

1 **Characterising and justifying sample size sufficiency in interview-based**
2 **studies: systematic analysis of qualitative health research over a 15-year**
3 **period**

4
5 **Authors:** Konstantina Vasileiou ^{1*}, Julie Barnett ¹, Susan Thorpe ², Terry Young ³

6
7 **Affiliations**

8 ¹ Department of Psychology, University of Bath, Building 10 West, Claverton Down, Bath, BA2 7AY,
9 United Kingdom. K.Vasileiou@bath.ac.uk; J.C.Barnett@bath.ac.uk

10 ² School of Psychology, Newcastle University, Ridley Building 1, Queen Victoria Road, Newcastle
11 upon Tyne, NE1 7RU, United Kingdom. susan.thorpe@newcastle.ac.uk

12 ³ Department of Computer Science, Brunel University London, Wilfred Brown Building 108, Uxbridge,
13 UB8 3PH, United Kingdom. terry.young@brunel.ac.uk

14
15 *** Corresponding author:** Building 10 West, Claverton Down, Bath, BA2 7AY, United Kingdom. E-mail
16 address: K.Vasileiou@bath.ac.uk; Phone: +44 (0) 1225 383167

17

18

19

20

21 **Abstract**

22 **Background:** Choosing a suitable sample size in qualitative research is an area of conceptual debate
23 and practical uncertainty. That sample size principles, guidelines and tools have been developed to
24 enable researchers to set, and justify the acceptability of, their sample size is an indication that the
25 issue constitutes an important marker of the quality of qualitative research. Nevertheless, research
26 shows that sample size sufficiency reporting is often poor, if not absent, across a range of
27 disciplinary fields.

28 **Methods:** A systematic analysis of single-interview-per-participant designs within three health-
29 related journals from the disciplines of psychology, sociology and medicine, over a 15-year period,
30 was conducted to examine whether and how sample sizes were justified and how sample size was
31 characterised and discussed by authors. Data pertinent to sample size were extracted and analysed
32 using qualitative and quantitative analytic techniques.

33 **Results:** Our findings demonstrate that provision of sample size justifications in qualitative health
34 research is limited; is not contingent on the number of interviews; and relates to the journal of
35 publication. Defence of sample size was most frequently supported across all three journals with
36 reference to the principle of saturation and to pragmatic considerations. Qualitative sample sizes
37 were predominantly – and often without justification – characterised as insufficient (i.e., ‘small’) and
38 discussed in the context of study limitations. Sample size insufficiency was seen to threaten the
39 validity and generalizability of studies’ results, with the latter being frequently conceived in
40 nomothetic terms.

41 **Conclusions:** We recommend, firstly, that qualitative health researchers be more transparent about
42 evaluations of their sample size sufficiency, situating these within broader and more encompassing
43 assessments of *data adequacy*. Secondly, we invite researchers critically to consider how saturation
44 parameters found in prior methodological studies and sample size community norms might best
45 inform, and apply to, their own project and encourage that data adequacy is best appraised with

46 reference to features that are *intrinsic* to the study at hand. Finally, those reviewing papers have a
47 vital role in supporting and encouraging transparent study-specific reporting.

48

49 **Keywords:** sample size; sample size justification; sample size characterisation; data adequacy;
50 qualitative health research; qualitative interviews; review; systematic analysis

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67 **Background**

68 Sample adequacy in qualitative inquiry pertains to the appropriateness of the sample *composition*
69 and *size*. It is an important consideration in evaluations of the quality and trustworthiness of much
70 qualitative research [1] and is implicated – particularly for research that is situated within a post-
71 positivist tradition and retains a degree of commitment to realist ontological premises – in appraisals
72 of validity and generalizability [2-5].

73 Samples in qualitative research tend to be small in order to support the depth of case-oriented
74 analysis that is fundamental to this mode of inquiry [5]. Additionally, qualitative samples are
75 purposive, that is, selected by virtue of their capacity to provide richly-textured information,
76 relevant to the phenomenon under investigation. As a result, purposive sampling [6,7] – as opposed
77 to probability sampling employed in quantitative research – selects ‘information-rich’ cases [8].
78 Indeed, recent research demonstrates the greater efficiency of purposive sampling compared to
79 random sampling in qualitative studies [9], supporting related assertions long put forward by
80 qualitative methodologists.

81 Sample size in qualitative research has been the subject of enduring discussions [4,10,11]. Whilst the
82 quantitative research community has established relatively straightforward statistics-based rules to
83 set sample sizes precisely, the intricacies of qualitative sample size determination and assessment
84 arise from the methodological, theoretical, epistemological, and ideological pluralism that
85 characterises qualitative inquiry (for a discussion focused on the discipline of psychology see [12]).
86 This mitigates against clear-cut guidelines, invariably applied. Despite these challenges, various
87 conceptual developments have sought to address this issue, with guidance and principles
88 [4,10,11,13-20], and more recently, an evidence-based approach to sample size determination seeks
89 to ground the discussion empirically [21-35].

90 Focusing on single-interview-per-participant qualitative designs, the present study aims to further
91 contribute to the dialogue of sample size in qualitative research by offering empirical evidence

92 around justification practices associated with sample size. We next review the existing conceptual
93 and empirical literature on sample size determination.

94 **Sample size in qualitative research: conceptual developments and empirical investigations**

95 Qualitative research experts argue that there is no straightforward answer to the question of ‘how
96 many’ and that sample size is contingent on a number of factors relating to epistemological,
97 methodological and practical issues [36]. Sandelowski [4] recommends that qualitative sample sizes
98 are large enough to allow the unfolding of a ‘new and richly textured understanding’ of the
99 phenomenon under study, but small enough so that the ‘deep, case-oriented analysis’ (p. 183) of
100 qualitative data is not precluded. Morse [11] posits that the more useable data are collected from
101 each person, the fewer participants are needed. She invites researchers to take into account
102 parameters, such as the scope of study, the nature of topic (i.e. complexity, accessibility), the quality
103 of data, and the study design. Indeed, the level of structure of questions in qualitative interviewing
104 has been found to influence the richness of data generated [37], and so, requires attention;
105 empirical research shows that open questions, which are asked later on in the interview, tend to
106 produce richer data [37].

107 Beyond such guidance, specific numerical recommendations have also been proffered, often based
108 on experts’ experience of qualitative research. For example, Green and Thorogood [38] maintain
109 that the experience of most qualitative researchers conducting an interview-based study with a fairly
110 specific research question is that little new information is generated after interviewing 20 people or
111 so belonging to one analytically relevant participant ‘category’ (pp. 102-104). Ritchie et al. [39]
112 suggest that studies employing individual interviews conduct no more than 50 interviews so that
113 researchers are able to manage the complexity of the analytic task. Similarly, Britten [40] notes that
114 large interview studies will often comprise of 50 to 60 people. Experts have also offered numerical
115 guidelines tailored to different theoretical and methodological traditions and specific research
116 approaches, e.g. grounded theory, phenomenology [11, 41]. More recently, a quantitative tool was

117 proposed [42] to support a priori sample size determination based on estimates of the prevalence of
118 themes in the population. Nevertheless, this more formulaic approach raised criticisms relating to
119 assumptions about the conceptual [43] and ontological status of 'themes' [44] and the linearity
120 ascribed to the processes of sampling, data collection and data analysis [45].

121 In terms of principles, Lincoln and Guba [17] proposed that sample size determination be guided by
122 the criterion of *informational redundancy*, that is, sampling can be terminated when no new
123 information is elicited by sampling more units. Following the logic of informational
124 comprehensiveness Malterud et al. [18] introduced the concept of *information power* as a pragmatic
125 guiding principle, suggesting that the more information power the sample provides, the smaller the
126 sample size needs to be, and vice versa.

