

1 **Human gene function publications that describe wrongly identified nucleotide sequence**
2 **reagents are unacceptably frequent within the genetics literature**

3 **Short title: Wrongly identified nucleotide sequences in gene function papers**

4 Yasunori Park¹, Rachael A West^{1,2}, Pranuajan Pathmendra¹, Bertrand Favier³, Thomas
5 Stoeger^{4,5,6}, Amanda Capes-Davis^{1,7}, Guillaume Cabanac⁸, Cyril Labbé⁹, Jennifer A
6 Byrne^{1,10,*}

7 ¹Faculty of Medicine and Health, The University of Sydney, NSW, Australia

8 ²Children's Cancer Research Unit, Kids Research, The Children's Hospital at Westmead,
9 Westmead, NSW, Australia

10 ³Univ. Grenoble Alpes, TIMC, Grenoble, France

11 ⁴Successful Clinical Response in Pneumonia Therapy (SCRIPT) Systems Biology Center,
12 Northwestern University, Evanston, United States.

13 ⁵Department of Chemical and Biological Engineering, Northwestern University, Evanston,
14 United States.

15 ⁶Center for Genetic Medicine, Northwestern University School of Medicine, Chicago, United
16 States

17 ⁷CellBank Australia, Children's Medical Research Institute, Westmead, New South Wales,
18 Australia

19 ⁸Computer Science Department, IRIT UMR 5505 CNRS, University of Toulouse, 118 route
20 de Narbonne, 31062 Toulouse Cedex 9, France

21 ⁹Univ. Grenoble Alpes, CNRS, Grenoble INP, LIG, Grenoble, France

22 ¹⁰NSW Health Statewide Biobank, NSW Health Pathology, Camperdown, NSW, Australia

23 *Corresponding author

24

25 **Keywords:** cancer; gene function; miRNA; non-coding RNA's; nucleotide sequence reagent;

26 paper mill; protein-coding gene

27

28 **Abstract**

29 Nucleotide sequence reagents underpin a range of molecular genetics techniques that have
30 been applied across hundreds of thousands of research publications. We have previously
31 reported wrongly identified nucleotide sequence reagents in human gene function
32 publications and described a semi-automated screening tool Seek & Blastn to fact-check the
33 targeting or non-targeting status of nucleotide sequence reagents. We applied Seek & Blastn
34 to screen 11,799 publications across 5 literature corpora, which included all original
35 publications in *Gene* from 2007-2018 and all original open-access publications in *Oncology*
36 *Reports* from 2014-2018. After manually checking the Seek & Blastn screening outputs for
37 over 3,400 human research papers, we identified 712 papers across 78 journals that described
38 at least one wrongly identified nucleotide sequence. Verifying the claimed identities of over
39 13,700 nucleotide sequences highlighted 1,535 wrongly identified sequences, most of which
40 were claimed targeting reagents for the analysis of 365 human protein-coding genes and 120
41 non-coding RNAs, respectively. The 712 problematic papers have received over 17,000
42 citations, which include citations by human clinical trials. Given our estimate that
43 approximately one quarter of problematic papers are likely to misinform or distract the future
44 development of therapies against human disease, urgent measures are required to address the
45 problem of unreliable gene function papers within the literature.

46

47 **Author summary**

48 This is the first study to have screened the gene function literature for nucleotide sequence
49 errors at the scale that we describe. The unacceptably high rates of human gene function
50 papers with incorrect nucleotide sequences that we have discovered represent a major
51 challenge to the research fields that aim to translate genomics investments to patients, and
52 that commonly rely upon reliable descriptions of gene function. Indeed, wrongly identified
53 nucleotide sequence reagents represent a double concern, as both the incorrect reagents
54 themselves and their associated results can mislead future research, both in terms of the
55 research directions that are chosen and the experiments that are undertaken. We hope that our
56 research will inspire researchers and journals to seek out other problematic human gene
57 function papers, as we are unfortunately concerned that our results represent the tip of a much
58 larger problem within the literature. We hope that our research will encourage more rigorous
59 reporting and peer review of gene function results, and we propose a series of responses for
60 the research and publishing communities.

61

62 **Introduction**

63 The promise of genomics to improve the health of cancer and other patients has resulted in
64 billions of dollars of research investment which have been accompanied by expectations of
65 similar quantum gains in health outcomes (1, 2). Since the first draft of the human genome
66 was reported (3, 4), a series of increasingly rapid technological advances has permitted the
67 routine sequencing of human genomes at scale (1, 2), and the increasing application of
68 genomics to inform clinical care (1, 2, 5). Despite the now routine capacity to sequence the
69 human genome, genomics research relies upon results produced by other research fields to
70 translate genome sequencing results to patients (5-7). For example, while whole genome
71 sequencing demonstrates that thousands of human genes are mutated or deregulated in human
72 cancers (1), knowledge of human gene function is required to prioritise individual gene
73 candidates for subsequent pre-clinical and translational studies (5-7).

74

75 A first step in triaging and prioritising gene candidates for further analysis is the
76 consideration of available knowledge of predicted and/or demonstrated gene functions (5-8).
77 High quality, reliable information about gene function is important to select the most
78 promising gene candidates and to then progress these candidates through pre-clinical and
79 translational research pipelines (8), which is supported by drug candidates with genetically
80 supported targets being significantly more likely to progress through phased clinical trials (9,
81 10). However, in contrast to the sophisticated platforms that produce genomic or
82 transcriptomic sequence data at scale, gene function experiments typically analyse single or
83 small numbers of genes through the application of more ubiquitous molecular techniques (6),
84 some of which have been in routine experimental use for 15-30 years. For example, gene
85 knockdown approaches have been widely employed to assess the consequences of reduced

86 gene expression in model systems (6). Similarly, RT-PCR is frequently used to analyse the
87 transcript levels of small groups of genes, either to confirm the effectiveness of gene
88 knockdown experiments or in association with other experimental techniques. The
89 widespread use and reporting of gene knockdown and PCR approaches reflect their low cost
90 and accessibility, in terms of the necessary reagents, laboratory equipment and facilities, and
91 the availability of technical expertise within the research community. As a consequence, the
92 results of experiments employing gene knockdown and/or PCR in the context of human
93 research have been described in hundreds of thousands of publications that are retrievable
94 through PubMed or Google Scholar.

95

96 Experiments that analyse the functions of individual genes typically require nucleotide
97 sequence reagents as either targeting and/or control reagents, with the sequences
98 corresponding to these reagents being disclosed to indicate the exact experiments that were
99 conducted (8, 11). As nucleotide sequence identities cannot be deduced by eye, DNA or RNA
100 reagent sequences must be paired with text descriptions of their genetic identities and
101 experimental use (8, 11-13). The integrity of reported experiments therefore requires both the
102 identities of nucleotide sequence reagents and their text descriptions to be correct (11, 12).
103 Accurate reporting of nucleotide sequence reagents is also critical to permit reagent reuse
104 across different experiments and publications (11).

105

106 The ubiquitous description of nucleotide sequence reagents within the biomedical and
107 genetics literature, combined with the routine pairing of nucleotide sequences and text
108 identifiers, are likely to contribute to tacit assumptions that reported nucleotide sequence
109 reagents are correctly identified. However, as nucleotide sequences cannot be understood by

110 eye, we have proposed that nucleotide sequence reagents are susceptible to different types of
111 errors (8, 11, 12). These error types represent the equivalent of spelling errors (12, 14, 15), as
112 well as identity errors, where a correct sequence is replaced by a different and possibly
113 genetically unrelated sequence (11-13, 16-21).

114

115 The problem of wrongly identified nucleotide sequences was recognised in the context of
116 DNA microarrays in the early 2000's, where wrongly identified sequence probes affected the
117 reliability and reproducibility of data from particular microarray platforms (22, 23). Our team
118 subsequently identified frequent wrongly identified sh/siRNA's and RT-PCR primers in
119 studies that reported the effects of knocking down single human genes in cancer cell lines
120 (12, 13). Some incorrect sh/siRNA's invalidated the experimental results reported, through
121 for example the repeated use of "non-targeting" shRNA's that were in fact indicated to target
122 specific human genes (11-13). Our discovery of incorrect and indeed impossible results
123 associated with wrongly identified nucleotide sequence reagents led us to develop a semi-
124 automated tool Seek & Blastn (S&B) to fact-check the reported identities of nucleotide
125 sequence reagents in human research publications (12, 24). The S&B tool scans text to
126 identify and extract nucleotide sequences and their associated text descriptors, submits
127 extracted nucleotide sequences to blastn analysis (25) to predict their genetic identities and
128 hence their targeting or non-targeting status, and then compares the predicted status of each
129 nucleotide sequence reagent with the claimed status within the text (12). The blastn results
130 for each extracted nucleotide sequence are then reported and any sequence whose text
131 identifier contradicts its blastn-predicted targeting/ non-targeting status is flagged as being
132 potentially incorrect (12). Flagged nucleotide sequences are then subjected to manual
133 verification, as described in our original publication (12) and an expanded online protocol
134 (26).

135

136 Our initial application of S&B identified 77/203 (38%) screened papers with incorrect
137 nucleotide sequence reagents, with our focus being the description of the S&B tool (12), as
138 opposed to its application. We have now employed S&B to screen original research papers
139 across 5 literature corpora, representing 3 targeted and two journal corpora. The 3 targeted
140 corpora included papers that were identified through literature searches that employed
141 specific keywords, and in some cases, PubMed similarity searches of index papers. The two
142 journal corpora included all original and original open-access papers in *Gene* and *Oncology*
143 *Reports*, respectively, as examples of journals that have published papers with wrongly
144 identified nucleotide sequences (11-13). S&B screening of 11,799 papers flagged 3,423
145 papers for further analysis, which required manually verifying the identities of over 13,700
146 nucleotide sequences. Our combined results provide very worrying evidence that human
147 research papers with wrongly identified nucleotide sequences are unacceptably frequent
148 within the literature.

149

150 **Results**

151 *Analysis of targeted publication corpora*

152 To extend our previous results from employing S&B to screen gene function papers (11, 12),
153 we used S&B to screen 3 targeted publication corpora (Fig 1) that were identified through
154 literature searches that employed specific keywords, and in some cases, PubMed similarity
155 searches of index papers (12, 13). In all cases, the keywords that were used to derive targeted
156 corpora did not refer to author affiliations, such as institution type or country of origin (see
157 Methods).

158

159 *Single gene knockdown (SGK) corpus*

160 We have previously identified frequent wrongly identified nucleotide sequences in papers
161 that reported the effects of knocking down a single human gene in cancer cell lines that
162 corresponded to a single cancer type, which we referred to as single gene knockdown (SGK)
163 papers (12, 13). As a subset of SGK papers analysed the same human genes across multiple
164 papers and cancer types (12, 13), we sought to identify all SGK papers for a set of 17 human
165 genes (*ADAM8*, *ANXA1*, *EAG1*, *GPR137*, *ICT1*, *KLF8*, *MACC1*, *MYO6*, *NOB1*, *PP4R1*, *PP5*,
166 *PPM1D*, *RPS15A*, *TCTN1*, *TPD52L2*, *USP39*, and *ZFX*), most of which had been analysed in
167 previously reported papers (12, 13).

168

169 By combining each of the 17 human gene identifiers with keywords previously used to
170 identify SGK papers (13), we identified 174 SGK papers published between 2006-2019
171 across 83 journals (Table 1). As most gene identifiers were selected from previously reported
172 SGK papers (12, 13), the SGK corpus consisted of 41 (24%) previously reported papers and

173 132 (76%) additional papers (Table 2). All 174 SGK papers analysed gene function in a
174 single human cancer type (Table 2). Across the 17 queried genes, we identified a median of 8
175 papers/ gene (range 3-20) that analysed a median of 8 human cancer types (range 3-11)
176 (Table 2, S1 Data). Most (136/174, 78%) SGK papers named a single queried gene in their
177 titles, with the remaining titles also referring to other human gene(s) and/or drugs, most
178 frequently cisplatin (S1 Data). Most (159/174, 91%) SGK papers were published by authors
179 from mainland China, almost all of which were affiliated with hospitals (147/159, 92%) (see
180 Methods) (Table 1). In contrast, less than half (6/15, 40%) SGK papers from 5 other countries
181 were affiliated with hospitals (Table 1, S1 Data).

182

183 S&B screening (24) flagged 144/174 (83%) SGK papers for further analysis (Fig 1). Manual
184 verification of the identities of all nucleotide sequences in flagged papers (see Methods)
185 confirmed that 75/174 (43%) SGK papers included 1-8 wrongly identified sequences/ paper
186 (Table 1, S1 Data). The 75 problematic SGK papers analysed 24 human cancer types, most
187 frequently brain cancer, where 1-9 problematic SGK papers were identified per queried gene
188 (Table 2). Whereas 31/75 (41%) problematic SGK papers have been reported in earlier
189 studies (12, 13), the remaining 45 SGK papers have not been previously analysed (Table 2).
190 The 75 problematic SGK papers were published across 42 journals, where Spandidos
191 Publications published the highest proportion (20/75, 27%) (S1 Data). Problematic SGK
192 papers described 115 wrongly identified sequences (Table 3), where half of these sequences
193 (57/115) targeted a gene or genomic sequence other than the claimed target, followed by
194 incorrect “non-targeting” reagents (44/115, 38%) (Fig 2, Table 3). The 71 incorrect targeting
195 sequences were claimed to interrogate 20 protein-coding genes (S1 Table). Most (67/115,
196 58%) incorrect sequences recurred across at least 2 SGK papers (Fig 3, S1 Table), where the
197 most frequent incorrect reagent was a previously described “non-targeting” shRNA that is

198 predicted to target *TPD52L2* (11-13). This shRNA or highly similar variants were employed
199 as “non-targeting” controls in 41% (31/75) problematic SGK papers (S1 Table).

200

201 *miR-145* corpus

202 PubMed similarity searches employing individual SGK papers identified numerous papers
203 that analysed the functions of different human miR’s in cancer cell lines. We selected *miR-*
204 *145* from many possible candidates to define a single miR corpus to be screened by S&B.
205 Papers that focussed upon *miR-145* were identified using PubMed similarity searches of
206 index papers (12, 13) and keyword searches of the Google Scholar database (see Methods). A
207 total of 163 *miR-145* papers were then screened by S&B to flag 50 *miR-145* papers for further
208 analysis (Fig 1). The 50 flagged *miR-145* papers were published between 2009-2019 across
209 35 journals (Table 1), and examined 18 human cancer types, where a single cancer type was
210 analysed in each paper (S2 Data). All flagged papers examined *miR-145* in combination with
211 1-5 other human gene(s) that were named in publication titles, with a minority of titles (5/50,
212 10%) also naming a single drug (S2 Data). Most (44/50, 88%) flagged *miR-145* papers were
213 published by authors from China, where most papers (40/44, 91%) were also affiliated with
214 hospitals (Table 1, S2 Data). The 6 *miR-145* papers from 5 other countries were affiliated
215 with institutions other than hospitals (Table 1, S2 Data).

