

1 **Development of a tool to accurately predict UK REF**
2 **funding allocation**

3

4 Shahd Al-Janabi¹, Lee Wei Lim², and Luca Aquili^{1,2*}

5 ¹College of Health & Human Sciences, Charles Darwin University, Darwin, Northern Territory,
6 Australia.

7 ²Neuromodulation Laboratory, School of Biomedical Sciences, Li Ka Shing Faculty of Medicine, The
8 University of Hong Kong, China.

9 *Corresponding author: Dr Luca Aquili. Email: luca.aquili@cdu.edu.au

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28 Keywords: REF; impact factor; metrics; funding

29

30 **Abstract**
31

32 Understanding the determinants of research funding allocation by funding bodies, such as the
33 Research Excellence Framework (REF) in the UK, is vital to help institutions prepare for
34 their research quality assessments. In these assessments, only publications ranked as 4* or 3*
35 (but not 2* or less) would receive funding. Correlational studies have shown that the impact
36 factor (IF) of a publication is associated with REF rankings. Yet, the precise IF boundaries
37 leading to each rank are unknown; for example, would a publication with an IF of 5 be
38 ranked 4* or less? Here, we provide a tool that **predicts the rank of each submitted**
39 **publication to (1) help researchers choose a publication outlet that would more likely lead to**
40 **the submission of their research output(s) by faculty heads in the next REF assessment,**
41 **thereby potentially improving their academic profile; and (2) help** faculty heads decide **which**
42 **outputs to submit for assessment, thereby** maximising their future REF scores **and ultimately**
43 **their research funding. Initially,** we applied our tool to the REF 2014 results for
44 Neuroscience, Psychiatry, and Psychology, which predicted publications ranked 4* with 95%
45 accuracy (IF ≥ 6.5), 3* with 98% accuracy (IF= 2.9-6.49), and 2* with 95% accuracy (IF=
46 1.3-2.89). **We then generalised these findings to another REF unit of assessment: Biological**
47 **Sciences to further demonstrate the predictive capacity of our tool.**

48

49

50

51

52

53

54

55

56 **Introduction**

57 The latest Research Excellence Framework (REF) exercise in 2014 distributed £1.6 billion
58 per year of research funding to universities and research institutions in the United Kingdom
59 (UK). In determining how funds are allocated, research outputs (65% weightage) submitted
60 to the REF are assessed by a panel of experts who assign a research quality grade or ranking
61 of 4* (world leading), 3* (internationally excellent), 2* (internationally recognised), 1*
62 (nationally recognised), or U/C (unclassified; below national standard) to each institution.
63 Importantly, only outputs judged to have a score of 4* or 3* are funded, with 4* (ranging
64 from £7504 to £14,639 depending on discipline) receiving four times as much money as 3*
65 (ranging from £1876 to £3659 depending on discipline) (Koya and Chowdhury 2017). The
66 method describe in this paper focused on the funding allocations based on submitted research
67 outputs. It should be noted that funding is also distributed on the basis of research impact
68 (20% weightage) defined by the REF as “an effect on, change or benefit to the economy,
69 society, culture, public policy or services, health, the environment or quality of life, beyond
70 academia”, which is a measure of Ph.D. completions, and also on the basis of research
71 infrastructure (i.e., laboratory facilities; cumulative 15% weightage). These weightages and
72 definition apply to both REF 2014 and the incoming REF 2021.

73

74 To gauge how well a university department might fare in the REF assessments, institutions
75 run mock REF exercises, in which research outputs are graded by other colleagues. These
76 mock assessments are time consuming (Farla and Simmonds 2015) and expensive (Jump
77 2015). Moreover, for academics who do not perform well in these internal mock assessments,
78 it can alarmingly lead to active research contracts turning into teaching-only or even potential
79 job losses (UCU 2013). As a result, there have been calls for the implementation of (costly)

80 alternatives to the peer-reviewed mock assessments (Stern and Nurse 2014), which are
81 primarily based on citation counts (Harnad 2009; Norris and Oppenheim 2003).

82
83 To date, some studies have shown that the H-index of a department correlates with the REF
84 rankings or REF GPA (grade point average) (Oppenheim 1995, 1997); however, other
85 investigations have not replicated these findings (Mryglod et al. 2015). Using a machine
86 learning approach, Balbuena (Balbuena 2018) demonstrated that the best predictors of REF
87 GPA were the number of Web Science documents, entry tariff, and the proportion of students
88 coming from independent schools. Additionally, Chowdhury et al. (Chowdhury et al. 2016;
89 Koya and Chowdhury 2017), reported a linear relationship between the number of **upper**
90 **quartile (Q1 or top 25%)** impact factor (IF) publications and the REF GPA. Although IF is
91 widely criticized by academics on numerous grounds (Callaway 2016; Hicks et al. 2015;
92 Paulus et al. 2018; Smaldino and McElreath 2016; Berenbaum 2019), unlike other metrics, it
93 has been shown to predict the future success of scientists (Acuna et al. 2012; Györfly et al.
94 2020; Van Dijk et al. 2014), and has been used by some universities to motivate staff to
95 publish in the most prestigious journals (e.g., through the provision of cash incentives for
96 publications with a high IF) (Abritis and McCook 2017; Quan et al. 2017). Furthermore, IF is
97 positively correlated with retraction rates (Fang and Casadevall 2011), possibly due to the
98 perceived benefits to one's academic standing driving publications in a high IF journal.
99 However, some of the above reviewed studies had limitations in that they examined variables
100 (e.g., H-index, entry tariff) that cannot be controlled by individual authors/faculties; hence,
101 knowledge of such relationships with REF GPA does not afford researchers any advantage.
102 Moreover, there is no correlational evidence between the controllable metrics (such as IF of
103 the journal) and REF GPA that allow the determination of the likely cut-offs leading to a 4*,
104 3*, or 2* rating.

105 In this research, we set out to describe a tool that predicts the rank of each submitted
106 publication. For authors, the purpose of this tool is to identify journals or, more precisely, the
107 impact factor associated with a journal that most likely leads to an award of a 4* or 3* rating
108 (for research outputs contributing 65% to the REF assessment). The importance of IF could
109 help individual authors to (1) plan research studies in such a way that they can publish in
110 outlets likely to result in a 4*/3* rating; and (2) be selected by their Faculty Head for the REF
111 assessment, thus enhancing their academic profile. For faculty heads, this tool would be most
112 beneficial in allowing them to strategically pick the research outputs of individual staff
113 members that have the best chance of being awarded a 4* or 3*, thereby potentially
114 increasing the funding allocation. Researchers and Faculty Heads already have a good sense
115 of what leads to a “good quality” paper; however, we outline here a novel, practical, and
116 accurate quantitative approach to define this “sense”. A discussion of the benefits and
117 limitations (methodological and philosophical) of this tool is also provided.

118 **Results**

119 **Accuracy of GPA prediction**

121 We assessed the accuracy of our tool by comparing our predicted GPAs with the REF 2014
122 GPAs for each institution. We compared the two sets of results for each institution by a
123 paired t-test (Fig. 1a). The analysis showed the actual REF 2014 GPAs ($M = 2.71$, $SD = 0.43$)
124 and our predicted GPAs ($M = 2.66$, $SD = 0.62$) were not statistically different ($t(80) = 1.43$, p
125 $= 0.16$). These results indicate our tool for assigning rank to each publication based on IF cut-
126 off values can accurately predict the ranking of publications, as our predicted GPAs for each
127 institution were comparable to the REF 2014 GPAs.

128

129

130

131 **Accuracy of the percentage of outputs per rank prediction**

132 To further check the accuracy of our tool, our predicted percentages of publications with
133 ranks ranging from 4* to U/C were compared with the percentage of publications awarded
134 each rank by REF 2014. First, we calculated the difference scores between each of the two
135 variables for each rank in each institution. Second, we conducted a one-way repeated
136 measures ANOVA with Bonferroni corrections on the mean difference scores of each rank
137 across institutions. Mauchly's Test of Sphericity was violated; hence, Lower-bound
138 corrections were used. We found no significant differences between rank categories ($F(1, 80)$
139 $= 2.93, p = 0.09$), which indicates the accuracy of our predicted percentage of outputs was
140 similar across ranks compared to the actual REF 2014 percentage of outputs (Fig. 1c). The
141 accuracy of our tool in predicting the percentage of outputs can be seen in Figure 1b. This
142 accuracy was calculated by summing the percentage of publications in each rank as identified
143 by the REF 2014 versus our tool. The error rate was calculated as the difference between the
144 two sums and was subtracted from the true positive (sum of the percentage of publications in
145 each rank identified by REF 2014) as an indication of our hit rate. The percentage accuracy
146 was then calculated using our hit rate divided by the true positive.

147

148 **Accuracy of rankings prediction**

149 We ordered the predicted GPAs to identify the predicted (output-based) ranking of each
150 institution. With this information, we conducted a Wilcoxon signed-rank test to compare the
151 REF 2014 output-based institution rankings with our predicted rankings to ascertain the
152 accuracy of our tool. The test showed our predicted institution rankings did not statistically
153 differ from that of the REF 2014 institution rankings ($Z = -0.082, p = 0.94$). The level of
154 agreement between the institution rankings derived from our IF-based tool versus the REF
155 2014 results was further confirmed by the Bland-Altman plot (Fig. 2). Specifically, the

156 Bland-Altman plot revealed no systematic differences or consistent bias between the
157 institution rankings derived using our tool versus the REF 2014 results. Indeed, few values
158 (<4%) lay outside of the limits of agreement (mean of differences ± 1.96 SD), indicating
159 agreement between the two measurements.

160

161 **Strength of GPA predictions**

162 We conducted two regression analyses to ascertain the extent to which our IF-based tool can
163 predict REF 2014 quality profile metrics; specifically, GPA. The first simple linear regression
164 analysis assessed whether our predicted GPA can be used to predict actual REF 2014 GPAs.
165 We found that our tool could accurately predict the GPA compared to the actual GPA
166 (unstandardized $\beta = 0.62$, $p < 0.0001$; overall model fit of $R^2 = 0.80$) (Fig. 3). This effect
167 remained even after conducting a hierarchical regression that accounted for each institution's
168 total output and 2013 University ranking. The model predicting GPA was the only
169 statistically significant model ($p < 0.0001$ vs. $p > 0.09$ for all other models), indicating our
170 tool based on IF cut-off values can predict GPA in the REF 2014.

