (accessed May 29, 2019).
- 804 Qualtrics. Available at <https://www.qualtrics.com/research-core/survey-software/> (accessed
805 February 19, 2023).
- 806 Raithel J. 2008. *Quantitative Forschung: Ein Praxiskurs*. Springer-Verlag.
- 807 Rubinstein A, Harner H, Krawczyk E, Simonson D, Katz G, Portner P. 2013. Toward fine-
808 grained annotation of modality in text. In: *Proceedings of the IWCS 2013 Workshop on*
809 *Annotation of Modal Meanings in Natural Language (WAMM)*. 38–46.
- 810 Saggion H, Ronzano F, Others. 2016. Trainable citation-enhanced summarization of scientific
811 articles. In: *Proceedings of the Joint Workshop on Bibliometric-enhanced Information*
812 *Retrieval and Natural Language Processing for Digital Libraries (BIRNDL)*. 175–186.
- 813 Snow R, O'Connor B, Jurafsky D, Ng AY. 2008. Cheap and Fast---but is It Good?: Evaluating
814 Non-expert Annotations for Natural Language Tasks. In: *Proceedings of the Conference on*
815 *Empirical Methods in Natural Language Processing*. EMNLP '08. Stroudsburg, PA, USA:
816 Association for Computational Linguistics, 254–263.
- 817 Thompson P, Nawaz R, McNaught J, Ananiadou S. 2011. Enriching a biomedical event corpus
818 with meta-knowledge annotation. *BMC bioinformatics* 12:393.
- 819 Sauvageot N, Alkerwi A, Adelin A, Guillaume M. 2013. Validation of the Food Frequency
820 Questionnaire Used to Assess the Association between Dietary Habits and Cardiovascular
821 Risk Factors in the NESCAV Study. 2013. *Journal of nutrition & food sciences* 3:1–8.

- 822 Vincze V, Szarvas G, Farkas R, Móra G, Csirik J. 2008. The BioScope corpus: biomedical texts
823 annotated for uncertainty, negation and their scopes. *BMC bioinformatics* 9 Suppl 11:S9.
- 824 Wang Y, Huang M, Zhao L, Others. 2016. Attention-based LSTM for aspect-level sentiment
825 classification. In: *Proceedings of the 2016 conference on empirical methods in natural*
826 *language processing*. 606–615.
- 827 Wilbur WJ, Rzhetsky A, Shatkey H. 2006. New directions in biomedical text annotation:
828 definitions, guidelines and corpus construction. *BMC bioinformatics* 7:356.
- 829 Wilkinson MD, Dumontier M, Aalbersberg IJJ, Appleton G, Axton M, Baak A, Blomberg N,
830 Boiten J-W, da Silva Santos LB, Bourne PE, Bouwman J, Brookes AJ, Clark T, Crosas M,
831 Dillo I, Dumon O, Edmunds S, Evelo CT, Finkers R, Gonzalez-Beltran A, Gray AJG, Groth
832 P, Goble C, Grethe JS, Heringa J, 't Hoen PAC, Hooft R, Kuhn T, Kok R, Kok J, Lusher SJ,
833 Martone ME, Mons A, Packer AL, Persson B, Rocca-Serra P, Roos M, van Schaik R,
834 Sansone S-A, Schultes E, Sengstag T, Slater T, Strawn G, Swertz MA, Thompson M, van
835 der Lei J, van Mulligen E, Velterop J, Waagmeester A, Wittenburg P, Wolstencroft K, Zhao
836 J, Mons B. 2016. The FAIR Guiding Principles for scientific data management and
837 stewardship. *Scientific data* 3:160018.
- 838 Xu S, Lorber MF. 2014. Interrater agreement statistics with skewed data: Evaluation of
839 alternatives to Cohen's kappa. *Journal of consulting and clinical psychology* 82:1219–1227.
- 840 Zerva C, Batista-Navarro R, Day P, Ananiadou S. 2017. Using uncertainty to link and rank
841 evidence from biomedical literature for model curation. *Bioinformatics* 33:3784–3792.
- 842

Figure 1(on next page)

How a claim becomes a fact.

