

24 **Abstract**

25 **(299 words, 300 words max.)**

26 Article-level measures of publication impact (alternative metrics or altmetrics) can help authors and
27 other stakeholders assess engagement with their research and the success of their communication
28 efforts. The wide variety of altmetrics can make interpretation and comparative assessment difficult;
29 available summary tools are either narrowly focused or do not reflect the differing values of metrics
30 from a stakeholder perspective. We created the EMPIRE (EMpirical Publication Impact and Reach
31 Evaluation) Index, a value-based, multi-component metric framework for medical publications.
32 Metric weighting and grouping were informed by a statistical analysis of 2891 Phase III clinical trial
33 publications and by a panel of stakeholders who provided value assessments. The EMPIRE Index
34 comprises three component scores (social, scholarly, and societal impact), each incorporating
35 related altmetrics indicating a different aspect of engagement with the publication. These are
36 averaged to provide a total impact score and benchmarked so that a score of 100 equals the mean
37 scores of Phase III clinical trial publications in the *New England Journal of Medicine* (NEJM) in 2016.
38 Predictor metrics are defined to estimate likely long-term impact. The social impact component
39 correlated strongly with the Altmetric Attention Score and the scholarly impact component
40 correlated modestly with CiteScore, with the societal impact component providing unique insights.
41 Analysis of fresh metrics collected 1 year after the initial dataset, including an independent sample,
42 showed that scholarly and societal impact scores continued to increase, whereas social impact
43 scores did not. Analysis of NEJM 'notable articles' showed that observational studies had the highest
44 total impact and component scores, except for societal impact, for which surgical studies had the
45 highest score. The EMPIRE Index provides a richer assessment of publication value than standalone
46 traditional and alternative metrics and may enable medical researchers to assess the impact of
47 publications easily and to understand what characterizes impactful research.

48 **Introduction**

49 The publication of clinical trial results and other medical advances is an ethical obligation and
50 benefits a variety of stakeholders. Published information can be used by physicians, other healthcare
51 practitioners, and patients to evaluate and understand potential treatments. Medical researchers
52 and academics can use published results to inform their own research endeavors and to advance
53 medical research. In addition, policymakers use published information to develop guidelines and
54 treatment protocols that help to guide changes to clinical practice.

55
56 Publications are therefore vehicles for communicating research insights for peer-to-peer validation
57 and discussion. Impact measurements aim to assess the utility of published research for its intended
58 audience as well as the effectiveness of the communication.

59
60 Objective measures of impact can support these endeavors by enabling comparative assessments to
61 be made. However, making such measurements is challenging owing to the lack of available data
62 and agreed definitions of impact. Historically, a common proxy for the publication impact of an
63 article has been the impact factor of the journal in which it is published. However, although the
64 journal impact factor (JIF) may help to identify journals with a high readership, it is widely recognized
65 to be a poor indicator of the quality or impact of individual research articles [1,2].

66
67 Article-level metrics avoid the category error of using JIF in this context. The number of citations is
68 the most well-known metric, but this reflects only scholarly activity and citations can take years to
69 accumulate [3]. Recently, the advent of alternative article-level metrics (altmetrics) has provided a
70 new way to evaluate the impact of scientific publications. A wide range of potential altmetrics exists,
71 signifying different interactions with the publication of interest but differing widely in quality and
72 representativeness [4]. The sheer volume of potential metrics is evident in the information gathered

73 by major aggregators including Altmetric, which collects nearly 20 different altmetrics, and PlumX,
74 which collects over 40 [5,6].

75

76 To make metrics easier to interpret, various approaches have been taken to distilling them into
77 simplified scores. The most well-known of these is the Altmetric Attention Score (AAS), which
78 weights a variety of individual metrics to reflect a subjective assessment of relative reach and
79 aggregates them into a single number. Attempts to reduce any complex set of metrics into one
80 linear scale have been criticized because they will tend to be driven by a single predictor, especially
81 when the variables included are correlated [7]. Indeed, the AAS is dominated by Twitter and, to a
82 lesser extent, news articles [7–9], so it does not reflect the impact of publications among researchers
83 or policy-makers.

84

85 The full range of altmetrics is, however, multifactorial because they have diverse origins and
86 represent different activities relating to publications [10–14]. The AAS is only weakly correlated with
87 the number of citations [15,16]. Among the most cited, downloaded, and mentioned articles
88 published in general medical journals, only 2.5% were found in all three lists [17]. This implies that
89 altmetrics cannot effectively be reduced to a single linear representation, and data reduction can, at
90 best, provide several scores that group together related metrics. As a result, a metric system with
91 summary scores designed on data-reduction principles must, if it includes diverse, weakly correlated
92 metrics, provide for several distinct factors [14,18].

93

94 We sought to develop a value-based, multi-component metric framework for medical publications,
95 the EMPIRE (EMpirical Publication Impact and Reach Evaluation) Index, that would allow authors and
96 other professionals within the medical and pharmaceutical fields to assess the impact of publications
97 in terms meaningful to them. The metric framework is also intended to monitor the long-term

98 impact of publications, predict the likely long-term impact using early indicators, and identify the
99 effectiveness of communication efforts surrounding publications.

100

101 Focusing on a single discipline, medicine, has several advantages when developing a metric
102 framework. First, value is inherently subjective and is likely to differ between disciplines. Similarly,
103 the relationship between metrics varies between scientific disciplines [10, 15, 19, 20]. Second, using
104 the number of citations alone is known to underestimate severely the impact of clinical intervention
105 research compared with basic and diagnostic medical research [21], underscoring the need for a
106 multivalent approach to impact assessment. Third, medicine and medical sciences is the scientific
107 discipline richest in metrics [15], providing a large dataset to examine.

108

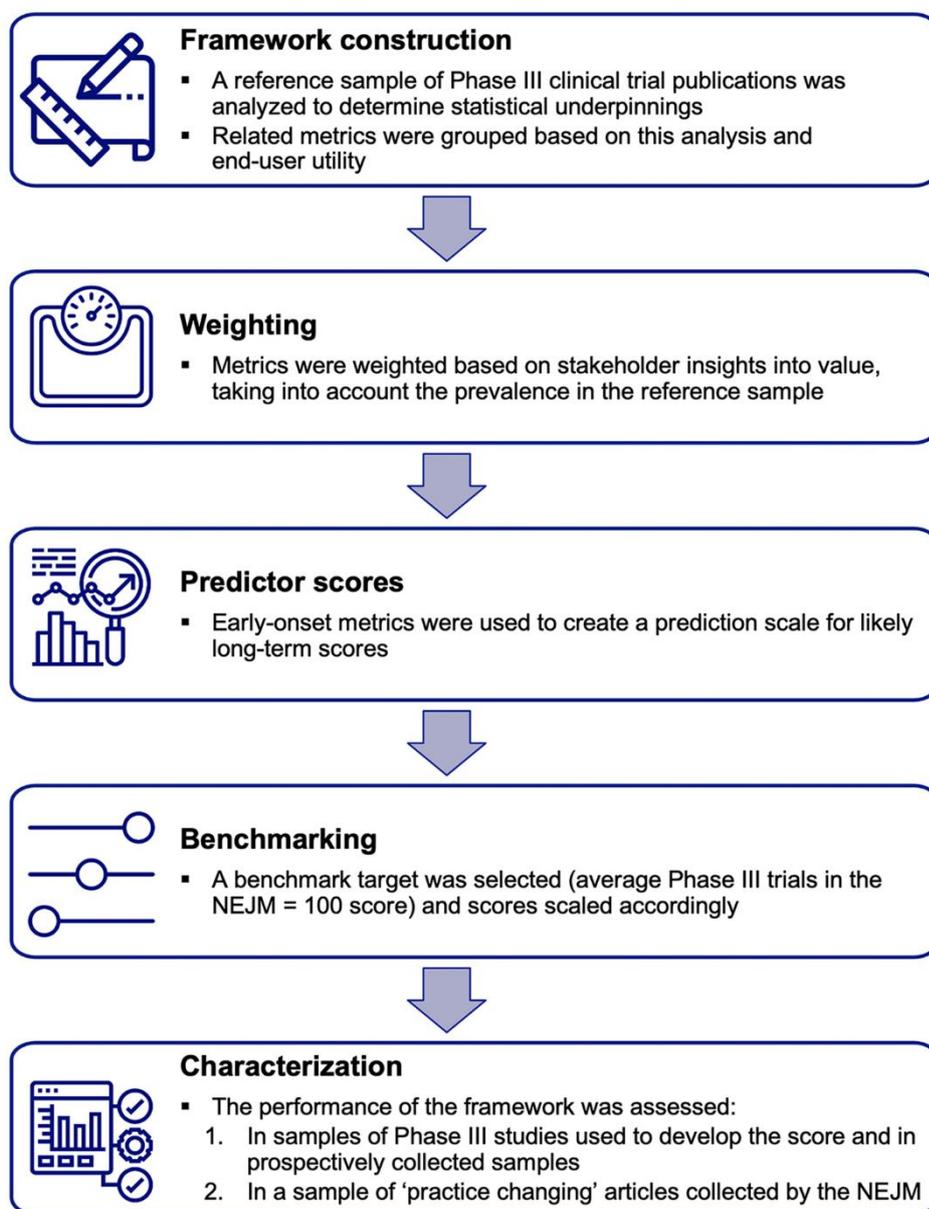
109 **Materials and methods**

110 **Approach to developing the scoring system**

111 Development of the scoring system for the EMPIRE Index proceeded through a series of stages,
112 outlined in Fig 1 and described in more detail in the sections below.

113

114 **Fig 1. Process for developing the scoring system.** NEJM, *New England Journal of Medicine*.



115

116

117 In summary, during **framework construction**, a large set of publications was generated to gain an in-
118 depth understanding of the statistical characteristics of altmetrics in a relevant sample. Publications
119 of Phase III clinical trials were chosen for analysis because these studies typically require a high
120 investment of resources and personnel and are most likely to have an impact on clinical practice. In
121 addition, they are likely to be rich in metrics – the mean number of metric counts has a substantial
122 effect on the size of the intercorrelation observed in a publication sample [22]. A series of statistical

123 analyses was then conducted to determine which metrics were comprehensive and provided useful
124 information, and how they were related to each other. The grouping and **weighting** of metrics was
125 informed by these analyses but was ultimately driven by an understanding of the type of interaction
126 each metric represented and by value judgments provided by a panel of stakeholders.

127

128 Once the structure and weighting of the metric system had been decided, **predictor scores** were
129 developed using altmetrics that accumulated rapidly. Scores for all components of the system were
130 then scaled to a **benchmark** representing a very high level of impact; for this, Phase III articles
131 published in the *New England Journal of Medicine* (NEJM) were chosen.

132

133 The last stage in development was to **characterize** the performance of the final scoring system. This
134 was carried out in three datasets: the original Phase III dataset, the Phase III dataset with metrics
135 updated after 1 year (and including 1 new year's worth of publications), and a dataset comprising
136 publications selected by NEJM editors that were likely to influence clinical practice.