216

217 Manual verification of S&B results revealed that most (31/50, 62%) flagged *miR-145* papers
218 described at least one wrongly identified sequence, with a median of one (range 1-5)
219 incorrect sequence/ paper (Table 1, S2 Data). The 31 problematic *miR-145* papers were
220 published from 2009-2019 across 25 journals, with the highest proportion published by Wiley
221 (S2 Data). The 31 problematic *miR-145* papers analysed 12 human cancer types, most

222 frequently colorectal or lung cancer, and described 49 wrongly identified sequences, most of
223 which (38/49, 78%) targeted a different gene or target from that claimed (Fig 2, Table 3). The
224 47 incorrect targeting sequences (Table 3) were claimed to interrogate 13 protein-coding and
225 4 non-coding RNA's (ncRNA's) (S1 Table). In contrast to SGK papers, most incorrect
226 sequences in *miR-145* papers were employed as (RT)-PCR primers (Table 3) and were
227 identified only once within the corpus (Fig 3). All problematic *miR-145* papers were
228 published by authors from China, where almost all papers (29/31, 94%) were affiliated with
229 hospitals (Table 1, S2 Data).

230

231 *Cisplatin and Gemcitabine (C+G) corpus*

232 We noted examples of SGK and *miR-145* papers that analysed the effects of drugs in cancer
233 cell lines (S1 Data, S2 Data). PubMed similarity searches of index papers (12, 13) combined
234 with keyword searches of the Google Scholar database were therefore employed to identify
235 papers that described either cisplatin or gemcitabine treatment of human cancer cell lines
236 and/or cancer patients (see Methods) (Fig 1). A total of 258 papers were screened by S&B to
237 flag 100 papers (n=50 papers for each drug) for further analysis as a combined cisplatin +
238 gemcitabine (C+G) corpus (Fig 1). The 100 flagged C+G papers were published between
239 2008-2019 across 48 journals (Table 1) and referred to a median of 2 (range 0-4) human
240 genes across their titles (S3 Data). The 100 flagged C+G papers examined 13 human cancer
241 types, where most (96/100, 96%) examined a single cancer type, typically pancreatic (35/100,
242 35%) or lung (22/100, 22%) cancer (S3 Data), reflecting the clinical use of cisplatin and
243 gemcitabine (27, 28). Most (90/100, 90%) C+G papers were published by authors from
244 China, where most (82/90, 91%) were also affiliated with hospitals (Table 1, S3 Data). In
245 contrast, 1/10 C+G papers from 8 other countries was hospital-affiliated (Table 1, S3 Data).

246

247 Approximately half (51/100, 51%) the flagged C+G papers were found to include a median
248 of 2 (range 1-8) wrongly identified sequences/ paper (Table 1). The 51 problematic C+G
249 papers were published between 2009-2019 across 31 journals (Table 1), where Springer
250 Nature published the highest proportion (15/51, 29%), followed by Elsevier (12/51, 24%) (S3
251 Data). The 51 problematic C+G papers examined 13 human cancer types, most frequently
252 pancreatic cancer, and described 109 wrongly identified nucleotide sequences, most of which
253 (79/109, 72%) targeted a gene or genomic sequence other than the claimed target (Fig 2,
254 Table 3, S1 Table). The 103 incorrect targeting sequences (Table 3) were claimed to
255 interrogate 31 protein-coding genes and 16 ncRNA's (S1 Table). As in *miR-145* papers, most
256 incorrect sequences in problematic C+G papers represented (RT)-PCR primers (Table 3) and
257 were identified once within the corpus (Fig 3). Almost all (50/51, 98%) problematic C+G
258 papers were published by authors from China, where almost all (48/50, 96%) were affiliated
259 with hospitals (Table 1, S3 Data).

260

261 ***Analysis of Gene and Oncology Reports corpora***

262 As papers in targeted corpora were selected using known features of human research papers
263 with incorrect nucleotide sequences, we complemented these analyses by screening original
264 research papers in the journals *Gene* and *Oncology Reports*. These journals were selected as
265 examples of journals that have published papers with incorrect nucleotide sequences (11-13),
266 where *Oncology Reports* also published the highest number of problematic SGK papers (S1
267 Data). S&B was employed to screen all original articles published in *Gene* from 2007-2018,
268 and all open-access articles published in *Oncology Reports* from 2014-2018 (Table 4).

269

270 Screening 7,399 original *Gene* articles from 2007-2018 flagged 742 (10%) papers for further
271 analysis (Fig 1) (see Methods). Manual verification of S&B outputs found that 17%
272 (128/742) flagged papers described a median of two (range 1-36) wrongly identified
273 sequences/ paper (Table 4, S4 Data). These 128 problematic papers referred to 186 human
274 genes (n=146 protein-coding, n=40 ncRNA's) across their publication titles (S4 Data).
275 Approximately half (65/128, 51%) the problematic *Gene* papers analysed gene function in
276 research contexts other than human cancer, most frequently by examining gene
277 polymorphisms in patient cohorts (16/65, 25%) (S4 Data). The remaining 60 papers analysed
278 17 different human cancer types, most frequently lung cancer (12/60, 20%) (S4 Data). A
279 minority of problematic *Gene* papers (7/128, 5%) referred to drugs within their titles. Manual
280 verification of over 5,200 sequences highlighted 284 wrongly identified sequences across the
281 128 problematic *Gene* papers. Almost all (275/284, 97%) incorrect sequences represented
282 targeting reagents (Fig 2, Table 3) for the analysis of 92 protein-coding genes and 24
283 ncRNA's (S1 Table). Most (261/279, 92%) incorrect sequences were described once within
284 the *Gene* corpus (Fig 3).

285

286 As *Oncology Reports* published many more papers per year than *Gene* from 2007-2018, we
287 employed S&B to screen open-access *Oncology Reports* articles from 2014-2018 (n=3,778
288 papers, 99% *Oncology Reports* articles) (Fig 1, Table 4). Almost half (1,709/ 3,778, 45%)
289 screened papers were flagged for further analysis (Fig 1), and over one quarter (436/1,709,
290 26%) of flagged papers were confirmed to describe a median of 2 (range 1-15) wrongly
291 identified sequences/ paper (Table 4, S5 Data). Almost all (432/436, 99%) problematic
292 *Oncology Reports* papers studied gene function in human cancer, most frequently lung
293 (54/432, 13%) or liver cancer (46/432, 11%). A subset (51/432, 12%) of problematic papers
294 referred to 42 different drugs across their titles, most frequently cisplatin or 5-fluorouracil (S5

295 Data). Manual verification of over 5,100 sequence identities confirmed 995 wrongly
296 identified sequences (S1 Table). Almost all (965/995, 97%) incorrect sequences represented
297 targeting reagents (Fig 2, Table 3) for the analysis of 262 protein-coding genes and 86
298 ncRNA's (S1 Table). Most (816/965, 85%) incorrect sequences were described once across
299 the *Oncology Reports* corpus (Fig 3).

300

301 *Geographic, institutional and temporal distributions of problematic Gene and Oncology*
302 *Reports papers*

303 The 128 problematic *Gene* papers were authored by teams from 19 countries (S1A Fig) (see
304 Methods). Just over half (69/128, 54%) problematic *Gene* papers were authored by teams
305 from China (Table 4, Fig 4, S1B Fig), followed by India (10/128, 8%) and Iran (9/128, 7%)
306 (S1A Fig). Similar results were obtained for the 95 problematic *Gene* papers from 2014-
307 2018, where teams from China authored 66% (63/95) papers (Fig 4). A significantly greater
308 proportion of problematic *Gene* papers from China were affiliated with hospitals (54/69,
309 78%), compared with papers from other countries (5/59, 8%) (Fisher's exact test, $p < 0.001$,
310 $n = 128$ papers) (Fig 4). This difference was also noted for problematic *Gene* papers from
311 2014-2018 (Fisher's exact test, $p < 0.001$, $n = 98$ papers) (Fig 4).

312

313 In contrast to the broader geographic distribution of problematic *Gene* papers, the 436
314 problematic *Oncology Reports* papers were authored by teams from just 13 countries (S2A
315 Fig). Most (393/436, 90%) problematic *Oncology Reports* papers were authored by teams
316 from China (Table 4, Fig 4), followed by much smaller proportions from South Korea
317 (14/436, 3%) and Japan (12/436, 3%) (S2A Fig). As noted for problematic *Gene* papers, a
318 significantly greater proportion of problematic *Oncology Reports* papers from China were

319 affiliated with hospitals (342/393, 87%) compared with papers from other countries (5/43,
320 12%) (Fisher's exact test, $p < 0.001$, $n = 436$ papers) (Fig 4).

321

322 We considered the distributions of problematic *Gene* and *Oncology Reports* papers according
323 to year of publication, country of origin, and affiliated institution type (Fig 5, S1 Fig, S2 Fig).
324 Problematic *Gene* papers were infrequent from 2007-2011 (1-4 papers/ year), rising to 8-38
325 papers/ year from 2012-2018, where the highest number of problematic papers was identified
326 in 2018 (Fig 5). These numbers correspond to 1.0%-4.2% of all original *Gene* papers
327 published per year from 2012-2018. Papers from China represented the majority of
328 problematic *Gene* papers from 2015-2018 (Fig 5), where most papers were also affiliated
329 with hospitals (S1B Fig). Compared with *Gene*, *Oncology Reports* published higher numbers
330 of problematic papers per year, corresponding to 8.3-12.6% original *Oncology Reports* papers
331 from 2014-2018 (Fig 5). Across all 5 years, most (87-93%) problematic *Oncology Reports*
332 papers were authored by teams from China, corresponding to 11% original *Oncology Reports*
333 articles in 2015-2017 (Fig 5), most of which were also affiliated with hospitals (S2B Fig).

334

335 ***Analysis of all problematic human gene function papers***

336 After adjusting for 9 duplicate papers across the 5 corpora, we identified 712 problematic
337 papers with wrongly identified sequences (Fig 1) that were published by 78 journals and 31
338 publishers (S6 Data). The 712 problematic papers included 1,535 wrongly identified
339 sequences, most of which were (RT-)PCR reagents (1,301/1,535, 85%), followed by
340 si/shRNA's (226/1,535, 15%) (Table 3).

341

342 As most incorrect reagents represented (RT-)PCR primers which are employed as paired
343 reagents, we considered the verified identities of primer pairs that were found to include at
344 least one wrongly identified primer (Fig 6). Problematic papers frequently paired one (RT-
345)PCR primer that targeted the claimed gene with a primer that was predicted to target a
346 different gene (n=237 papers), or to be non-targeting in human (n=118 papers) (Fig 6). Many
347 problematic papers (n=192) described primer pairs that were predicted to target the same
348 incorrect gene (Fig 6). Problematic papers also combined one non-targeting primer with
349 another that was predicted to target an incorrect gene (n=70 papers), two non-targeting
350 primers (n=63 papers), and/or primers that were predicted to target two different incorrect
351 genes (n=42 papers) (Fig 6). Notably, 21% (276/1,301) incorrect (RT-)PCR primers were
352 predicted to target an orthologue of the claimed gene, typically in rat or mouse (S1 Table).

353

354 *Bibliometric analysis of human genes analysed in problematic papers*

355 Almost all (1,442/1,535, 94%) incorrect sequences represented targeting reagents that were
356 claimed to target 365 protein-coding genes and 120 ncRNA's (S1 Table). The remaining 88
357 sequences represented incorrect “non-targeting” sequences that were instead predicted to
358 target 35 genes, most of which (28/35, 80%) were protein-coding genes.

359

360 To count the numbers of papers in PubMed that are associated with protein-coding genes in
361 n=709 problematic papers (Fig 7), we used the gene2pubmed service of the National Center
362 for Biotechnology Information (29), restricting these analyses to protein-coding gene
363 identifiers that mapped to official gene names. Primary protein-coding genes, which
364 represented the first-listed genes in publication titles or abstracts, tended to be associated with
365 more papers in PubMed than a randomly chosen human protein-coding gene (median

366 publication numbers: 167 vs 31, $P < 10^{-109}$, two-sided Mann-Whitney U test) (S3A Fig). Only
367 two genes that were the primary focus of least two problematic papers (*TCTN1* and *GPR137*)
368 (Table 2) have appeared in fewer publications in PubMed than a randomly chosen human
369 protein-coding gene (Fig 7A). We repeated these analyses to examine the protein-coding
370 genes that were claimed as targets by wrongly identified reagents. Again, most wrongly
371 identified target genes have appeared in more papers than a randomly chosen protein-coding
372 gene (median publication numbers: 238 vs 31, $P < 10^{-94}$, two-sided Mann-Whitney U test)
373 (S3B Fig). The most frequent wrongly claimed gene targets were *GAPDH* and *ACTB* (Fig
374 7B), reflecting their widespread use as RT-PCR control genes. In summary, these analyses
375 demonstrate that problematic papers can focus upon and/or employ reagents that are wrongly
376 claimed to target highly-investigated human genes such as *BCL2*, *EGFR*, *PTEN*, *STAT3*, and
377 *CCND1* (Fig 7).

378

379 *Post-publication correction, citation and curation of problematic papers*

380 We considered whether any problematic papers have been the subject of post-publication
381 notices, such as retractions, expressions of concern or corrections (11). Only 2% (11/712) of
382 problematic papers have been retracted, where most (8/11) retraction notices did not refer to
383 wrongly identified sequence(s), and 3 problematic papers have been subject to expressions of
384 concern (S2 Table). Although we excluded papers in which incorrect sequences had been
385 subsequently corrected (see Methods), we noted 5 corrections to problematic papers that
386 addressed issues other than incorrect sequences (S2 Table). Almost all (693/712, 97%)
387 problematic papers therefore remain uncorrected within the literature.

388

389 We then considered how problematic papers have been curated within gene knowledgebases
390 and cited within the literature. Between 1-207 problematic papers were found within 5 gene
391 knowledgebases that rely upon text mining (30-34), where knowledgebases of miR functions
392 contained the most problematic papers (S3 Table). In March 2021, the 712 problematic
393 papers had been cited 17,183 times according to Google Scholar. Subsets of problematic
394 C+G, *Gene* and *Oncology Reports* papers have also been cited by one or more clinical trials
395 (Fig 8A). Given expected publication delays between pre-clinical and clinical research, we
396 extended these data by considering the approximate potential to translate (APT) for
397 problematic papers (35) according to publication corpus (Fig 8B). The APT metric uses the
398 combination of concepts contained within a paper to infer the probability that the paper will
399 be cited by future clinical trials or guidelines (35). The average APT for problematic papers
400 in the 5 corpora ranged from 15-35% (Fig 8B), indicating that 15-35% of problematic papers
401 in each corpus resemble papers that will be cited by clinical research. Without timely
402 interventions, or without heuristics not captured by the models underlying the APT (35),
403 around one quarter of problematic papers are likely to directly misinform or distract the
404 future development of therapies against human disease.

405

406 **Discussion**

407 Experimental analyses of gene function require nucleotide sequence reagent identities to
408 precisely match their published descriptions. Wrongly identified nucleotide sequence
409 reagents therefore represent a threat to the continuum of genetics research, from population-
410 based genomic sequencing to pre-clinical functional analyses, and the translation of these
411 results to patients. While there have been previous reports of wrongly identified PCR and
412 gene knockdown reagents in single papers or small cohorts (11-21), the present study is the
413 first to systematically fact-check the identities of nucleotide sequences in over 3,400 papers.
414 Our supported application of S&B (12, 24, 26) for both targeted corpora and journal
415 screening identified 712 papers in 78 journals that describe 1-36 wrongly identified
416 nucleotide sequences/ paper. These 712 papers described over 1,500 wrongly identified
417 sequences, most of which were claimed targeting reagents for the study of over 480 human
418 genes. The sheer number of problematic papers and incorrect reagents that we have
419 uncovered from screening a very small fraction of the human gene function literature predicts
420 a problem of alarming proportions that requires urgent and co-ordinated action.