171

172 **Generalisability of tool**

173 We examined the generalisability of our tool by applying it to a different area evaluated in the
174 REF 2014 Unit of Assessment: Biological Sciences. Such an investigation should confirm the
175 predictive value of our tool. First, we assessed the accuracy of our tool by comparing our
176 predicted GPAs with the REF 2014 GPAs for each institution who submitted publications to
177 the Biological Sciences Unit of Assessment. As above, we compared the two sets of results
178 for each institution by a paired t-test (Fig. 5a). The analysis showed the actual REF 2014
179 GPAs ($M = 2.89$, $SD = 0.40$) and our predicted GPAs ($M = 2.93$, $SD = 0.43$) were not
180 statistically different ($t(43) = 1.50$, $p = 0.14$). The accuracy of our tool in predicting the

181 percentage of outputs in Biological Sciences can be seen in Figure 5b. Second, we conducted
182 a simple linear regression analysis to probe whether or not our tool could accurately predict
183 the GPAs compared to the actual GPAs of institutions who submitted to the Biological
184 Sciences Unit of Assessment. This analysis showed our method was accurate (unstandardized
185 $\beta = 0.86, p < 0.0001$; overall model fit of $R^2 = 0.85$) (Fig. 6). This effect remained even after
186 conducting a hierarchical regression that accounted for each institution's total output and
187 2013 ranking. As before, the model predicting GPA was the only statistically significant
188 model ($p < 0.0001$ vs. $p > 0.14$ for all other models). These results collectively indicate that
189 our tool can accurately predict GPA in the REF 2014 for another Units of Assessment besides
190 Neuroscience, Psychiatry, and Psychology.

191 **Discussion**

192
193

194 The primary goal of this study was to identify the IF threshold values for predicting the
195 different REF rankings and GPAs. These findings have the potential to be used by individual
196 authors and faculty heads to make informed decisions on (1) where to submit and publish
197 articles; and (2) which research outputs should be submitted to the next REF assessment to
198 maximise funding allocation. We first briefly summarise the key findings and then provide an
199 in-depth interpretation of these results and their implications.

200

201 Our analysis of the REF 2014 Unit of Assessment: Neuroscience, Psychiatry, and Psychology
202 demonstrated accuracy scores in excess of 95% for predicting the percent of publications
203 ranked 4*, 3*, or 2*. Specifically, we revealed that to receive the highest possible score of
204 4*, a publication submitted to REF must have an IF equal to or greater than 6.5. We
205 estimated with 98% accuracy that a publication with an IF of 2.9 or less receiving a rank of
206 2* or below was unlikely to receive funding. We replicated these findings in the REF unit of

207 assessment: Biological Sciences. Although the IF boundaries changed, the accuracy of the
208 model predictions for publications ranked 4*, 3*, or 2* remained largely similar. Our
209 findings significantly expand on previous correlational studies by showing that the IF of
210 publications can be used to predict REF ranking and GPA (Chowdhury et al. 2016; Koya and
211 Chowdhury 2017). We have developed a precise numerical tool to predict the research output
212 quality score based on REF 2014 results.

213

214 Our estimates for the IF cut-offs are based on the total GPA achieved by each institution and
215 the corresponding percentages of 4*, 3*, 2*, 1*, and U/C as assigned in the REF 2014
216 published by the *Times Higher Education*. Figure 1a shows the similarity of our mean
217 estimated REF GPAs for all universities (2.66) compared to the actual REF results (2.71);
218 Figure 1b shows the similarities of our estimated percentages of research outputs with 4*, 3*,
219 2*, 1*, and U/C scores versus the assigned scores, and Figure 4 shows the similarities
220 between the university REF GPA rankings and our predictions, overall confirming the
221 accuracy of our estimated IF cut-offs used in the REF 2014 Unit of Assessment. It is
222 important to note, however, that our IF cut-offs are not predictions at the individual output
223 level, as this information is not publicly available. The report commissioned by the Higher
224 Education Funding Council for England (HEFCE) is the only study that has attempted to
225 match individual research outputs to REF scores (Wouters et al. 2015). The authors
226 concluded that “The statistics presented do not overwhelmingly support the use of metrics as
227 a replacement for a peer-review-driven model of research quality assessment”. However, it is
228 impossible to further comment on their findings, as the methodological details in their report
229 were scarce and IF was not used as a metric to derive prediction cut-offs.

230

231 Interestingly, our hierarchical regression analyses showed that both the total number of
232 research outputs and university rankings (aka: league tables, as determined by The Guardian
233 University guide and ranking for 2013) were not significant predictors of REF GPA. These
234 results suggest that the evaluators in the assessments, such as REF 2014, primarily (if not
235 solely) rely on the journal IF to make their decisions on the rank (4* to U/C) of a submitted
236 publication.

237

238 It is important to acknowledge the current study has limitations and the results should be
239 considered with a note of **methodological and philosophical** caution. The specific IF
240 boundaries established in this study are limited to predicting the REF rankings from **two**
241 **Units of Assessment: (1) Neuroscience, Psychiatry, and Psychology, and (2) Biological**
242 **Sciences; thus**, future studies are needed to confirm the **accuracy** of our tool when applied to
243 other fields. **Nevertheless**, given that other disciplines, such as Clinical Medicine,
244 Agriculture, Veterinary and Food Science, Chemistry, Economics, and Econometrics, have
245 been reported to have similar correlations between REF quality and SCImago (an alternative
246 measure of The Journal Citation Reports IF, which is highly correlated with IF; Rocha-e-
247 Silva 2010) to that in Neuroscience, Psychiatry, and Psychology **and Biological Sciences** (R
248 ≥ 0.437 ; (Wouters et al. 2015), **we expect our tool will have predictive capacity when also**
249 **applied to these fields. It is important to stress that our tool may be less precise for some**
250 **Units of Assessment (e.g., Law, Business and Management Studies, Arts and Design, Social**
251 **Work and Social Policy, etc.) where a predominant proportion of the assessed research**
252 **outputs were not from journal articles (e.g., from chapters in a book, conference**
253 **contributions, patents, exhibitions, etc.), which do not have an IF.**

254 It should also be noted that our IF cut-offs were determined post-hoc (i.e., after the REF 2014
255 data was made available). The next REF assessment is due to take place this year (2021), and

256 our IF cut-offs will need to be adjusted to take into account any increases in journal IF values
257 since 2013 (Althouse et al. 2009). The median IF for Neuroscience-related journals increased
258 by 22% from 2007 to 2013, whereas the increase from 2013 to 2020 has so far been about
259 5%. Additionally, individual researchers/faculties using our tool need to take the 2020 IF
260 values corresponding to the bottom publications in each rank (4* - U/C) and recreate the
261 confidence intervals (CI) around the new mean IF of each rank across institutions. It is likely
262 that these changes will not grossly affect the accuracy of our model¹. The expectation that our
263 tool will continue to be generally accurate relies on two assumptions. First, there should be a
264 high correlation in the proportions of 4* and 3* assigned to each institution between REF
265 2014 and REF 2021, as we showed in our analyses of the Research Assessment Exercise
266 (RAE) 2008 and REF 2014². The pattern of results indicates that the number of outputs
267 receiving each rank may be stable over time. Second, the distribution of REF quality ratings
268 should remain largely similar to that in 2021. This was not the case between RAE 2008 and
269 REF 2014, where the number of 4* and 3* ratings doubled in 2014; however, this difference
270 has been widely acknowledged as an inflation of grades (Marginson 2014). Our IF cut-offs
271 for the next REF assessment (2021) also need to take into account any new REF 2021
272 submission rules different from REF 2014 (e.g., number of outputs per individual (REF

1 To test the accuracy of our model with different IF cut-offs, under the assumption that the IFs of publications submitted to exercises, such as the REF 2014, may change over time, we examined the accuracy of our model with cut-offs that are 10% greater and 10% lower than the cut-offs (lower-bound confidence interval) used in this study. The variations in our IF cut-offs indicate that our model remained accurate, particularly for the first two ranks. For 10% greater cut-offs: 84% accuracy for 4*, 91% for 3*, 89% for 2*, 86% for 1*, and 21% for U/C; and for 10% lower cut-offs: 75% accuracy for 4*, 95% for 3*, 78% for 2*, 50% for 1*, and 24% for U/C. This analysis indicated our tool can predict possible REF rankings and GPAs with changing IF cut-offs.

2 We extracted the percentage of publications awarded a 4* and 3* during RAE 2008 for the Unit of Assessment: Psychology (note that Psychiatry and Neuroscience were measured separately in very few entries). With this information, we conducted a correlational analysis between the percentage outputs for RAE 2008 and REF 2014 for each rank. This analysis revealed a medium-to-large correlation between the percentage outputs receiving a 4* ($r = 0.78, p < 0.0001$) and 3* ($r = 0.50, p < 0.0001$) in the 2008 versus 2014 exercises, suggesting the number of outputs receiving each rank may be stable over time.

273 2019)). As our tool was based on REF 2014, any predictions of results not yet published are
274 bound to be less precise.