137

138 **Sample acquisition**

139 **Reference Phase III sample**

140 We identified a sample of publications (the reference Phase III sample) that was representative of
141 the primary output of clinical medicine (Phase III clinical trials) as well as being sufficiently large to
142 permit statistical and longitudinal analysis. Data were obtained across 3 years of publications to
143 ensure the sample was large enough for analysis and included publications old enough to have
144 accumulated citations in guidelines and policy documents, while minimizing the impact of
145 confounding factors related to the change in use of publications over time (in particular, changes in
146 social media mentions). Non-English publications were excluded because the distribution of
147 altmetrics for these was likely to differ substantially from that of publications in English (e.g. news

148 coverage). The search was conducted on May 23, 2019 in PubMed, using the search term: ("clinical
149 trial, phase iii"[Publication Type]) AND (("2016/05/01"[Date - Publication] : "2019/05/01"[Date -
150 Publication]) AND Clinical Trial[ptyp] AND English[lang]).

151

152 Altmetrics for this sample were obtained on May 27, 2019. Article publication dates were obtained
153 from Altmetric Explorer and were used to split the sample into two subsamples – the older 50% (1H)
154 and the younger 50% (2H) – to assess the effect of temporal change in altmetrics.

155

156 **Benchmark NEJM Phase III sample**

157 The benchmark sample provided a ‘target’ against which to calibrate metrics achieved by other
158 publications. For this reason, a sample was chosen from a journal widely considered the ‘gold
159 standard’ for clinical trial publications, the NEJM, which has the highest JIF of all general medical
160 journals and describes itself as “the world’s leading medical journal and website” [23]. The
161 benchmark sample comprised all Phase III clinical trial articles published in the NEJM in 2016
162 (manually identified from a sample of all clinical trials obtained via a PubMed search). The year was
163 selected to allow the accumulation of metrics such as article or guideline citations, and to match the
164 base year in the reference Phase III sample. Altmetrics for the benchmark sample were obtained on
165 July 31, 2019.

166

167 **1-year update Phase III sample**

168 An independent sample was obtained to assess the metric framework for consistency. This sample
169 was identified on June 6, 2020 using the same search terms as the reference Phase III sample but for
170 the consecutive 12-month period (i.e. ("clinical trial, phase iii"[Publication Type]) AND
171 (("2019/05/01"[Date - Publication] : "2020/05/01"[Date - Publication]) AND Clinical Trial[ptyp] AND
172 English[lang])). Metrics for this 1-year update Phase III sample as well as for the original reference

173 Phase III sample were acquired on June 7, 2020 (approximately 1 year after the original metrics were
174 acquired).

175

176 To enable analysis of temporal changes, both the updated reference sample and the prospective
177 Phase III sample were divided into 12-month subsamples (May 1 to April 31) based on publication
178 dates provided by Altmetric Explorer. Publications with a publication date before May 1, 2016
179 according to Altmetric Explorer were excluded.

180

181 **NEJM notable articles sample**

182 An additional independent sample was identified with which to assess framework performance in
183 other types of clinical research, especially the utility of the societal impact component. Annually, the
184 editor of the NEJM curates a selection of articles published in the journal that year that they believe
185 have practice-changing potential ('notable articles'). We identified all of these articles for the years
186 2016, 2017, 2018, and 2019 [24–27], and obtained altmetrics for them on January 8, 2020. Articles
187 were classified by the authors under a broad typology: interventional (studies describing an
188 intervention with a medical treatment intended for clinical practice), observational (prospective and
189 retrospective non-interventional studies), innovative (publications describing novel techniques or
190 assays), and surgical.

191

192 **Acquisition of altmetrics and other metrics**

193 Data for all publications were obtained from the five sources listed below.

- 194 • Altmetric Explorer [6]: This was the primary source for altmetrics data as well as publication
195 dates).

- 196 • PlumX [5]: In addition to a wide range of metrics similar to those provided by Altmetric
197 Explorer, PlumX provided some unique metrics such as citations in articles classified by
198 Medline’s indexers as ‘clinical practice guideline’ (PubMed guidelines).
- 199 • Pubstrat Journal Database [28]: This was scraped to determine JIFs for journals identified by
200 Altmetric Explorer in the acquired datasets.
- 201 • CiteScore [29]: A journal-level, citation-based metric, similar to JIF. CiteScore was
202 downloaded for all journals on August 7, 2019 and CiteScore values for 2016 were used.
- 203 • Scimago Journal Ranking [30]: A journal-level, citation-based metric that used a PageRank
204 algorithm.

205

206 In addition to these standard metrics, original tweets and retweets (provided by Altmetric.com)
207 were obtained for the reference Phase III sample.

208

209 In a similar way to the exploratory analysis of Costas et al. (2015), an ‘altmetrics-driven’ universe of
210 publications was created in which all publications had at least one altmetric or citation (via Altmetric
211 Explorer) [12]. Costas et al. noted that this analysis did not result in a meaningful impact on the
212 precision of altmetrics as predictive tools for citations, but did reduce the zero inflation that can
213 confound statistical analysis.

214

215 **Statistical analysis**

216 Analyses were conducted in Microsoft Excel using the Analyse-it plugin (Analyse-it Software, Ltd.,
217 Leeds, United Kingdom). Descriptive statistics were obtained and Spearman rank correlations
218 between individual altmetrics were calculated. In addition, exploratory factor analysis was used to
219 provide insights into how best to group similar metrics. Factor analysis assumes that latent or
220 underlying factors exist that causally influence the observations. For the purposes of metric

221 development, we wanted to explore the hypothesis that publications have an intrinsic ‘social’
222 interest leading to social media mentions that is fundamentally different from an intrinsic ‘scholarly’
223 interest leading to citations. An alternative data-reduction technique, principal component analysis,
224 simply creates one or more index variables explaining as much statistical variance as possible
225 without regard to theoretical differences in the metrics. In practice, the two approaches yield similar
226 results.

227
228 We used maximum likelihood factor analysis with oblique (oblimin) rotation. Because altmetrics
229 follow a power-law distribution [31], data were log-transformed before factor analysis. All data were
230 increased by 1, which allows the discretized lognormal distribution to be fitted to the full range of
231 data [32]. Adding a positive constant to the dependent variable is a common solution to the problem
232 of log-transformation of datasets containing zeros, although it does introduce a small distortion to
233 the data [33]. Regression analyses were conducted using multiple linear regression on the
234 untransformed data.

235
236 EMPIRE Index scores for the NEJM notable articles were averaged over the different years (2016–
237 2019). To control for the impact of time on the accumulation of altmetrics, EMPIRE Index scores for
238 articles in each year were expressed as a percentage of the average score of observational studies
239 (the highest-scoring article type), and the average of these yearly percentages was taken.

240

241 **Value assessment**

242 An internal Novartis cross-functional stakeholder panel meeting was convened on July 9, 2019,
243 comprising representatives from scientific communications, medical, commercial, launch strategy,
244 and medical analytics departments. Participants reviewed information on the analyses conducted as
245 well as background information on metrics, and provided qualitative insights into the interpretation
246 and importance of key metrics. Quantitative value assessments were obtained through points

247 allocation (i.e. participants were given a fixed number of points to distribute among metrics
248 according to the value seen in them). Points were summed and the proportion of points allocated to
249 each metric was calculated.

250

251 **Predictor scores**

252 Two predictor scores were developed based on metrics that accumulate rapidly after publication.
253 The early predictor score included altmetrics that accumulated most rapidly (Twitter, Facebook, and
254 news mentions) [3, 34, 35], and included CiteScore, used here as a proxy for the readership and
255 interest in a journal. The intermediate predictor score included blog mentions, F1000Prime
256 mentions, and Mendeley readers – altmetrics that accumulate more slowly, but still faster than
257 metrics with high lag, such as citations.

258

259 The basis of each predictor score was a multiple linear regression of the altmetrics included in the
260 predictor against the total impact score in the reference Phase III sample. Weightings for each metric
261 were calculated as follows:

$$262 \quad \text{weighting} = \beta \frac{\text{sum}_m}{\int \text{sum}_{m1,m2\dots}}$$

263 where β is β from linear regression, sum_m is the sum total of the incidence of the target metric in the
264 reference sample, and $\text{sum}_{m1,m2\dots}$ is the sum total of all metrics included in the predictor score.

265

266 **Results**

267 **Framework construction**

268 The initial search found 3498 Phase III clinical publications, of which altmetrics for 3450 were
269 identifiable by PlumX and 2891 by Altmetric Explorer. The analysis set comprised 2891 articles with
270 at least one metric identified by Altmetric Explorer, of which eight were unavailable in the PlumX

271 dataset. Publication metric characteristics of this sample are shown in S1 Table. Several altmetrics
272 had a very low density so were discarded for further analysis (e.g. Weibo, LinkedIn, Google+,
273 Pinterest, Q&A, peer review, video, and syllabi mentions). Some altmetrics were retained despite a
274 low density as they were thought to provide unique insights relevant to the objectives (policy,
275 patent, F1000Prime, Wikipedia, and guideline [from PlumX] mentions). Some metrics of high
276 relevance (abstract and publication views and downloads) were discarded because the quality of the
277 data was inconsistent – in particular, many papers had numerous citations and Mendeley readers
278 without recorded views or downloads, suggesting that coverage was incomplete.

279
280 Journal-level metrics were not included in the EMPIRE Index total impact score or component
281 scores, but they were considered potential components of predictor scores. Given that the coverage
282 obtained with CiteScore was higher than with the other two journal-level metrics examined (JIF and
283 Scimago Journal Ranking – S1 Table part C), CiteScore was selected for further analyses.

284
285 Pairwise Spearman correlations between the altmetrics included are shown in S2 Table. The most
286 common metrics were strongly correlated (news, blog, and Twitter mentions, Mendeley readers,
287 and Dimensions citations). The strongest correlation was seen between Mendeley readers and
288 Dimensions citations, although Facebook mentions and tweets were also strongly correlated. In
289 addition, original tweets and retweets were highly correlated with each other and with total tweets,
290 suggesting that a single measure (total tweets) is sufficient. Other metrics showed only weak
291 correlations with each other.

292
293 Dividing the reference Phase III sample into two subsamples according to the publication date
294 provided by Altmetric Explorer revealed important differences (S1 Fig). The more recent half of the
295 publications (2H, after May 21, 2017) had higher mean Twitter, Wikipedia, and other counts, but
296 lower Dimensions citations, than the older half (1H).

297

298 Three-factor analysis was conducted on the full range of metrics selected for inclusion (S3 Table).

299 Two-factor analysis was also carried out on a subset of metrics excluding those with low incidence

300 (policy document, PubMed guideline, and patent mentions) (S4 Table). These analyses revealed

301 consistent groupings, such as Mendeley readers with Dimensions citations, and news, blog, and

302 Wikipedia mentions.

303

304 **Weighting**

305 Based on the results of these analyses and considerations, a framework for grouping metrics was

306 developed comprising three component scores: social impact (news, blog, Twitter, Facebook, and

307 Wikipedia mentions), scholarly impact (Mendeley readers, Dimensions citations, and F1000Prime

308 posts), and societal impact (mentions in policy documents, PubMed guidelines, and patents). An

309 initial statistical estimate of weightings was calculated as the inverse proportion of counts of each

310 altmetric in the reference Phase III sample relative to the total number of all altmetric counts.

311

312 Discussions during the stakeholder panel meeting revealed the central importance given to guideline

313 and policy document citations as a measure of article impact. This was also reflected in the

314 quantitative session, in which guidelines and policy documents were allocated over one-third of the

315 total points (Table 1).