These sentences represent a series of scholarly assertions about the same biological phenomenon, revealing that the core assertion transforms from a hedging sentence into statements resembling fact through several steps, but without additional evidence. (de Waard, 2012)

How a claim becomes a fact

*“These miRNAs neutralize p53- mediated CDK inhibition, **possibly** through direct inhibition of the expression of the tumor suppressor LATS2.” (Voorhoeve et al. 2007)*

*“In a genetic screen, miR-372 and miR-373 **were found to** allow proliferation of primary human cells that express oncogenic RAS and active p53, **possibly** by inhibiting the tumor suppressor LATS2 (Voorhoeve et al., 2006).” (Kloosterman and Plasterk 2006)*

*“[On the other hand,] two miRNAs, miRNA-372 and-373, function as **potential** novel oncogenes in testicular germ cell tumors by inhibition of LATS2 expression, **which suggests that** Lats2 is an important tumor suppressor (Voorhoeve et al., 2006).” (Yabuta et al. 2007)*

*“Two oncogenic miRNAs, miR-372 and miR-373, **directly inhibit** the expression of Lats2, **thereby** allowing tumorigenic growth in the presence of p53 (Voorhoeve et al., 2006).” (Okada et al. 2011)*

Figure 2

Example of the Survey 1 questionnaire interface.

A scholarly assertion is highlighted in blue, in its original context. Participants are asked to characterize the blue assertion, using one of 4 levels of certainty (High, Medium High, Medium Low or Low).



CENTRO DE BIOTECNOLOGÍA
Y GENÓMICA DE PLANTAS



Scientific Statement:

Our observations raise the important prediction that many malignancies considered to be non-TK driven because of the absence of a dominant TK mutation may indeed be dependent on TK signaling. It is likely that in different cell types, different PTPs may play roles similar to PTPN12 in suppressing tumorigenesis, possibly by antagonizing different combinations of TKs.

Forget what you know about biology... What do you think is the certainty level expressed by the authors in the statement highlighted in blue?

High

Medium High

Medium Low

Low

Figure 3

Spearman Rank Correlation and hierarchically-clustered heatmap comparing the statements assigned to the Certainty Categories among all three questionnaires.

The clustering tree and heatmap are based on participants' responses from questionnaires S1, S2 and S3. Certainty categories from surveys S1 S2 and S3 were combined and clustered based on the similarity of the statements appearing in those categories (top and left hierarchical axes). Colored boxes show the Spearman's rank-order correlation for each pairwise category as both a heatmap and its corresponding numerical coefficient.