275 Whilst our tool was highly accurate (i.e., above 90%) in identifying the ranks of 2*, 3*, and
276 4*, it was less accurate for U/C and 1* ranks. We can provide some plausible explanations as
277 to why this pattern of results may have occurred. First, we automatically assigned a rank of
278 U/C to any research article which we could not derive an IF (see Methods for the derivation
279 of an IF for each publication). This strategy likely resulted in over assigning U/C ranks when
280 a panel of reviewers would have judged the publications differently. Take for example the
281 case of a newly formed journal from an established publishing group (e.g., *Nature*,
282 *Frontiers*). Because these publishing groups are well-known by academics, it is foreseeable
283 that a panel of reviewers would have assigned a higher score than U/C to a publication e.g., in
284 *Nature Human Behaviour* and/or *Frontiers in Physiology*, despite them not yet having an IF
285 or a stable IF pattern. Second, the predicted IF range for an output classified as U/C and 1*
286 was much smaller than for 2*, 3* and 4* ranks (see Fig. 1 and Fig. 5). This pattern of results
287 may lead to minor IF deviations that result in underclassifying and/or overclassifying a
288 research output to one rank versus another. Third and related to the second point, we can
289 speculate that when the IF of a journal is very low, reviewers who may not be familiar with
290 these journals may use other heuristics that are as salient as IF for determining the rank of a
291 paper, such as the institutional ranking of the corresponding author. Lastly, the sample size
292 for U/C and 1* ranks was between one tenth and one twentieth of 2*, 3*, and 4* ranks, which
293 likely increased the error rate. Indeed a previous study on the Finnish ranking of research
294 publications found that misclassification/model errors occurred even when citation-based
295 metrics predicted most of the expert-based rankings for papers (Saarela et al. 2016). Despite
296 these methodological limitations, it is important to note that the practical utility of our tool

297 remains largely unaffected, given that we identified that only the outputs assigned a value of
298 3* and 4* will receive funding under the REF scheme.

299 Our analyses indicate that funding bodies, such as REF 2014 in the UK, may base their
300 allocation decisions on the IF of the outputs submitted for assessment by each institution.
301 Given the speed in which these decisions must be made, this strategy makes sense insofar as
302 journals with higher IFs are more visible, hence, one can use this metric as a ‘quick and dirty’
303 index for (a) the importance of each output to a wide readership; and (b) the (world-class)
304 success of the researchers/institution of each output. If our assumption are indeed true then an
305 obvious outcome, as previously stated, is that it may be advantageous for researchers to
306 publish in higher IF journals and/or for Faculty Heads to submit only high IF outputs for
307 assessment in order to increase funding allocation to the institution. Although this strategy is
308 *implied* by our findings, our view is that researchers, Faculty Heads, and funding bodies
309 should consider the ethics of using IF as a prime measure of deciding where to publish or for
310 allocating funding, primarily because IF was conceived as a metric of journal usage not
311 author scholarship. The formula for IF is the number of cited articles published within a given
312 period divided by the total number of citable outputs published within that period (Jones
313 2013; Marson 2020). This formula (a) does not index the way in which an article was read or
314 used following publication, thus obscuring the assessment of the publication’s impact or
315 importance; (b) does not index the actual merit of a publication in that journal because the
316 increased citation could be due to weaknesses and flaws in the publication as opposed to its
317 strengths, thus obscuring the assessment of the researcher’s success; (c) does not index
318 publications that have been cited in text books and not in or in addition to an indexable
319 journal, thus obscuring both the above assessments; and (d) does not have high reliability
320 (Greenwood 2007), though it may have high validity in some fields (Saha et al. 2003, also see
321 Jarwal et al. 2009; Law and Leung 2020)). In listing the limitations of applying IF as a

322 metric, we seek to highlight the ethical issue(s) that may arise from its use in strategy,
323 planning, and decision-making, as per the implications of our findings. In the future, the
324 increasing number of easily accessible tools that chart bibliometrics may encourage funding
325 bodies to make more use of an output's citation count as opposed to IF in decision-making
326 (Marson 2020).

327

328 In summary, our results provide a simple practical tool that can help individual researchers
329 and/or heads of departments/faculties in the Unit of Assessment: Neuroscience, Psychiatry,
330 and Psychology and Biological Sciences to accurately estimate the UK REF quality score
331 based on the IF of their publications. Although our findings were based on REF 2014, our
332 additional analyses ¹ demonstrate that the tool likely has predictive capabilities for REF 2021
333 (despite the approaching REF 2021 submission deadline). Our findings largely agree with
334 previous literature on the association between various publication metrics and scientific
335 success/research funding, but further expands on these studies by providing a more cost and
336 time effective approach for evaluating research papers in the context of REF compared to
337 mock peer-reviews (Farla and Simmonds 2015; Jump 2015). It is possible that this tool could
338 be applied to other research funding allocation exercises (e.g., ERA Research Assessment in
339 Australia) with similar utility after some modifications.

340

341 ***Methods***

342 **Data sources and variables**

343 The primary data for this study was the REF 2014 results (Research Excellence Framework
344 2014) in the UK. As proof of concept, we focused on the Units of Assessment: Neuroscience,
345 Psychiatry, and Psychology, and Biological Sciences to assess the accuracy of our tool in
346 predicting REF rankings. These units were chosen because 99% of their outputs were

347 research articles (and hence likely to have an associated IF). There were 81 institutions who
348 submitted publications for the Unit of Assessment: Neuroscience, Psychiatry, and Psychology
349 totalling a combined 9,121 publications. On the other hand, there were 44 institutions who
350 submitted publications for the Unit of Assessment: Biological Sciences totally a combined
351 8,608 publications. For each Unit of Assessment, we extrapolated the output and quality
352 profiles from the REF 2014 dataset: each institution name, the journals of the submitted
353 publications, and the evaluations for each institution, including the percentage of publications
354 awarded each rank (4* to U/C) at the culmination of the REF 2014 assessment. We added
355 four variables to this extrapolated dataset. The first variable was the total number of outputs
356 (i.e., publications) submitted by each institution to the Unit of Assessment for evaluation. The
357 second variable was the IF of each journal as of 2013 (sourced from Journal Citation Reports
358 2014). Publications for which a 2013 IF could not be sourced were either assigned the IF for
359 the preceding or following year. We were unable to source an IF in 1.29% of publications for
360 the Unit of Assessment: Neuroscience, Psychiatry, and Psychology, and 1.20% of
361 publications for the Unit of Assessment: Biological Science. The third variable added was the
362 Grade Point Average (GPA) awarded to each institution for each Unit of Assessment
363 (Research Excellence Framework 2014: Institutions Ranked by Subject 2014), which is
364 calculated by multiplying the percentage of publications in each rank by its rating (adding the
365 total across all ranks and dividing by 100), giving a GPA index of an institution's overall
366 quality of research. Finally, the fourth variable added was the ranking of each institution in
367 2013 (University guide 2013: University league table 2013).

368 **Calculation of the prediction cut-off values**

369 We first defined the IF boundaries or cut-off values for each of the REF ranks in each Unit of
370 Assessment using the 2014 dataset. We calculated the number of publications assigned each
371 rank of 4*, 3*, 2*, 1*, and U/C for each institution in each Unit of Assessment. Next, we

372 ordered the publication ranks by their IF for each institution **in each Unit of Assessment** and
373 identified the IF value for the publication that was likely to be at the bottom of each rank in
374 each institution. The IF cut-offs were then identified by calculating the number of
375 publications assigned a rank from 4* to U/C in each institution. For example, for a given
376 institution with six publications assigned a 4* by REF 2014, we ranked these submitted
377 publications by their IF and identified the publication with the lowest IF for that institution **in**
378 **that rank**. We then calculated the mean IF for each rank across institutions **in each Unit of**
379 **Assessment** with their corresponding confidence intervals. The lower bound confidence
380 interval served as our prediction cut-off **for each rank in each Unit of Assessment** (see Fig. 1b
381 and Fig. 5b).

382 **Predicting the ranking of publications and GPA**

383
384 Using IF cut-off values **for each Unit of Assessment**, we aimed to predict the rank of each
385 publication submitted to REF 2014. For example, all publications with an impact factor equal
386 to or above the lower bound confidence interval for a 4* (**e.g., 6.54 for Neuroscience,**
387 **Psychiatry, and Psychology**) were assigned a rank of 4*, whereas all publications with an
388 impact factor less than the lower bound confidence interval for a 1* (**e.g., 0.71 for**
389 **Neuroscience, Psychiatry, and Psychology**) were assigned a rank of U/C. Next, we calculated
390 the percentage of publications in an institution that we predicted to be awarded each rank.
391 With that information, we also calculated the predicted GPA for each institution. This
392 calculation used the same formula described above (see Fig. 4 for comparisons).

393

394

395

396

397

398

399 **Figure 1. Panel A** shows the mean GPA awarded by REF 2014 versus the predicted mean GPA determined by
400 our tool **for the Unit of Assessment: Neuroscience, Psychiatry, and Psychology**. **Panel B** shows the percentage
401 of publications awarded each rank across institutions by REF 2014 versus our predicted percentage of
402 publications for **the Unit of Assessment: Neuroscience, Psychiatry, and Psychology** using our tool based on
403 impact factor cut-off values. We used those cut-off values to predict the number of publications receiving each
404 rank from 4* to U/C according to the REF 2014 Unit of Assessment: Neuroscience, Psychiatry, and Psychology.
405 **Panel C** shows the difference score between REF 2014 and our predicted percentage of publications in each
406 rank for **the Unit of Assessment: Neuroscience, Psychiatry, and Psychology** to validate the accuracy of the tool.
407 All error bars reflect standard error of the mean.

408

409 **Figure 2.** Bland-Altman plot for institution rankings derived from our tool versus the REF 2014 results **for the**
410 **Unit of Assessment: Neuroscience, Psychiatry, and Psychology**. The centre line indicates the mean difference,
411 whereas the lines on either side are the limits of agreement ($\text{mean} \pm 1.96 \text{ SD}$).

412

413 **Figure 3.** The relationship between the predicted GPA using our tool and REF 2014 GPA for each institution
414 **that submitted publications to the Unit of Assessment: Neuroscience, Psychiatry, and Psychology**. The blue line
415 is the best fit ($r = 0.90$) and green dotted lines are the upper and lower confidence interval boundaries.

416

417 **Figure 4.** The GPA awarded by REF 2014 (in green) and the predicted GPA by our tool (in blue) for each
418 institution **that submitted publications to the Unit of Assessment: Neuroscience, Psychiatry, and Psychology**.
419 The GPA is used to assess the quality of research published in the fields of Neuroscience, Psychiatry, and
420 Psychology that takes into account the percentage of outputs in each institution assigned a rank from 4* to U/C.
421 Panel A represents institutions that are in the upper GPA quartile of REF 2014, panels B and C represent
422 institutions in the middle quartiles (2nd and 3rd, respectively), and panel D represents institutions in the lower
423 GPA quartile.