316

317 **Table 1. Value accorded to metrics by the stakeholder panel (quantitative scoring).**

Metric	Allocation of points	Percentage of points
Twitter mentions	21	10
Facebook mentions	21	10
Blog mentions	14	6
News mentions	19	9
Wikipedia mentions	7	3
Dimensions citations	32	15
Mendeley readers	14	6
F1000Prime mentions	12	6
Guideline mentions	47	22
Policy mentions	26	12
Patent mentions	3	1

318

319 Weightings derived from the statistical approach were revised to reflect findings from stakeholder
320 value assessments. The selected weightings and their contribution to the total impact score based
321 on the sample are shown in Table 2. In general, the approach taken was to balance the weighting
322 such that the percentage contribution to scores in publications in the reference Phase III sample
323 resembled the stakeholder value, while acknowledging relative importance (e.g. of news articles vs
324 blogs) and prevalence (e.g. when Wikipedia entries were too infrequent to make a meaningful
325 contribution without greatly inflated weighting relative to the value accorded by the stakeholder
326 panel). To combine statistical and value-based weighting effectively, some related metrics were
327 considered as combined entities (i.e. Twitter and Facebook mentions were allocated a combined
328 20% of points by stakeholders, and contributed a combined 17.7% to the total impact score in the
329 reference Phase III sample).

330

331 **Table 2. Weighting assigned to metrics included in the social, scholarly, and societal impact scores,**
332 **along with their contribution to total impact scores in the reference sample.**

Metric	Total in reference sample	Percentage of all metrics in reference sample	Social weighting	Scholarly weighting	Societal weighting	Percentage contribution to total in reference sample
Twitter mentions	94,235	29.25	3	–	–	17.0
Facebook mentions	3821	1.19	3	–	–	0.7
Blog mentions	1086	0.34	10	–	–	0.7
News mentions	18,539	5.75	15	–	–	16.7
Wikipedia mentions	70	0.02	5	–	–	0.0
Dimensions citations	78,785	24.45	–	4	–	18.9
Mendeley readers	124,866	38.75	–	1	–	7.5
F1000Prime mentions	252	0.08	–	15	–	0.2
Guideline mentions	183	0.06	–	–	1800	19.8
Policy mentions	321	0.10	–	–	900	17.4
Patent mentions	59	0.02	–	–	300	1.1
Percentage contribution to total	–	–	35.1	26.7	38.2	100

334 **Predictor scores**

335 The variance in total impact scores explained by each predictor score was moderately high (early
336 predictor vs total impact score, $r^2 = 0.56$; intermediate predictor vs total impact score, $r^2 = 0.65$, S2
337 Fig). An overall predictor score can be calculated as the average of early and intermediate predictor
338 scores. The variance in total impact scores explained by the overall predictor score was also
339 moderate (overall predictor vs total impact score, $r^2 = 0.69$). Weightings calculated for each of the
340 variables in the predictor score are shown in Table 3.

341

342 **Table 3. Weightings assigned to metrics included in the early and intermediate predictor scores.**

Metric	Total in reference sample	Early predictor score	Intermediate predictor score	Percentage contribution to overall predictor score in reference sample
CiteScore	15,384	57	–	26.2
News mentions	18,539	22	–	12.2
Twitter mentions	94,235	3	–	8.2
Facebook mentions	3821	30	–	3.5
Mendeley readers	124,866	–	12	44.8
Blog mentions	1086	–	125	4.1
F1000Prime mentions	252	–	146	1.1

343

344 **Benchmarking**

345 In total, 74 Phase III publications from the NEJM published in 2016 were identified for the
346 benchmark sample. The non-adjusted, non-adjusted overall predictor score was selected as the

347 benchmark for predictor scores, and the non-adjusted total impact score was selected for total,
 348 social, scholarly, and societal impact scores (Table 4).

349

350 **Table 4. Scores in the benchmark sample before and after benchmark adjustment.** Non-adjusted

351 scores chosen as benchmarks are shown in bold.

	Early predictor score	Intermediate predictor score	Overall predictor score^a	Social	Scholarly	Societal	Total^b
Mean non-adjusted score	3218	4414	3816	1622	1593	2854	6068
Benchmark value	3816	3816	3816	2023	2023	2023	6068
Mean benchmarked score	84	116	100	80	79	141	100

352 ^aThe overall predictor score is the average of early and intermediate predictor scores.

353 ^bThe non-adjusted total impact score is the sum of the social, scholarly, and societal impact scores.

354 The adjusted total impact score is the average of the adjusted component scores.

355

356 Dividing the non-adjusted total benchmark by 3 before applying it to the component scores had the

357 effect of upscaling them so that the adjusted total impact score represents the mean of the

358 components (rather than the sum, as in the non-adjusted total impact score). EMPIRE Index scores

359 are calculated by dividing the unadjusted score of interest by the appropriate benchmark and

360 multiplying by 100.

361

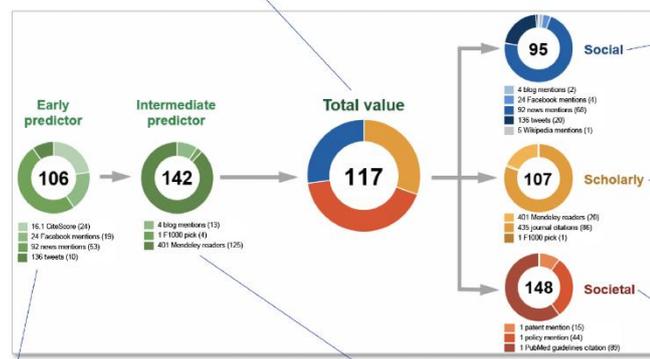
362 Final scoring framework

363 An overview of the final EMPIRE Index framework is shown in Fig 2. The framework comprises the
 364 three component scores (social, scholarly, and societal impact), which are averaged to provide a
 365 total impact score. Each component score incorporates a separate type of altmetric, indicating a
 366 different aspect of engagement with the publication. The framework also includes the two predictor
 367 scores.

368

369 **Fig 2. Example of the EMPIRE Index score for a single publication.** HCP, healthcare provider; NEJM,
 370 *New England Journal of Medicine*.

The **total** impact score represents a weighted average of the social, scholarly, and societal impact scores
 A **total** impact score of 100 is equivalent to the average score of Phase III articles published in the NEJM in 2016



The **social** score represents the impact of the article in public domains such as social media and news
Audience: HCPs, non-specialists, healthcare support staff, patients, and other members of the public

The **scholarly** score represents the impact of the article in academic domains such as journal citations and scholarly reference libraries
Audience: specialists, experts, scientists, and academics

The **societal** score represents the impact of the article in treatment guidelines, policy documents, and patents
Audience: healthcare and policy decision makers, disease management bodies

The **early predictor** score uses metrics that accumulate quickly to estimate future total impact

The **intermediate predictor** score uses metrics that are intermediate between early metrics and late metrics (citations and societal metrics) to provide an additional estimate total impact

371

372

373 **Characterization of the EMPIRE Index**

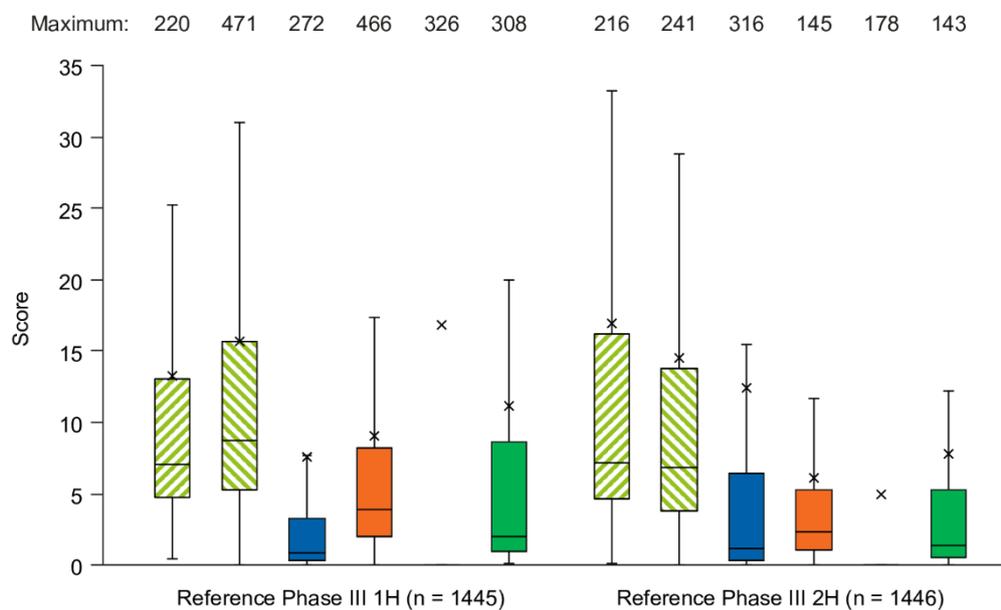
374 **Characterization in samples used in development**

375 The distributions of scores in the reference sample 1H and 2H, and in the benchmark sample, are
376 shown in Fig 3. Of note, social impact scores were lower and societal impact scores were higher in
377 1H than in 2H. Predictor scores were higher than total impact scores in the reference Phase III
378 sample but not in the benchmark NEJM Phase III sample, and median social impact scores were
379 closer to median total impact scores in the benchmark NEJM Phase III sample than in the reference
380 Phase III sample.

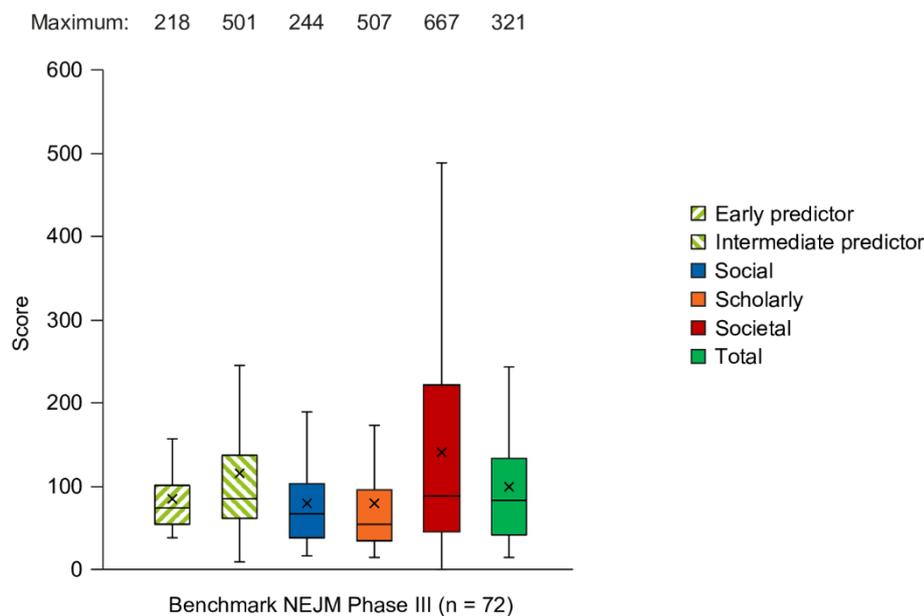
381

382 **Fig 3. Distribution of scores.** Distribution of scores in (A) the reference Phase III sample and (B) the
 383 benchmark NEJM Phase III sample. Box = 1Q–2Q, whiskers = $1.5 \times$ interquartile range, X = mean. 1H,
 384 older 50%; 2H, younger 50%. NEJM, *New England Journal of Medicine*.