421

422 Before outlining the possible significance of our results, we must recognise the limitations of
423 the present study. Firstly, we recognise that by verifying the identities of over 13,700
424 sequences, we are likely to have committed some errors of our own. The most challenging
425 incorrect reagents that we encountered are claimed human targeting sequences that appear to
426 have no human target. As blastn is indicated to have a very low but measurable false-negative
427 rate (36, 37), we acknowledge that a small fraction of what appear to be non-targeting
428 reagents may in fact be targeting reagents, as claimed. Nonetheless, as we have taken
429 extensive steps to verify the identities of individual nucleotide sequences (see Methods), we

430 reasonably expect that false-positives at both the sequence and publication level will be rare
431 within our dataset. Many papers also described at least two wrongly identified nucleotide
432 sequences, building in some protection against false-positives at the publication level.

433

434 We also recognise that we have screened three targeted corpora, of which two (the *miR-145*
435 and C+G corpora) did not include all available papers. We therefore recognise that the rates
436 of problematic papers within these corpora may not reflect the rates of problematic papers in
437 other targeted corpora or in the wider literature. Nonetheless, the relative proportions of
438 nucleotide sequence error types in problematic *miR-145* and C+G papers were similar to
439 those identified by screening papers in both *Gene* and *Oncology Reports* (Fig 2). This
440 indicates that features of problematic papers identified within the targeted corpora may
441 extend to papers elsewhere, particularly as problematic *miR-145* and C+G papers were
442 published across 47 journals.

443

444 Finally, we recognize that by only verifying the S&B outputs for a minority of papers
445 published in *Gene* and *Oncology Reports*, we will not have described all papers with wrongly
446 identified nucleotide sequences in these journals. Indeed, the identified numbers and
447 proportions of problematic papers are very likely to represent underestimates, as these
448 proportions were derived from verifying sequences in just 10% *Gene* papers from 2007-2018
449 and 45% *Oncology Reports* papers from 2014-2018. As S&B does not analyze papers where
450 most species identifiers correspond to species other than human (26), and we did not verify
451 the identities of sequences that were claimed to target non-human species, we also recognize
452 that our results do not indicate whether wrongly identified nucleotide sequences occur in

453 papers that examine other species. Analyses of targeted corpora identified around reference
454 papers may begin to answer this question.

455

456 Despite these limitations, our results paint a worrying picture of the reliability of some
457 sections of the human gene-focussed literature. The consequences of many incorrect gene
458 function papers therefore deserve consideration. As previously discussed, papers that describe
459 incorrect nucleotide sequences could encourage the incorrect selection of genes for further
460 experimentation, possibly at the expense of more productive candidates (8). This could be
461 exacerbated where multiple problematic papers report similar results for the analysis of the
462 same gene (8, 13). For example, series of 3-20 single gene knockdown papers that universally
463 claim that the gene target plays a causal role in 3-11 different cancer types could collectively
464 encourage further research. Many incorrect gene function papers could also lead to the
465 overestimation of knowledge of gene function from text mining approaches (30-34),
466 particularly given the assumed reliability of published experimental results (38). Our results
467 demonstrate that problematic papers are already indexed within gene function knowledge
468 bases (30-34).

469

470 As experimental reagents, wrongly identified nucleotide sequences carry the additional risk
471 of being wrongly used in other studies (8, 11). Although siRNA's and shRNA's are
472 increasingly purchased from external companies as preformulated reagents, many researchers
473 continue to order custom-made PCR primers. The most frequent incorrect reagent type that
474 we have identified can therefore easily be reused from the literature. Experiments that either
475 attempt to replicate published results associated with incorrect reagents and/or unknowingly
476 reuse incorrect reagents are likely to generate unexpected results that may then remain

477 unpublished. As possible evidence of this, we could only identify one study that reported
478 incorrect sequences based on the results of follow-up experiments (20) that our analyses also
479 identified. Given that over 17,000 citations have been accumulated by the problematic papers
480 that we have identified, it seems inevitable that unreliable gene function papers are already
481 wasting time and resources.

482

483 Given the numbers of problematic papers and incorrect reagents that we have identified, we
484 must consider possible explanations for both incorrect nucleotide sequences and the papers
485 that describe them. As nucleotide sequences are prone to being affected by different error
486 types, we believe that wrongly identified sequences represent unintended errors that persist
487 when authors, editors and peer reviewers do not check sequence identities during manuscript
488 preparation or review (8, 11, 12). We recognize that wrongly identified sequences in papers
489 could also reflect some form of research sabotage. However, as workplace sabotage is
490 typically directed towards known individuals (39, 40), this seems an unlikely explanation for
491 wrongly identified sequences across hundreds of gene function papers published by many
492 different authors.

493

494 We also recognize that some of the wrongly identified sequences that we have described have
495 almost certainly arisen in the context of genuine research. Nonetheless, some wrongly
496 identified sequences were associated with results that appeared highly implausible. For
497 example, the majority of incorrect (RT-)PCR primer pairs that we identified should have
498 uniformly failed to generate (RT-)PCR products across all templates, and yet appeared to
499 generate results that were consistent with primers targeting the claimed genes. As previously
500 reported (11-13), we also identified papers where the non-targeting si/shRNA was verified to

501 target the gene of interest and yet still generated the expected negative or baseline results.
502 Based on the many discrepancies between verified sequence identities and their associated
503 results, it would appear that many experiments across hundreds of gene function papers were
504 not performed as described.

505

506 Although most analyses of experimental errors consider the context of genuine research (41-
507 43), errors can also arise within the context of fraudulent or fabricated research (43, 44). We
508 have previously proposed that incorrect nucleotide sequences could flag that some affected
509 papers are fraudulent (8, 11-13, 45, 46). Due to numerous unexpected similarities between
510 single gene knockdown papers with incorrect sequences, we initially proposed that these
511 papers may reflect the undeclared involvement of organisations such as paper mills (13)
512 which have been alleged to mass-produce fraudulent manuscripts for publication (47). The
513 production of fraudulent manuscripts at scale may involve the use of manuscript templates,
514 leading to unusual degrees of similarity between the resulting publications (8, 13, 46). Papers
515 with incorrect nucleotide sequences showed features consistent with the possible use of
516 manuscript templates, such as similarities in textual and figure organisation, outlier levels of
517 textual similarity, superficial explanations for the analysis of particular genes, and generic
518 experimental approaches regardless of predicted gene function (8, 13, 45, 46). In addition to
519 the use of manuscript templates, we have proposed that producing many gene-focused
520 manuscripts at minimal cost could involve writers with either an incomplete understanding of
521 the experiments that they are describing and/or limited time for quality control (8, 46). These
522 conditions could lead to incorrect nucleotide sequences being a feature of gene function
523 manuscripts from paper mills (8, 46) where incorrect sequences could also be reused across
524 different manuscripts (11-13). The description of errors which are incompatible with reported

525 results and the presence of multiple incorrect sequences across different techniques within
526 single papers could therefore be features of experiments that were never actually performed.

527

528 We have proposed that human genes represent attractive publication targets for paper mills,
529 with many under-studied human genes (48-50) that can be targeted in different cancer or
530 disease types that can then be distributed across different authors and journals over many
531 years (8, 45, 46). However, whereas we had previously hypothesised that under-studied genes
532 might be preferentially targeted by paper mills (8), in the present analysis, problematic papers
533 rarely focussed on under-studied human protein-coding genes, and instead either focussed on
534 or employed incorrect reagents that were claimed to target protein-coding genes that had
535 received prior attention. While under-studied genes may present more individual publication
536 opportunities, papers that describe human genes with known functions and significance could
537 carry more editorial and reader interest, increasing the likelihood of problematic manuscripts
538 being accepted for publication and then cited by future research.

539

540 Our analyses of both targeted and journal corpora indicate that ncRNA's may provide a
541 further layer of possibilities for the fabrication of gene-focussed papers. While we recognise
542 that studying genes in different diseases and analysing the functions of ncRNA's in
543 combination with other genes are features of genuine research, our results suggest that a
544 focus upon ncRNA's such as miR's could allow the inclusion of more topic variables within
545 manuscripts from paper mills, such as ncRNA's and protein-coding genes that are studied
546 across different disease types, with or without drug or natural product treatments. Examining
547 ncRNA's in combination with other gene(s) could allow larger and more diverse publication
548 series to be created, compared with those that focus on single genes. Furthermore, as

549 ncRNA's possess largely numeric identifiers that may be more difficult to recognise and
550 recall than alphanumeric gene identifiers, any focus upon ncRNA's could contribute to large
551 publication series being less visible within the literature. As papers that describe miR
552 functions have also been shown to be highly cited (51), miR's and other ncRNA's could
553 represent attractive target genes for paper mills.

554

555 In both targeted and journal corpora, most problematic papers were authored by teams from
556 mainland China that were also overwhelmingly affiliated with hospitals. While researchers in
557 different countries may use paper mills to meet publication targets or quotas (49), paper mills
558 have been most widely discussed in the context of academics and medical doctors in China
559 (47, 53-59). Stringent publication requirements may represent a particular challenge for
560 hospital doctors in China, where some hospital doctors have described limited time, training
561 and/or opportunities to undertake research (47, 53-59). Large numbers of human gene
562 function papers with incorrect nucleotide sequences that list hospital affiliations in China
563 could reflect hospital doctors turning to paper mills to meet publication requirements,
564 whereas the contrasting institutional profiles of problematic papers from other countries
565 could highlight different publication pressures elsewhere. In summary, our results combined
566 with previous reports (47, 53-59) indicate that publication pressure upon hospital doctors in
567 China may be exerting measurable effects upon the human gene function literature.

568

569 In summary, we are concerned that the sheer number of human genes that are available for
570 analysis, combined with research drivers that favour the continued investigation of genes of
571 known function (48-50), are unwittingly providing an extensive source of topics around
572 which gene function papers can be fraudulently created. Furthermore, since genuine pre-

573 clinical gene research requires specialised expertise, time and material resources (13), the
574 mass production of fraudulent gene function papers could be quicker and cheaper by orders
575 of magnitude (8). Given the number of human genes which can be studied either singly or in
576 combination with other genes and topics such as drugs and analysed across different cancer
577 types or other diseases, combined with unrealistic demands for research productivity (46, 53-
578 59), the publication of fraudulent gene function papers could potentially outstrip the
579 publication of genuine gene function research.

580

581 *Future directions*

582 Our results indicate that the problem of incorrect gene function papers requires urgent action.
583 Within the research community, this can take place in several ways. As the validation of
584 nucleotide sequence identities using algorithms such as blastn (12, 25) represents a routine
585 activity for teams that investigate gene function, we hope that our results will encourage
586 researchers to unfailingly check the identities of published nucleotide sequences, both in the
587 context of their own research and during peer review. Researchers encountering wrongly
588 identified sequences can describe these to authors, journal(s) and/or PubPeer (60) using the
589 reporting fields that we have proposed (11). Researchers can also compare the claims of
590 gene-focused papers with those from high-throughput experimental studies (61) and/or
591 predictive algorithms (38, 62). Professional societies can reinforce the importance of reagent
592 verification through conference presentations, education programs, and journal editorials, and
593 can advocate for tangible incentives to encourage further fact-checking of the genetics
594 literature.

595

596 Our analysis of only a small proportion of human gene function papers, combined with the
597 discovery that most incorrect nucleotide sequences are unique within screened literature
598 corpora, highlights the need for further screening studies to identify other problematic gene
599 function papers. Journals that published problematic papers in targeted corpora represent
600 possible targets for future screening approaches. As incorrect nucleotide sequences are
601 unlikely to be found in all problematic gene function papers, future research should also
602 combine analyses of nucleotide sequences with other features of concern such as manipulated
603 or recurring experimental images and data (46, 63, 64).

604

605 Unfortunately, efforts from the research community alone will not solve the problem that we
606 have described. Similarly, recent changes to researcher assessment (65, 66) will not address
607 problematic papers that have already been published. The ability to alter the published record
608 relies upon the engagement and co-operation of journals and publishers (11). Over the past
609 year, growing numbers of journals have begun to recognize the issue of manuscripts and
610 publications from paper mills (67-77), including papers that analyse gene function and drug
611 treatment in cell lines (71, 74). While the described efforts to screen incoming manuscripts
612 are welcome and should be extended to all journals that publish gene function research,
613 screening incoming manuscripts must be coupled with addressing problematic papers that are
614 already embedded in the literature (71, 74-76). These efforts could be supported by gene
615 function experts who could explain the significance of incorrect nucleotide sequences and/or
616 provide training for editorial staff, particularly as the necessary researcher skills are already
617 widely available. To overcome the protracted timeframes that can be associated with journal
618 investigations of incorrect sequences (11), we have proposed the rapid publication of editorial
619 notes to transparently flag papers with verifiable errors while journal and institutional
620 investigations proceed (46).

621

622 *Summary and conclusions*

623 To fully extend the benefits of genomics towards patients and broader populations, it is
624 widely recognised that we must understand the functions of every human gene (1, 2).
625 However, the availability of many human genes for experimental analysis, combined with
626 research drivers that favour the continued investigation of genes of known function (46-49),
627 may unwittingly provide an extensive source of topics around which gene function papers
628 can be fraudulently created. Whereas genuine gene research requires time, expertise, and
629 material resources, the mass production of fraudulent gene function papers by paper mills
630 could be quicker and cheaper by orders of magnitude (8). Given the number of human genes
631 whose functions can analysed singly and/or in combination with other genes and/or drugs
632 across different cancer types or other diseases, combined with acute demands for research
633 productivity that may not always be matched by researcher capacity and training (78),
634 fraudulent gene function papers could unfortunately outstrip the publication of genuine gene
635 function research. Indeed, the possible extent of the problem of unreliable human gene
636 function papers is indicated by the lack of overlap between the problematic papers that we
637 have reported, and other papers of concern reported elsewhere (71, 74, 79). While publishers
638 and journals decide how to address this urgent problem, laboratory scientists, text miners and
639 clinical researchers must approach the human gene function literature with a critical mindset,
640 and carefully evaluate the merits of individual papers before acting upon their results.

641

642 **Materials and Methods**

643 *Identification of literature corpora*

644 *Single Gene Knockdown (SGK) corpus*

645 Single gene knockdown (SGK) papers were identified by combining each of 17 human gene
646 identifiers (*ADAM8*, *ANXA1*, *EAG1*, *GPR137*, *ICT1*, *KLF8*, *MACC1*, *MYO6*, *NOBI*, *PP4R1*,
647 *PP5*, *PPM1D*, *RPS15A*, *TCTN1*, *TPD52L2*, *USP39*, and *ZFX*) with the search string "cancer
648 AND/OR knockdown AND/OR lentivirus" (13), to search PubMed and Google Scholar
649 databases in June 2019 using the "allintext:" function for Google Scholar searches. No
650 publication date ranges, country-specific or journal-based search terms were used to limit
651 search results. Papers were visually inspected to confirm that papers described gene
652 knockdown experiments that targeted one of the 17 human genes in human cancer cell lines.