424

425 **Figure 5.** shows the mean GPA awarded by REF 2014 versus the predicted mean GPA determined by our tool
426 **for the Unit of Assessment: Biological Sciences**. **Panel B** shows the percentage of publications awarded each
427 rank across institutions by REF 2014 versus our predicted percentage of publications in the Unit of Assessment:

428 Biological Sciences using our tool based on impact factor cut-off values. We used those cut-off values to predict
429 the number of publications receiving each rank from 4* to U/C according to the REF 2014 Unit of Assessment:
430 Biological Sciences.

431

432 **Figure 6.** The relationship between the predicted GPA using our tool and REF 2014 GPA for each institution
433 that submitted publications to the Unit of Assessment: Biological Sciences. The blue line is the best fit ($r = 0.92$)
434 and green dotted lines are the upper and lower confidence interval boundaries.

435

436

437 ***Declarations***

438 ***Funding***

439 The study received no funding.

440 ***Competing interests***

441 The authors declare no competing interests.

442

443 ***Availability of data and material***

444 Data available from public resources (REF 2014, The Guardian University Guide) and the

445 Journal of Citation Reports.

446 ***Code availability***

447 N/A

448 ***Author contributions***

449 LA conceived the study; SA and LA developed the methodology and analysed the data; SA,

450 LWL and LA wrote the manuscript.

451

452

453

454 **REFERENCES**

455

456 Abritis, A., & McCook, A. (2017). Cash incentives for papers go global. *Science*,
457 357(6351), 541-541, doi:10.1126/science.357.6351.541.

458 Acuna, D. E., Allesina, S., & Kording, K. P. (2012). Predicting scientific success.
459 *Nature*, 489(7415), 201-202, doi:10.1038/489201a.

460 Althouse, B., West, J., Bergstrom, T., & Bergstrom, C. (2009). Differences in Impact
461 Factor Across Fields and Over Time. *Journal of the American Society for*
462 *Information Science and Technology*, 60, doi:10.1002/asi.20936.

463 Balbuena, L. D. (2018). The UK Research Excellence Framework and the Matthew
464 effect: Insights from machine learning. *PLOS ONE*, 13(11), e0207919,
465 doi:10.1371/journal.pone.0207919.

466 Berenbaum, M. R. (2019). Impact factor impacts on early-career scientist careers.
467 *Proceedings of the National Academy of Sciences*, 116(34), 16659-16662,
468 doi:10.1073/pnas.1911911116.

469 Callaway, E. (2016). Beat it, impact factor! Publishing elite turns against
470 controversial metric. *Nature News*, 535(7611), 210.

471 Chowdhury, G., Koya, K., & Philipson, P. (2016). Measuring the impact of research:
472 lessons from the UK's Research Excellence Framework 2014. *PLOS ONE*,
473 11(6), e0156978.

474 Fang, F. C., & Casadevall, A. (2011). Retracted Science and the Retraction Index.
475 *Infection and Immunity*, 79(10), 3855-3859, doi:10.1128/iai.05661-11.

476 Farla, K., & Simmonds, P. (2015). REF 2014 accountability review: Costs, benefits
477 and burden. *Brighton, UK: Technopolis/ group*. Retrieved October, 1, 2019.

478 Greenwood, D. C. (2007). Reliability of journal impact factor rankings. *BMC medical*
479 *research methodology*, 7(1), 1-6.

480 Györfly, B., Herman, P., & Szabó, I. (2020). Research funding: past performance is a
481 stronger predictor of future scientific output than reviewer scores. *Journal of*
482 *Informetrics*, 14(3), 101050, doi:<https://doi.org/10.1016/j.joi.2020.101050>.

483 Harnad, S. (2009). Open access scientometrics and the UK Research Assessment
484 Exercise. *Scientometrics*, 79(1), 147-156.

485 Hicks, D., Wouters, P., Waltman, L., De Rijcke, S., & Rafols, I. (2015).
486 Bibliometrics: the Leiden Manifesto for research metrics. *Nature*, 520(7548),
487 429-431.

488 Jarwal, S. D., Brion, A. M., & King, M. L. (2009). Measuring research quality using
489 the journal impact factor, citations and 'Ranked Journals': Blunt instruments
490 or inspired metrics? *Journal of Higher Education Policy and Management*,
491 31(4), 289-300.

492 Jones, J. (2013). The impact of impact factors and the ethics of publication. Springer.
493 Journal Citation Reports (R) (2014). In T. Reuters (Ed.).

494 Jump, P. (2015). Winners and losers in HEFCE funding allocations. *Times Higher*
495 *Education*, 26, 6-9.

496 Koya, K., & Chowdhury, G. (2017). Metric-based vs peer-reviewed evaluation of a
497 research output: Lesson learnt from UK's national research assessment
498 exercise. *PLOS ONE*, 12(7), e0179722.

499 Law, R., & Leung, D. (2020). Journal impact factor: A valid symbol of journal
500 quality? *Tourism Economics*, 26(5), 734-742.

501 Marginson, S. (2014). UK research is getting better all the time—or is it? *The*
502 *Guardian*.

503 Marson, S. (2020). Editorial: Is the Impact Factor (IF) Ethical To Use for

504 Promotion and Tenure Decisions? *Journal of Social Work Values and Ethics*, 17, 2-5.
505 Mryglod, O., Kenna, R., Holovatch, Y., & Berche, B. (2015). Predicting results of the
506 Research Excellence Framework using departmental h-index. *Scientometrics*,
507 102(3), 2165-2180.

508 Norris, M., & Oppenheim, C. (2003). Citation counts and the research assessment
509 exercise V. *Journal of Documentation*.

510 Oppenheim, C. (1995). The correlation between citation counts and the 1992
511 Research Assessment Exercise Ratings for British library and information
512 science university departments. *Journal of Documentation*.

513 Oppenheim, C. (1997). The correlation between citation counts and the 1992 research
514 assessment exercise ratings for British research in genetics, anatomy and
515 archaeology. *Journal of Documentation*.

516 Paulus, F. M., Cruz, N., & Krach, S. (2018). The impact factor fallacy. *Frontiers in*
517 *psychology*, 9, 1487.

518 Quan, W., Chen, B., & Shu, F. (2017). Publish or impoverish. *Aslib Journal of*
519 *Information Management*.

520 REF (2019). Guidance on Submissions. Research England Bristol.
521 Research Excellence Framework (2014). <http://www.ref.ac.uk/2014/>.

522 Research Excellence Framework 2014: Institutions Ranked by Subject (2014).
523 [https://www.timeshighereducation.com/sites/default/files/Attachments/2014/1](https://www.timeshighereducation.com/sites/default/files/Attachments/2014/12/17/g/o/1/sub-14-01.pdf)
524 [2/17/g/o/1/sub-14-01.pdf](https://www.timeshighereducation.com/sites/default/files/Attachments/2014/12/17/g/o/1/sub-14-01.pdf).

525 Rocha-e-Silva, M. (2010). Impact factor, scimago indexes and the brazilian journal
526 rating sytem: where do we go from here? *Clinics*, 65(4), 351-355.

527 Saarela, M., Kärkkäinen, T., Lahtonen, T., & Rossi, T. (2016). Expert-based versus
528 citation-based ranking of scholarly and scientific publication channels.
529 *Journal of Informetrics*, 10(3), 693-718.

530 Saha, S., Saint, S., & Christakis, D. A. (2003). Impact factor: a valid measure of
531 journal quality? *Journal of the Medical Library Association*, 91(1), 42.

532 Smaldino, P. E., & McElreath, R. (2016). The natural selection of bad science. *Royal*
533 *Society open science*, 3(9), 160384.

534 Stern, N., & Nurse, P. (2014). It's our duty to assess the costs of the REF. *Times*
535 *Higher Education*. <https://www.timeshighereducation.com/comment/letters/its-our-duty-to-assess-the-costs-of-the-ref/2017479>.
536 *article*.