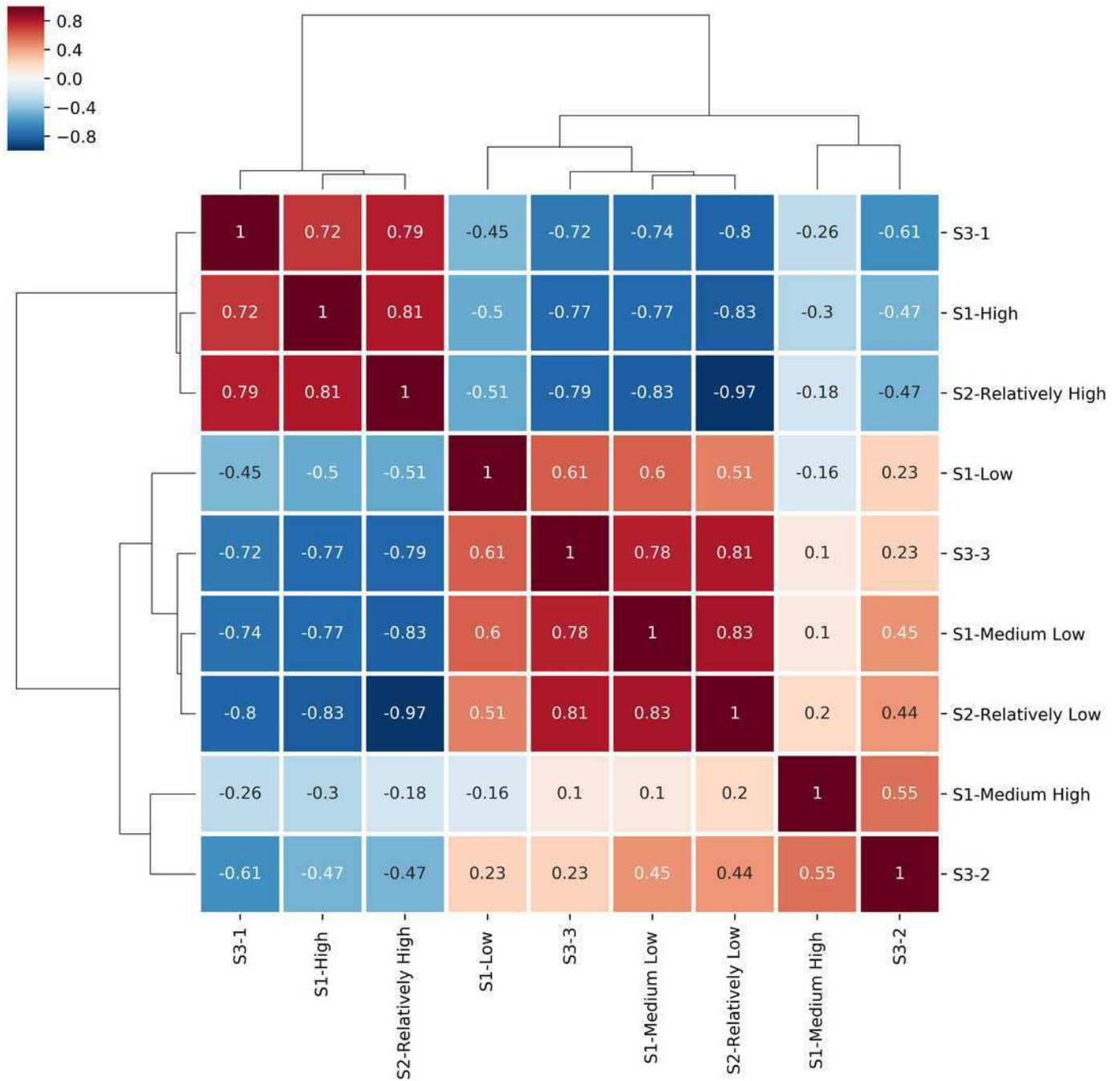


Figure 4

Majority rule output for deciding optimal number of clusters (k) in the three surveys.

(A) Majority rule indicates three clusters for Survey 1. (B) Majority rule indicates two clusters for Survey 2, though there is notable support for three clusters. (C) Majority rule indicates three clusters for Survey 3, with notable support for two clusters.