653

654 *miR-145 corpus*

655 The *miR-145* corpus included papers that analysed human *miR-145* function in human cell
656 lines. Two reference papers PMID 29749434 and PMID 29217166, where PMID 29217166
657 was verified to describe incorrect nucleotide sequence reagents (24, 26, see below), were
658 used in PubMed similarity searches conducted in September 2019 and October 2020.
659 Additional papers were identified through Google Scholar searches employing the keywords
660 "gene" + "miR-145" + "cancer" conducted in April 2019 and September 2020. Publication
661 dates were limited to 2019 to broadly align with the SGK corpus. All identified papers were
662 visually inspected to confirm the analysis of human *miR-145* function in human cell lines.

663

664 *Cisplatin + Gemcitabine (C+G) corpus*

665 The cisplatin + gemcitabine (C+G) corpus included papers that described either cisplatin or
666 gemcitabine treatment of human cancer cell lines and/or biospecimens from cisplatin or
667 gemcitabine-treated cancer patients, where most papers also investigated human gene
668 function. Two reference papers PMID's 30250547 and 26852750 that both described
669 incorrect sequences (24, 26) were used as reference papers in PubMed similarity searches
670 conducted in September 2019 and October 2020. Additional papers were identified using
671 Google Scholar searches with the search string "gene", "cancer", "cisplatin" +/- "miR"
672 conducted in September 2019 and October 2020. A PubMed similarity search for PMID
673 26852750 conducted in September 2019 also identified 5 papers that referred to gemcitabine
674 treatment (PMID's 18636187, 26758190, 28492560, 30117016, 31272718). Four of these
675 papers (PMID's 26758190, 28492560, 30117016, 31272718) described incorrect sequences
676 and were used as reference papers for PubMed similarity searches conducted in September
677 2019 and October 2020. Additional publications were identified through Google Scholar
678 searches with the query "gene", "cancer", "gemcitabine" +/- "miR" between September 2019
679 and October 2020. In all cases, publication dates were limited to 2019 to align with other
680 targeted corpora. Papers were visually inspected to confirm that they studied either cisplatin
681 or gemcitabine treatment in the context of human cancer cell lines or biospecimens, and to
682 exclude papers from other targeted corpora.

683

684 *Journal corpora*

685 *Gene* and *Oncology Reports* were selected for S&B screening as examples of journals that
686 have published papers with incorrect nucleotide sequences (11-13), where *Oncology Reports*
687 also published the highest number of problematic SGK papers (S1 file). Whereas *Gene*

688 (published by Elsevier) encompasses a broad range of gene research across different species,
689 *Oncology Reports* (published by Spandidos Publications) focusses on human cancer research.

690

691 *Gene* papers from January 2007 to December 2018 were retrieved using the Web of Science
692 search criteria: PY= "2007-2018" AND SO= "GENE" AND DT= ("Article" OR "Review").

693 *Oncology Reports* papers from January 2014 to December 2018 were retrieved using the Web
694 of Science search criteria: PY="2014-2018" AND SO= "ONCOLOGY REPORTS" AND

695 DT= ("Article" OR "Review"). In the case of *Gene* papers, DOI's were retrieved, and PDF
696 files were downloaded using the Elsevier Application Programming Interface with Crossref

697 Content negotiation (<http://tdmsupport.crossref.org>), whereas open-access *Oncology Reports*
698 articles were directly downloaded from www.spandidos-publications.com.

699

700 ***Seek & Blastn screening***

701 SGK, *miR-145* and C+G papers were named using PMID's or journal identifiers and
702 screened by S&B as described (12, 24, 26). All SGK papers identified for the 17 selected

703 human genes were screened by S&B. In the case of *miR-145* and C+G papers, S&B
704 screening was conducted until 50 *miR-145*, cisplatin and gemcitabine papers were flagged for

705 further analysis, either because S&B had flagged at least one wrongly identified reagent or
706 had failed to extract any sequences from the text. This required S&B screening of 163 *miR-*

707 *145* papers and 258 C+G papers. S&B screening was conducted in 2019 and/or 2020, with all
708 papers flagged by S&B in 2019 being rescreened by S&B in 2020.

709

710 *Gene* papers were labelled with PMID's, and batched pdf files were zipped into two
711 compressed files according to publication dates (2007-2013 and 2014-2018). *Oncology*
712 *Reports* papers were labelled by PMID's and journal identifiers. S&B screening was
713 conducted between July-October 2019, with all papers rescreened in November 2020-
714 February 2021. *Gene* and *Oncology Reports* papers were flagged for further analysis where
715 S&B had either flagged at least one nucleotide sequence or had failed to extract any
716 sequences from the text.

717

718 ***Visual inspection of papers following S&B screening***

719 Papers were visually inspected to determine the claimed genetic and/or experimental identity
720 of each sequence. If the claimed target or experimental use of any sequence was not evident,
721 or if a sequence was claimed to target a species other than human, the sequence was excluded
722 from further analysis. Papers that had been subject to post-publication corrections where
723 wrongly identified nucleotide sequences had been corrected were also excluded. We included
724 retracted papers, to align with previous descriptions of SGK papers (11-13), and in
725 recognition of the possibility of retracted papers continuing to be cited (80).

726

727 ***Manual verification of nucleotide sequence reagent identities***

728 Nucleotide sequence identities were manually confirmed for all sequences that were not
729 (correctly) extracted and/or flagged as being possibly incorrect by S&B, as described (26).
730 For the *Oncology Reports* and *Gene* corpora, this involved checking at least 34% and 54% of
731 all sequences, respectively. Further verification steps were performed for particular reagents,
732 as follows:

733 (i) For reagents that were claimed to target specific gene polymorphisms or mutant sequences
734 and for which no sequence match could be identified by either blastn or blat (25, 80), manual
735 sequence alignments were performed in Word with the query sequence in forward, forward
736 complement, reverse and reverse complement orientations, against either the sequence
737 corresponding to the accession number provided within the text, or to the most relevant
738 genomic sequence found in NCBI GenBank, according to the text claim. Sequences were
739 delineated using the R studio “stringr” library and accepted as targeting if the specified
740 mutated base(s), when reverted to their original base(s) as described in NCBI dbSNP
741 <https://www.ncbi.nlm.nih.gov/snp/>, allowed the reagent to target the wild-type sequence
742 according to previously published targeting criteria (12).

743 (ii) If no significant matches were identified for reagents specified for the analysis of mutant
744 or variant targets, mismatches within the nucleotide sequence were converted to the wild type
745 sequence, either as described in the publication or according to dbSNP and reanalysed as
746 described (26). Reagents that were indicated to target the claimed wild-type sequence were
747 accepted as correct targeting reagents.

748 (iii) All flagged incorrect targeting sequences were double-checked through additional blastn
749 searches against the database: “Homo sapiens (taxid:9606)”, optimized for “Somewhat
750 similar sequences (blastn)”, using an expect threshold 1000, in February 2021.

751 Nucleotide sequence reagents that were verified to have been wrongly identified were
752 assigned to one of 3 previously described error categories (11, 12):

753 (i) Reagents claimed to represent targeting reagents but verified to target a human gene or
754 target other than that claimed within the text. This error category included miR-targeting
755 reverse RT-PCR primers with incorrect gene targeting descriptions, as supported by sequence
756 verification (26), and by having been employed to analyse gene(s) other than the claimed

757 miR and/or as a claimed universal primer according to Google Scholar searches (11, 13).

758 While we recognise that some of these reagents could amplify the claimed miR target as

759 described, their descriptions as specific targeting reagents were incorrect and could lead to

760 incorrect RT-PCR primer reuse.

761 (ii) Reagents claimed to target a human gene or genomic sequence but verified to be non-

762 targeting in human. These reagents included RT-PCR primers that targeted introns or other

763 non-transcribed regions within claimed genes.

764 (iii) Reagents claimed to represent non-targeting reagents in human but verified to target a

765 human gene or genomic sequence.

766

767 ***Additional publication analyses***

768 For all papers subjected to S&B analysis, publication titles were visually inspected to identify

769 human gene identifiers, human cancer types, and drug identifiers which were confirmed

770 through Google searches. Human genes were categorized as either protein-coding or

771 ncRNA's (miRs, lncRNA or circRNA) according to GeneCards

772 (<https://www.genecards.org/>).

773

774 Journal publishers were identified via the SCImago database (<https://www.scimagojr.com/>).

775 The Journal Impact Factor corresponding to the (closest) publication year of each paper was

776 obtained from the Clarivate Analytics Journal Citation Reports database

777 (<https://jcr.clarivate.com/>). Numbers of original articles published per year by *Gene* and

778 *Oncology Reports* were obtained from Clarivate InCites (<https://incites.clarivate.com/>), under

779 Entity type= "Publication Sources", Publication Date= 2007-2018, DT= include only

780 "Article". The country of origin of each paper was assigned according to the affiliations of at
781 least half of the listed authors. Papers were considered to be affiliated with hospital(s) if the
782 institutional affiliations of at least half of the listed authors were associated with one or more
783 of the keywords: "clinic", "health cent", "hosp", "hospital", "infirmary", "sanatorium",
784 "surgery". Papers not meeting this criterion were considered to be affiliated with institutions
785 other than hospitals. Proportions of problematic papers (from China versus all other
786 countries, hospitals versus other institutions) were compared using the Fisher's Exact test
787 (SPSS statistics 27).

788

789 *Bibliometric analysis of human genes in problematic papers*

790 Linkages of protein-coding genes to publications were obtained via gene2pubmed from
791 NCBI NIH (<https://ftp.ncbi.nlm.nih.gov/gene/DATA/gene2pubmed.gz>) on 15 July, 2021 as
792 described (29, 49). Two-sided Mann Whitney U tests were performed using SciPy (82). Post-
793 publication notices linked with problematic papers were identified through PubMed and
794 Google Scholar searches. PubMed ID's or other publication identifiers were employed as
795 search queries of gene knowledgebases in May 2021 (30-34). Publication citation counts are
796 those reported by Google Scholar in March, 2021. Problematic papers cited by clinical trials
797 were cited by at least one publication within the NIH Open Citation collection (83), which in
798 MedLine carried the annotation of a publication type of any "clinical trial" (without
799 distinguishing clinical trial stage). The approximate potential to translate (APT) for
800 problematic papers in each corpus was calculated as described (35) and obtained from iCite
801 (83).

802

803

804 **Acknowledgements**

805 **Funding:** JAB and CL gratefully acknowledge funding from the US Office of Research
806 Integrity, grant ID ORIIR180038-01-00. JAB, CL and ACD gratefully acknowledge grant
807 funding from the National Health and Medical Research Council of Australia, Ideas grant ID
808 APP1184263. TS gratefully acknowledges funding from the National Science Foundation,
809 1956338, SCISIPBIO: A data-science approach to evaluating the likelihood of fraud and
810 error in published studies; K99AG068544, National Institutes on Aging, Integrative Multi-
811 Scale Systems Analysis of Gene-Expression-Driven Aging Morbidity; National Institute of
812 Allergy and Infectious Diseases, AI135964, Successful Clinical Response In Pneumonia
813 Therapy (SCRIPT) Systems Biology Center. The authors thank journal editorial staff for
814 discussions and support of this study.

815

816 **Authors' contributions (alphabetical order): Conceptualization:** Jennifer Byrne,
817 Guillaume Cabanac, Cyril Labbé, Thomas Stoeger; **Methodology:** Jennifer Byrne, Guillaume
818 Cabanac, Cyril Labbé, Thomas Stoeger; **Formal analysis:** Jennifer Byrne, Bertrand Favier,
819 Yasunori Park, Pranujan Pathmendra, Thomas Stoeger, Rachael West; **Writing - original**
820 **draft preparation:** Jennifer Byrne, Yasunori Park, Thomas Stoeger, Rachael West; **Writing**
821 **- review and editing:** Jennifer Byrne, Guillaume Cabanac, Amanda Capes-Davis, Bertrand
822 Favier, Cyril Labbé, Yasunori Park, Pranujan Pathmendra, Thomas Stoeger, Rachael West;
823 **Funding acquisition:** Jennifer Byrne, Amanda Capes Davis, Cyril Labbé, Thomas Stoeger;
824 **Supervision:** Jennifer Byrne

825

826 **Conflicts of Interest:** The authors have no relevant financial or non-financial interests to
827 disclose.

828 **References**

- 829 1. Shendure J, Findlay GM, Snyder MW. Genomic Medicine-Progress, Pitfalls, and
830 Promise. *Cell*. 2019;177(1):45-57. PMID: 30901547
- 831 2. Green ED, Gunter C, Biesecker LG, Di Francesco V, Easter CL, Feingold EA, et al.
832 Strategic vision for improving human health at The Forefront of Genomics. *Nature*.
833 2020;586(7831):683-692. PMID: 33116284
- 834 3. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial
835 sequencing and analysis of the human genome. *Nature*. 2001;409(6822):860-921. PMID:
836 11237011.
- 837 4. Venter JC, Adams MD, Myers EW, Li PW, Mural RJ, Sutton GG, et al. The sequence of
838 the human genome. *Science*. 2001;291(5507):1304-13051. PMID: 11181995.
- 839 5. Boutros PC. The path to routine use of genomic biomarkers in the cancer clinic. *Genome*
840 *Res*. 2015(10);25:1508-1513. PMID: 26430161
- 841 6. Kaelin WG Jr. Common pitfalls in preclinical cancer target validation. *Nat Rev Cancer*.
842 2017;17(7):425-440. PMID: 28524181
- 843 7. Hahn WC, Bader JS, Braun TP, Califano A, Clemons PA, Druker BJ, et al. An expanded
844 universe of cancer targets. *Cell*. 2021;184(5):1142-1155. PMID: 33667368
- 845 8. Byrne JA, Grima N, Capes-Davis A, Labbé C. The possibility of systematic research
846 fraud targeting under-studied human genes: causes, consequences and potential
847 solutions. *Biomarker Insights* 2019;14:1177271919829162. PMID: 30783377
- 848 9. Nelson MR, Tipney H, Painter JL, Shen J, Nicoletti P, Shen Y, et al. The support of
849 human genetic evidence for approved drug indications. *Nat Genet*. 2015;47(8):856-860.
850 PMID: 26121088

- 851 10. King EA, Davis JW, Degner JF. Are drug targets with genetic support twice as likely to
852 be approved? Revised estimates of the impact of genetic support for drug mechanisms on
853 the probability of drug approval. *PLoS Genet.* 2019;15(12):e1008489. PMID: 31830040
- 854 11. Byrne JA, Park Y, West RA, Capes-Davis A, Cabanac G, Labbé C. The thin ret(raction)
855 line: biomedical journal responses to reports of incorrect non-targeting nucleotide
856 sequence reagents in human gene knockdown publications. *Scientometrics.*
857 2021;126:3513-3534.
- 858 12. Labbé C, Grima N, Gautier T, Favier B, Byrne JA. Semi-automated fact-checking of
859 nucleotide sequence reagents in biomedical research publications: the Seek & Blastn
860 tool. *PLoS ONE.* 2019;14(3):e0213266. PMID: 30822319
- 861 13. Byrne JA, Labbé C. Striking similarities between publications from China describing
862 single gene knockdown experiments in human cancer cell lines. *Scientometrics.*
863 2017;110:1471-1493.
- 864 14. Habbal W, Monem F, Gärtner BC. Errors in published sequences of human
865 cytomegalovirus primers and probes: do we need more quality control? *J Clin Microbiol.*
866 2005;43:5408-5409. PMID: 16208034
- 867 15. Shannon BA, Cohen RJ, Garrett KL. Influence of 16S rDNA primer sequence
868 mismatches on the spectrum of bacterial genera detected in prostate tissue by universal
869 eubacterial PCR. *The Prostate.* 2008;68(14):1487-1491. PMID: 18651564
- 870 16. Chiarella P, Carbonari D, Iavicoli S. Utility of checklist to describe experimental
871 methods for investigating molecular biomarkers. *Biomarkers Med.* 2015;9(10):989-995.
872 PMID: 26439607
- 873 17. Katavetin P, Nangaku M, Fujita T. Wrong primer for rat angiotensinogen mRNA. *Am J*
874 *Physiol Renal Physiol.* 2005;288(5):F1078. PMID: 15821251