127 Undoubtedly, the most widely used principle for determining sample size and evaluating its
128 sufficiency is that of *saturation*. The notion of saturation originates in grounded theory [15] – a
129 qualitative methodological approach explicitly concerned with empirically-derived theory
130 development – and is inextricably linked to theoretical sampling. Theoretical sampling describes an
131 iterative process of data collection, data analysis and theory development whereby data collection is
132 governed by emerging theory rather than predefined characteristics of the population. Grounded
133 theory saturation (often called theoretical saturation) concerns the theoretical categories – as
134 opposed to data – that are being developed and becomes evident when 'gathering fresh data no
135 longer sparks new theoretical insights, nor reveals new properties of your core theoretical
136 categories' [46, p. 113]. Saturation in grounded theory, therefore, does not equate to the more
137 common focus on data repetition and moves beyond a singular focus on sample size as the
138 justification of sampling adequacy [46, 47]. Sample size in grounded theory cannot be determined a
139 priori as it is contingent on the evolving theoretical categories.

140 Saturation – often under the terms of 'data' or 'thematic' saturation – has diffused into several
141 qualitative communities beyond its origins in grounded theory. Alongside the expansion of its

142 meaning, being variously equated with ‘no new data’, ‘no new themes’, and ‘no new codes’,
143 saturation has emerged as the ‘gold standard’ in qualitative inquiry [2, 26]. Nevertheless, and as
144 Morse [49] asserts, whilst saturation is the most frequently invoked ‘guarantee of qualitative rigor’,
145 ‘it is the one we know least about’ (p. 587). Certainly researchers caution that saturation is less
146 applicable to, or appropriate for, particular types of qualitative research (e.g. conversation analysis,
147 [48]; phenomenological research, [50]) whilst others reject the concept altogether [19, 51].

148 Methodological studies in this area aim to provide guidance about saturation and develop a practical
149 application of processes that ‘operationalise’ and evidence saturation. Guest, Bunce, and Johnson
150 [26] analysed 60 interviews and found that saturation of themes was reached by the twelfth
151 interview. They noted that their sample was relatively homogeneous, their research aims focused,
152 so studies of more heterogeneous samples and with a broader scope would be likely to need a larger
153 size to achieve saturation. Extending the enquiry to multi-site, cross-cultural research, Hagaman and
154 Wutich [28] showed that sample sizes of 20 to 40 interviews were required to achieve data
155 saturation of meta-themes that cut across research sites. In a theory-driven content analysis, Francis
156 et al. [25] reached data saturation at the 17th interview for all their pre-determined theoretical
157 constructs. The authors further proposed two main principles upon which specification of saturation
158 be based: (a) researchers should a priori specify an *initial analysis sample* (e.g. 10 interviews) which
159 will be used for the first round of analysis and (b) a *stopping criterion*, that is, a number of interviews
160 (e.g. 3) that needs to be further conducted, the analysis of which will not yield any new themes or
161 ideas. For greater transparency, Francis et al. [25] recommend that researchers present cumulative
162 frequency graphs supporting their judgment that saturation was achieved. A comparative method
163 for themes saturation (CoMeTS) has also been suggested [23] whereby the findings of each new
164 interview are compared with those that have already emerged and if it does not yield any new
165 theme, the ‘saturated terrain’ is assumed to have been established. Because the order in which
166 interviews are analysed can influence saturation thresholds depending on the richness of the data,
167 Constantinou et al. [23] recommend reordering and re-analysing interviews to confirm saturation.

168 Hennink, Kaiser and Marconi's [29] methodological study sheds further light on the problem of
169 specifying and demonstrating saturation. Their analysis of interview data showed that *code*
170 *saturation* (i.e. the point at which no additional issues are identified) was achieved at 9 interviews,
171 but *meaning saturation* (i.e. the point at which no further dimensions, nuances, or insights of issues
172 are identified) required 16-24 interviews. Although *breadth* can be achieved relatively soon,
173 especially for high-prevalence and concrete codes, *depth* requires additional data, especially for
174 codes of a more conceptual nature.

175 Critiquing the concept of saturation, Nelson [19] proposes five conceptual depth criteria in grounded
176 theory projects to assess the robustness of the developing theory: (a) theoretical concepts should be
177 supported by a wide range of evidence drawn from the data; (b) be demonstrably part of a network
178 of inter-connected concepts; (c) demonstrate subtlety; (d) resonate with existing literature; and (e)
179 can be successfully submitted to tests of external validity.

180 Other work has sought to examine practices of sample size reporting and sufficiency assessment
181 across a range of disciplinary fields and research domains, from nutrition [34] and health education
182 [32], to education and the health sciences [22, 27], information systems [30], organisation and
183 workplace studies [33], human computer interaction [21], and accounting studies [24]. Others
184 investigated PhD qualitative studies [31] and grounded theory studies [35]. Incomplete and
185 imprecise sample size reporting is commonly pinpointed by these investigations whilst assessment
186 and justifications of sample size sufficiency are even more sporadic.

187 Sobal [34] examined the sample size of qualitative studies published in the Journal of Nutrition
188 Education over a period of 30 years. Studies that employed individual interviews ($n = 30$) had an
189 average sample size of 45 individuals and none of these explicitly reported whether their sample size
190 sought and/or attained saturation. A minority of articles discussed how sample-related limitations,
191 (with the latter most often concerning the type of sample, rather than the size) limited
192 generalizability. A further systematic analysis [32] of health education research over 20 years

193 demonstrated that interview-based studies averaged 104 participants (range 2 to 720 interviewees).
194 However, 40% did not report the number of participants. An examination of 83 qualitative interview
195 studies in leading information systems journals [30] indicated little defence of sample sizes on the
196 basis of recommendations by qualitative methodologists, prior relevant work, or the criterion of
197 saturation. Rather, sample size seemed to correlate with factors such as the journal of publication or
198 the region of study (US vs Europe vs Asia). These results led the authors to call for more rigor in
199 determining and reporting sample size in qualitative information systems research and to
200 recommend optimal sample size ranges for grounded theory (i.e. 20-30 interviews) and single case
201 (i.e. 15-30 interviews) projects.

202 Similarly, fewer than 10% of articles in organisation and workplace studies provided a sample size
203 justification relating to existing recommendations by methodologists, prior relevant work, or
204 saturation [33], whilst only 17% of focus groups studies in health-related journals provided an
205 explanation of sample size (i.e. number of focus groups), with saturation being the most frequently
206 invoked argument, followed by published sample size recommendations and practical reasons [22].
207 The notion of saturation was also invoked by 11 out of the 51 most highly cited studies that
208 Guetterman [27] reviewed in the fields of education and health sciences, of which six were grounded
209 theory studies, four phenomenological and one a narrative inquiry. Finally, analysing 641 interview-
210 based articles in accounting, Dai et al. [24] called for more rigor since a significant minority of studies
211 did not report precise sample size.

212 Despite increasing attention to rigor in qualitative research (e.g. [52]) and more extensive
213 methodological and analytical disclosures that seek to validate qualitative work [24], sample size
214 reporting and sufficiency assessment remain inconsistent and partial, if not absent, across a range of
215 research domains.

216 **Objectives of the present study**

217 The present study sought to enrich existing systematic analyses of the customs and practices of
218 sample size reporting and justification by focusing on qualitative research relating to health.
219 Additionally, this study attempted to expand previous empirical investigations by examining how
220 qualitative sample sizes are characterised and discussed in academic narratives. Qualitative health
221 research is an inter-disciplinary field that due to its affiliation with medical sciences, often faces
222 views and positions reflective of a quantitative ethos. Thus qualitative health research constitutes an
223 *emblematic case* that may help to unfold underlying philosophical and methodological differences
224 across the scientific community that are crystallised in considerations of sample size. The present
225 research, therefore, incorporates a comparative element on the basis of three different disciplines
226 engaging with qualitative health research: medicine, psychology, and sociology. We chose to focus
227 our analysis on single-per-participant-interview designs as this not only presents a popular and
228 widespread methodological choice in qualitative health research, but also as the method where
229 consideration of sample size – defined as the number of interviewees – is particularly salient.