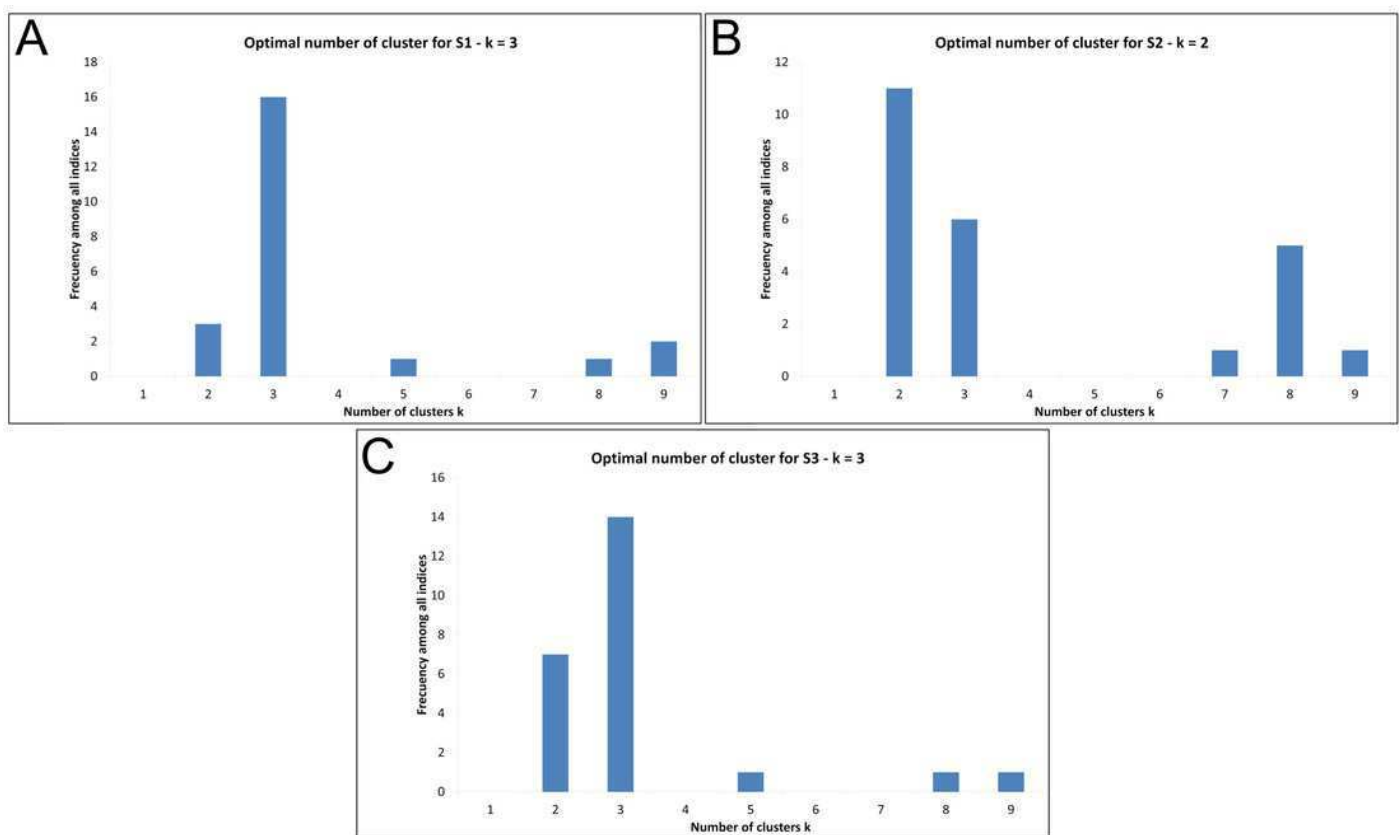


Figure 5

Principal component analysis of questionnaire responses in the three surveys.

Bi-plot of certainty level distribution over results from k-means clustering (colors) for: Survey 1 (A), Survey 2 with three clusters (B), Survey 2 with two clusters (C) and Survey 3 (D). Each dot represents a statement. Red lines are the eigenvectors for each component.