(A)



(B)



385

386

387 The correlations between component scores, the AAS, and CiteScore are shown in Table 5.

388 Correlations between component scores were relatively low, the greatest being between social and

389 scholarly impact scores. The social impact score correlated strongly with AAS, and both social and
 390 scholarly impact scores correlated modestly with CiteScore. However, the societal impact score is
 391 quite distinct from AAS, CiteScore, and the other component scores. Although predictor scores were
 392 moderately successful at predicting the total impact score, they were only weakly related to the
 393 societal impact score.

394

395 **Table 5. Correlations (Spearman r) between component scores, AAS, and CiteScore in the**
 396 **reference Phase III sample. Correlations > 0.6 are shown in bold.**

Score	Early predictor	Intermediate predictor	Social	Scholarly	Societal	Total	AAS	CiteScore
Early predictor	–	0.61	0.79	0.63	0.19	0.70	0.76	0.91
Intermediate predictor	0.61	–	0.59	0.87	0.26	0.76	0.59	0.55
Social	0.79	0.59	–	0.59	0.18	0.74	0.95	0.58
Scholarly	0.63	0.87	0.59	–	0.32	0.84	0.59	0.58
Societal	0.19	0.26	0.18	0.32	–	0.55	0.27	0.16
Total	0.70	0.76	0.74	0.84	0.55	–	0.78	0.57
AAS	0.76	0.59	0.95	0.59	0.27	0.78	–	0.56
CiteScore	0.91	0.55	0.58	0.58	0.16	0.57	0.56	–

397 AAS, Altmetric Attention Score.

398

399 **Characterization in prospectively collected samples**

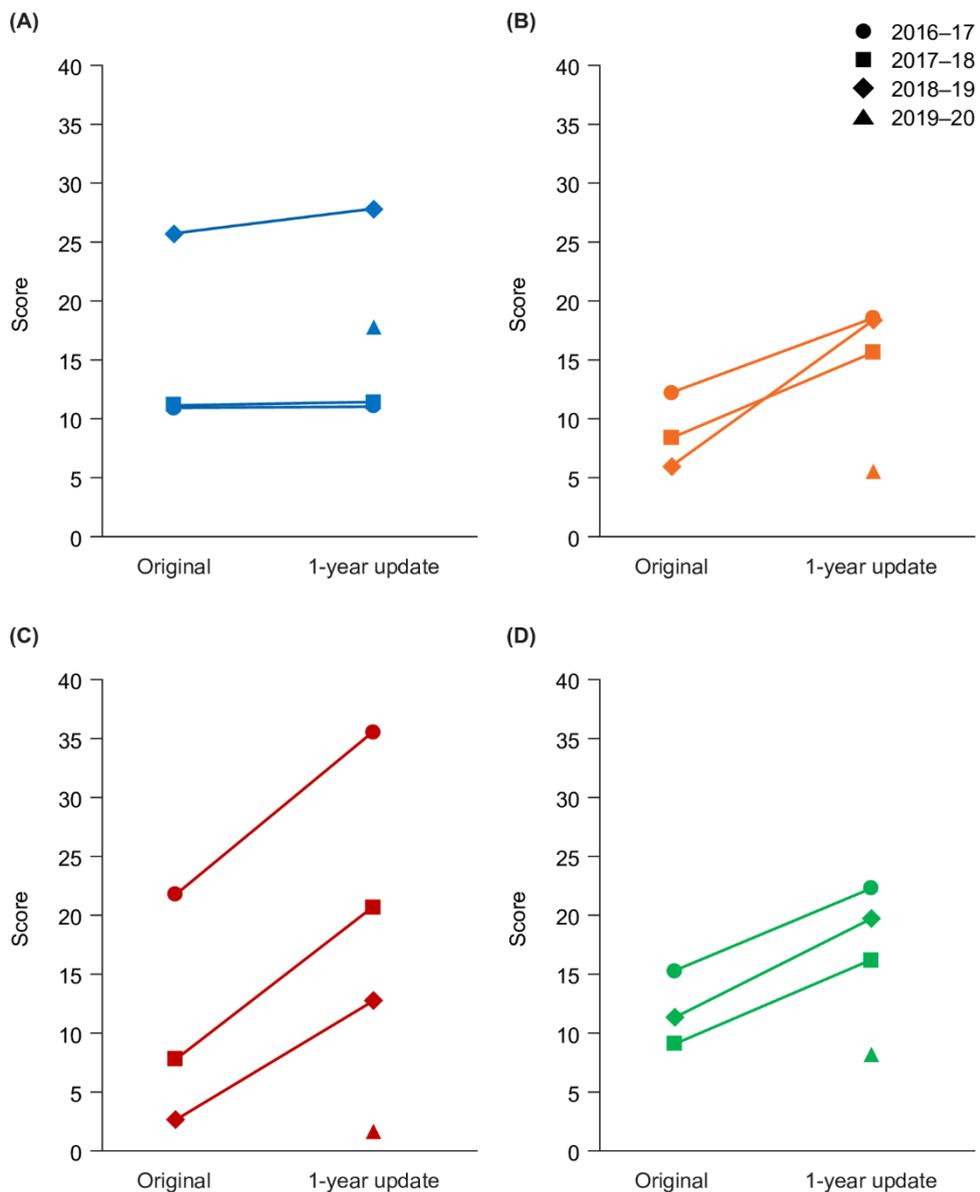
400 **1-year update Phase III sample characterization**

401 Publication dates obtained from Altmetric Explorer indicated that 194 articles were published prior
402 to May 1, 2015; 1173 from May 1, 2016 to April 30, 2017; 1101 from May 1, 2017 to April 30, 2018;
403 and 423 from May 1, 2018 to April 30, 2019. The drop in publication numbers in the latter period
404 most likely reflects a lag in MEDLINE indexing. The 2019–2020 search identified 503 publications, of
405 which 435 met the date criteria based on publication dates obtained from Altmetric Explorer. Mean
406 EMPIRE Index scores in these year groups in both the original altmetric acquisition and the 1-year
407 update are shown in Fig 4. Little change was found in the social impact component. Scholarly impact
408 and, especially, societal impact continued to accumulate. The greatest increase in scholarly impact
409 was seen in the most recent publications, while societal impact scores increased similarly across all 3
410 years sampled.

411

412 **Fig 4. Mean impact scores in the original reference Phase III sample and the 1-year update Phase**

413 **III sample. (A) Social, (B) scholarly, (C) societal, and (D) total mean impact scores.**



414

415

416 **NEJM notable articles characterization**

417 In total, 48 notable articles were identified by NEJM editors from 2016 to 2019. Mean impact scores

418 from the 2016 subset are shown in Fig 5, with mean scores from the 2016 benchmark NEJM Phase III

419 sample for comparison. Notable articles had higher social and societal impact than benchmark

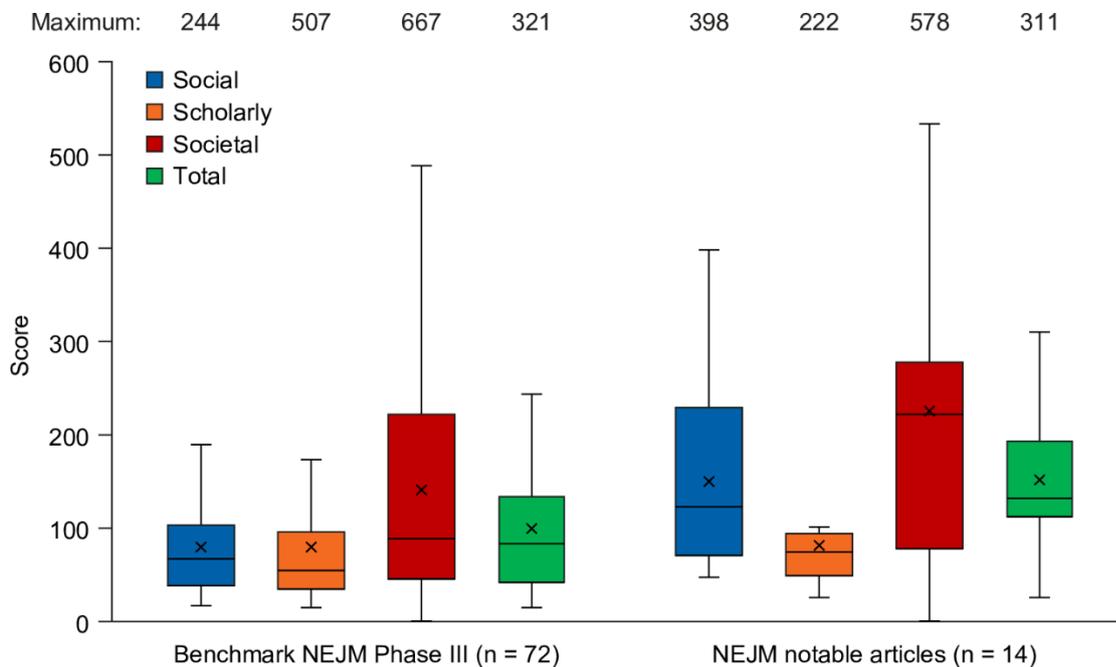
420 articles.

421

422 **Fig 5. Mean impact scores for NEJM notable articles from 2016 compared with benchmark scores.**

423 Box = 1Q–2Q; whiskers = $1.5 \times$ interquartile range; X = mean. NEJM, *New England Journal of*

424 *Medicine*.



425

426

427 Of the 48 articles, the focus was assessed to be interventional in 24 cases, observational in 10 cases,

428 innovative in 6 cases, and surgical in 8 cases. After adjusting for publication year, observational

429 studies were found to have the highest total impact, with other publication types having lower

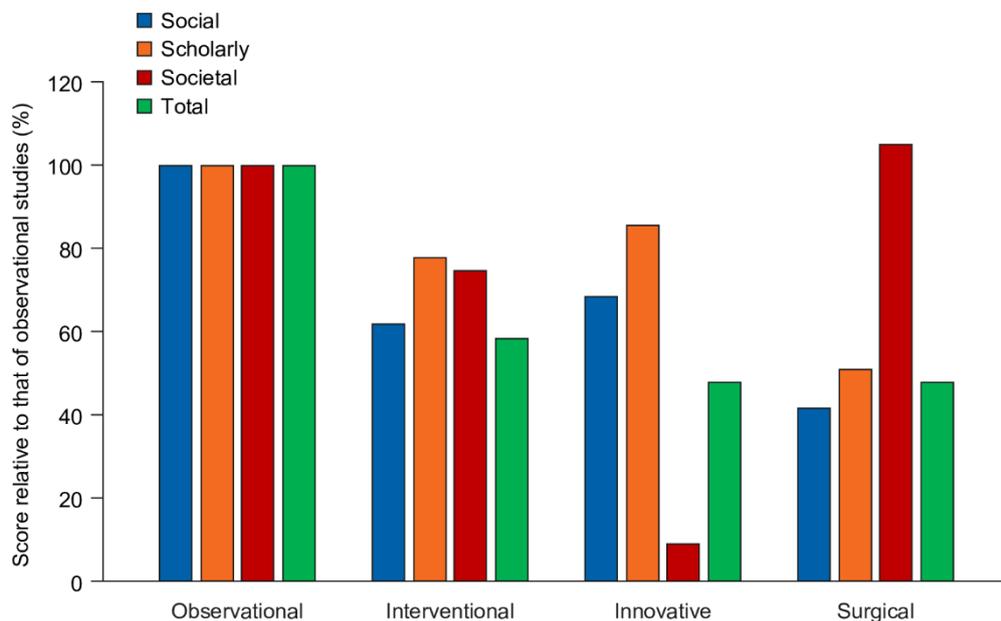
430 impact scores across all component scores except for the societal impact of surgical studies.