- 875 18. Kocemba KA, Dudzik P, Ostrowska B, Laidler P. Incorrect Analysis of MCAM Gene
876 Promoter Methylation in Prostate Cancer. *The Prostate*. 2016;76(15):1464-1465. PMID:
877 27418327
- 878 19. Tamm A. Incorrect primer sequences in the article on methylprednisolone treatment.
879 *Acta Neurol Scand*. 2016;134(1):90. PMID: 27251935
- 880 20. Khoshi A, Sirghani A, Ghazisaeedi M, Mahmudabadi AZ, Azimian A. Association
881 between TPO Asn698Thr and Thr725Pro gene polymorphisms and serum anti-TPO
882 levels in Iranian patients with subclinical hypothyroidism. *Hormones*. 2017;16:75-83.
883 PMID: 28500830
- 884 21. Bustin S, Nolan T. Talking the talk, but not walking the walk: RT-qPCR as a paradigm
885 for the lack of reproducibility in molecular research. *Eur J Clin Invest*. 2017;47(10):756-
886 774. PMID: 28796277.
- 887 22. Harbig J, Sprinkle R, Enkemann SA. A sequence-based identification of the genes
888 detected by probesets on the Affymetrix U133 plus 2.0 array. *Nucleic Acids Res*.
889 2005;33(3):e31. PMID: 15722477
- 890 23. Draghici S, Khatri P, Eklund AC, Szallasi Z. Reliability and reproducibility issues in
891 DNA microarray measurements. *Trends Genet*. 2006;22(2):101-109. PMID: 16380191
- 892 24. Labbé C, Cabanac G, West RA, Gautier T, Favier B, Byrne JA. Flagging errors in
893 biomedical papers: to what extent does the leading publication format impede automatic
894 error detection? *Scientometrics*. 2020;124:1139-1156.
- 895 25. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search
896 tool. *J Mol Biol*. 1990;215:403-410. PMID: 2231712
- 897 26. <https://www.protocols.io/view/seek-amp-blastn-standard-operating-procedure-bjhpkj5n>
- 898 27. Dasari S, Tchounwou PB. Cisplatin in cancer therapy: molecular mechanisms of action.
899 *Eur J Pharmacol*. 2014;740:364-378. PMID: 25058905

- 900 28. de Sousa Cavalcante L, Monteiro G. Gemcitabine: metabolism and molecular
901 mechanisms of action, sensitivity and chemoresistance in pancreatic cancer. *Eur J*
902 *Pharmacol.* 2014;741:8-16. PMID: 25084222
- 903 29. Maglott D, Ostell J, Pruitt KD, Tatusova T. Entrez Gene: gene-centered information at
904 NCBI. *Nucleic Acids Res.* 2011;39:D52-D57. PMID: 2111545
- 905 30. Li G, Ross KE, Arighi CN, Peng Y, Wu CH, Vijay-Shanker K. miRTex: A Text Mining
906 System for miRNA-Gene Relation Extraction. *PLoS Comput Biol.*
907 2015;11(9):e1004391. PMID: 26407127
- 908 31. Vergoulis T, Kanellos I, Kostoulas N, Georgakilas G, Sellis T, Hatzigeorgiou A, et al.
909 mirPub: a database for searching microRNA publications. *Bioinformatics.*
910 2015;31(9):1502-1504. PMID: 25527833
- 911 32. Lever J, Zhao EY, Grewal J, Jones MR, Jones SJM. CancerMine: a literature-mined
912 resource for drivers, oncogenes and tumor suppressors in cancer. *Nat Methods.*
913 2019;16(6):505-507. PMID: 31110280
- 914 33. Jang YE, Jang I, Kim S, Cho S, Kim D, Kim K, et al. ChimerDB 4.0: an updated and
915 expanded database of fusion genes. *Nucleic Acids Res.* 2020;48(D1):D817-D824. PMID:
916 31680157
- 917 34. Roychowdhury D, Gupta S, Qin X, Arighi CN, Vijay-Shanker K. emiRIT: A text-mining
918 based resource for microRNA information. *bioRxiv.* 2020.11.05.370593; doi:
919 <https://doi.org/10.1101/2020.11.05.370593>
- 920 35. Hutchins BI, Davis MT, Meseroll RA, Santangelo GM. Predicting translational progress
921 in biomedical research. *PLoS Biol.* 2019;17(10):e3000416. PMID: 31600189.
- 922 36. Koonin EV, Galperin MY. Principles and Methods of Sequence Analysis. In: *Sequence-*
923 *Evolution- Function: Computational Approaches in Comparative Genomics.* Kluwer
924 Academic, Boston, 2003.

- 925 37. Buhler JD, Lancaster JM, Jacob AC, Chamberlain RD. Mercury BLASTN: Faster DNA
926 Sequence Comparison Using a Streaming Hardware Architecture. In: Proceedings of
927 Reconfigurable Systems Summer Institute, July 2007.
- 928 38. Pinkney HR, Wright BM, Diermeier SD. The lncRNA Toolkit: Databases and In Silico
929 Tools for lncRNA Analysis. *Noncoding RNA*. 2020;6(4):49. PMID: 33339309
- 930 39. Anderson MS, Ronning EA, De Vries R, Martinson BC. The perverse effects of
931 competition on scientists' work and relationships. *Sci Eng Ethics*. 2007;13:437-461.
932 PMID: 18030595
- 933 40. Serenko A. Knowledge sabotage as an extreme form of counterproductive behaviour:
934 conceptualisation, typology and empirical demonstration. *J Knowledge Production*.
935 2019;23(7):1260-1288.
- 936 41. Brown AW, Kaiser KA, Allison DB. Issues with data and analyses: Errors, underlying
937 themes, and potential solutions. *Proc Natl Acad Sci USA*. 2018;115(11):2563-2570.
938 PMID: 29531079
- 939 42. Wren JD. Algorithmically outsourcing the detection of statistical errors and other
940 problems. *EMBO J*. 2018;37(12):e99651. PMID: 29794111
- 941 43. Allchin D. Error types. *Perspectives on Science*. 2001;9(1):38-58.
- 942 44. Stroebe W, Postmes T, Spears R. Scientific misconduct and the myth of self-correction
943 in science. *Perspect Psychol Sci*. 2012;7:670-688. PMID: 26168129
- 944 45. Byrne J. We need to talk about systematic fraud. *Nature*. 2019;566(7742):9. PMID:
945 30728525
- 946 46. Byrne JA, Christopher J. Digital magic, or the dark arts of the 21st century-how can
947 journals and peer reviewers detect manuscripts and publications from paper mills? *FEBS*
948 *Lett*. 2020;594(4):583-589. PMID: 32067229

- 949 47. Else H, Van Noorden R. The fight against fake-paper factories that churn out sham
950 science. *Nature*. 2021;591(7851):516-519. PMID: 33758408
- 951 48. Pfeiffer T, Hoffmann R. Temporal patterns of genes in scientific publications. *Proc Natl*
952 *Acad Sci U S A*. 2007;104(29):12052-12056. PMID: 17620606
- 953 49. Stoeger T, Gerlach M, Morimoto RI, Nunes Amaral LA. Large-scale investigation of the
954 reasons why potentially important genes are ignored. *PLoS Biol*. 2018;16(9):e2006643.
955 PMID: 30226837
- 956 50. Gates AJ, Gysi DM, Kellis M, Barabási AL. A wealth of discovery built on the Human
957 Genome Project - by the numbers. *Nature*. 2021;590(7845):212-215. PMID: 33568828
- 958 51. Fire M, Guestrin C. Over-optimization of academic publishing metrics: observing
959 Goodhart's Law in action. *Gigascience*. 2019;8(6):giz053. PMID: 31144712
- 960 52. Hu Z, Wu Y. An empirical analysis on number and monetary value of ghostwritten
961 papers in China. *Curr Sci*. 2013;105:1230-1234.
- 962 53. Liu X, Chen X. Journal Retractions: Some Unique Features of Research Misconduct in
963 China. *J Schol Pub*. 2018;49(3):305-319.
- 964 54. Han J, Li Z. How Metrics-Based Academic Evaluation Could Systematically Induce
965 Academic Misconduct: A Case Study. *East Asian Sci Tech Soc*. 2018;12(2):165-179.
- 966 55. Zhao T, Dai T, Lun Z, Gao Y. An Analysis of Recently Retracted Articles by Authors
967 Affiliated with Hospitals in Mainland China. *J Schol Pub*. 2021;52(2):107-122.
- 968 56. Li Y. Chinese doctors connecting to the English publishing world: literature access,
969 editorial services, and training in publication skills. *Publications*. 2014;2:1-13.
- 970 57. Li Y. Chinese medical doctors negotiating the pressure of the publication requirement.
971 *Iberica*. 2014;28:107-126.
- 972 58. Tian M, Su Y, Ru X. Perish or publish in China: pressures on young Chinese scholars to
973 publish in internationally indexed journals. *Publications*. 2016;4:9.

- 974 59. Yi N, Nemery B, Dierickx K. Perceptions of research integrity and the Chinese situation:
975 In-depth interviews with Chinese biomedical researchers in Europe. *Account Res.*
976 2019;28:405-426. PMID: 31379202.
- 977 60. Barbour B, Stell BM. PubPeer: Scientific assessment without metrics. In: Biagioli M,
978 Lippman A, editors. *Gaming the metrics: Misconduct and manipulation in academic*
979 *research*. MIT Press; 2020. pp. 149-155.
- 980 61. Tsherniak A, Vazquez F, Montgomery PG, Weir BA, Kryukov G, Cowley GS, et al.
981 *Defining a Cancer Dependency Map*. *Cell*. 2017;170(3):564-576. PMID: 28753430
- 982 62. Chen L, Heikkinen L, Wang C, Yang Y, Sun H, Wong G. Trends in the development of
983 miRNA bioinformatics tools. *Brief Bioinform*. 2019;20(5):1836-1852. PMID: 29982332
- 984 63. Fanelli D, Costas R, Fang FC, Casadevall A, Bik EM. Testing Hypotheses on Risk
985 Factors for Scientific Misconduct via Matched-Control Analysis of Papers Containing
986 Problematic Image Duplications. *Sci Eng Ethics*. 2019;25(3):771-789. PMID: 29460082
- 987 64. Acuna DE, Brookes PS, Kording KP. Bioscience-scale automated detection of figure
988 element reuse. *bioRxiv* 2018;269415. doi: <https://doi.org/10.1101/269415>
- 989 65. Wang L, Liu Z. Keeping a clean research environment: Addressing research misconduct
990 and improving scientific integrity in China. *Cancer Lett*. 2019;464:1-4. PMID: 31430529
- 991 66. Zhang L, Sivertsen G. The New Research Assessment Reform in China and Its
992 Implementation. *Schol Assess Reports*. 2020;2(1):3.
- 993 67. Miyakawa T. No raw data, no science: another possible source of the reproducibility
994 crisis. *Mol Brain*. 2020;13(1):24. PMID: 32079532
- 995 68. Pinna N, Clavel G, Roco MC. The *Journal of Nanoparticle Research* victim of an
996 organized rogue editor network! *J Nanopart Res*. 2020;22:376.
- 997 69. Pines J. Image integrity and standards. *Open Biol*. 2020;10(6):200165. PMID: 32574547

- 998 70. Hackett R, Kelly S. Publishing ethics in the era of paper mills. *Biol Open*.
999 2020;9(10):bio056556. PMID: 33115756
- 1000 71. Seifert R. How Naunyn-Schmiedeberg's Archives of Pharmacology deals with fraudulent
1001 papers from paper mills. *Naunyn Schmiedebergs Arch Pharmacol*. 2021;394(3):431-436.
1002 PMID: 33547901
- 1003 72. Frederickson RM, Herzog RW. Keeping Them Honest: Fighting Fraud in Academic
1004 Publishing. *Mol Ther*. 2021;29(3):889-890. PMID: 33581045
- 1005 73. Calver M. Combatting the rise of paper mills. *Pacific Conserv Biol*. 2021;27:1-2.
- 1006 74. Fisher L, Cox R. *RSC Advances* Editorial: retraction of falsified manuscripts. *RSC Adv*.
1007 2021;11:4194.
- 1008 75. Behl C. Science integrity has been never more important: It's all about trust. *J Cell*
1009 *Biochem*. 2021;22(7):694-695. PMID: 33559144
- 1010 76. Heck S, Bianchini F, Souren NY, Wilhelm C, Ohl Y, Plass C. Fake data, paper mills, and
1011 their authors: The International Journal of Cancer reacts to this threat to scientific
1012 integrity. *Int J Cancer*. 2021;149(3):492-493. PMID: 33905542
- 1013 77. Cooper CDO, Han W. A new chapter for a better Bioscience Reports. *Biosci Rep*.
1014 2021;41(5):BSR20211016. PMID: 33942870
- 1015 78. Memon AR, Rathore FA. The rising menace of scholarly black market Challenges and
1016 solutions for improving research in low- and middle-income countries. *J Pak Med Assoc*.
1017 2021;71(6):1523-1526. PMID: 34111064
- 1018 79. Bik E. The Tadpole paper mill. [https://scienceintegritydigest.com/2020/02/21/the-](https://scienceintegritydigest.com/2020/02/21/the-tadpole-paper-mill/)
1019 [tadpole-paper-mill/](https://scienceintegritydigest.com/2020/02/21/the-tadpole-paper-mill/)
- 1020 80. Schneider J, Ye D, Hill AM, Whitehorn AS. Continued post-retraction citation of a
1021 fraudulent clinical trial report, 11 years after it was retracted for falsifying data.
1022 *Scientometrics*. 2020;125:2877-2913.

- 1023 81. Lee CM, Barber GP, Casper J, Clawson H, Diekhans M, Gonzalez JN, et al. UCSC
1024 Genome Browser enters 20th year. *Nucleic Acids Res.* 2020;48(D1):D756-D761. PMID:
1025 31691824
- 1026 82. Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al.
1027 SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat Methods.*
1028 2020;17(3):261-272. PMID: 32015543
- 1029 83. Hutchins BI, Baker KL, Davis MT, Diwersy MA, Haque E, Harriman RM, et al. The
1030 NIH Open Citation Collection: A public access, broad coverage resource. *PLoS Biol.*
1031 2019;17(10):e3000385. PMID: 31600197
- 1032

1033 **Figure and Supporting Information Captions**

1034 **Fig 1. Diagram describing the five literature corpora screened by S&B.** For each corpus
1035 (top row), the diagram shows the numbers of articles that were (i) screened by S&B (white),
1036 (ii) flagged by S&B with sequences manually verified (grey), and (iii) found to be
1037 problematic by describing at least one wrongly identified nucleotide sequence (dark grey).
1038 Total numbers of problematic papers and wrongly identified sequences are indicated below
1039 the diagram, corrected for duplicate papers between the corpora.