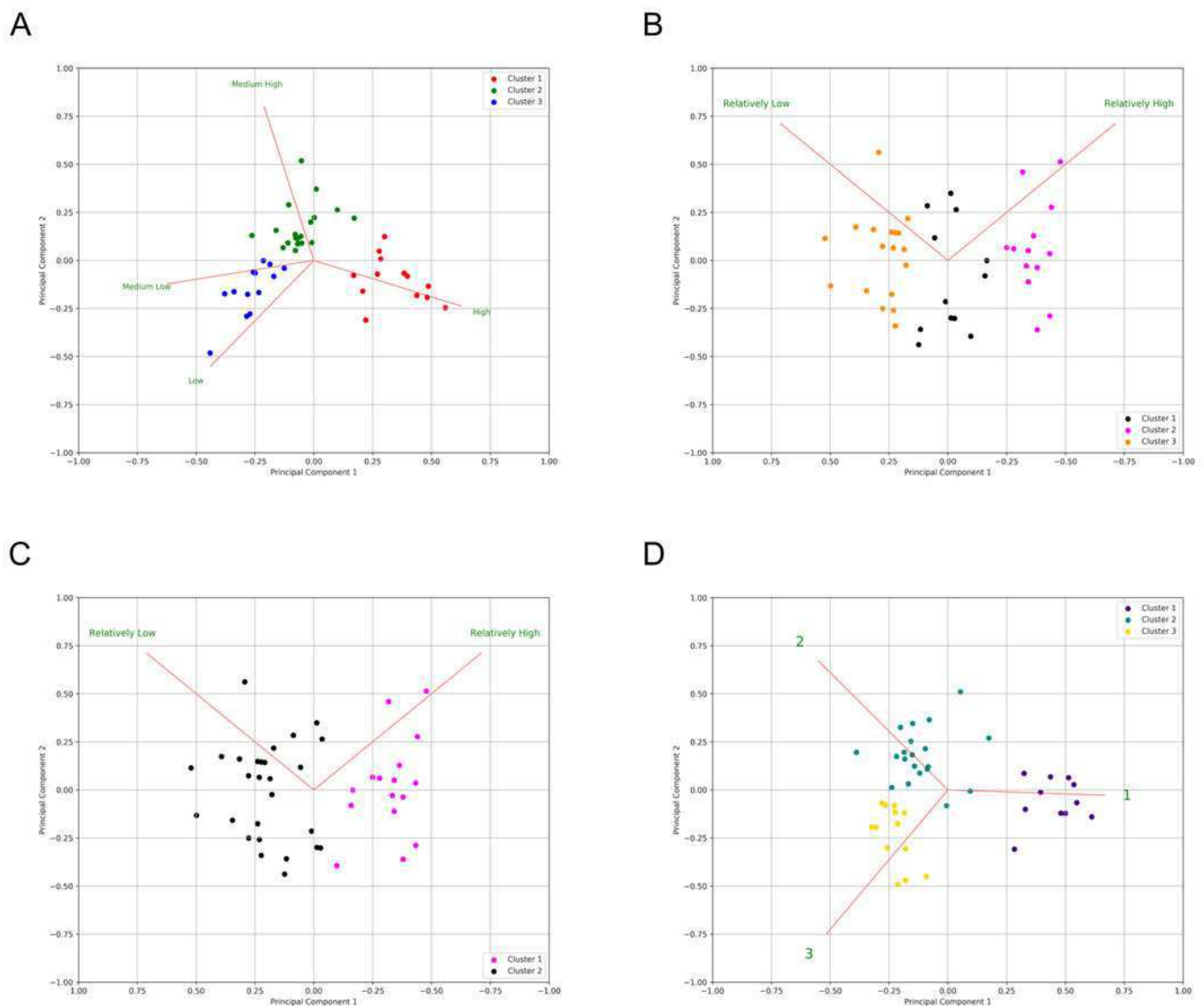


Figure 6(on next page)

Automated classification of scholarly assertions related to the accumulation of beta-APP protein in muscle fibres, color coded as green (Category A - highest certainty), orange (Category B - medium certainty) and red (Category C - lowest certainty).

(A and B) Two citation chains showing that the degree of certainty expressed in the most recent statement is higher than that in the cited text. (C) A selection of statements identified by Greenberg, 2009, as being potentially indicative of 'citation distortion'. In this panel, there is a general trend to higher certainty over time, with the exception of an early high-certainty statement by Mastaglia in 2003 (second row from the bottom).

A

"We have previously demonstrated that accumulation of A β PP epitopes precedes other abnormalities in IBM muscle fibers" (Askanas et al., 2000b)

" β APP accumulation is considered to play a major role in the pathogenesis of IBM and AD and is thought to precede other changes in both diseases"(Askanas et al., 1996)

"Those muscle fibers, widely prevalent in our one case of hereditary IBM, may represent early changes of IBM and therefore be analogous to the finding in AD brains where PAP accumulations in the "diffuse" Congoed-negative plaques seem to represent early changes"(Askanas, Engel & Alvarez, 1992)

B

"We have previously demonstrated that accumulation of A β PP epitopes precedes other abnormalities in IBM muscle fibers" (Askanas et al., 2000b)

"Increased β APP-mRNA and increased accumulation of β APP epitopes appear to precede other abnormalities in IBM muscle fiber" (Askanas et al., 1997a)

"One possibility is that one protein is accumulated first, due to excessive synthesis, e.g., excessive transcription of mRNA in the IBMs is known for beta APP"(Askanas & Engel, 1995)