230 **Methods**

231 **Study design**

232 A structured search for articles reporting cross-sectional, interview-based qualitative studies was
233 carried out and eligible reports were systematically reviewed and analysed employing both
234 quantitative and qualitative analytic techniques.

235 We selected journals which (a) follow a peer review process, (b) are considered high quality and
236 influential in their field as reflected in journal metrics, and (c) are receptive to, and publish,
237 qualitative research (Additional File 1 presents the journals' editorial positions in relation to
238 qualitative research and sample considerations where available). Three health-related journals were
239 chosen, each representing a different disciplinary field; the *British Medical Journal* (BMJ)
240 representing medicine, the *British Journal of Health Psychology* (BJHP) representing psychology, and
241 the *Sociology of Health & Illness* (SHI) representing sociology.

242 **Search strategy to identify studies**

243 Employing the search function of each individual journal, we used the terms ‘interview*’ AND
244 ‘qualitative’ and limited the results to articles published between 1 January 2003 and 22 September
245 2017 (i.e. a 15-year review period).

246 **Eligibility criteria**

247 To be eligible for inclusion in the review, the article had to report a cross-sectional study design.
248 Longitudinal studies were thus excluded whilst studies conducted within a broader research
249 programme (e.g. interview studies nested in a trial, as part of a broader ethnography, as part of a
250 longitudinal research) were included if they reported only single-time qualitative interviews. The
251 method of data collection had to be individual, synchronous qualitative interviews (i.e. group
252 interviews, structured interviews and e-mail interviews over a period of time were excluded), and
253 the data had to be analysed qualitatively (i.e. studies that quantified their qualitative data were
254 excluded). Mixed method studies and articles reporting more than one qualitative method of data
255 collection (e.g. individual interviews and focus groups) were excluded. Figure 1, a PRISMA flow
256 diagram [53], shows the number of: articles obtained from the searches and screened; papers
257 assessed for eligibility; and articles included in the review (Additional File 2 provides the full list of
258 articles included in the review and their unique identifying code – e.g. BMJ01, BJHP02, SHI03). One
259 review author (KV) assessed the eligibility of all papers identified from the searches. When in doubt,
260 discussions about retaining or excluding articles were held between KV and JB in regular meetings,
261 and decisions were jointly made.

262 - Insert *Figure 1* here -

263 **Data extraction and analysis**

264 A data extraction form was developed (see Additional File 3) recording three areas of information:
265 (a) information about the article (e.g. authors, title, journal, year of publication etc.); (b) information

266 about the aims of the study, the sample size and any justification for this, the participant
267 characteristics, the sampling technique and any sample-related observations or comments made by
268 the authors; and (c) information about the method or technique(s) of data analysis, the number of
269 researchers involved in the analysis, the potential use of software, and any discussion around
270 epistemological considerations. The Abstract, Methods and Discussion (and/or Conclusion) sections
271 of each article were examined by one author (KV) who extracted all the relevant information. This
272 was directly copied from the articles and, when appropriate, comments, notes and initial thoughts
273 were written down.

274 To examine the kinds of sample size justifications provided by articles, an inductive content analysis
275 [54] was initially conducted. On the basis of this analysis, the categories that expressed qualitatively
276 different sample size justifications were developed.

277 We also extracted or coded quantitative data regarding the following aspects:

- 278 - Journal and year of publication
- 279 - Number of interviews
- 280 - Number of participants
- 281 - Presence of sample size justification(s) (Yes/No)
- 282 - Presence of a particular sample size justification category (Yes/No), and
- 283 - Number of sample size justifications provided

284 Descriptive and inferential statistical analyses were used to explore these data.

285 A thematic analysis [55] was then performed on all scientific narratives that discussed or
286 commented on the sample size of the study. These narratives were evident both in papers that
287 justified their sample size and those that did not. To identify these narratives, in addition to the
288 methods sections, the discussion sections of the reviewed articles were also examined and relevant
289 data were extracted and analysed.

290 **Results**

291 In total, 214 articles – 21 in the BMJ, 53 in the BJHP and 140 in the SHI – were eligible for inclusion in
 292 the review. Table 1 provides basic information about the sample sizes – measured in number of
 293 interviews – of the studies reviewed across the three journals. Figure 2 depicts the number of
 294 eligible articles published each year per journal.

295 Table 1

296 *Descriptive statistics of the sample sizes of eligible articles across the three journals*

Sample size of studies	BMJ (n = 21)	BJHP (n = 53)	SHI (n = 140)
<i>Mean (SD) number of interviews</i>	44.5 (29.3)	18.1 (10.4)	37.4 (28)
<i>Min number of interviews</i>	19	6	7
<i>Max number of interviews</i>	128	55	197
<i>Median</i>	31	15	30.5

297

298 - Insert *Figure 2* here -

299 Pairwise comparisons following a significant Kruskal-Wallis¹ test indicated that the studies published
 300 in the BJHP had significantly ($p < .001$) smaller samples sizes than those published either in the BMJ
 301 or the SHI. Sample sizes of BMJ and SHI articles did not differ significantly from each other.

302 **Sample size justifications: results from the quantitative and qualitative content analysis**

303 Ten (47.6%) of the 21 BMJ studies, 26 (49.1%) of the 53 BJHP papers and 24 (17.1%) of the 140 SHI
 304 articles provided some sort of sample size justification. As shown in Table 2, the majority of articles
 305 which justified their sample size provided one justification (70% of articles); fourteen studies (25%)
 306 provided two distinct justifications; one study (1.7%) gave three justifications and two studies (3.3%)
 307 expressed four distinct justifications.

308 Table 2

309 *Number and percentage of ‘justifying’ articles and number of justifications stated by ‘justifying’*
 310 *articles*

How many justifications were provided by the ‘justifying’ articles?	BMJ	BJHP	SHI	Total
One justification	6	17	19	42 (70%)
Two justifications	2	8	5	15 (25%)
Three justifications	1	0	0	1 (1.7%)
Four justifications	1	1	0	2 (3.3%)
Total N of ‘justifying’ articles	10	26	24	60
(out of eligible articles)	(21)	(53)	(140)	(214)
% of ‘justifying’ articles	47.6	49.1	17.1	28

311

312 There was no association between the number of interviews (i.e. sample size) conducted and the
 313 provision of a justification ($r_{pb} = .054$, $p = .433$). Within journals, Mann-Whitney tests indicated that
 314 sample sizes of ‘justifying’ and ‘non-justifying’ articles in the BMJ and SHI did not differ significantly
 315 from each other. In the BJHP, ‘justifying’ articles (*Mean rank* = 31.3) had significantly larger sample
 316 sizes than ‘non-justifying’ studies (*Mean rank* = 22.7; $U = 237.000$, $p < .05$).

317 There was a significant association between the journal a paper was published in and the provision
 318 of a justification ($\chi^2 (2) = 23.83$, $p < .001$). BJHP studies provided a sample size justification
 319 significantly more often than would be expected ($z = 2.9$); SHI studies significantly less often ($z = -$
 320 2.4). If an article was published in the BJHP, the odds of providing a justification were 4.8 times
 321 higher than if published in the SHI. Similarly if published in the BMJ, the odds of a study justifying its
 322 sample size were 4.5 times higher than in the SHI.

323 The qualitative content analysis of the scientific narratives identified eleven different sample size
 324 justifications. These are described below and illustrated with excerpts from relevant articles. By way

325 of a summary, the frequency with which these were deployed across the three journals is indicated
326 in Table 3.