1040

1041 **Fig 2. Percentages of sequence identity error types in each corpus.** Percentages of
1042 wrongly identified nucleotide sequence reagents that correspond to the 3 identity error types
1043 (Y axis) in each corpus (X axis). Percentages corresponding to each error type are indicated,
1044 rounded to the nearest single digit. The numbers of incorrect sequences in each corpus are
1045 shown below the X axis.

1046

1047 **Fig 3. Percentages of wrongly identified nucleotide sequences that were either unique or**
1048 **repeated within each corpus.** Percentages of wrongly identified sequences that were
1049 identified at least twice in any single corpus (black) are shown above each image, rounded to
1050 the nearest single digit. All other wrongly identified sequences were unique in the indicated
1051 corpus (grey). Numbers of wrongly identified sequences identified in each corpus are shown
1052 below each image.

1053

1054 **Fig 4. Percentages of problematic *Gene* and *Oncology Reports* papers according to**
1055 **hospital affiliation status and country of origin.** Percentages of problematic *Gene* and

1056 *Oncology Reports* papers according to hospital affiliation status (Y axis) from either China or
1057 all other countries (X axis). The journal and relevant date ranges of problematic papers are
1058 shown above each panel. Problematic papers that were (not) affiliated with hospitals are
1059 shown in blue (grey), respectively. Percentages shown have been rounded to the nearest
1060 single digit. Numbers of problematic papers from China or all other countries are indicated
1061 below the X axis. For the comparisons shown in each panel, significantly higher proportions
1062 of problematic papers from China were affiliated with hospitals versus problematic papers
1063 from other countries (Fisher's Exact test, $p < 0.001$).

1064

1065 **Fig 5. Percentages and numbers of problematic *Gene* and *Oncology Reports* papers per**
1066 **year.** Percentages of all *Gene* or *Oncology Reports* papers that were found to be problematic
1067 (Y axis) per publication year (X axis). The journal and relevant publication year ranges are
1068 shown above each panel. Problematic papers from China or all other countries are shown in
1069 orange or grey, respectively. Percentages shown are rounded to one decimal place. Total
1070 numbers of problematic papers per year are shown below each graph.

1071

1072 **Fig 6. Summary of (RT-)PCR primer pairings that involved at least one wrongly**
1073 **identified primer.** For $n=851$ primer pairs that were claimed to target particular genes/
1074 sequences (gene X) (left panel), one or both primers were predicted to be incorrect (right
1075 panel), either by targeting unrelated genes or sequences (gene Y or gene Z), or by having no
1076 predicted human target (no target). Numbers of primer pairs and affected papers are indicated
1077 below each incorrect primer pair category. Some problematic papers described more than one
1078 (category of) incorrect primer pairing. Left- or right-hand primers are not intended to indicate
1079 forward or reverse primer orientations.

1080

1081 **Fig 7. Numbers of past research papers that have studied human protein-coding genes**
1082 **in problematic papers.** Numbers (log base 10) of problematic papers (Y axis) versus past
1083 research papers (X axis) for A: primary protein-coding genes in problematic papers and B:
1084 claimed protein-coding gene targets of wrongly identified reagents. Vertical dashed lines
1085 indicate the median number of research papers for protein-coding genes, with the associated
1086 interquartile range shown in grey. Subsets of protein-coding genes are highlighted in each
1087 panel.

1088

1089 **Fig 8. Clinical trial citations and approximate potential to translate for problematic**
1090 **papers.** A: Percentages of problematic papers that are cited at least once according to the
1091 NIH Open Citation Collection (Y axis), according to publication corpus (X axis). Error bars
1092 indicate 95% confidence intervals of bootstrapped estimates of percentages. Numbers of
1093 problematic papers with at least one clinical citation are shown below the X axis for each
1094 corpus. B: Average approximate potential to translate for problematic papers (Y axis)
1095 according to publication corpus (X axis). Error bars indicate bootstrapped 95% confidence
1096 intervals. Numbers of problematic papers for which the approximate potential to translate
1097 was computed by iCite are shown below the X axis for each corpus.

1098

1099 **S1 Fig. Problematic *Gene* papers per year, according to country of origin and**
1100 **institutional affiliation type.** Total numbers of problematic papers for each country/ group
1101 of countries are shown in the upper left corner of each panel. A: Numbers of problematic
1102 *Gene* papers (Y axes) per publication year (X axes) according to country of origin, shown
1103 above each graph. Countries are shown in alphabetical order, from top left. B: Numbers of

1104 problematic *Gene* papers (Y axis) per publication year (X axis) from China (right panel) or
1105 all other countries (left panel). Papers affiliated with hospitals or other institution types are
1106 shown in blue or grey, respectively. Numbers of problematic papers per year are shown
1107 below each stacked bar graph.

1108

1109 **S2 Fig. Problematic *Oncology Reports* papers per year, according to country of origin**
1110 **and institutional affiliation type.** Total numbers of problematic papers for each country/
1111 group of countries are shown in the upper left corner of each panel. A: Numbers of
1112 problematic *Oncology Reports* papers (Y axes) per publication year (X axes) according to
1113 country of origin, shown above each graph. Countries are shown in alphabetical order, from
1114 top left. B: Numbers of problematic *Oncology Reports* papers (Y axis) per publication year
1115 (X axis) from China (right panel) or all other countries (left panel). Papers affiliated with
1116 hospitals or other institution types are shown in blue or grey, respectively. Numbers of
1117 problematic papers per year are shown below each stacked bar graph.

1118

1119 **S3 Fig. Human protein-coding genes in problematic papers appear frequently in**
1120 **PubMed.** Numbers (log base 10) of PubMed papers (Y axis) for primary protein-coding
1121 genes in problematic papers (A) (green) or claimed protein-coding gene targets of wrongly
1122 identified reagents (B) (green), versus all other human protein-coding genes (A, B) (grey).
1123 Horizontal lines within box plots indicate median values, with box plots showing percentiles
1124 according to letter-proportions, ie +/- 25% percentile, +/- (25 + 25/2)% percentile, +/- (25 +
1125 25/2 + 25/4)% percentile.

1126

1127 **S1 Table. Wrongly identified nucleotide sequence reagents.** Wrongly identified sequences
1128 are listed for each corpus in separate tabs, as well as a combined list of all wrongly identified
1129 sequences. Primary human protein-coding and ncRNA's refer to the first-mentioned gene in
1130 each category in each publication title.

1131 **S2 Table. Author corrections, expressions of concern and retractions associated with**
1132 **problematic papers.**

1133 **S3 Table. Problematic papers curated within gene knowledgebases.** Problematic papers
1134 identified within CancerMine, ChimerDB, emiRIT, miRText, and mirPub are shown in
1135 separate tabs.

1136

1137 **S1 Data. Single gene knockdown (SGK) corpus.** All SGK papers screened by S&B and
1138 problematic SGK papers are shown in separate tabs. Primary human genes refer to the first-
1139 mentioned gene in the publication title or abstract.

1140 **S2 Data. *miR-145* corpus.** All *miR-145* papers flagged by S&B screening and problematic
1141 *miR-145* papers are shown in separate tabs. Primary human genes refer to the first-mentioned
1142 gene(s) in the publication title or abstract.

1143 **S3 Data. Cisplatin + gemcitabine (C+G) corpus.** All C+G papers flagged by S&B
1144 screening and problematic C+G papers are shown in separate tabs. Primary human genes
1145 refer to the first-mentioned gene(s) in the publication title or abstract.

1146 **S4 Data. Problematic *Gene* papers.** Primary human genes refer to the first-mentioned
1147 gene(s) in the publication title or abstract.

1148 **S5 Data. Problematic *Oncology Reports* corpus.** Primary human genes refer to the first-
1149 mentioned gene(s) in the publication title or abstract.

1150 **S6 Data. Journals and publishers that have published problematic papers.** Journals and
1151 publishers of problematic SGK, *miR-145* and C+G papers and all problematic papers are
1152 shown in separate tabs.

Table 1: Descriptions of the targeted corpora screened by Seek & Blastn with manual verification of nucleotide sequence reagent identities

	Single gene knockdown (SGK)		<i>miR-145</i>		Cisplatin + Gemcitabine (C+G)	
	Corpus	Problematic	Corpus	Problematic	Corpus	Problematic
Number of papers (% of corpus)	174 (100%)	75 (43%)	50 (100%)	31 (62%)	100 (100%)	51 (50%)
Number of journals	83	42	35	25	48	31
Publication year median (range)	2015 (2006-2019) ^a	2015 (2010-2019)	2017 (2009-2019)	2017 (2009-2019)	2017 (2008-2019)	2017 (2009-2019)
Journal Impact Factor at publication year median (range)	2.204 (0.098-8.459)	1.778 (0.098-5.712)	3.34 (0.700-9.050)	3.23 (0.700-8.278)	3.571 (1.099-10.391)	3.041 (1.099-8.579)
Number of sequences/ paper median (range)	6 (0-24)	6 (1-24)	11 (4-46)	10 (4-46)	11 (2-71)	12 (2-70)
Number of incorrect sequences/ paper median (range)	ND	1 (1-8)	ND	1 (1-5)	ND	2 (1-8)
Papers from China proportion (%)	159/174 (91%)	73/75 (97%)	44/50 (88%)	31/31 (100%)	90/100 (90%)	50/51 (98%)
Papers from China affiliated with hospitals proportion (%)	147/159 (92%)	68/73 (93%)	40/44 (91%)	28/31 (90%)	82/90 (91%)	48/50 (96%)
Papers from all other countries affiliated with hospitals proportion (%)	6/15 (40%)	0/2 (0%)	0/6 (0%)	0/3 (0%)	1/10 (10%)	0/1 (0%)
Papers with post-publication notices ^b proportion (%)	20/174 (12%)	13/75 (17%)	1/50 (2%)	1/31 (3%)	1/100 (1%)	1/51 (2%)

^aSGK papers were published until June 2019

^bPost-publication notices include retractions, expressions of concern and corrections

Table 2: Cancer types studied in the Single Gene Knockdown (SGK) corpus, where each cancer type corresponds to a single paper

Gene	Previously reported SGK papers	New SGK papers
<i>ADAM8</i>	Liver ^a	Breast, Breast, Colorectal, Gastric, Liver, Lung, Pancreatic
<i>ANXA1</i>	N/A	Breast , Breast, Breast, Esophageal, Leukemia, Liver, Lung, Prostate
<i>EAG1</i>	Liposarcoma, Osteosarcoma	Brain, Osteosarcoma, Osteosarcoma, Ovarian, Sarcoma
<i>GPR137</i>	Bladder, Brain, <u>Colorectal</u>^b, Pancreatic	Brain, Gastric, Leukemia , Liver, Osteosarcoma, Ovarian , Prostate
<i>ICT1</i>	<u>Brain</u>	Breast, Gastric, Leukemia, Lung, Lymphoma, Prostate
<i>KLF8</i>	Osteosarcoma	Bladder, Brain, Brain, Brain, Breast, Colorectal, Colorectal, Gastric, Gastric, Gastric, Gastric, Liver, Liver, Nasopharyngeal, Oral, Ovarian, Pancreatic, Renal
<i>MACC1</i>	Ovarian	Bladder, Brain, Cervical, Cervical, Colorectal, Colorectal, Esophageal, Esophageal, Gallbladder, Gastric, Liver, Liver, Lung, Oral, Oral, Nasopharyngeal, Ovarian, Ovarian, Skin
<i>MYO6</i>	<u>Brain, Colorectal, Liver, Lung</u>	Breast, Gastric, Oral, Prostate
<i>NOB1</i>	Brain, Breast , Colorectal, Liver, Osteosarcoma , Ovarian, Prostate	Laryngeal, Lung, Lung , Oral, Osteosarcoma, Renal, Thyroid, Thyroid
<i>PP4R1</i>	<u>Breast, Liver</u>	Lung
<i>PP5</i>	Colorectal, Ovarian	Bladder, Brain, Leukemia , Liver, Osteosarcoma, Pancreatic, Prostate
<i>PPM1D</i>	Bladder, Lung	Brain, Brain, Breast, Breast, Liver, Pancreatic
<i>RPS15A</i>	Brain, Lung	Brain, Gastric, Leukemia , Liver, Lung, Osteosarcoma, Renal, Thyroid
<i>TCTN1</i>	<u>Brain</u> , Brain, Pancreatic	Brain, Colorectal, Gastric, Thyroid
<i>TPD52L2</i>	<u>Brain, Breast, Gastric, Liver, Oral</u>	Brain
<i>USP39</i>	<u>Liver, Thyroid</u>	Breast, Colorectal, Colorectal, Gastric, Liver, Liver, Lung, Oral, Osteosarcoma, Renal, Skin
<i>ZFX</i>	Brain, Breast	Brain, Brain, Gallbladder, Laryngeal, Leukemia, Lung, Lung, Oral, Oral, Osteosarcoma, Pancreatic, Prostate, Renal

^aCancer types shown in bold correspond to problematic papers with wrongly identified nucleotide sequence(s)

^bUnderlined cancer types correspond to papers that have been retracted or assigned an expression of concern

Table 3: Wrongly identified nucleotide sequences summarized according to experimental technique and identity error type

Corpus	Technique	“Non-targeting” yet targeting proportion (%)	“Targeting” yet non-targeting proportion (%)	Targeting wrong gene/ sequence proportion (%)	Total per corpus proportion (%)
SGK (n = 115 reagents in n = 75 papers)	PCR ^a	0/45 (0)	7/14 (50)	45/57 (79)	52/115 (45)
	Gene knockdown ^b	44/44 (100)	7/14 (50)	12/57 (21)	63/115 (55)
	Other ^c	0/45 (0)	0/14 (0)	0/57 (0)	0/115 (0)
	Total (Error type)	44/44 (100)	14/14 (100)	57/57 (100)	115/115 (100)
miR-145 (n = 49 reagents in n = 31 papers)	PCR	0/2 (0)	8/9 (89)	33/38 (87)	41/49 (84)
	Gene knockdown	2/2 (100)	1/9 (11)	5/38 (13)	8/49 (16)
	Other	0/2 (0)	0/9 (0)	0/38 (0)	0/49 (0)
	Total (Error type)	2/2 (100)	9/9 (100)	38/38 (100)	49/49 (100)
C+G (n = 109 reagents in n = 51 papers)	PCR	0/6 (0)	23/24 (96)	73/79 (93)	96/109 (88)
	Gene knockdown	4/6 (67)	1/24 (4)	5/79 (6)	10/109 (9)
	Other	2/6 (33)	0/24 (0)	1/79 (1)	3/109 (3)
	Total (Error type)	6/6 (100)	24/24 (100)	79/79 (100)	109/109 (100)
Gene (n = 284 reagents in n = 128 papers)	PCR	0/9 (0)	35/42 (83)	218/233 (94)	253/284 (88)
	Gene knockdown	9/9 (100)	7/42 (17)	15/233 (6)	31/284 (11)
	Other	0/9 (0)	0/42 (0)	0/233 (0)	0/284 (0)
	Total (Error type)	9/9 (100)	42/42 (100)	233/233 (100)	284/284 (100)
Oncology Reports (n = 995 reagents in n = 436 papers)	PCR	0/36 (0)	296/335 (88)	573/630 (91)	869/995 (87)
	Gene knockdown	30/30 (100)	37/335 (11)	54/630 (8)	121/995 (12)
	Other	0/36 (0)	2/335 (1)	3/630 (1)	5/995 (1)
	Total (Error type)	30/30 (100)	335/335 (100)	630/630 (100)	995/995 (100)
Total (n = 1,535 reagents in n = 712 papers)	PCR	0/89 (0)	364/416 (87)	937/1,030 (90)	1,301/1,535 (84)
	Gene knockdown	87/89 (98)	50/416 (12)	89/1,030 (9)	226/1,535 (15)
	Other	2/89 (2)	2/416 (1)	4/1,030 (1)	8/1,535 (1)
	Total (Error type)	89/89 (100)	416/416 (100)	1,030/1,030 (100)	1,535/1,535 (100)