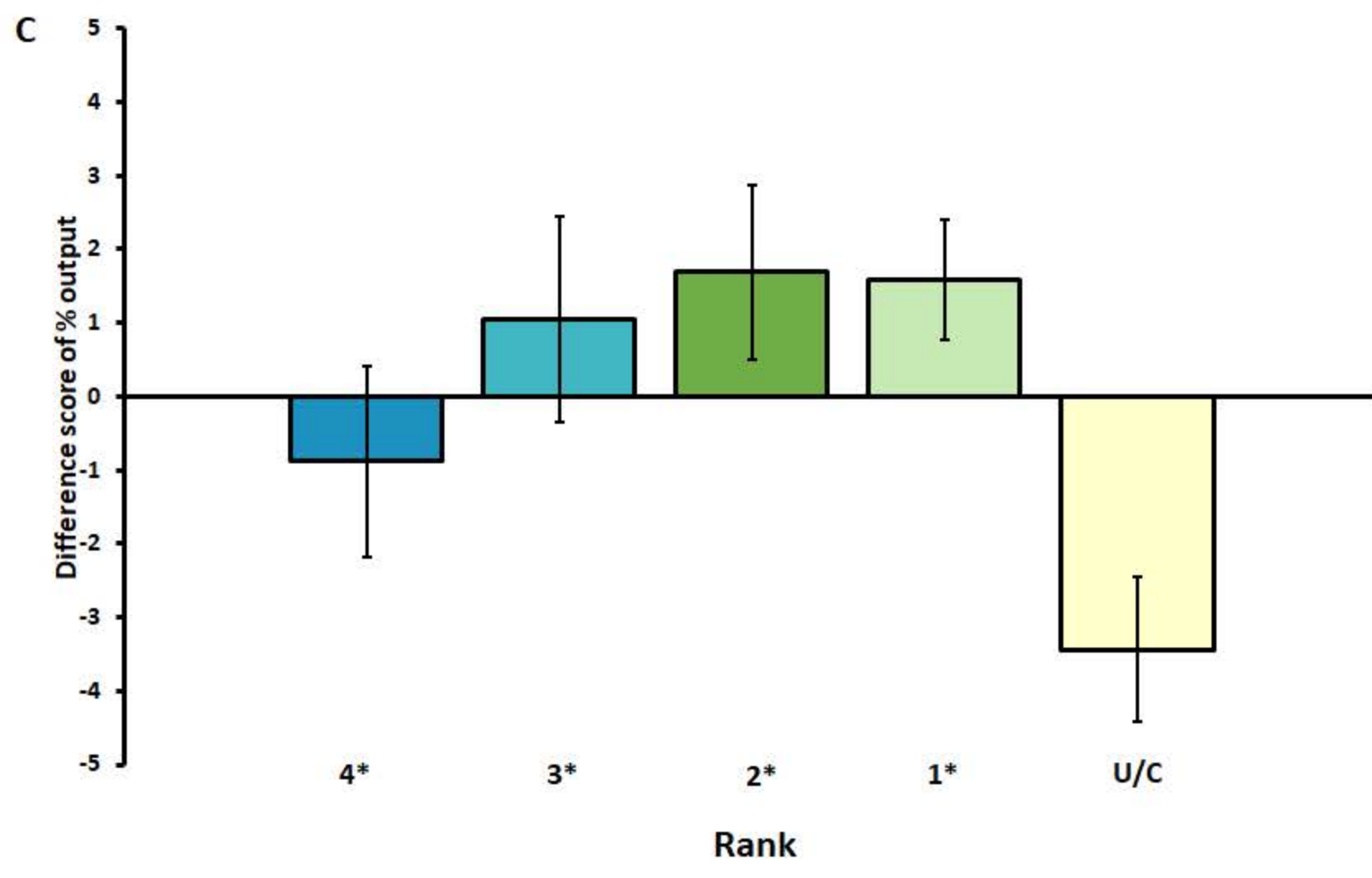
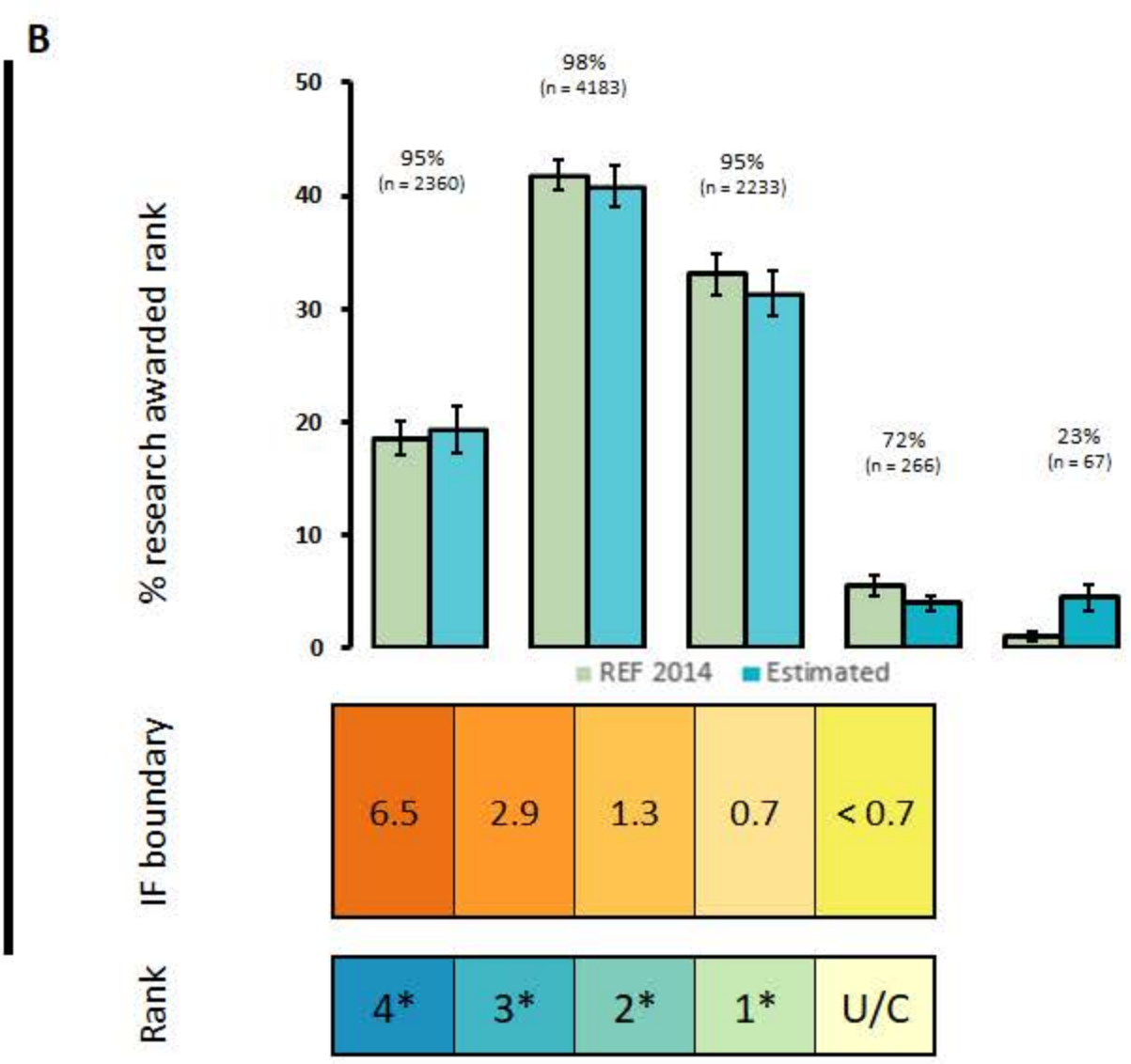
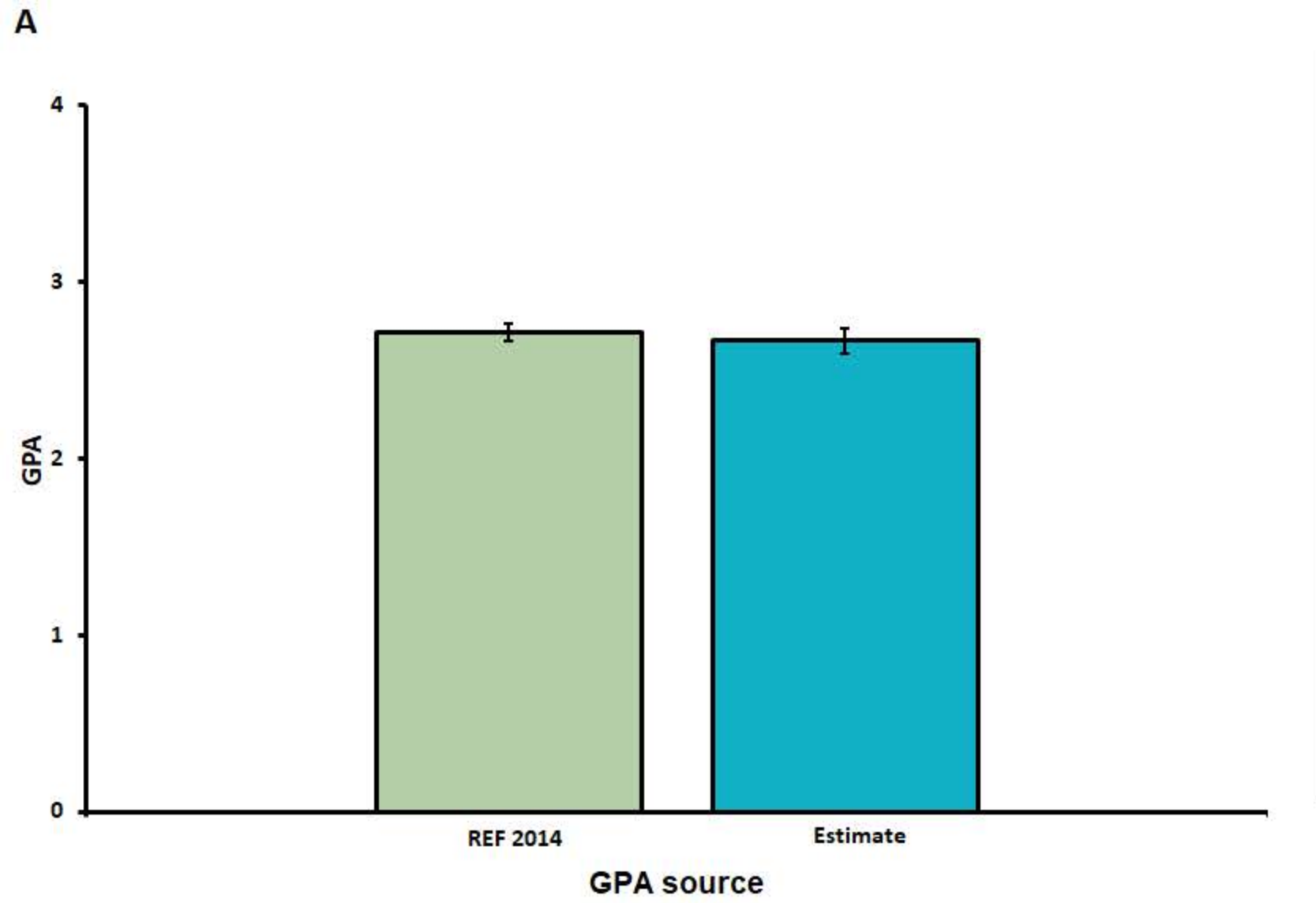
537 UCU (2013). The Research Excellence Framework UCU Survey Report. UCU
538 London.