327 Table 3

328 *Commonality, type and counts of sample size justifications across journals*

Commonality of justifications across journals	Qualitatively different justifications	BMJ	BJHP	SHI	Total
Justifications shared by all 3 journals	1. Saturation	7	20	19	46
	2. Pragmatic considerations	1	4	3	8
	3. Qualities of the analysis	1	6	0	7
Justifications shared by 2 journals	4. Meet sampling requirements	2	0	4	6
	5. Sample size guidelines	0	5	1	6
	6. In line with existing research	2	1	0	3
	7. Richness and volume of data	1	0	1	2
Justifications found in 1 journal only	8. Meet research design requirements	2	0	0	2
	9. Researchers' previous experience	1	0	0	1
	10. Nature of study	0	1	0	1
	11. Further sampling to check findings consistency	0	0	1	1
	<i>Total</i>	<i>17</i>	<i>37</i>	<i>29</i>	<i>83</i>

329

330 **Saturation**

331 Saturation was the most commonly invoked principle (55.4% of all justifications) deployed by studies
332 across all three journals to justify the sufficiency of their sample size. In the BMJ, two studies
333 claimed that they achieved *data saturation* (BMJ17; BMJ18) and one article referred descriptively to

334 achieving saturation without explicitly using the term (BMJ13). Interestingly, BMJ13 included data in
335 the analysis beyond the point of saturation in search of ‘unusual/deviant observations’ and with a
336 view to establishing findings consistency.

337 *Thirty three women were approached to take part in the interview study. Twenty seven*
338 *agreed and 21 (aged 21-64, median 40) were interviewed before data saturation was*
339 *reached (one tape failure meant that 20 interviews were available for analysis). (BMJ17)*

340 *No new topics were identified following analysis of approximately two thirds of the*
341 *interviews; however, all interviews were coded in order to develop a better understanding of*
342 *how characteristic the views and reported behaviours were, and also to collect further*
343 *examples of unusual/deviant observations. (BMJ13)*

344 Two articles reported pre-determining their sample size with a view to achieving data saturation
345 (BMJ08 – see extract in section *In line with existing research*; BMJ15 – see extract in section
346 *Pragmatic considerations*) without further specifying if this was achieved. One paper claimed
347 *theoretical saturation* (BMJ06) conceived as being when “no further recurring themes emerging
348 from the analysis” whilst another study argued that although the analytic categories were highly
349 saturated, it was not possible to determine whether theoretical saturation had been achieved
350 (BMJ04). One article (BMJ18) cited a reference to support its position on saturation.

351 In the BJHP, six articles claimed that they achieved *data saturation* (BJHP21; BJHP32; BJHP39;
352 BJHP48; BJHP49; BJHP52) and one article stated that, given their sample size and the guidelines for
353 achieving data saturation, it anticipated that saturation would be attained (BJHP50).

354 *Recruitment continued until data saturation was reached, defined as the point at which no*
355 *new themes emerged. (BJHP48)*

356 *It has previously been recommended that qualitative studies require a minimum sample size*
357 *of at least 12 to reach data saturation (Clarke & Braun, 2013; Fugard & Potts, 2014; Guest,*

358 *Bunce, & Johnson, 2006) Therefore, a sample of 13 was deemed sufficient for the qualitative*
359 *analysis and scale of this study. (BJHP50)*

360 Two studies argued that they achieved *thematic saturation* (BJHP28 – see extract in section *Sample*
361 *size guidelines*; BJHP31) and one (BJHP30) article, explicitly concerned with theory development and
362 deploying theoretical sampling, claimed both theoretical and data saturation.

363 *The final sample size was determined by thematic saturation, the point at which new data*
364 *appears to no longer contribute to the findings due to repetition of themes and comments by*
365 *participants (Morse, 1995). At this point, data generation was terminated. (BJHP31)*

366 Five studies argued that they achieved (BJHP05; BJHP33; BJHP40; BJHP13 – see extract in section
367 *Pragmatic considerations*) or anticipated (BJHP46) saturation without any further specification of the
368 term. BJHP17 referred descriptively to a state of achieved saturation without specifically using the
369 term. *Saturation of coding*, but not saturation of themes, was claimed to have been reached by one
370 article (BJHP18). Two articles explicitly stated that they did not achieve saturation; instead claiming a
371 level of *theme completeness* (BJHP27) or that themes being replicated (BJHP53) were arguments for
372 sufficiency of their sample size.

373 *Furthermore, data collection ceased on pragmatic grounds rather than at the point when*
374 *saturation point was reached. Despite this, although nuances within sub-themes were still*
375 *emerging towards the end of data analysis, the themes themselves were being replicated*
376 *indicating a level of completeness. (BJHP27)*

377 Finally, one article criticised and explicitly renounced the notion of data saturation claiming that, on
378 the contrary, the criterion of *theoretical sufficiency* determined its sample size (BJHP16).

379 *According to the original Grounded Theory texts, data collection should continue until there*
380 *are no new discoveries (i.e., 'data saturation'; Glaser & Strauss, 1967). However, recent*
381 *revisions of this process have discussed how it is rare that data collection is an exhaustive*

382 *process and researchers should rely on how well their data are able to create a sufficient*
383 *theoretical account or ‘theoretical sufficiency’ (Dey, 1999). For this study, it was decided that*
384 *theoretical sufficiency would guide recruitment, rather than looking for data saturation.*

385 (BJHP16)

386 Ten out of the 20 BJHP articles that employed the argument of saturation used one or more citations
387 relating to this principle.

388 In the SHI, one article (SHI01) claimed that it achieved *category saturation* based on authors’
389 judgment.

390 *This number was not fixed in advance, but was guided by the sampling strategy and the*
391 *judgement, based on the analysis of the data, of the point at which ‘category saturation’ was*
392 *achieved. (SHI01)*

393 Three articles described a state of achieved saturation without using the term or specifying what
394 sort of saturation they had achieved (i.e. data, theoretical, thematic saturation) (SHI04; SHI13;
395 SHI30) whilst another four articles explicitly stated that they achieved saturation (SHI100; SHI125;
396 SHI136; SHI137). Two papers stated that they achieved *data saturation* (SHI73 – see extract in
397 section *Sample size guidelines*; SHI113), two claimed *theoretical saturation* (SHI78; SHI115) and two
398 referred to achieving *thematic saturation* (SHI87; SHI139) or to *saturated themes* (SHI29; SHI50).

399 *Recruitment and analysis ceased once theoretical saturation was reached in the categories*
400 *described below (Lincoln and Guba 1985). (SHI115)*

401 *The respondents’ quotes drawn on below were chosen as representative, and illustrate*
402 *saturated themes. (SHI50)*

403 One article stated that thematic saturation was anticipated with its sample size (SHI94). Briefly
404 referring to the difficulty in pinpointing achievement of theoretical saturation, SHI32 (see extract in
405 section *Richness and volume of data*) defended the sufficiency of its sample size on the basis of “the

406 high degree of consensus [that] had begun to emerge among those interviewed”, suggesting that
407 information from interviews was being replicated. Finally, SHI112 (see extract in section *Further*
408 *sampling to check findings consistency*) argued that it achieved *saturation of discursive patterns*.
409 Seven of the 19 SHI articles cited references to support their position on saturation (see Additional
410 File 4 for the full list of citations used by articles to support their position on saturation across the
411 three journals).

412 Overall, it is clear that the concept of saturation encompassed a wide range of variants expressed in
413 terms such as saturation, data saturation, thematic saturation, theoretical saturation, category
414 saturation, saturation of coding, saturation of discursive themes, theme completeness. It is
415 noteworthy, however, that although these various claims were sometimes supported with reference
416 to the literature, they were not evidenced in relation to the study at hand.

417 ***Pragmatic considerations***

418 The determination of sample size on the basis of pragmatic considerations was the second most
419 frequently invoked argument (9.6% of all justifications) appearing in all three journals. In the BMJ,
420 one article (BMJ15) appealed to pragmatic reasons, relating to time constraints and the difficulty to
421 access certain study populations, to justify the determination of its sample size.