431 Innovative studies had notably low societal impact, indicating that they were infrequently

432 referenced in guidelines or policy documents (Fig 6).

433

434 **Fig 6. EMPIRE Index scores for NEJM Notable articles.** EMPIRE Index scores are expressed as
435 percentages of the scores achieved by observational studies.



436

437

438 Discussion

439 We have developed the EMPIRE Index, a metric framework to assess the multidimensional impact of
440 medical publications, including the potential impact on clinical practice. It avoids the pitfalls of JIF-
441 based research assessment and unidimensional scoring systems. It also fulfills the Leiden criteria of
442 being open, transparent, and simple [36].

443

444 The EMPIRE Index aggregates selected article-level metrics into meaningful component scores and
445 weights them according to the value placed on them by members of a stakeholder panel and
446 statistical analysis of a representative sample of articles. It differs conceptually from both the AAS
447 and other, recently developed scoring systems: the #SoME_Score [37], the Weighted Altmetric
448 Impact, and the Inverse Altmetric Impact [18,38]. First, the value-based approach to the weighting
449 and grouping of metrics recognizes that simple statistical associations may be sample-dependent
450 and may not relate to underlying conceptual underpinnings. Second, the EMPIRE Index is specifically

451 designed for medical publications. Many studies have documented different scales and relationships
452 between metrics in various disciplines [10,15,19,20] and, given that the value of each metric is
453 inherently subjective, this value is unlikely to be consistent across scholarly disciplines. Third, the
454 EMPIRE Index is scaled against a clearly defined, relevant benchmark, because interpretation of a
455 novel composite metric is difficult without such a reference point.

456

457 Such are the potential advantages of the EMPIRE Index. However, its utility is dependent on the
458 robustness of the selection grouping, weighting of metrics, and benchmarking, as well as its
459 performance in the evaluation of suitable publications. In the process of investigating these factors,
460 a series of results of broad interest to the altmetrics community were generated. These will be
461 discussed in the sections that follow.

462

463 **Metric selection**

464 Suitable metrics were identified for inclusion by reviewing the coverage and density of metrics
465 obtained for the reference Phase III sample through the two most established metric providers:
466 Altmetric and PlumX. Previous research has shown significant differences between these providers
467 in terms of Mendeley readers and Twitter coverage, as a result of different approaches to collecting,
468 tracking, and updating metrics [39–41]. Furthermore, the approach to covering news and blog posts
469 differs greatly between PlumX and Altmetric Explorer [42]. We found broadly similar metrics
470 between the two providers except for news articles, with Altmetric reporting twice as many for our
471 sample as PlumX, and Facebook mentions (because Altmetric extracts only mentions on Facebook
472 pages whereas PlumX also extracts ‘likes’).

473

474 The reference dataset was selected to provide a sample rich in altmetrics. PlumX identified at least
475 one metric for 99% of our sample, while the figure was 83% for Altmetric. This result compares
476 favorably with that of previous work [12,15,19,20,43,44], most likely indicating the increasing

477 volume of altmetric activity. One important metric not included was article views and downloads.
478 Although these data were provided by PlumX, we found them to be patchy, with many articles
479 reporting metrics such as tweets or Mendeley readers but no page views or downloads on the
480 EBSCO information service. This resulted in weak and spurious correlations (data not shown), similar
481 to the findings of Maggio et al. (2018) [45].

482
483 Similar to previous investigators, we found that news, blog, Twitter, and Facebook mentions,
484 Mendeley readers, and Dimensions citations were the most common metrics in our sample. These
485 metrics were included in our analysis, as well as additional metrics that, although rare, provided
486 valuable insights into article impact: citations in policy documents, guidelines, patents, Wikipedia,
487 and F1000Prime.

488

489 **Rationale for the three component scores**

490 The EMPIRE Index comprises three component scores, each representing a different factor
491 underlying the observed patterns of metrics seen in the reference Phase III sample. The social impact
492 component represents actions that involve or are accessible to the general public as well as
493 healthcare professionals and academics. The scholarly impact component represents actions with an
494 academic focus. The societal impact component represents actions in which the publication has
495 been used to inform decisions around optimal care or, in the case of patents, medical advances.

496

497 To inform our metric grouping, correlation analyses were performed and exploratory factor analysis
498 was used. These revealed a close connection between Mendeley readers and Dimensions citations,
499 in line with findings from previous research [46–48]. Correlations were also found between Twitter
500 and Facebook mentions, and news/blog and social media mentions, which again aligns with previous
501 observations [47, 48].

502

503 No meaningful correlations were found between mentions in F1000Prime articles, policy documents,
504 guidelines, patents, or Wikipedia articles and other metrics. These metrics have not previously been
505 widely studied, and the low correlations observed may reflect their very small coverage – over 90%
506 of publications score zero on these metrics. However, Bornmann and Haunschild (2018) reported
507 that F1000Prime recommendations were more closely correlated with Mendeley readers and
508 Dimensions citations than with Twitter mentions [49].

509

510 Pairwise correlations can give useful insights into relationships between different metrics, but for
511 the purposes of reducing data into composite scores it is helpful to understand the shared variance
512 between multiple metrics. The exploratory factor analysis in our study produced findings consistent
513 with those reported in previous literature [10–13]. Separating articles into those that were older
514 (1H) and younger (2H) showed that citations (including policy and guideline mentions) and Mendeley
515 readers consistently grouped into one factor; news, blogs, Wikipedia, and F1000Prime mentions
516 grouped into a second factor; and Twitter (and, usually, Facebook) mentions comprised a third
517 factor. A two-factor analysis excluding policy document, guideline, and patent mentions confirmed
518 that Mendeley readers and Dimensions citations formed a separate group from the remaining
519 metrics.

520

521 Each altmetric represents a different action on the part of an audience; this has implications for how
522 we understand the meaning of individual metrics [4] and whether these statistical associations
523 represent meaningful groupings. For example, much remains unknown about the motivation for
524 tweeting, given that most tweets are empty of context [50] and content [51]. Often all that is certain
525 is that the tweeter felt the research interesting enough to broadcast. Social media platforms are
526 known to be used mostly by the general public, so a central motivation for scholars to tweet is likely
527 to be to communicate and explain their work to lay people [52]. This may be particularly true of
528 publications in biomedical sciences, which attain greater Twitter interest than those in other

529 scholarly disciplines [52]. Twitter communities linked through publication tweets tend to be led by
530 organizational accounts associated with well-known journals or leading scholars [53], although at
531 least half of sharing on social media is likely to be non-academic [54,55].

532

533 Reference manager data have been suggested as an alternative to download counts as a source of
534 readership evidence [3]. Although Mendeley users often add articles to their library with the
535 intention of citing them, many also add these for professional or teaching purposes, which may
536 explain why some articles have many readers but few citations.

537

538 Interestingly, articles rated in F1000Prime reviews as ‘good for teaching’ received higher Twitter
539 scores, but not higher Mendeley scores, than those that were not rated this way. The reverse was
540 true for articles considered a ‘technical advance’ [56].

541

542 **Weighting**

543 The weighting of metrics in the EMPIRE Index was based on three considerations: the prevalence of
544 metrics in the reference sample (highly prevalent metrics were weighted less), the need for each
545 component to make a substantial contribution to the total impact score, and the value given to each
546 metric as an indicator of impact. As a result, the weighting is quite different from other approaches
547 based on purely statistical considerations.

548

549 Several approaches have determined weighting by regressing altmetrics on citations. These typically
550 result in, for example, higher weighting given to blog posts and Mendeley readers than to news
551 articles (because blog posts are relatively uncommon) [15,37,45,57]. Because the target variable is
552 journal citations, each Mendeley save or F1000 citation may be weighted in a similar way to or
553 higher than a policy document citation [37, 57]. Ortega has developed weightings based on principal
554 component analysis and also on inverse prevalence (so that the rarest metrics receive the highest

555 weighting). The two approaches create very different weightings – for example, a news article
556 carries half the weight of a publication citation in the Weighted Altmetric Impact, but eight times the
557 weight of a publication citation in the Inverse Altmetric Impact [18,38]. These statistical approaches
558 give very different results from the weighting developed for the EMPIRE Index.

559

560 **Predictor scores**

561 Given that some altmetrics accumulate early, there is long-standing interest in the use of a limited
562 selection of rapidly accumulating altmetrics to identify publications likely to have high long-term
563 impact. Earlier work has employed multivariate regression with citations as a measure of long-term
564 impact [9,15,37,45,57,58] but, as we have seen, citations are only one of several measures of long-
565 term impact.

566

567 Among common metrics, tweets and news articles accumulate most rapidly after publication, while
568 Mendeley readers, blogs, and F1000Prime articles increase more gradually [3,34,35]. Wikipedia and
569 policy document mentions can, like article citations, take well over a year to accumulate [34,59]. The
570 EMPIRE Index addresses this by using two predictor scores – early and intermediate. The early
571 predictor score also uses CiteScore, a journal-based metric. CiteScore, in this context, can be thought
572 of as a proxy for the exposure an article is likely to have; it has previously been shown that
573 combining citations over the first year with JIFs accurately predicts future citations [59,60].

574

575 Predictor scores are a purely statistical construct so the weighting is quite different from the EMPIRE
576 Index itself; however, the weighting is also different from methods employed in previous work using
577 citations as a target. Compared with studies mentioned earlier that used statistically based
578 weighting with only citations as a target, in the EMPIRE Index predictor scores, Mendeley readers
579 carry less weight relative to news article citations. This most likely reflects the broader basis of the
580 EMPIRE Index compared with citation-only targets.

581

582 The reasonably strong relationship between predictor scores and the total impact score in the
583 reference Phase III sample is to be expected, given that they share many of the same metrics.
584 However, the weak correlation with the societal impact score indicates that the predictor scores will
585 lack precision in identifying high-impact publications (given the importance of the contribution of
586 societal impact to the total impact). Further work using longitudinal datasets is required to improve
587 these predictor scores.

588

589 **Responsiveness and characterization**

590 The responsiveness and utility of the EMPIRE Index was evaluated in several ways. Averages and
591 distributions of scores in the reference Phase III sample and the benchmark NEJM sample were
592 explored, showing that both samples had similar social and scholarly metrics and the latter had far
593 higher societal metrics. Because the scores were scaled to the benchmark NEJM sample, this
594 resulted in predictor scores lacking sensitivity for lower-impact publications (i.e. although they
595 retained precision for identifying higher-impact articles, they tended to overpredict the impact of
596 lower-impact articles uniformly).

597

598 The social score was shown to be closely correlated with the AAS. The AAS weights metrics in a way
599 that is not possible for users of the Altmetric Explorer dashboard – news outlets are weighted in a
600 proprietary (and undisclosed) tier system, while retweets are assigned only 75% of the weight of
601 original tweets [6]. The high correlation between the social score and the AAS thus reassures users
602 that these nuances make little difference.

603

604 Changes over time were evaluated in a 1-year follow-up of the reference Phase III sample. The
605 minimal change in the social impact component further underlines the similarity of this component
606 to the AAS, and supports the notion that news article and tweet metrics accrue soon after

607 publication. Both scholarly and societal impact scores continued to increase, and further follow-up is
608 needed to identify the point at which these scores plateau.