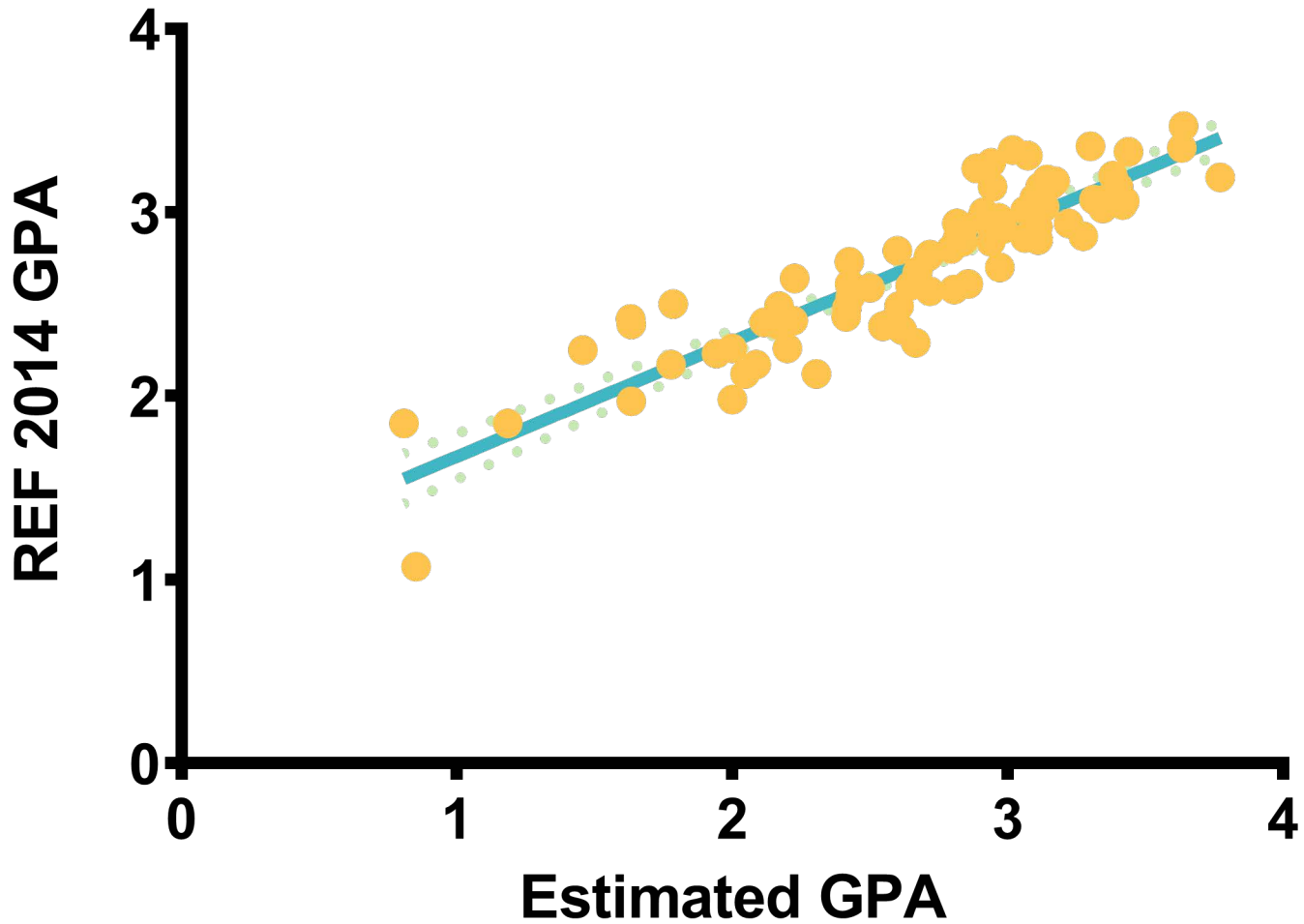
539 University guide 2013: University league table (2013).
540 [https://www.theguardian.com/education/table/2012/may/21/university-league-](https://www.theguardian.com/education/table/2012/may/21/university-league-table-2013)
541 [table-2013](https://www.theguardian.com/education/table/2012/may/21/university-league-table-2013).

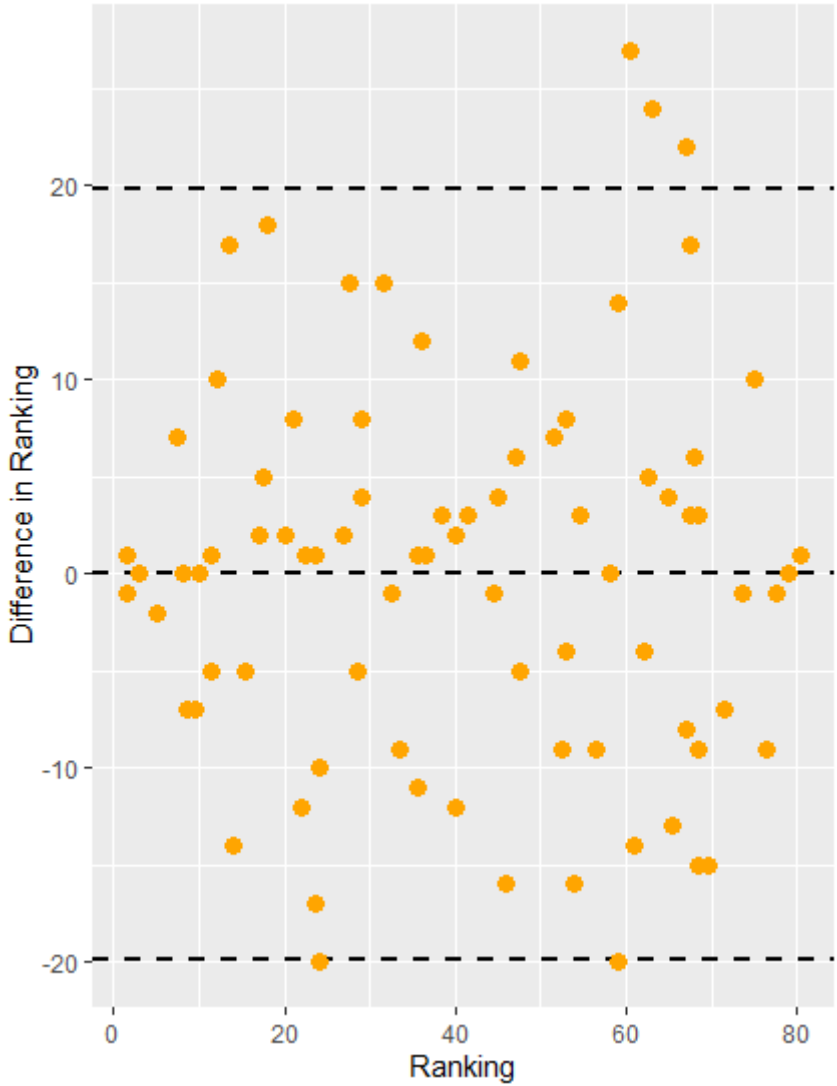
542 Van Dijk, D., Manor, O., & Carey, L. B. (2014). Publication metrics and success on
543 the academic job market. *Current Biology*, 24(11), R516-R517.