^aPCR = Human gene or genomic targeting primers for PCR, RT-PCR or methylation-specific PCR

^bGene knockdown = siRNA or shRNA

^cOther = Claimed Ribozyme, Talen, mimic sequences, other oligonucleotide sequences

Table 4: Summary of features of *Gene* and *Oncology Reports* journals and problematic papers

Feature	<i>Gene</i>	<i>Oncology Reports</i>
Publication years screened by Seek & Blastn	2007-2018	2014-2018
Journal Impact Factor (range during years screened)	2.082-2.871	2.301-3.041
Flagged/ screened papers proportion (%)	742/7,399 (10%)	1,709/3,778 (45%)
Problematic/ flagged papers proportion (%)	128/742 (17%)	436/1,709 (26%)
Incorrect sequences/ problematic paper median (range)	2 (1-36)	2 (1-15)
Problematic papers from China proportion (%)	69/128 (54%)	393/436 (90%)
Problematic papers from all other countries proportion (%)	59/128 (46%)	43/436 (10%)
Problematic papers from China affiliated with hospitals proportion (%)	54/69 (78%)	342/393 (87%)
Problematic papers from all other countries affiliated with hospitals proportion (%)	5/59 (9%)	5/43 (12%)
Retracted or corrected problematic papers proportion (%)	2/128 (2%)	2/436 (0.5%)

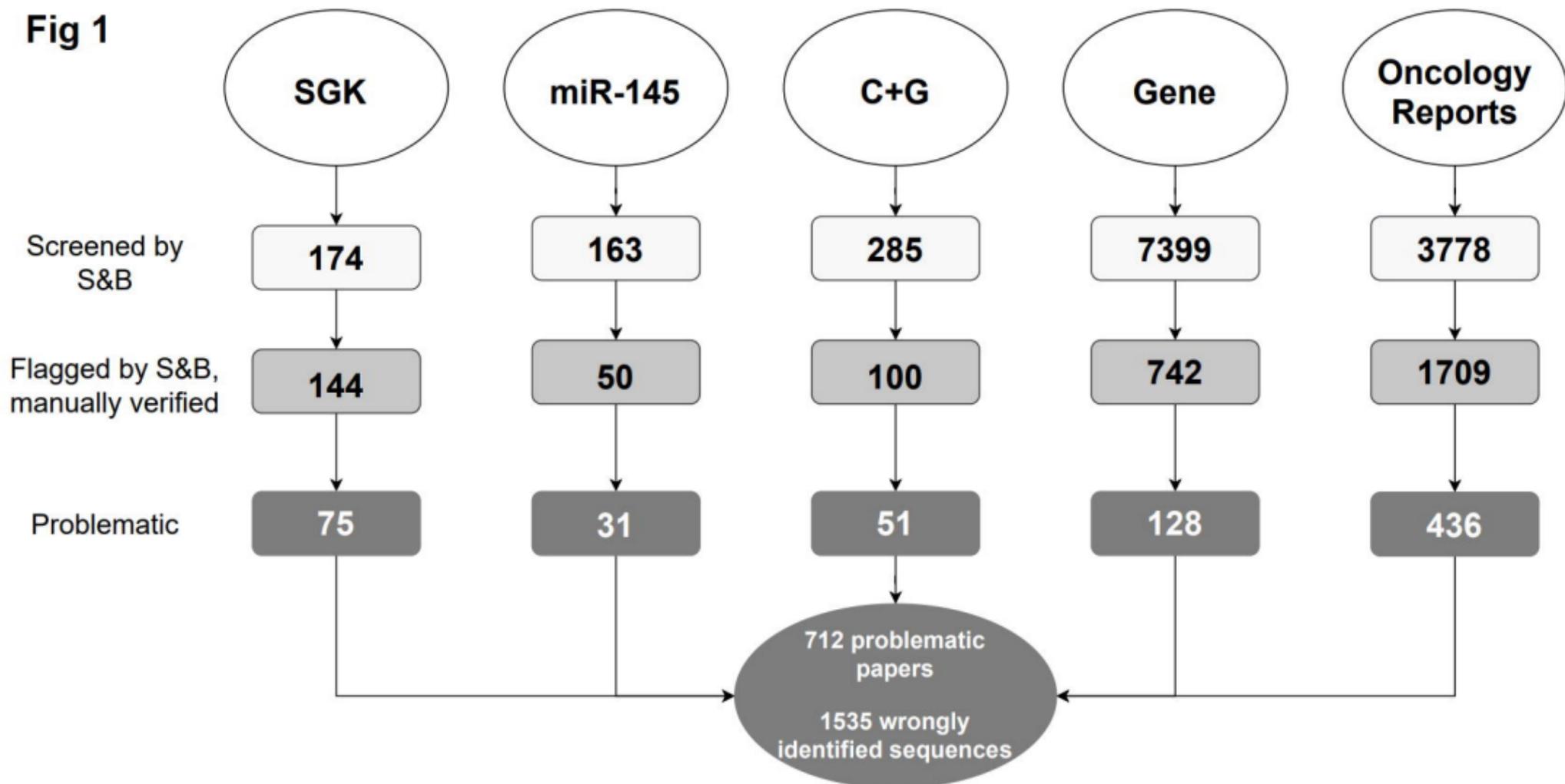
Fig 1

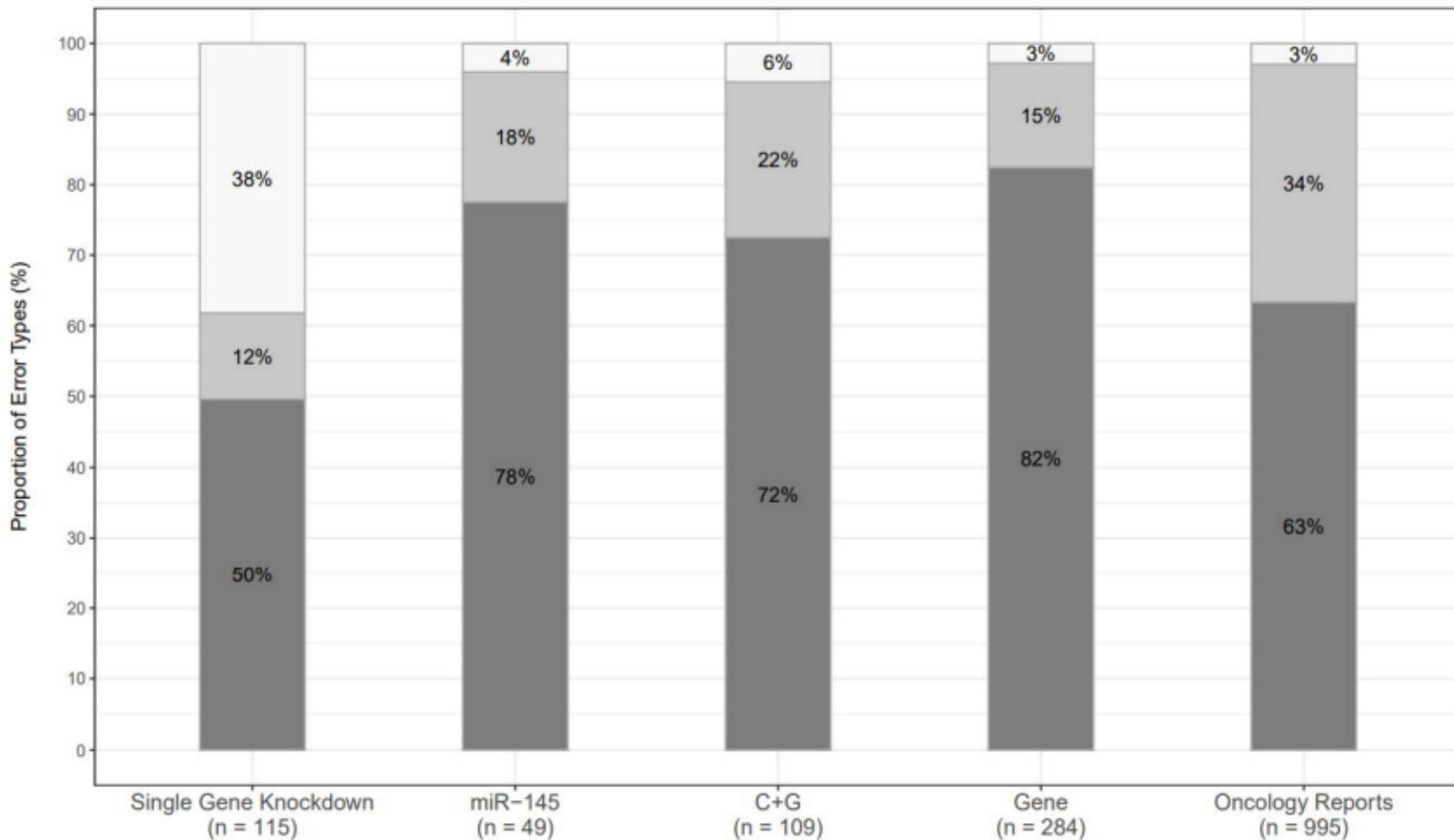
Fig 2Error Types  'Non-targeting' yet targeting  'Targeting' yet non-targeting  Targeting wrong gene or genomic sequence

Fig 3

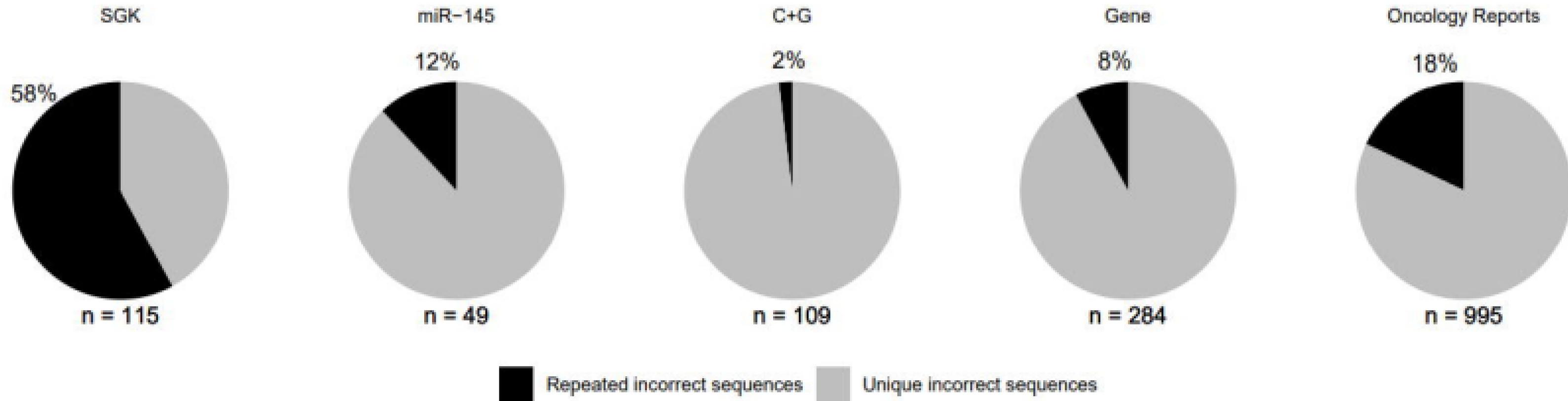


Fig 4

Publications not from hospitals (grey) Publications from hospitals (blue)