C

Recently it was reported that s-IBM vacuolated muscle fibers, and those in some other vacuolar myopathies, contain a marker of autophagosomes, but only in s-IBM is it colocalized with A β PP[18]. (Askanas and Engel 2007)

Overexpression of amyloid precursor protein (APP) and subsequent accumulation of cleaved fragments including β -amyloid in vacuolated muscle fibers is considered a central mechanism in the pathogenesis of s-IBM.[2] (Lünemann et al. 2007)

it is now established that A β /A β PP is also abnormally accumulated in muscle fibers of s-IBM patients, where they are considered to play an important pathogenetic role[4,5,6,7] (Askanas and Engel 2006)

A possibility that excessive accumulation of A β PP/A β induces inflammation has been proposed by us and by others.[1-3,7,10] (Askanas and Engel 2003)

Deposition of the A β fragment of the amyloid precursor protein is a feature of affected muscle in IBM (see below) and it has been shown that muscle cells can secrete A β . [10] Interaction of A β with muscle cells in turn can stimulate IL-6 production by these cells [19]... (Mastaglia et al. 2003)

However, in some abnormal muscle fibers in IBM, the accumulation of β APP appears to extend outside the muscle fiber boundary. This may have been attributable to a fragility of the fiber's surface membrane, which could have been transiently broken.[25] (Baron et al. 2001)

Figure 7 (on next page)

An exemplar prototype NanoPublication including certainty annotations.

The figure shows how certainty classifications could be used as additional, and important metadata when added to text-mining pipelines. A NanoPublication is a machine-readable representation of a scholarly assertion, carrying with it all of its provenance. In this exemplar (hypothetical) NanoPublication for statement #29 in this study, the concept being asserted (that microRNA mir-155 has the function of a Tumor Suppressor) is captured using ontologically-based concepts in the “assertion” block of the NanoPublication (red text). The “provenance” block then carries a variety of information about the original text, including our proposed annotation of the category of certainty that was detected in that statement (blue text) as being a part of Certainty Category A, which could be used to filter assertions based on the degree of certainty they express. The final block, “pubinfo”, contains authorship, license, and citation information for the NanoPublication itself, expressing the terms of usage of this metadata, and who to cite (green text). This entire structure can be interpreted by automated agents, and fully complies with the FAIR Data Principles.

```

@prefix this: <http://linkeddata.systems/nanopubs/ABC123> .
@prefix sub: <http://linkeddata.systems/nanopubs/ABC123#> .
@prefix void: <http://rdfs.org/ns/void#> .
@prefix dcterms: <http://purl.org/dc/terms/> .
@prefix dcelem: <http://purl.org/dc/elements/1.1/> .
@prefix np: <http://www.nanopub.org/nschema#> .
@prefix sio: <http://semanticscience.org/resource/> .
@prefix pav: <http://swan.mindinformatics.org/ontologies/1.2/pav/> .
@prefix prov: <http://www.w3.org/ns/prov#> .
@prefix xsd: <http://www.w3.org/2001/XMLSchema#> .
@prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
@prefix dcat: <http://www.w3.org/ns/dcat#> .
@prefix schema: <https://schema.org/> .
@prefix orca: <https://vocab.deri.ie/orca#/> .
@prefix cancer: <http://purl.obolibrary.org/obo/NCIT_> .
@prefix mir: <http://www.mirbase.org/cgi-bin/mirna_entry.pl?acc=> .
@prefix certainty: <http://w3id.org/orca-x#> .

sub:Head {
  this: np:hasAssertion sub:assertion ;
  np:hasProvenance sub:provenance ;
  np:hasPublicationInfo sub:pubinfo ;
  a np:Nanopublication .
}

sub:assertion {
  mir:MI0000681 sio:has-function cancer:C17362 .
  mir:MI0000681 rdfs:label 'MicroRNA miR-34'@en .
  cancer:C17362 rdfs:label 'Tumor Suppressor'@en .
}

sub:provenance {
  sub:assertion dcterms:author "Certainty Classifier" ;
  dcterms:title "Automated Certainty Classification of Statement
from doi:10.1371/journal.pone.0073940 page 2, line 19, character 29" ;
  dcat:distribution sub:_1 ;
  prov:wasGeneratedBy "Mario Prieto's Certainty Classifier" ;
  orca:hasConfidenceLevel certainty:CategoryA .

  sub:_1 dcelem:format "application/pdf" ;
  a void:Dataset , dcat:Distribution ;
  schema:identifier "10.1371/journal.pone.0073940" ;
  dcat:downloadURL <https://journals.plos.org/plosone/article/file?id=10.1371/journal.pone.0073940> .
}

sub:pubinfo {
  this: dcterms:created "2018-11-21T00:00:00.0Z"^^xsd:dateTime ;
  dcterms:rights <https://creativecommons.org/publicdomain/zero/1.0> ;
  dcterms:rightsHolder <https://orcid.org/0000-0001-9833-8031> ;
  pav:authoredBy "Mario Prieto" , <https://orcid.org/0000-0001-9833-8031> ;
  pav:versionNumber "1" .
}