422 *On the basis of the researchers’ previous experience and the literature,[30, 31] we estimated*
423 *that recruitment of 15-20 patients at each site would achieve data saturation when data*
424 *from each site were analysed separately. We set a target of seven to 10 caregivers per site*
425 *because of time constraints and the anticipated difficulty of accessing caregivers at some*
426 *home based care services. This gave a target sample of 75-100 patients and 35-50 caregivers*
427 *overall. (BMJ15)*

428 In the BJHP, four articles mentioned pragmatic considerations relating to time or financial
429 constraints (BJHP27 – see extract in section *Saturation*; BJHP53), the participant response rate

430 (BJHP13), and the fixed (and thus limited) size of the participant pool from which interviewees were
431 sampled (BJHP18).

432 *We had aimed to continue interviewing until we had reached saturation, a point whereby*
433 *further data collection would yield no further themes. In practice, the number of individuals*
434 *volunteering to participate dictated when recruitment into the study ceased (15 young*
435 *people, 15 parents). Nonetheless, by the last few interviews, significant repetition of*
436 *concepts was occurring, suggesting ample sampling. (BJHP13)*

437 Finally, three SHI articles explained their sample size with reference to practical aspects: time
438 constraints and project manageability (SHI56), limited availability of respondents and project
439 resources (SHI131), and time constraints (SHI113).

440 *The size of the sample was largely determined by the availability of respondents and*
441 *resources to complete the study. Its composition reflected, as far as practicable, our interest*
442 *in how contextual factors (for example, gender relations and ethnicity) mediated the illness*
443 *experience. (SHI131)*

444 **Qualities of the analysis**

445 This sample size justification (8.4% of all justifications) was mainly employed by BJHP articles and
446 referred to an intensive, idiographic and/or latently focused analysis, i.e. that moved beyond
447 description. More specifically, six articles defended their sample size on the basis of an intensive
448 analysis of transcripts and/or the idiographic focus of the study/analysis. Four of these papers
449 (BJHP02; BJHP19; BJHP24; BJHP47) adopted an Interpretative Phenomenological Analysis (IPA)
450 approach.

451 *The current study employed a sample of 10 in keeping with the aim of exploring each*
452 *participant's account (Smith et al., 1999). (BJHP19)*

453 BJHP47 explicitly renounced the notion of saturation within an IPA approach. The other two BJHP
454 articles conducted thematic analysis (BJHP34; BJHP38). The level of analysis – i.e. latent as opposed
455 to a more superficial descriptive analysis – was also invoked as a justification by BJHP38 alongside
456 the argument of an intensive analysis of individual transcripts.

457 *The resulting sample size was at the lower end of the range of sample sizes employed in*
458 *thematic analysis (Braun & Clarke, 2013). This was in order to enable significant reflection,*
459 *dialogue, and time on each transcript and was in line with the more latent level of analysis*
460 *employed, to identify underlying ideas, rather than a more superficial descriptive analysis*
461 *(Braun & Clarke, 2006). (BJHP38)*

462 Finally, one BMJ paper (BMJ21) defended its sample size with reference to the complexity of the
463 analytic task.

464 *We stopped recruitment when we reached 30-35 interviews, owing to the depth and*
465 *duration of interviews, richness of data, and complexity of the analytical task. (BMJ21)*

466 ***Meet sampling requirements***

467 Meeting sampling requirements (7.2% of all justifications) was another argument employed by two
468 BMJ and four SHI articles to explain their sample size. Achieving maximum variation sampling in
469 terms of specific interviewee characteristics determined and explained the sample size of two BMJ
470 studies (BMJ02; BMJ16 – see extract in section *Meet research design requirements*).

471 *Recruitment continued until sampling frame requirements were met for diversity in age, sex,*
472 *ethnicity, frequency of attendance, and health status. (BMJ02)*

473 Regarding the SHI articles, two papers explained their numbers on the basis of their sampling
474 strategy (SHI01- see extract in section *Saturation*; SHI23) whilst sampling requirements that would
475 help attain sample heterogeneity in terms of a particular characteristic of interest was cited by one
476 paper (SHI127).

477 *The combination of matching the recruitment sites for the quantitative research and the*
478 *additional purposive criteria led to 104 phase 2 interviews (Internet (OLC): 21; Internet (FTF):*
479 *20); Gyms (FTF): 23; HIV testing (FTF): 20; HIV treatment (FTF): 20.) (SHI23)*

480 *Of the fifty interviews conducted, thirty were translated from Spanish into English. These*
481 *thirty, from which we draw our findings, were chosen for translation based on heterogeneity*
482 *in depressive symptomology and educational attainment. (SHI127)*

483 Finally, the pre-determination of sample size on the basis of sampling requirements was stated by
484 one article though this was not used to justify the number of interviews (SHI10).

485 **Sample size guidelines**

486 Five BJHP articles (BJHP28; BJHP38 – see extract in section *Qualities of the analysis*; BJHP46; BJHP47;
487 BJHP50 – see extract in section *Saturation*) and one SHI paper (SHI73) relied on citing existing sample
488 size guidelines or norms within research traditions to determine and subsequently defend their
489 sample size (7.2% of all justifications).

490 *Sample size guidelines suggested a range between 20 and 30 interviews to be adequate*
491 *(Creswell, 1998). Interviewer and note taker agreed that thematic saturation, the point at*
492 *which no new concepts emerge from subsequent interviews (Patton, 2002), was achieved*
493 *following completion of 20 interviews. (BJHP28)*

494 *Interviewing continued until we deemed data saturation to have been reached (the point at*
495 *which no new themes were emerging). Researchers have proposed 30 as an approximate or*
496 *working number of interviews at which one could expect to be reaching theoretical*
497 *saturation when using a semi-structured interview approach (Morse 2000), although this can*
498 *vary depending on the heterogeneity of respondents interviewed and complexity of the*
499 *issues explored. (SHI73)*

500 **In line with existing research**

501 Sample sizes of published literature in the area of the subject matter under investigation (3.5% of all
502 justifications) were used by 2 BMJ articles as guidance and a precedent for determining and
503 defending their own sample size (BMJ08; BMJ15 – see extract in section *Pragmatic considerations*).

504 *We drew participants from a list of prisoners who were scheduled for release each week,*
505 *sampling them until we reached the target of 35 cases, with a view to achieving data*
506 *saturation within the scope of the study and sufficient follow-up interviews and in line with*
507 *recent studies [8-10]. (BMJ08)*

508 Similarly, BJHP38 (see extract in section *Qualities of the analysis*) claimed that its sample size was
509 within the range of sample sizes of published studies that use its analytic approach.

510 ***Richness and volume of data***

511 BMJ21 (see extract in section *Qualities of the analysis*) and SHI32 referred to the richness, detailed
512 nature, and volume of data collected (2.3% of all justifications) to justify the sufficiency of their
513 sample size.

514 *Although there were more potential interviewees from those contacted by postcode*
515 *selection, it was decided to stop recruitment after the 10th interview and focus on analysis of*
516 *this sample. The material collected was considerable and, given the focused nature of the*
517 *study, extremely detailed. Moreover, a high degree of consensus had begun to emerge*
518 *among those interviewed, and while it is always difficult to judge at what point ‘theoretical*
519 *saturation’ has been reached, or how many interviews would be required to uncover*
520 *exception(s), it was felt the number was sufficient to satisfy the aims of this small in-depth*
521 *investigation (Strauss and Corbin 1990). (SHI32)*

522 ***Meet research design requirements***

523 Determination of sample size so that it is in line with, and serves the requirements of, the research
524 design (2.3% of all justifications) that the study adopted was another justification used by 2 BMJ
525 papers (BMJ16; BMJ08 – see extract in section *In line with existing research*).

526 *We aimed for diverse, maximum variation samples [20] totalling 80 respondents from*
527 *different social backgrounds and ethnic groups and those bereaved due to different types of*
528 *suicide and traumatic death. We could have interviewed a smaller sample at different points*
529 *in time (a qualitative longitudinal study) but chose instead to seek a broad range of*
530 *experiences by interviewing those bereaved many years ago and others bereaved more*
531 *recently; those bereaved in different circumstances and with different relations to the*
532 *deceased; and people who lived in different parts of the UK; with different support systems*
533 *and coroners' procedures (see tables 1 and 2 for more details). (BMJ16)*

534 **Researchers' previous experience**

535 The researchers' previous experience (possibly referring to experience with qualitative research) was
536 invoked by BMJ15 (see extract in section *Pragmatic considerations*) as a justification for the
537 determination of sample size.