609

610 Finally, an independent dataset was investigated: articles selected by NEJM editors for their practice-
611 changing potential. These papers had substantially higher societal impact than the benchmark set of
612 NEJM Phase III articles, supporting the sensitivity of the societal impact component in identifying
613 practice-changing publications. Furthermore, innovative articles were found to have relatively low
614 societal impact, indicating that although these are of interest to scholars and wider society, they do
615 not directly feed into clinical practice changes. Conversely, articles on surgery had a high impact on
616 practice even though social and academic interest was low.

617

618 **Weaknesses**

619 Although the EMPIRE Index provides advantages over existing metric approaches, it has some
620 potential weaknesses. For example, grouping and value weighting have a large subjective
621 component that may not reflect the value assigned to metrics by others. However, the transparent
622 nature of the approach will hopefully stimulate further debate and discussion around the inherent
623 subjectivity and allow for future refinements.

624

625 The analyses conducted were based on a closely defined subset of medical publications, in terms of
626 both content (Phase III trials) and publication date. As metrics evolve over time owing to changes in
627 the way audiences engage with publications or technical advances in the way metrics are recorded,
628 these original analyses and assumptions may not apply. They may also not apply to other publication
629 types or study designs, and may vary across disease areas. Predictor scores are based on results of
630 cross-sectional, rather than longitudinal, analyses; further follow-up will allow these scores to be
631 refined and improved. Furthermore, benchmarking to very high-impact articles results in predictor
632 scores that tend to overestimate the final impact of more usual articles.

633

634 Lastly, although the scoring system is transparent and reproducible, it depends on metrics
635 aggregated by two different proprietary systems. These metrics may not be available to all intended
636 users of the index.

637

638 **Conclusions**

639 The EMPIRE Index is a novel metric framework incorporating three component scores that respond
640 to different types of publication impact: social, scholarly, and societal. Whereas the social impact
641 score is similar to the AAS and the scholarly impact score is closely linked to (but broader than)
642 article citations, the societal impact score reflects a key and distinct aspect of publication impact. In
643 a similar way to the AAS, the EMPIRE Index weights metrics subjectively to reflect their value from
644 the user's perspective as well as by prevalence. Unlike the AAS, it is designed for a limited subject
645 area (medicine) and weights and benchmarks the metrics accordingly. It also has a clear, transparent
646 explanation of the scoring system, and provides predictor scores to give an early estimate of likely
647 future impact.

648

649 Several potential uses are envisaged for the EMPIRE Index. Because it provides a richer assessment
650 of publication value than standalone traditional and alternative metrics, it will enable individuals
651 involved in medical research to assess the impact of related publications easily and to understand
652 what characterizes impactful research. It can also be used to assess the effectiveness of
653 communications around publications and publication enhancements such as infographics and
654 explanatory videos. Fuller validation of the EMPIRE Index requires additional prospective and cross-
655 sectional studies, which are ongoing.

656

657 **Acknowledgments**

658 The authors thank Heather Lang of Oxford PharmaGenesis, Oxford, UK for providing valuable
659 insights in the initial design phases of the project and for co-facilitating the stakeholder insights
660 workshop. Editorial support (manuscript proofreading, figure drawing, and project management)
661 was provided by Oxford PharmaGenesis.

662

663 **References**

- 664 1. Haustein S, Larivière V. The use of bibliometrics for assessing research: possibilities,
665 limitations and adverse effects. In: Welpel I, Wollersheim J, Ringelhan S, Osterloh M, editors.
666 Incentives and performance: governance of research organizations. Springer, Cham; 2015. pp.
667 121–139. doi:10.1007/978-3-319-09785-5_8
- 668 2. Raff JW. The San Francisco declaration on research assessment. *Biology Open*. 2013. pp. 533–
669 534. doi:10.1242/bio.20135330
- 670 3. Mohammadi E, Thelwall M. Readership Data and Research Impact. 2019 [cited 9 Feb 2019].
671 Available: <https://arxiv.org/ftp/arxiv/papers/1901/1901.08593.pdf>
- 672 4. Tahamtan I, Bornmann L. Altmetrics and societal impact measurements: Match or mismatch?
673 A literature review. *Prof la Inf*. 2020;29. doi:10.3145/epi.2020.ene.02
- 674 5. Plum Analytics. About PlumX Metrics. 2020 [cited 4 Jan 2021]. Available:
675 <https://plumanalytics.com/learn/about-metrics/>
- 676 6. Altmetric.com. How is the Altmetric Attention Score calculated? In: Altmetric.com [Internet].
677 2020 [cited 1 Jan 2021]. Available:
678 [https://help.altmetric.com/support/solutions/articles/6000060969-how-is-the-altmetric-](https://help.altmetric.com/support/solutions/articles/6000060969-how-is-the-altmetric-attention-score-calculated-)
679 [attention-score-calculated-](https://help.altmetric.com/support/solutions/articles/6000060969-how-is-the-altmetric-attention-score-calculated-)
- 680 7. Copiello S. Multi-criteria altmetric scores are likely to be redundant with respect to a subset
681 of the underlying information. *Scientometrics*. 2020;124: 819–824. doi:10.1007/s11192-020-
682 03491-9
- 683 8. Luc JGY, Percy E, Hirji S, Vervoort D, Mann GK, Phan K, et al. Predictors of High-Impact Articles
684 in The Annals of Thoracic Surgery. *Ann Thorac Surg*. 2020;110: 2096–2103.
685 doi:10.1016/j.athoracsur.2020.04.102
- 686 9. Studenic P, Ospelt C. Do you tweet?: Trailing the connection between Altmetric and research
687 impact! *RMD Open*. 2020;6: e001034. doi:10.1136/rmdopen-2019-001034

- 688 10. Thelwall M, Nevill T. Could scientists use Altmetric.com scores to predict longer term citation
689 counts? *J Informetr.* 2018;12: 237–248. doi:10.0.3.248/j.joi.2018.01.008
- 690 11. Bornmann L, Haunschild R. Normalization of zero-inflated data: An empirical analysis of a new
691 indicator family and its use with altmetrics data. *J Informetr.* 2018;12: 998–1011.
692 doi:10.1016/j.joi.2018.01.010
- 693 12. Costas R, Zahedi Z, Wouters P. Do “altmetrics” correlate with citations? Extensive comparison
694 of altmetric indicators with citations from a multidisciplinary perspective. *J Assoc Inf Sci
695 Technol.* 2015;66: 2003–2019. doi:10.1002/asi.23309
- 696 13. Mukherjee B, Subotić S, Chaubey AK. And now for something completely different: the
697 congruence of the Altmetric Attention Score’s structure between different article groups.
698 *Scientometrics.* 2018;114: 253–275. doi:10.1007/s11192-017-2559-8
- 699 14. Buttlere B, Buder J. Personalizing papers using Altmetrics: comparing paper ‘Quality’ or
700 ‘Impact’ to person ‘Intelligence’ or ‘Personality.’ *Scientometrics.* 2017;111: 219–239.
701 doi:10.1007/s11192-017-2246-9
- 702 15. Hassan SU, Imran M, Gillani U, Aljohani NR, Bowman TD, Didegah F. Measuring social media
703 activity of scientific literature: an exhaustive comparison of scopus and novel altmetrics big
704 data. *Scientometrics.* 2017;113: 1037–1057. doi:10.1007/s11192-017-2512-x
- 705 16. Nocera AP, Boyd CJ, Boudreau H, Hakim O, Rais-Bahrami S. Examining the Correlation
706 Between Altmetric Score and Citations in the Urology Literature. *Urology.* 2019;134: 45–50.
707 doi:10.1016/j.urology.2019.09.014
- 708 17. Hong JH, Yoon DY, Lim KJ, Moon JY, Baek S, Seo YL, et al. Characteristics of the most cited,
709 most downloaded, and most mentioned articles in general medical journals: a comparative
710 bibliometric analysis. *Healthcare.* 2020;8: 492. doi:10.3390/healthcare8040492
- 711 18. Ortega JL. Proposal of composed altmetric indicators based on prevalence and impact
712 dimensions. *J Informetr.* 2020;14: 1–18. doi:10.1016/j.joi.2020.101071
- 713 19. Haustein S, Costas R, Larivière V. Characterizing social media metrics of scholarly papers: The

- 714 effect of document properties and collaboration patterns. PLoS One. 2015;10: e0120495.
715 doi:10.1371/journal.pone.0120495
- 716 20. Zahedi Z, Haustein S. On the relationships between bibliographic characteristics of scientific
717 documents and citation and Mendeley readership counts: A large-scale analysis of Web of
718 Science publications. J Informetr. 2018;12: 191–202. doi:10.1016/j.joi.2017.12.005
- 719 21. van Eck NJ, Waltman L, van Raan AFJ, Klautz RJM, Peul WC. Citation Analysis May Severely
720 Underestimate the Impact of Clinical Research as Compared to Basic Research. PLoS One.
721 2013;8: e62395. doi:10.1371/journal.pone.0062395
- 722 22. Thelwall M. Interpreting correlations between citation counts and other indicators.
723 Scientometrics. 2016;108: 337–347. doi:10.1007/s11192-016-1973-7
- 724 23. About NEJM. [cited 4 Feb 2021]. Available: <https://www.nejm.org/about-nejm/about-nejm>
- 725 24. Drazan JM. Notable articles of 2016. A collection of articles selected by NEJM editors. 2016.
726 Available: <http://cdn.nejm.org/pdf/Notable-Articles-2016.pdf>
- 727 25. Drazan JM. Notable articles of 2016. A collection of articles selected by NEJM editors. 2016.
728 Available: <http://cdn.nejm.org/pdf/NEJM-Notable-Articles-2017.pdf>
- 729 26. Drazan JM. Notable articles of 2018. A collection of articles selected by NEJM editors. 2018.
730 Available: <http://cdn.nejm.org/pdf/Notable-Articles-2018.pdf>
- 731 27. Rubin EJ. Notable Articles of 2019. A collection of articles selected by NEJM editors. 2019.
732 Available: <https://cdn.nejm.org/pdf/Notable-Articles-of-2019.pdf>
- 733 28. Pubstrat, Journal Database | Anju Software. [cited 13 Jan 2021]. Available:
734 <https://www.anjusoftware.com/solutions/medical-affairs/pubstrat/pubstrat-journalselector>
- 735 29. Scopus Citescore. [cited 13 Jan 2021]. Available: <https://www.scopus.com/sources>
- 736 30. SJR : Scientific Journal Rankings. [cited 13 Jan 2021]. Available:
737 <https://www.scimagojr.com/journalrank.php>
- 738 31. Kumar BS, Basu A, Singh VK, Gupta S, Muhuri PK. Do “altmetric mentions” follow Power
739 Laws? Evidence from social media mention data in Altmetric.com. [cited 18 Dec 2020].