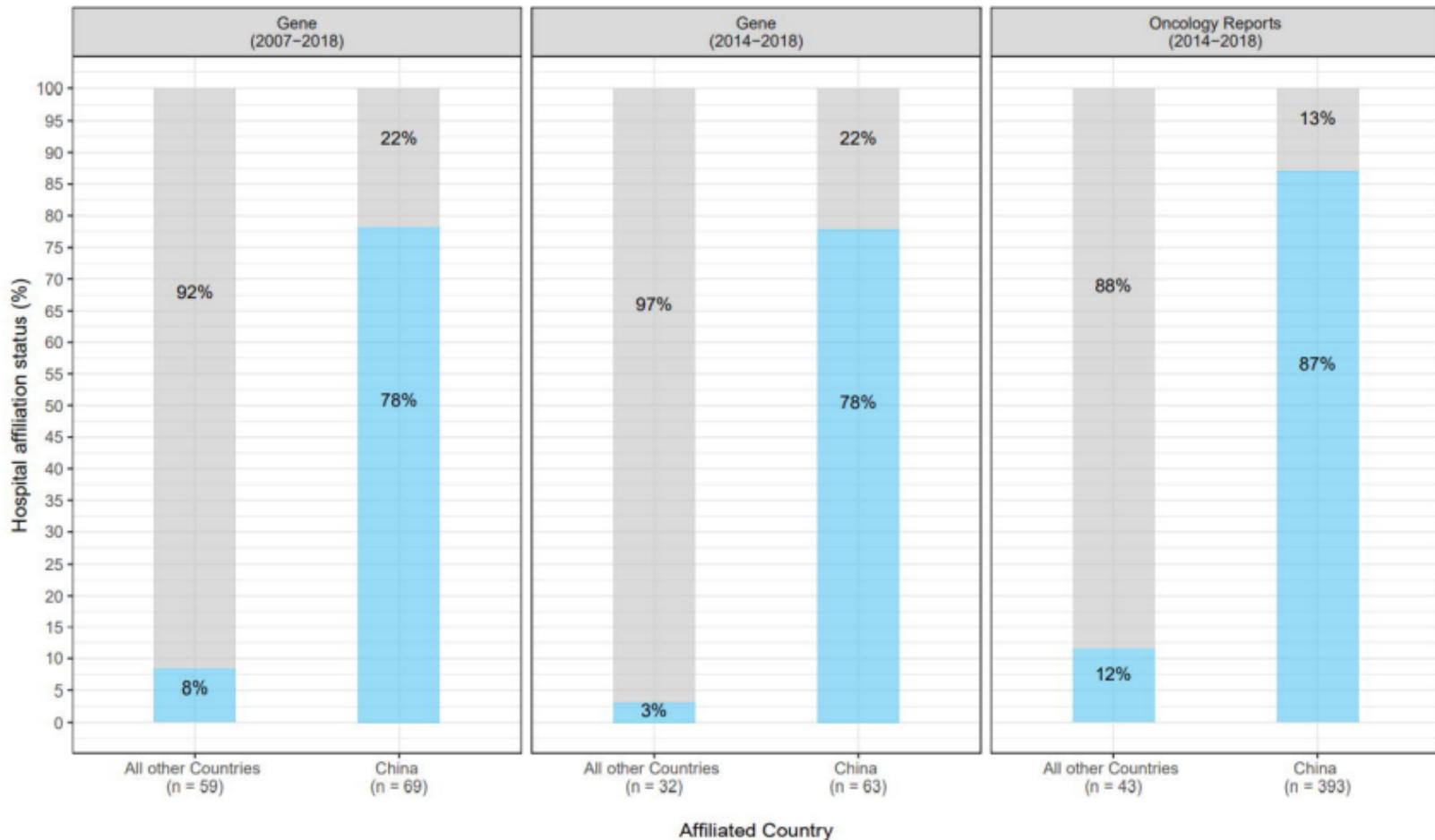


Fig 5

All other Countries China

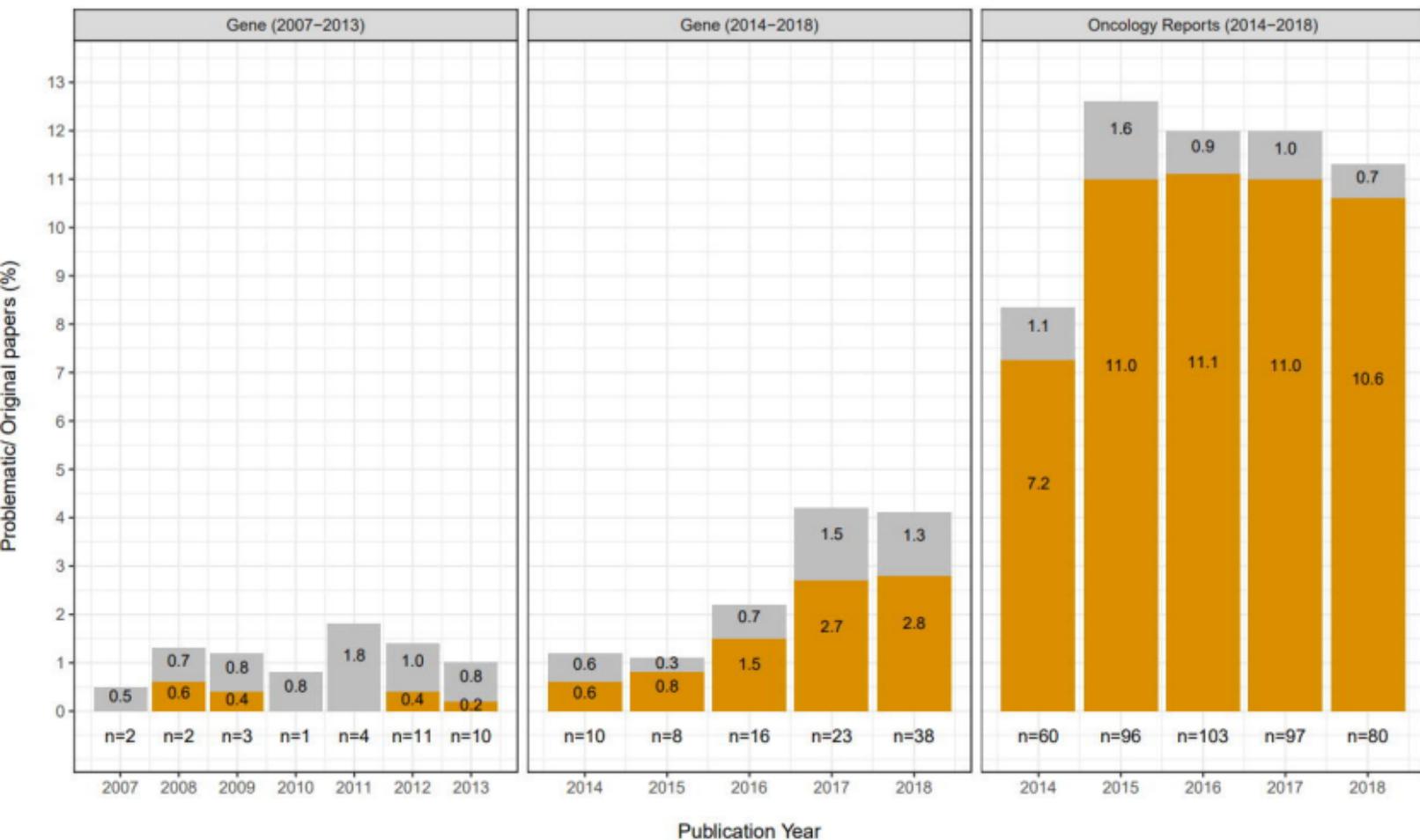


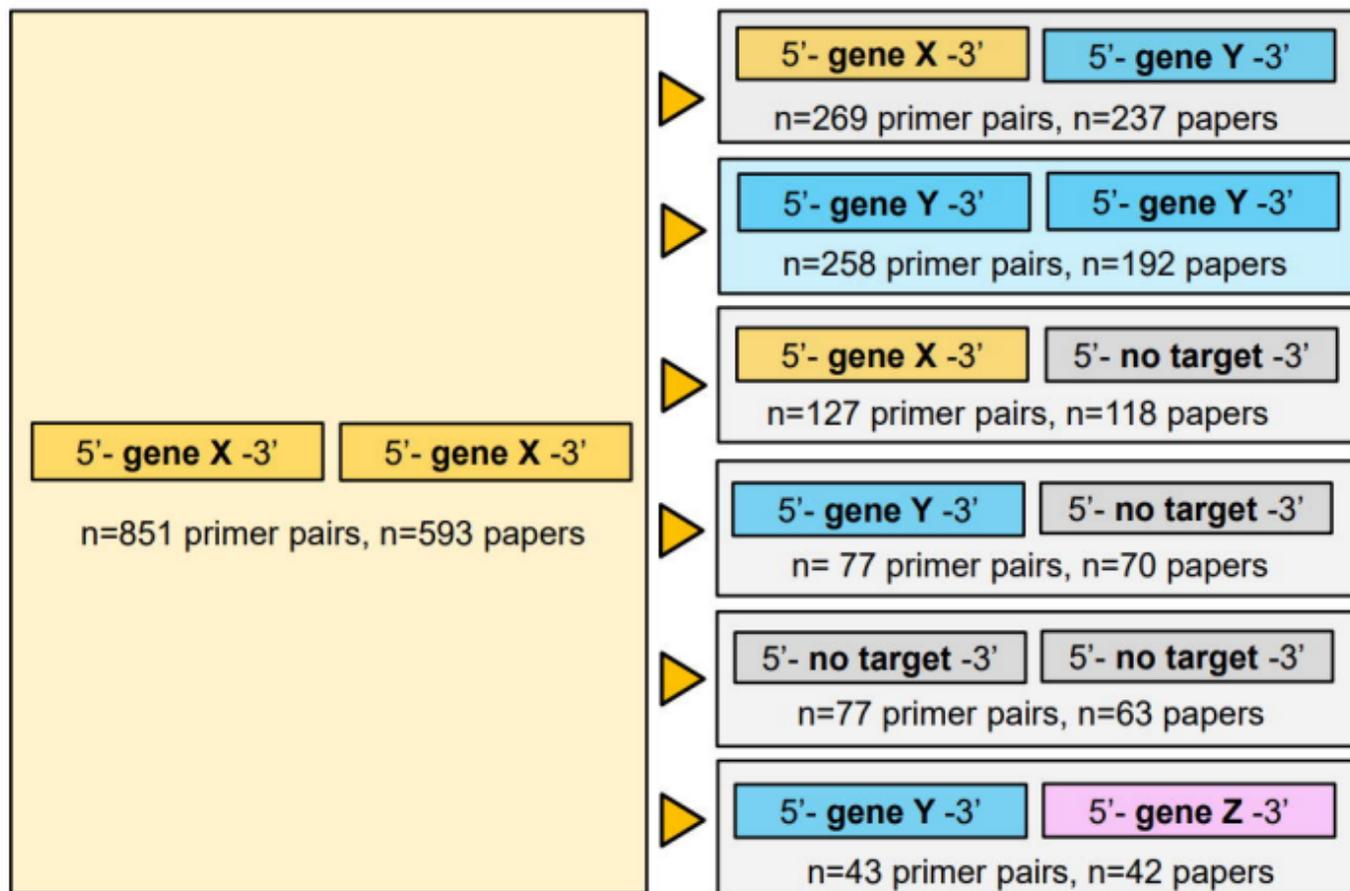
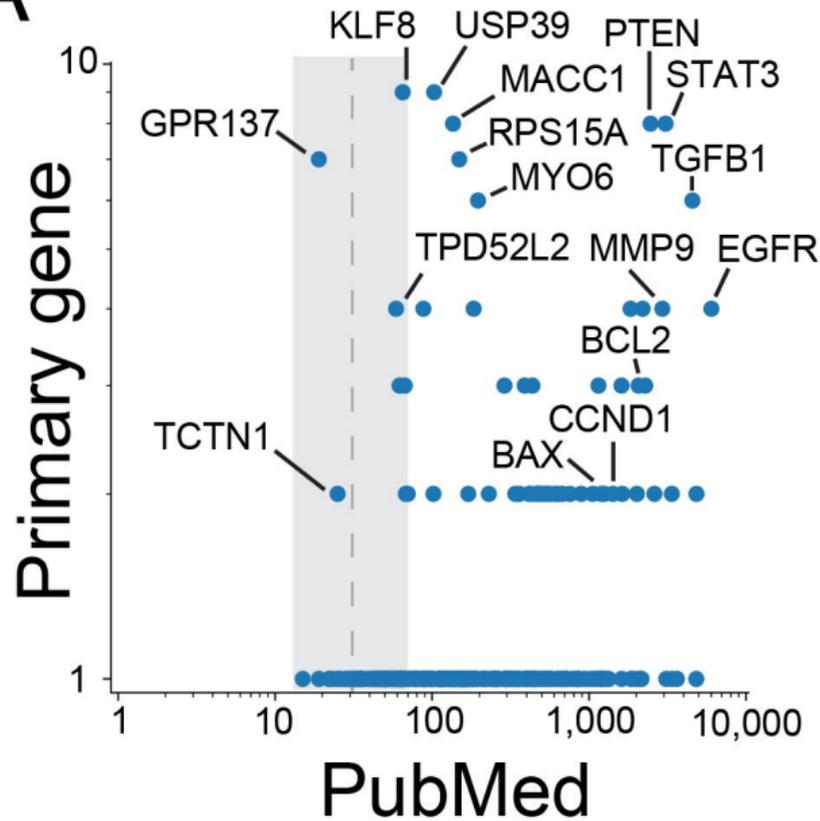
Fig 6**Claimed primer identities****Verified primer identities**

Fig 7

A



B

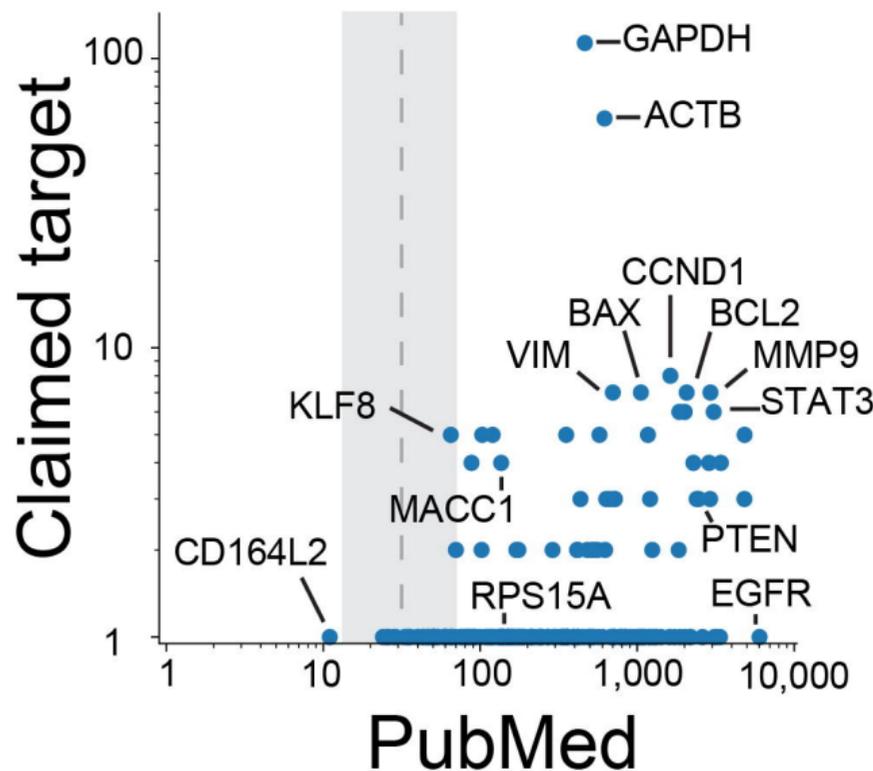
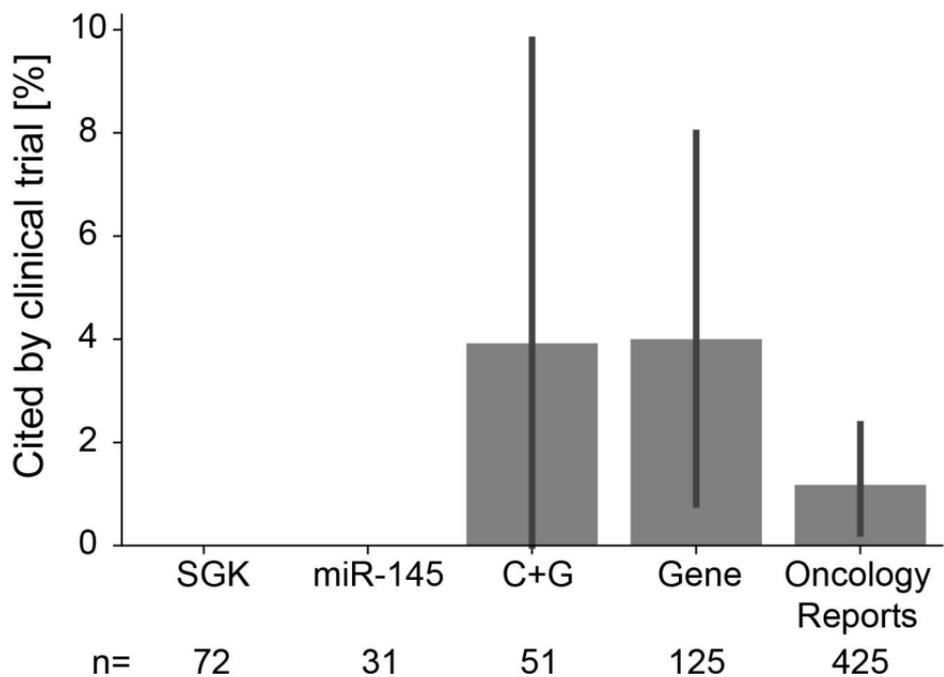


Fig 8

A



B