538 **Nature of study**

539 One BJHP paper argued that the sample size was appropriate for the exploratory nature of the study
540 (BJHP38).

541 *A sample of eight participants was deemed appropriate because of the exploratory nature of*
542 *this research and the focus on identifying underlying ideas about the topic. (BJHP38)*

543 **Further sampling to check findings consistency**

544 Finally, SHI112 argued that once it had achieved saturation of discursive patterns, further sampling
545 was decided and conducted to check for consistency of the findings.

546 *Within each of the age-stratified groups, interviews were randomly sampled until saturation*
547 *of discursive patterns was achieved. This resulted in a sample of 67 interviews. Once this*
548 *sample had been analysed, one further interview from each age-stratified group was*
549 *randomly chosen to check for consistency of the findings. Using this approach it was possible*
550 *to more carefully explore children’s discourse about the ‘I’, agency, relationality and power in*
551 *the thematic areas, revealing the subtle discursive variations described in this article.*
552 (SHI112)

553 **Thematic analysis of passages discussing sample size**

554 This analysis resulted in two overarching thematic areas; the first concerned the variation in the
555 characterisation of sample size sufficiency, and the second related to the perceived threats deriving
556 from sample size insufficiency.

557 ***Characterisations of sample size sufficiency***

558 The analysis showed that there were three main characterisations of the sample size in the articles
559 that provided relevant comments and discussion: (a) the vast majority of these qualitative studies (n
560 = 42) considered their sample size as ‘small’ and this was seen and discussed as a limitation; only two
561 articles viewed their small sample size as desirable and appropriate (b) a minority of articles ($n = 4$)
562 proclaimed that their achieved sample size was ‘sufficient’; and (c) finally, a small group of studies (n
563 = 5) characterised their sample size as ‘large’. Whilst achieving a ‘large’ sample size was sometimes
564 viewed positively because it led to richer results, there were also occasions when a large sample size
565 was problematic rather than desirable.

566 ***‘Small’ but why and for whom?*** A number of articles which characterised their sample size as ‘small’
567 did so against an implicit or explicit quantitative framework of reference. Interestingly, three studies
568 that claimed to have achieved data saturation or ‘theoretical sufficiency’ with their sample size,
569 discussed or noted as a limitation in their discussion their ‘small’ sample size, raising the question of

570 why, or for whom, the sample size was considered small given that the qualitative criterion of
571 saturation had been satisfied.

572 *The current study has a number of limitations. The sample size was small (n = 11) and,*
573 *however, large enough for no new themes to emerge. (BJHP39)*

574 *The study has two principal limitations. The first of these relates to the small number of*
575 *respondents who took part in the study. (SHI73)*

576 Other articles appeared to accept and acknowledge that their sample was flawed because of its
577 small size (as well as other compositional ‘deficits’ e.g. non-representativeness, biases, self-
578 selection) or anticipated that they might be criticized for their small sample size. It seemed that the
579 imagined audience – perhaps reviewer or reader – was one inclined to hold the tenets of
580 quantitative research, and certainly one to whom it was important to indicate the recognition that
581 small samples were likely to be problematic. That one’s sample might be thought small was often
582 construed as a limitation couched in a discourse of regret or apology.

583 Very occasionally, the articulation of the small size as a limitation was explicitly aligned against an
584 espoused positivist framework and quantitative research.

585 *This study has some limitations. Firstly, the 100 incidents sample represents a small number*
586 *of the total number of serious incidents that occurs every year.²⁶ We sent out a nationwide*
587 *invitation and do not know why more people did not volunteer for the study. Our lack of*
588 *epidemiological knowledge about healthcare incidents, however, means that determining an*
589 *appropriate sample size continues to be difficult. (BMJ20)*

590 Indicative of an apparent oscillation of qualitative researchers between the different requirements
591 and protocols demarcating the quantitative and qualitative worlds, there were a few instances of
592 articles which briefly recognised their ‘small’ sample size as a limitation, but then defended their

593 study on more qualitative grounds, such as their ability and success at capturing the complexity of
594 experience and delving into the idiographic, and at generating particularly rich data.

595 *This research, while limited in size, has sought to capture some of the complexity attached to*
596 *men's attitudes and experiences concerning incomes and material circumstances. (SHI35)*

597 *Our numbers are small because negotiating access to social networks was slow and labour*
598 *intensive, but our methods generated exceptionally rich data. (BMJ21)*

599 *This study could be criticised for using a small and unrepresentative sample. Given that older*
600 *adults have been ignored in the research concerning suntanning, fair-skinned older adults are*
601 *the most likely to experience skin cancer, and women privilege appearance over health when*
602 *it comes to sunbathing practices, our study offers depth and richness of data in a*
603 *demographic group much in need of research attention. (SHI57)*

604 **'Good enough' sample sizes:** Only four articles expressed some degree of confidence that their
605 achieved sample size was sufficient. For example, SHI139, in line with the justification of thematic
606 saturation that it offered, expressed trust in its sample size sufficiency despite the poor response
607 rate. Similarly, BJHP04, which did not provide a sample size justification, argued that it targeted a
608 larger sample size in order to eventually recruit a sufficient number of interviewees, due to
609 anticipated low response rate.

610 *Twenty-three people with type I diabetes from the target population of 133 (i.e. 17.3%)*
611 *consented to participate but four did not then respond to further contacts (total N = 19). The*
612 *relatively low response rate was anticipated, due to the busy life-styles of young people in*
613 *the age range, the geographical constraints, and the time required to participate in a semi-*
614 *structured interview, so a larger target sample allowed a sufficient number of participants to*
615 *be recruited. (BJHP04)*

616 Two other articles (BJHP35; SHI32) linked the claimed sufficiency to the scope (i.e. ‘small, in-depth
617 investigation’), aims and nature (i.e. ‘exploratory’) of their studies, thus anchoring their numbers to
618 the particular context of their research. Nevertheless, claims of sample size sufficiency were
619 sometimes undermined when they were juxtaposed with an acknowledgement that a larger sample
620 size would be more scientifically productive.

621 *Although our sample size was sufficient for this exploratory study, a more diverse sample*
622 *including participants with lower socioeconomic status and more ethnic variation would be*
623 *informative. A larger sample could also ensure inclusion of a more representative range of*
624 *apps operating on a wider range of platforms. (BJHP35)*

625 **‘Large’ sample sizes - Promise or peril?** Three articles (BMJ13; BJHP05; BJHP48) which all provided
626 the justification of saturation, characterised their sample size as ‘large’ and narrated this
627 oversufficiency in positive terms as it allowed richer data and findings and enhanced the potential
628 for generalisation. The type of generalisation aspired to (BJHP48) was not further specified however.

629 *This study used rich data provided by a relatively large sample of expert informants on an*
630 *important but under-researched topic. (BMJ13)*

631 *Qualitative research provides a unique opportunity to understand a clinical problem from the*
632 *patient’s perspective. This study had a large diverse sample, recruited through a range of*
633 *locations and used in-depth interviews which enhance the richness and generalizability of the*
634 *results. (BJHP48)*

635 And whilst a ‘large’ sample size was endorsed and valued by some qualitative researchers, within the
636 psychological tradition of IPA, a ‘large’ sample size was counter-normative and therefore needed to
637 be justified. Four BJHP studies, all adopting IPA, expressed the appropriateness or desirability of
638 ‘small’ sample sizes (BJHP41; BJHP45) or hastened to explain why they included a larger than typical
639 sample size (BJHP32; BJHP47). For example, BJHP32 below provides a rationale for how an IPA study

640 can accommodate a large sample size and how this was indeed suitable for the purposes of the
641 particular research. To strengthen the explanation for choosing a non-normative sample size,
642 previous IPA research citing a similar sample size approach is used as a precedent.