- 740 Available: <https://arxiv.org/abs/2011.09079>
- 741 32. Thelwall M. The discretised lognormal and hooked power law distributions for complete
742 citation data: Best options for modelling and regression. *J Informetr.* 2016;10: 336–346.
743 doi:10.1016/j.joi.2015.12.007
- 744 33. Bellego C, Pape L-D. Dealing with the log of zero in regression models. *Série des Doc Trav.*
745 2019; 16. doi:10.2139/ssrn.3444996
- 746 34. Fang Z, Costas R. Studying the accumulation velocity of altmetric data tracked by
747 Altmetric.com. *Scientometrics.* 2020;123: 1077–1101. doi:10.1007/s11192-020-03405-9
- 748 35. Ortega JL. The life cycle of altmetric impact: A longitudinal study of six metrics from PlumX. *J*
749 *Informetr.* 2018;12: 579–589. doi:10.1016/j.joi.2018.06.001
- 750 36. Hicks D, Wouters P, Waltman L, de Rijcke S, Rafols I. The Leiden Manifesto for research
751 metrics. *Nature.* 2015;520: 429–31. doi:10.1038/520429a
- 752 37. Sathianathen NJ, Iii RL, Murphy DG, Loeb S, Bakker C, Lamb AD, et al. Social Media Coverage
753 of Scientific Articles Immediately after Publication Predicts Subsequent Citations-#SoME-
754 Impact Score: Observational Analysis. *J Med Internet Res.* 2020;22: e12288.
755 doi:10.2196/12288
- 756 38. Ortega JL. Disciplinary differences of the impact of altmetric. *FEMS Microbiol Lett.* 2018;365:
757 49. doi:10.1093/femsle/fny049
- 758 39. Zahedi Z, Costas R. General discussion of data quality challenges in social media metrics:
759 Extensive comparison of four major altmetric data aggregators. *PLoS One.* 2018;13:
760 e0197326. doi:10.1371/journal.pone.0197326
- 761 40. Ortega JL. Reliability and accuracy of altmetric providers: a comparison among Altmetric.com,
762 PlumX and Crossref Event Data. *Scientometrics.* 2018;116: 2123–2138. doi:10.1007/s11192-
763 018-2838-z
- 764 41. Meschede C, Siebenlist T. Cross-metric compatibility and inconsistencies of altmetrics.
765 *Scientometrics.* 2018;115: 283–297. doi:10.1007/s11192-018-2674-1

- 766 42. Ortega JL. Blogs and news sources coverage in altmetrics data providers: a comparative
767 analysis by country, language, and subject. *Scientometrics*. 2020; 555–572.
768 doi:10.1007/s11192-019-03299-2
- 769 43. Haneef R, Ravaud P, Baron G, Ghosn L, Boutron I. Factors associated with online media
770 attention to research: a cohort study of articles evaluating cancer treatments. *Res Integr Peer*
771 *Rev*. 2017;2: 9. doi:10.1186/s41073-017-0033-z
- 772 44. Robinson-García N, Torres-Salinas D, Zahedi Z, Costas R. New data, new possibilities:
773 Exploring the insides of altmetric.com. *Prof la Inf*. 2014;23: 359–366.
774 doi:10.3145/epi.2014.jul.03
- 775 45. Maggio LA, Leroux TC, Meyer HS, Artino AR. #MedEd: exploring the relationship between
776 altmetrics and traditional measures of dissemination in health professions education.
777 *Perspect Med Educ*. 2018;7: 239–247. doi:10.1007/s40037-018-0438-5
- 778 46. Sheikh AM, Brown SK. Tracing the path from social attention to scientific impact. *Cardiovasc*
779 *Res*. 2019;115: e172–e176. doi:10.1093/cvr/cvz276
- 780 47. Ebrahimi S, Setareh F. Direct and indirect influence of altmetrics on citation in social systems:
781 Assessing a new conceptual model. *Int J Inf Sci Manag*. 2018;16: 161–173.
- 782 48. Nuzzolese AG, Ciancarini P, Gangemi A, Peroni S, Poggi F, Presutti V. Do altmetrics work for
783 assessing research quality? *Scientometrics*. 2019;118: 539–562. doi:10.1007/s11192-018-
784 2988-z
- 785 49. Bornmann L, Haunschild R. Do altmetrics correlate with the quality of papers? A large-scale
786 empirical study based on F1000Prime data. *PLoS One*. 2018;13: e0197133.
787 doi:10.1371/journal.pone.0197133
- 788 50. Pulido CM, Redondo-Sama G, Sordé-Martí T, Flecha R. Social impact in social media: A new
789 method to evaluate the social impact of research. *PLoS One*. 2018;13: e0203117.
790 doi:10.1371/journal.pone.0203117
- 791 51. Robinson-Garcia N, Costas R, Isett K, Melkers J, Hicks D. The unbearable emptiness of

- 792 tweeting — about journal articles. PLoS One. 2017;12: e0183551.
793 doi:10.1371/journal.pone.0183551
- 794 52. Haustein S. Scholarly twitter metrics. In: Glänzel W, Moed H, Schmoch U, Thelwall M, editors.
795 Handbook of quantitative science and technology research. Springer; 2019. pp. 729–760.
796 doi:10.1007/978-3-030-02511-3_28
- 797 53. Said A, Bowman TD, Abbasi RA, Aljohani NR, Hassan SU, Nawaz R. Mining network-level
798 properties of Twitter altmetrics data. Scientometrics. 2019;120: 217–235.
799 doi:10.1007/s11192-019-03112-0
- 800 54. Côté IM, Darling ES. Scientists on Twitter: Preaching to the choir or singing from the
801 rooftops? Facets. 2018;3: 682–694. doi:10.1139/facets-2018-0002
- 802 55. Mohammadi E, Barahmand N, Thelwall M. Who shares health and medical scholarly articles
803 on Facebook? Learn Publ. 2019; 1–9. doi:10.1002/leap.1271
- 804 56. Bornmann L. Usefulness of altmetrics for measuring the broader impact of research: A case
805 study using data from PLOS and F1000Prime. Aslib J Inf Manag. 2015;67: 305–319.
806 doi:10.1108/AJIM-09-2014-0115
- 807 57. Smith ZL, Chiang AL, Bowman D, Wallace MB. Longitudinal relationship between social media
808 activity and article citations in the journal Gastrointestinal Endoscopy. Gastrointest Endosc.
809 2019;90: 77–83. doi:10.1016/j.gie.2019.03.028
- 810 58. Thelwall M. Early Mendeley readers correlate with later citation counts. Scientometrics.
811 2018;115: 1231–1240. doi:10.1007/s11192-018-2715-9
- 812 59. Abramo G, D’Angelo CA, Felici G. Predicting publication long-term impact through a
813 combination of early citations and journal impact factor. J Informetr. 2019;13: 32–49.
814 doi:10.1016/j.joi.2018.11.003
- 815 60. Stegehuis C, Litvak N, Waltman L. Predicting the long-term citation impact of recent
816 publications. J Informetr. 2015;9: 642–657. doi:10.1016/j.joi.2015.06.005
817

818 **Supporting information**

819 **S1 Table. Summary statistics for metrics obtained via (A) Altmetric, (B) PlumX, and (C) journal-**
 820 **level, citation-based indices.** Coverage is the proportion of articles with > 0 on that metric.

821 **(A)**

Metric	Coverage	Mean	Median	Maximum
News mentions	0.35	6.4	0	365
Blog mentions	0.14	0.4	0	27
Policy mentions	0.08	0.1	0	4
Patent mentions	0.01	0.0	0	8
Twitter mentions	0.85	32.6	5	1701
Original tweets	0.85	10.5	3	347
Retweets	0.59	22.2	1	1533
Peer review mentions	0.00	0.0	0	2
Weibo mentions	0.00	0.0	0	1
Facebook mentions	0.33	1.3	0	55
Wikipedia mentions	0.02	0.0	0	8
Google+ mentions	0.08	0.1	0	10
LinkedIn mentions	0.00	0.0	0	0
Reddit mentions	0.04	0.0	0	6
Pinterest mentions	0.00	0.0	0	0
F1000Prime mentions	0.07	0.1	0	3
Q&A mentions	0.00	0.0	0	1
Video mentions	0.01	0.0	0	4
Syllabi mentions	0.00	0.0	0	0
Mendeley readers	0.99	43.2	24	1325
Dimensions citations	0.94	27.3	9	2021

822

823 **(B)**

Metric	Coverage	Mean	Median	Maximum
Captures:Exports-Saves:EBSCO	0.47	6.0	0	456
Captures:Readers:Mendeley	0.80	41.6	19	1286
Captures:Readers:CiteULike	0.01	0.0	0	3
Citations:Clinical Citations:PubMed Guidelines	0.05	0.1	0	3
Citations:Clinical Citations:DynaMed Plus	0.16	0.2	0	5
Citations:Citation Indexes:Scopus	0.92	24.9	8	1971
Citations:Citation Indexes:PubMed	0.23	3.9	0	795
Citations:Citation Indexes:CrossRef	0.73	9.3	2	454
Social Media:Tweets:Twitter	0.78	28.0	3	1892
Social Media:Shares, Likes & Comments:Facebook	0.21	36.5	0	39,422
Mentions:References:Wikipedia	0.03	0.0	0	7

NOMF003

EMPIRE Index development

6 April 2020

Mentions:Blog	0.11	0.2	0	13
Mentions:Comments:Reddit	0.00	0.0	0	17
Mentions:News Mentions:News	0.29	3.1	0	261
Usage:Views:Figshare	0.01	0.3	0	68
Usage:Abstract Views:DSpace	0.00	0.5	0	315
Usage:Abstract Views:SciELO	0.00	0.1	0	87
Usage:Abstract Views:Digital Commons	0.02	0.2	0	58
Usage:Abstract Views:Expert Gallery Suite	0.00	0.0	0	8
Usage:Abstract Views:EBSCO	0.71	68.9	16	3124
Usage:Clicks:Bitly	0.13	7.6	0	2862
Usage:Downloads:Figshare	0.01	0.2	0	76
Usage:Downloads:UWA Research Repository	0.00	0.0	0	6
Usage:Downloads:Digital Commons	0.01	0.5	0	220
Usage:Downloads:Expert Gallery Suite	0.00	0.0	0	7
Usage:Downloads:EBSCO	0.00	0.0	0	1
Usage:Full Text Views:PubMedCentral	0.02	9.9	0	2119
Usage:Full Text Views:SciELO	0.00	0.8	0	1064
Usage:Full Text Views:PLoS	0.02	27.5	0	7710
Usage:Full Text Views:EBSCO	0.13	3.7	0	2663
Usage:Link-outs:EBSCO	0.56	7.1	1	374

824

825 (C)

Metric	Density	Mean	Median	Maximum
CiteScore	1.00	5.3	4.03	19.14
Journal Impact Factor	0.91	14.1	5.231	79.258
Scimago Journal Ranking	0.79	4.6	2.339	19.476

826

827 **S2 Table. Correlations (Spearman’s r) between investigational metrics in the sample of Phase III clinical trial publications.** Correlations > 0.5 are
 828 shown in bold.