```

Table 1 (on next page)

Comparison of corpora and approaches used in prior investigations into scholarly certainty

1

Table 1: Comparison of corpora and approaches used in prior investigations into scholarly certainty

	N° of annotators	Annotator expertise	Text provenance	Discourse segment source	Approach to automated detection	Number of certainty classification classes	Corpus Size	Meta knowledge examined
Light, Qiu, and Srinivasan (2004)	4	following annotation guidelines	Medline	Abstract	SVM	3	2,093 statements	certainty
Malhotra et al. (2013)	3	following annotation guidelines	Medline	Abstract	Maximum Entropy	4	350 abstracts	certainty
Zerva et al. (2017)	7+2	biomedical	GENIA-MK, BioNLP-ST	Abstract, Text Event	Random Forest classifier + Rule Induction	up to 5	652 passages	certainty
A. De Waard and Maat (2012)	2	publishing	2 articles (Voorhoeve et al. 2006), (Zimmermann et al. 2005)	Full text	N/A	4	812 clauses	certainty, basis, source
Friedman et al. (1994)	3	physics	Columbia Presbyterian Medical Database	Free text	Natural Language Processor	4	230 reports	certainty, degree, change, status, quantity, descriptor
Wilbur, Rzhetsky, and Shatkay (2006)	3+9	following annotation guidelines	Ten research articles	Full text	N/A	4	101 sentences	focus, polarity, certainty, evidence, and directionality
Vincze et al. (2008)	3	linguistics	Clinical, FlyBase, BMC Bioinfo	Free Text, Full Text, Abstract	N/A	2	20,924 statements	certainty, negation
Thompson et al. (2011)	2	following annotation guidelines	Medline	Abstract	N/A	3	36,858 events	manner, source, polarity, certainty, knowledge type
This	375	biomedical	TAC 2014	Full Text	Neural	3	45	certainty

Manuscript					Network		statements	
------------	--	--	--	--	---------	--	------------	--

2

3

Table 2 (on next page)

Jaccard similarity clusters resulting from K-Means applied to questionnaire results.

Jaccard similarity index on k-means results from scaled questionnaire responses. The score is the result from statements' labels pairwise comparison. A dash indicates that it is not possible to compare due to differing cluster size.

1

Table 2: Jaccard similarity clusters resulting from K-Means applied to questionnaire results.

Jaccard similarity index on k-means results from scaled questionnaire responses. The score is the result from statements' labels pairwise comparison. A dash indicates that it is not possible to compare due to differing cluster size.