643 *Small scale IPA studies allow in-depth analysis which would not be possible with larger*
644 *samples (Smith et al., 2009). (BJHP41)*

645 *Although IPA generally involves intense scrutiny of a small number of transcripts, it was*
646 *decided to recruit a larger diverse sample as this is the first qualitative study of this*
647 *population in the United Kingdom (as far as we know) and we wanted to gain an overview.*
648 *Indeed, Smith, Flowers, and Larkin (2009) agree that IPA is suitable for larger groups.*
649 *However, the emphasis changes from an in-depth individualistic analysis to one in which*
650 *common themes from shared experiences of a group of people can be elicited and used to*
651 *understand the network of relationships between themes that emerge from the interviews.*
652 *This large-scale format of IPA has been used by other researchers in the field of false-positive*
653 *research. Baillie, Smith, Hewison, and Mason (2000) conducted an IPA study, with 24*
654 *participants, of ultrasound screening for chromosomal abnormality; they found that this*
655 *larger number of participants enabled them to produce a more refined and cohesive account.*
656 *(BJHP32)*

657 The IPA articles found in the BJHP were the only instances where a 'small' sample size was
658 advocated and a 'large' sample size problematized and defended. These IPA studies illustrate that
659 the characterisation of sample size sufficiency can be a function of researchers' theoretical and
660 epistemological commitments rather than the result of an 'objective' sample size assessment.

661 ***Threats from sample size insufficiency***

662 As shown above, the majority of articles that commented on their sample size, simultaneously
663 characterized it as small and problematic. On those occasions that authors did not simply cite their

664 'small' sample size as a study limitation but rather continued and provided an account of how and
665 why a small sample size was problematic, two important scientific qualities of the research seemed
666 to be threatened: the generalizability and validity of results.

667 **Generalizability:** Those who characterised their sample as 'small' connected this to the limited
668 potential for generalization of the results. Other features related to the sample – often some kind of
669 compositional particularity – were also linked to limited potential for generalisation. Though not
670 always explicitly articulated to what form of generalisation the articles referred to (see BJHP09),
671 generalisation was mostly conceived in nomothetic terms, that is, it concerned the potential to draw
672 inferences from the sample to the broader study population ('representational generalisation' – see
673 BJHP31) and less often to other populations or cultures.

674 *It must be noted that samples are small and whilst in both groups the majority of those*
675 *women eligible participated, generalizability cannot be assumed. (BJHP09)*

676 *The study's limitations should be acknowledged: Data are presented from interviews with a*
677 *relatively small group of participants, and thus, the views are not necessarily generalizable to*
678 *all patients and clinicians. In particular, patients were only recruited from secondary care*
679 *services where COFP diagnoses are typically confirmed. The sample therefore is unlikely to*
680 *represent the full spectrum of patients, particularly those who are not referred to, or who*
681 *have been discharged from dental services. (BJHP31)*

682 Without explicitly using the term generalisation, two SHI articles noted how their 'small' sample size
683 imposed limits on 'the extent that we can extrapolate from these participants' accounts' (SHI114) or
684 to the possibility 'to draw far-reaching conclusions from the results' (SHI124).

685 Interestingly, only a minority of articles alluded to, or invoked, a type of generalisation that is aligned
686 with qualitative research, that is, idiographic generalisation (i.e. generalisation that can be made
687 *from and about cases* [5]). These articles, all published in the discipline of sociology, defended their

688 findings in terms of the possibility of drawing logical and conceptual inferences to other contexts
689 and of generating understanding that has the potential to advance knowledge, despite their ‘small’
690 size. One article (SHI139) clearly contrasted nomothetic (statistical) generalisation to idiographic
691 generalisation, arguing that the lack of statistical generalizability does not nullify the ability of
692 qualitative research to still be relevant beyond the sample studied.

693 *Further, these data do not need to be statistically generalisable for us to draw inferences*
694 *that may advance medicalisation analyses (Charmaz 2014). These data may be seen as an*
695 *opportunity to generate further hypotheses and are a unique application of the*
696 *medicalisation framework. (SHI139)*

697 *Although a small-scale qualitative study related to school counselling, this analysis can be*
698 *usefully regarded as a case study of the successful utilisation of mental health-related*
699 *resources by adolescents. As many of the issues explored are of relevance to mental health*
700 *stigma more generally, it may also provide insights into adult engagement in services. It*
701 *shows how a sociological analysis, which uses positioning theory to examine how people*
702 *negotiate, partially accept and simultaneously resist stigmatisation in relation to mental*
703 *health concerns, can contribute to an elucidation of the social processes and narrative*
704 *constructions which may maintain as well as bridge the mental health service gap. (SHI103)*

705 Only one article (SHI30) used the term *transferability* to argue for the potential of wider relevance of
706 the results which was thought to be more the product of the composition of the sample (i.e. diverse
707 sample), rather than the sample size.

708 **Validity:** The second major concern that arose from a ‘small’ sample size pertained to the internal
709 validity of findings (i.e. here the term is used to denote the ‘truth’ or credibility of research findings).
710 Authors expressed uncertainty about the degree of confidence in particular aspects or patterns of
711 their results, primarily those that concerned some form of differentiation on the basis of relevant
712 participant characteristics.

713 *The information source preferred seemed to vary according to parents' education; however,*
714 *the sample size is too small to draw conclusions about such patterns. (SHI80)*

715 *Although our numbers were too small to demonstrate gender differences with any certainty,*
716 *it does seem that the biomedical and erotic scripts may be more common in the accounts of*
717 *men and the relational script more common in the accounts of women. (SHI81)*

718 In other instances, articles expressed uncertainty about whether their results accounted for the full
719 spectrum and variation of the phenomenon under investigation. In other words, a 'small' sample size
720 (alongside compositional 'deficits' such as a not statistically representative sample) was seen to
721 threaten the 'content validity' of the results which in turn led to constructions of the study
722 conclusions as tentative.

723 *Data collection ceased on pragmatic grounds rather than when no new information*
724 *appeared to be obtained (i.e., saturation point). As such, care should be taken not to*
725 *overstate the findings. Whilst the themes from the initial interviews seemed to be replicated*
726 *in the later interviews, further interviews may have identified additional themes or provided*
727 *more nuanced explanations. (BJHP53)*

728 *...it should be acknowledged that this study was based on a small sample of self-selected*
729 *couples in enduring marriages who were not broadly representative of the population. Thus,*
730 *participants may not be representative of couples that experience postnatal PTSD. It is*
731 *therefore unlikely that all the key themes have been identified and explored. For example,*
732 *couples who were excluded from the study because the male partner declined to participate*
733 *may have been experiencing greater interpersonal difficulties. (BJHP03)*

734 In other instances, articles attempted to preserve a degree of credibility of their results, despite the
735 recognition that the sample size was 'small'. Clarity and sharpness of emerging themes and

736 alignment with previous relevant work were the arguments employed to warrant the validity of the
737 results.

738 *This study focused on British Chinese carers of patients with affective disorders, using a*
739 *qualitative methodology to synthesise the sociocultural representations of illness within this*
740 *community. Despite the small sample size, clear themes emerged from the narratives that*
741 *were sufficient for this exploratory investigation. (SHI98)*

742 **Discussion**

743 The present study sought to examine how qualitative sample sizes in health-related research are
744 characterised and justified. In line with previous studies [22,30,33,34] the findings demonstrate that
745 reporting of sample size sufficiency is limited; just over 50% of articles in the BMJ and BJHP and 82%
746 in the SHI did not provide any sample size justification. Providing a sample size justification was not
747 related to the number of interviews conducted, but it was associated with the journal that the article
748 was published in, indicating the influence of disciplinary or publishing norms, also reported in prior
749 research [30]. This lack of transparency about sample size sufficiency is problematic given that most
750 qualitative researchers would agree that it is an important marker of quality [56,57]. Moreover, and
751 with the rise of qualitative research in social sciences, efforts to synthesise existing evidence and
752 assess its quality are obstructed by poor reporting [58,59].