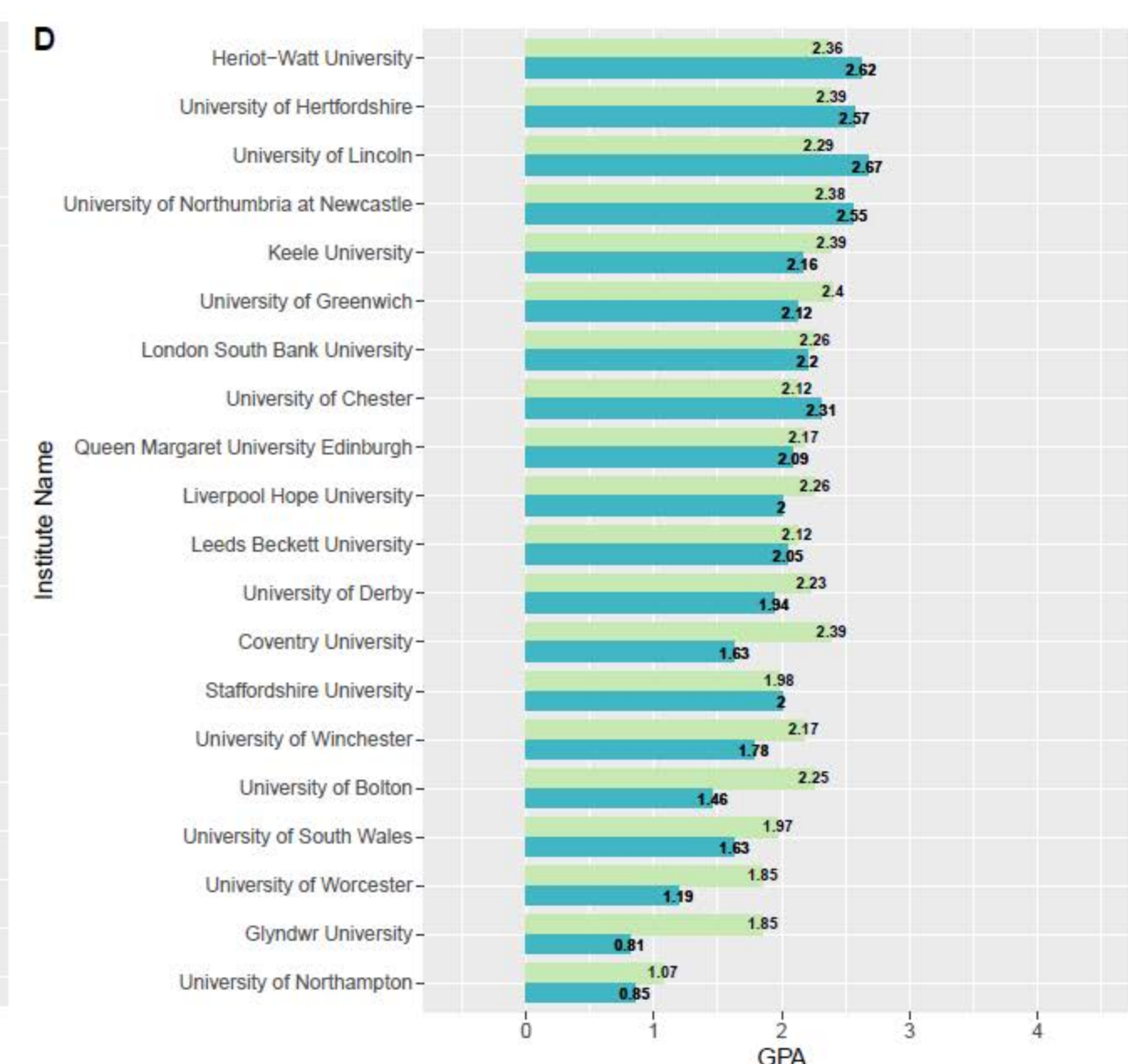
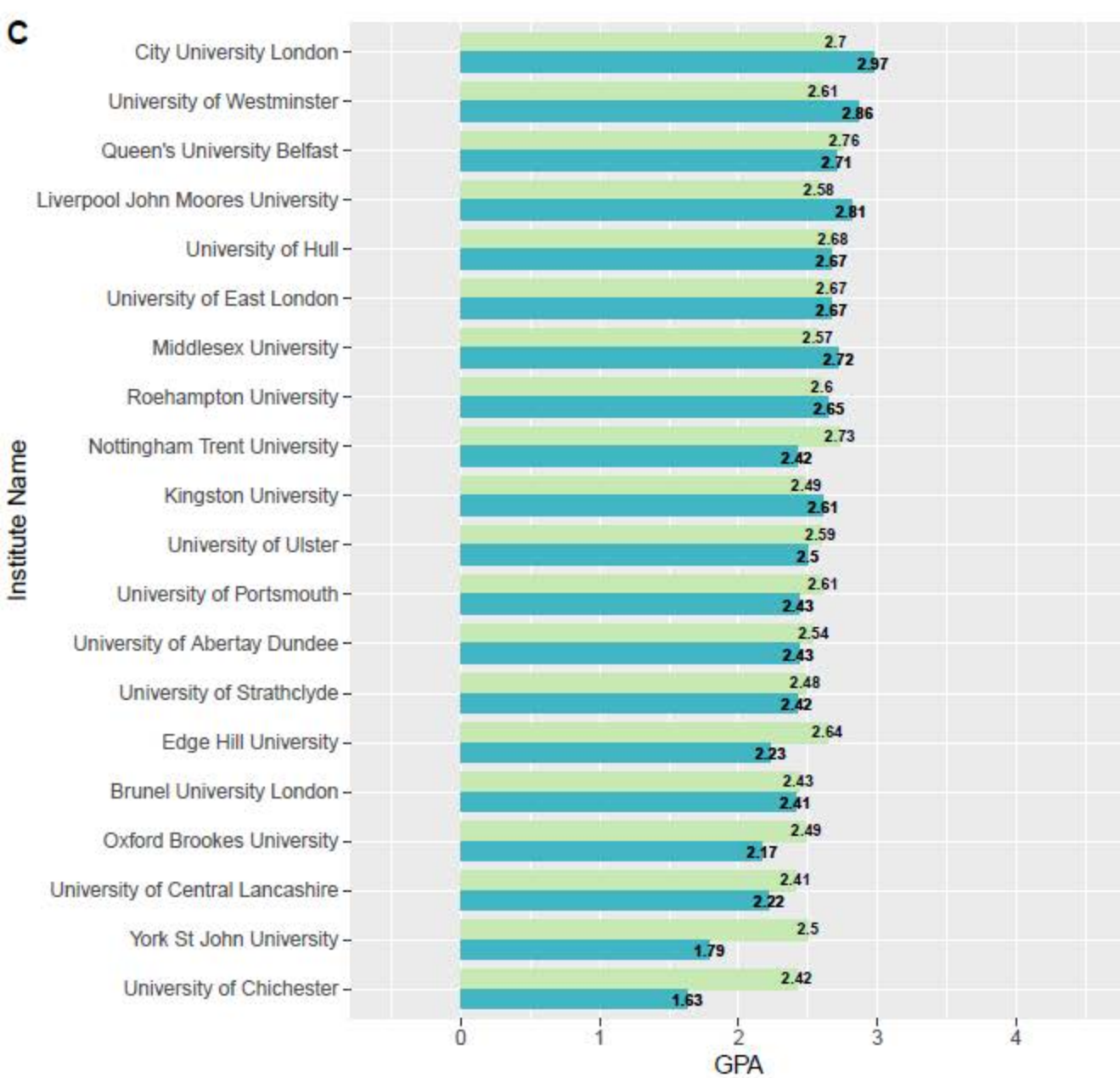
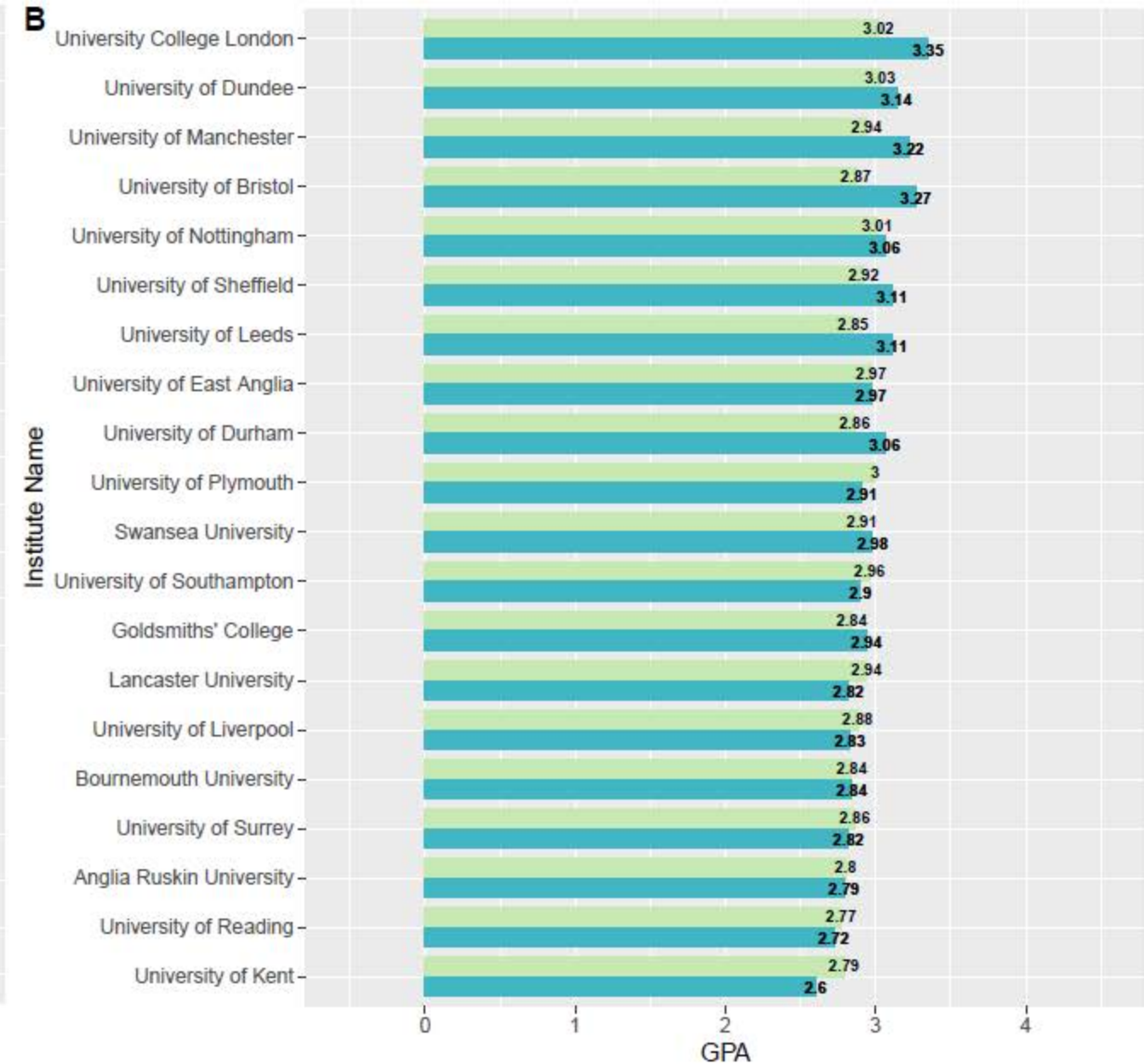
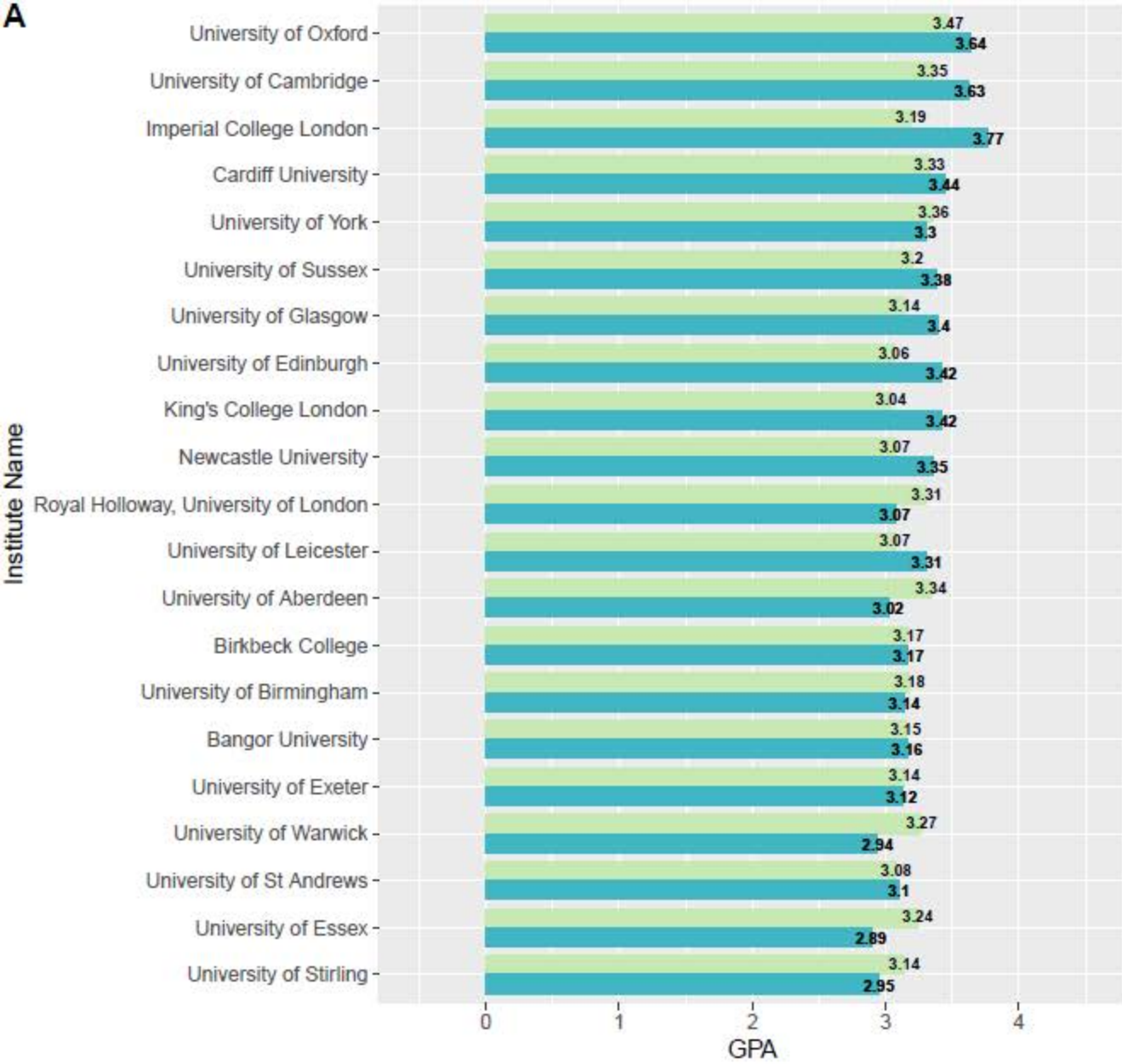
544 Wouters, P., Thelwall, M., Kousha, K., Waltman, L., De Rijcke, S., Rushforth, A., et
545 al. (2015). The metric tide: Correlation analysis of REF2014 scores and
546 metrics (Supplementary Report II to the Independent Review of the Role of
547 Metrics in Research Assessment and Management). *London: Higher*
548 *Education Funding Council for England (HEFCE)*.