	S1-S2	S1-S3	S2-S3
Cluster 1-1	0.923	0.923	0.786
Cluster 1-2	-	-	-
Cluster 1-3	0	0	-
Cluster 2-1	-	-	-
Cluster 2-2	0.474	0.737	0.833
Cluster 2-3	-	-	-
Cluster 3-1	0	0	-
Cluster 3-2	-	-	-
Cluster 3-3	0.846	0.692	0.684

2

Table 3 (on next page)

Categorization consistency of statements (by statement number) for survey S1

1

Table 3: Categorization consistency of statements (by statement number) for survey S1

Agreement Level	High	% of Corpus	Medium High	% of Corpus	Medium Low	% of Corpus	Low	% of Corpus
Almost Perfect [0.81-1.00]	29	2.2%	0	0%	0	0%	0	0%
Substantial [0.61-0.8]	25, 27, 30	6.6%	5	2.2%	0	0%	0	0%
Moderate [0.41-0.6]	4, 28, 42	6.6%	19, 35, 37, 40, 45	11.11%	21, 36, 44	6.6%	0	0%
Fair [0.21-0.4]	3, 15, 22, 38	8.9%	2, 8, 9, 16, 17, 20, 34, 39	17.7%	1, 6, 10, 11, 12, 14, 18, 26, 33	20%	0	0%
Poor [0.2]	13	2.2%						
Double-Classified	41, 43	4.4%	7, 23, 24, 31, 32, 41, 43	15.5%	7, 23, 24, 31, 32,	11.11%		

2

Table 4 (on next page)

Categorization consistency of statements (by statement number) for survey S2

1

Table 4: Categorization consistency of statements (by statement number) for survey S2

Agreement Level	Relatively High	% of Corpus	Relatively Low	% of Corpus
Almost Perfect [0.81-1.00]	25, 27, 28, 29, 30, 41	13.32%	36, 44	4.4%
Substantial [0.61-0.8]	3, 15, 22, 38, 40, 42, 43	15.5%	0	0%
Moderate [0.41-0.6]	5, 6, 9,	6.6%	10, 11, 14, 18, 31, 33, 39	15.5%
Fair [0.21-0.4]	4, 37, 45	6.6%	1, 12, 13, 16, 19, 21, 23, 24, 32, 34	22.2%
Poor [0.2]	2, 7, 8, 17, 20, 26, 35	15.5%		
Double Classified	0	0%		

2

Table 5 (on next page)

Categorization consistency of statements (by statement number) for survey S3

1
2**Table 5: Categorization consistency of statements (by statement number) for survey S3**

Agreement Level	Category 1	% of Corpus	Category 2	% of Corpus	Category 3	% of Corpus
Almost Perfect [0.81-1.00]	3, 15	4.4%	0	0%	0	0%
Substantial [0.61-0.8]	27, 28, 29, 38, 42	11.11%	0	0%	0	0%
Moderate [0.41-0.6]	4, 25, 30, 41, 43	11.11%	2, 16, 17, 23, 26, 33, 34, 35, 37, 40	22.22%	0	0%
Fair [0.21-0.4]	22	2.2%	1, 6, 8, 9, 10, 11, 12, 13, 18, 19, 20, 31, 32, 45	31.11%	21, 24, 36, 44	8.8%
Poor [0.2]	5, 7, 14, 39	8.8%				
Double Classified	0	0%				

3

Table 6 (on next page)

Analysis of Principal Components of survey S1

1
2

Table 6: Analysis of Principal Components of survey S1

Principal Components:	Comp.1	Comp.2	Comp.3	Comp.4
High	0.620	-0.235	0.185	0.725
Medium High	-0.210	0.795	0.473	0.317
Medium Low	-0.618	-0.121	-0.480	0.611
Low	-0.436	-0.546	0.716	0.013
Component variances	2.238	1.275	0.382	0.105
Standard deviation	1.496	1.129	0.618	0.324
Proportion of Variance	0.560	0.319	0.095	0.026
Cumulative Proportion	0.560	0.878	0.974	1.000

3
4

Table 7 (on next page)

Performance of the neural network model on the 45 publicly-annotated statements

1

Table 7: Performance of the neural network model on the 45 publicly-annotated statements

S3 Majority Rule vs. Author's Classification				S3 Majority Rule vs. Model's Classification		
	Precision	Recall	Overall accuracy	Precision	Recall	Overall accuracy
Category 1	0.857	0.923	0.933	0.786	0.786	0.867
Category 2	0.692	0.947	0.800	0.778	0.808	0.756
Category 3	1.000	0.385	0.822	0.250	0.200	0.844
Average	0.849	0.751	0.851	0.604	0.598	0.822
Kappa	0.649			0.512		

2