829

Metric	News mentions	Blog mentions	Policy mentions	Patent mentions	Twitter mentions	Original tweets	Retweets	Facebook mentions	Wikipedia mentions	F1000 mentions	Mendeley readers	Dimensions citations	PubMed guidelines
News mentions	–	0.529	0.165	0.066	0.530	0.535	0.519	0.479	0.160	0.302	0.504	0.518	0.187
Blog mentions	0.529	–	0.136	0.062	0.455	0.461	0.459	0.476	0.191	0.327	0.406	0.373	0.127
Policy mentions	0.165	0.136	–	0.044	0.116	0.121	0.120	0.130	0.086	0.101	0.185	0.220	0.134
Patent mentions	0.066	0.062	0.044	–	0.048	0.050	0.045	0.048	0.073	0.038	0.073	0.110	0.098
Twitter mentions	0.530	0.455	0.116	0.048	–	0.958	0.932	0.555	0.140	0.270	0.539	0.521	0.140
Original tweets	0.535	0.461	0.121	0.050	0.958	–	0.819	0.566	0.140	0.266	0.552	0.522	0.145
Retweets	0.519	0.459	0.120	0.045	0.932	0.819	–	0.534	0.145	0.283	0.509	0.505	0.138
Facebook mentions	0.479	0.476	0.130	0.048	0.555	0.566	0.534	–	0.155	0.288	0.448	0.417	0.130
Wikipedia mentions	0.160	0.191	0.086	0.073	0.140	0.140	0.145	0.155	–	0.163	0.130	0.138	0.057
F1000Prime mentions	0.302	0.327	0.101	0.038	0.270	0.266	0.283	0.288	0.163	–	0.264	0.284	0.118
Mendeley readers	0.504	0.406	0.185	0.073	0.539	0.552	0.509	0.448	0.130	0.264	–	0.737	0.221
Dimensions citations	0.518	0.373	0.220	0.110	0.521	0.522	0.505	0.417	0.138	0.284	0.737	–	0.274
PubMed guidelines mentions	0.187	0.127	0.134	0.098	0.140	0.145	0.138	0.130	0.057	0.118	0.221	0.274	–

830

831 **S3 Table. Three-factor analysis of included metrics in (A) the full sample, (B) older papers (1H), and**
 832 **(C) younger papers (2H). Highest loadings for each metric are shown in bold.**

833 **(A)**

Metric	1	2	3
News mentions	-0.54	0.19	0.14
Blog mentions	-0.88	-0.02	-0.01
Policy mentions	-0.18	0.32	-0.17
Patent mentions	-0.12	0.14	-0.14
Twitter mentions	-0.02	0.06	0.84
Facebook mentions	-0.43	0.01	0.43
Wikipedia mentions	-0.33	0.05	-0.09
F1000Prime mentions	-0.42	0.11	0.00
Mendeley readers	-0.06	0.65	0.17
Dimensions citations	0.02	0.96	-0.01
PubMed guidelines mentions	-0.07	0.40	-0.17

834 **(B)**

Metric	1	2	3
News mentions	-0.19	0.04	-0.58
Blog mentions	0.00	-0.01	-0.84
Policy mentions	-0.29	-0.06	-0.19
Patent mentions	-0.13	-0.05	-0.09
Twitter mentions	-0.08	0.87	0.00
Facebook mentions	0.09	0.42	-0.48
Wikipedia mentions	-0.06	-0.06	-0.31
F1000Prime mentions	-0.13	-0.04	-0.45
Mendeley readers	-0.65	0.13	-0.05
Dimensions citations	-0.93	0.01	0.00
PubMed guidelines mentions	-0.41	-0.10	-0.08

835 **(C)**

Metric	1	2	3
News mentions	-0.11	-0.25	0.52
Blog mentions	0.03	0.04	0.92
Policy mentions	0.16	-0.26	0.12
Patent mentions	0.11	-0.03	0.15
Twitter mentions	-0.98	0.00	0.01
Facebook mentions	-0.26	-0.12	0.54
Wikipedia mentions	0.09	-0.05	0.30
F1000Prime mentions	-0.03	-0.09	0.39
Mendeley readers	-0.07	-0.71	0.11
Dimensions citations	0.02	-0.96	-0.04
PubMed guidelines mentions	0.04	-0.25	-0.03

836

837 **S4 Table. Two-factor analysis of metrics excluding citations in policy documents, PubMed**
838 **guidelines, and patents in (A) the full sample, (B) older papers (1H), and (C) younger papers (2H).**
839 Highest loadings for each metric are shown in bold.

840 **(A)**

Metric	1	2
News mentions	0.68	-0.13
Blog mentions	0.86	0.09
Twitter mentions	0.65	-0.19
Facebook mentions	0.84	0.02
Wikipedia mentions	0.27	0.01
F1000Prime mentions	0.44	-0.05
Mendeley readers	0.14	-0.73
Dimensions citations	-0.04	-0.93

841

842 **(B)**

Metric	1	2
News mentions	0.15	0.62
Blog mentions	-0.04	0.81
Twitter mentions	0.34	0.49
Facebook mentions	-0.05	0.81
Wikipedia mentions	0.01	0.28
F1000Prime mentions	0.07	0.44
Mendeley readers	0.77	0.06
Dimensions citations	0.91	-0.02

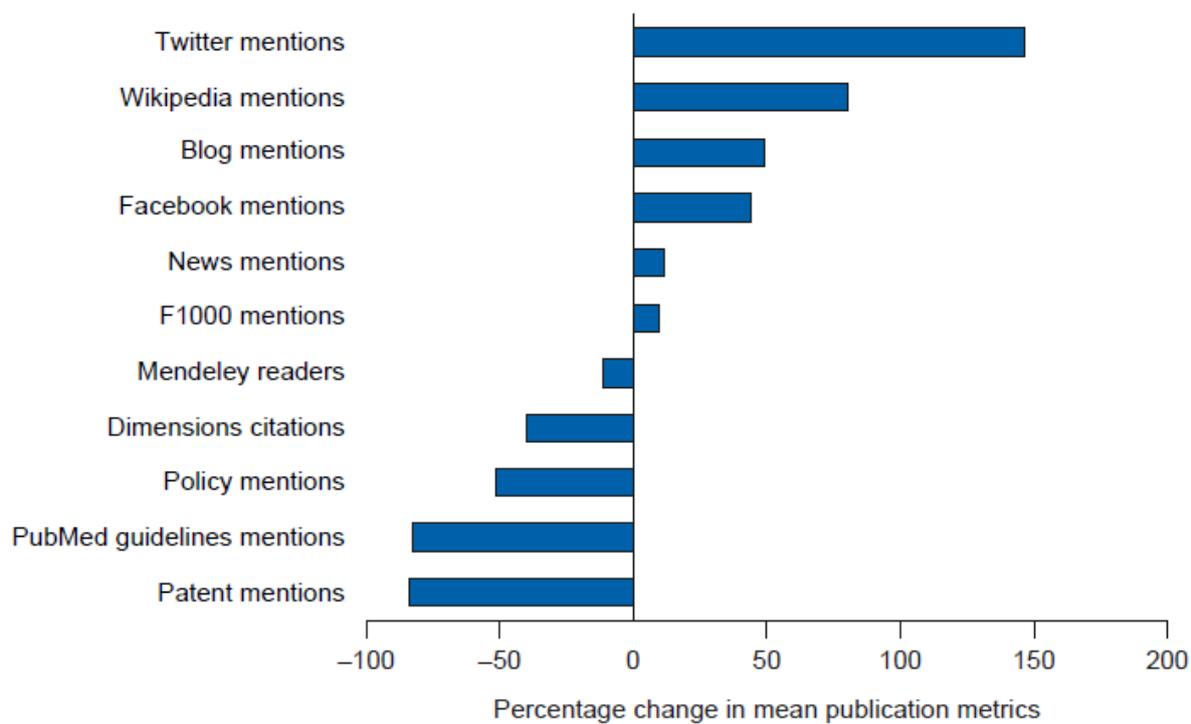
843

844 **(C)**

Metric	1	2
News mentions	-0.66	-0.18
Blog mentions	-0.91	0.12
Twitter mentions	-0.63	-0.24
Facebook mentions	-0.83	-0.04
Wikipedia mentions	-0.25	0.01
F1000 mentions	-0.46	-0.02
Mendeley readers	-0.17	-0.71
Dimensions citations	0.03	-0.93

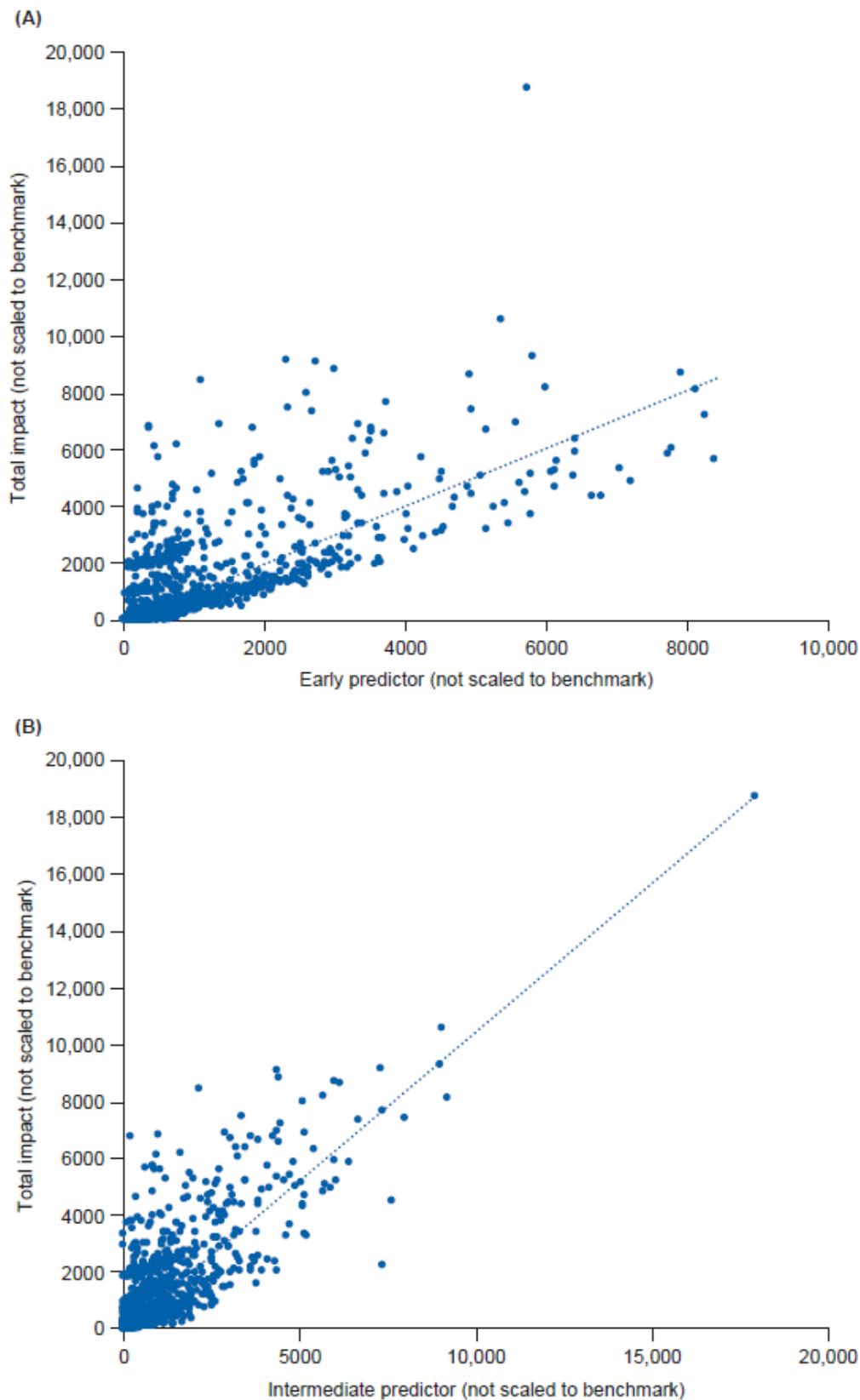
845

846 **S1 Fig. Percentage change in mean publication metrics in the more recent half of the publications**
847 **(2H) versus the older half (1H).**



848

849 **S2 Fig. Correlation of total impact scores with (A) early predictor and (B) intermediate predictor**
850 **scores. Scores shown are not adjusted to the benchmark.**



851