549
550

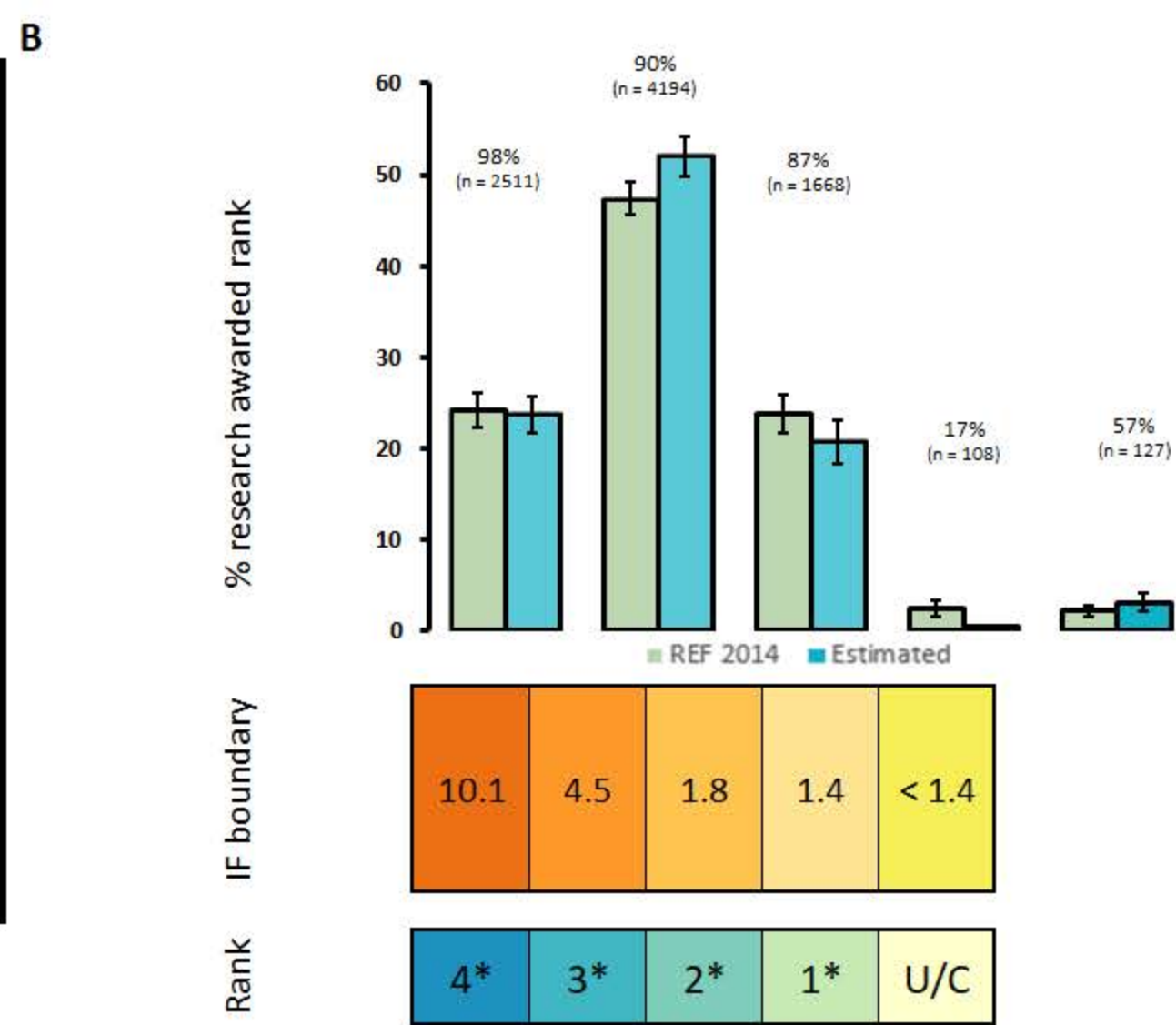
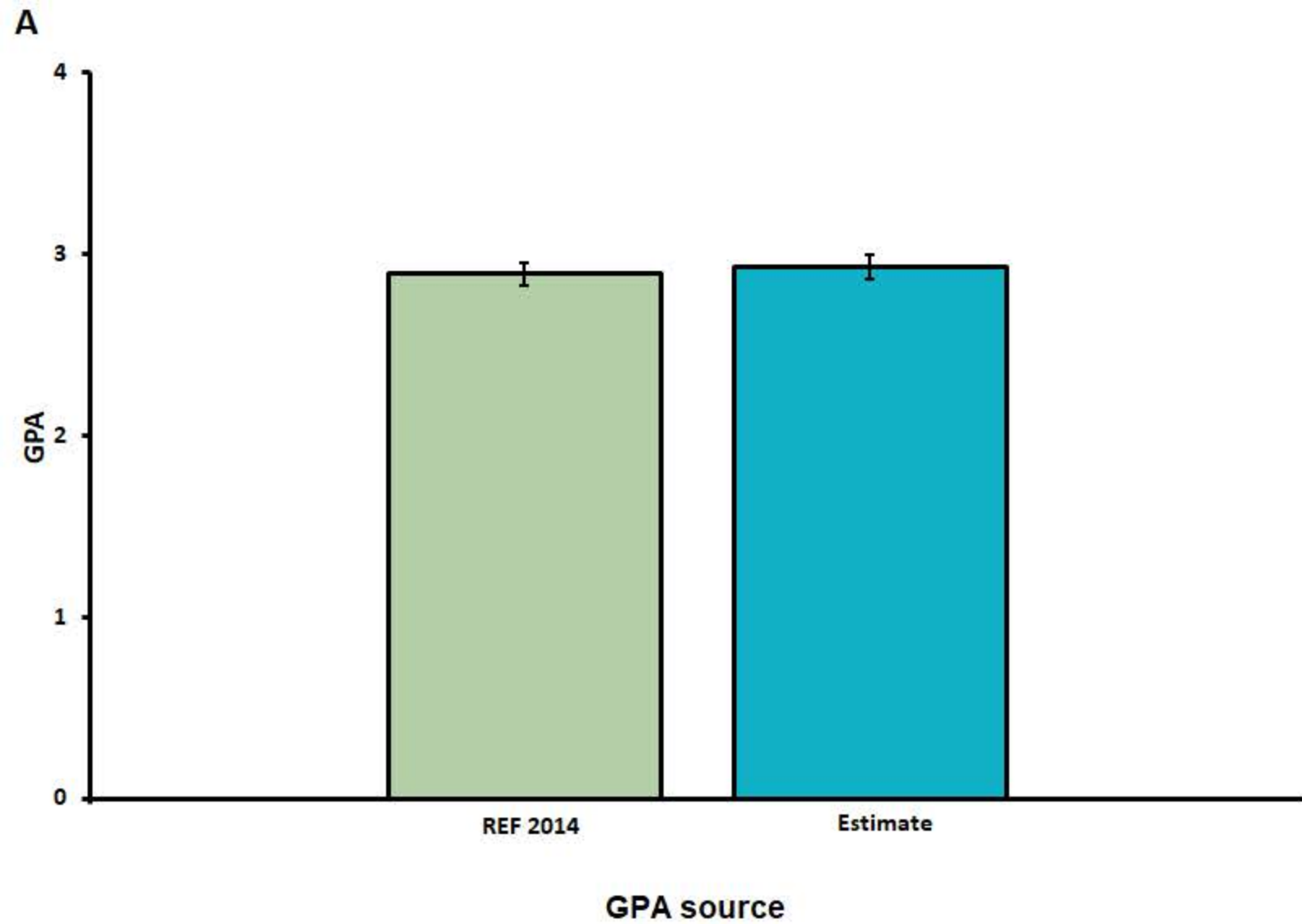








Type Estimate REF 2014



REF 2014 GPA