753 When authors justified their sample size, our findings indicate that sufficiency was mostly appraised
754 with reference to features that were intrinsic to the study, in agreement with general advice on
755 sample size determination [4,11,36]. The principle of saturation was the most commonly invoked
756 argument [22] accounting for 55% of all justifications. A wide range of variants of saturation was
757 evident corroborating the proliferation of the meaning of the term [48] and reflecting different
758 underlying conceptualisations or models of saturation [20]. Nevertheless, claims of saturation were
759 never substantiated in relation to procedures conducted in the study itself, endorsing similar
760 observations in the literature [25,30,47]. Claims of saturation were sometimes supported with

761 citations of other literature, suggesting a removal of the concept away from the characteristics of
762 the study at hand. Pragmatic considerations, such as resource constraints or participant response
763 rate and availability, was the second most frequently used argument accounting for approximately
764 10% of justifications and another 23% of justifications also represented intrinsic-to-the-study
765 characteristics (i.e. qualities of the analysis, meeting sampling or research design requirements,
766 richness and volume of the data obtained, nature of study, further sampling to check findings
767 consistency).

768 Only, 12% of mentions of sample size justification pertained to arguments that were external to the
769 study at hand, in the form of existing sample size guidelines and prior research that sets precedents.
770 Whilst community norms and prior research can establish useful rules of thumb for estimating
771 sample sizes [60] – and reveal what sizes are more likely to be acceptable within research
772 communities – researchers should avoid adopting these norms uncritically, especially when such
773 guidelines [e.g. 30,35], might be based on research that does not provide adequate evidence of
774 sample size sufficiency. Similarly, whilst methodological research that seeks to demonstrate the
775 achievement of saturation is invaluable since it explicates the parameters upon which saturation is
776 contingent and indicates when a research project is likely to require a smaller or a larger sample [e.g.
777 29], specific numbers at which saturation was achieved within these projects cannot be routinely
778 extrapolated for other projects. We concur with existing views [11,36] that the consideration of the
779 characteristics of the study at hand, such as the epistemological and theoretical approach, the
780 nature of the phenomenon under investigation, the aims and scope of the study, the quality and
781 richness of data, or the researcher’s experience and skills of conducting qualitative research, should
782 be the primary guide in determining sample size and assessing its sufficiency.

783 Moreover, although numbers in qualitative research are not unimportant [61], sample size should
784 not be considered alone but be embedded in the more encompassing examination of *data adequacy*
785 [56,57]. Erickson’s [62] dimensions of ‘evidentiary adequacy’ are useful here. He explains the

786 concept in terms of adequate amounts of evidence, adequate variety in kinds of evidence, adequate
787 interpretive status of evidence, adequate disconfirming evidence, and adequate discrepant case
788 analysis. All dimensions might not be relevant across all qualitative research designs, but this
789 illustrates the thickness of the concept of data adequacy, taking it beyond sample size.

790 The present research also demonstrated that sample sizes were commonly seen as 'small' and
791 insufficient and discussed as limitation. Often unjustified (and in two cases incongruent with their
792 own claims of saturation) these findings imply that sample size in qualitative health research is often
793 adversely judged (or expected to be judged) against an implicit, yet omnipresent, quasi-quantitative
794 standpoint. Indeed there were a few instances in our data where authors appeared, possibly in
795 response to reviewers, to resist to some sort of quantification of their results. This implicit reference
796 point became more apparent when authors discussed the threats deriving from an insufficient
797 sample size. Whilst the concerns about internal validity might be legitimate to the extent that
798 qualitative research projects, which are broadly related to realism, are set to examine phenomena in
799 sufficient breadth and depth, the concerns around generalizability revealed a conceptualisation that
800 is not compatible with purposive sampling. The limited potential for generalisation, as a result of a
801 small sample size, was often discussed in nomothetic, statistical terms. Only occasionally was
802 analytic or idiographic generalisation invoked to warrant the value of the study's findings [5,17].

803 ***Strengths and limitations of the present study***

804 We note, first, the limited number of health-related journals reviewed, so that only a 'snapshot' of
805 qualitative health research has been captured. Examining additional disciplines (e.g. nursing
806 sciences) as well as inter-disciplinary journals would add to the findings of this analysis.

807 Nevertheless, our study is the first to provide some comparative insights on the basis of disciplines
808 that are differently attached to the legacy of positivism and analysed literature published over a
809 lengthy period of time (15 years). Guetterman [27] also examined health-related literature but this
810 analysis was restricted to 26 most highly cited articles published over a period of five years whilst

811 Carlsen and Glenton's [22] study concentrated on focus groups health research. Moreover, although
812 it was our intention to examine sample size justification in relation to the epistemological and
813 theoretical positions of articles, this proved to be challenging largely due to absence of relevant
814 information, or the difficulty into discerning clearly articles' positions [63] and classifying them under
815 specific approaches (e.g. studies often combined elements from different theoretical and
816 epistemological traditions). We believe that such an analysis would yield useful insights as it links the
817 methodological issue of sample size to the broader philosophical stance of the research. Despite
818 these limitations, the analysis of the characterisation of sample size and of the threats seen to
819 accrue from insufficient sample size, enriches our understanding of sample size (in)sufficiency
820 argumentation by linking it to other features of the research. As the peer-review process becomes
821 increasingly public, future research could usefully examine how reporting around sample size
822 sufficiency and data adequacy might be influenced by the interactions between authors and
823 reviewers.

824 **Conclusions**

825 The past decade has seen a growing appetite in qualitative research for an evidence-based approach
826 to sample size determination and to evaluations of the sufficiency of sample size. Despite the
827 conceptual and methodological developments in the area, the findings of the present study confirm
828 previous studies in concluding that appraisals of sample size sufficiency are either absent or poorly
829 substantiated. To ensure and maintain high quality research that will encourage greater appreciation
830 of qualitative work in health-related sciences [64], we argue that qualitative researchers should be
831 more transparent and thorough in their evaluation of sample size as part of their appraisal of data
832 adequacy. We would encourage the practice of appraising sample size sufficiency with close
833 reference to the study at hand and would thus caution against responding to the growing
834 methodological research in this area with a decontextualised application of sample size numerical
835 guidelines, norms and principles. Although researchers might find sample size community norms

836 serve as useful rules of thumb, we recommend methodological knowledge is used to critically
837 consider how saturation and other parameters that affect sample size sufficiency pertain to the
838 specifics of the particular project. Those reviewing papers have a vital role in encouraging
839 transparent study-specific reporting. The review process should support authors to exercise nuanced
840 judgments in decisions about sample size determination in the context of the range of factors that
841 influence sample size sufficiency and the specifics of a particular study. In light of the growing
842 methodological evidence in the area, transparent presentation of such evidence-based judgement is
843 crucial and in time should surely obviate the seemingly routine practice of citing the ‘small’ size of
844 qualitative samples among the study limitations.

845

846 **Abbreviations:** BMJ: British Medical Journal; BJHP: British Journal of Health Psychology; SHI:
847 Sociology of Health & Illness; IPA: Interpretative Phenomenological Analysis.

848 **Ethics approval and consent to participate:** Not applicable

849 **Availability of data and materials:** Supporting data can be accessed in the original publications.
850 Additional File 2 lists all eligible studies that were included in the present analysis.

851 **Competing interests:** Terry Young is an academic who undertakes research and occasional
852 consultancy in the areas of health technology assessment, information systems, and service design.
853 He is unaware of any direct conflict of interest with respect to this paper. All other authors have no
854 competing interests to declare.

855 **Consent for publication:** Not applicable

856 **Funding:** This research was initially conceived of and partly conducted with financial support from
857 the *Multidisciplinary Assessment of Technology Centre for Healthcare (MATCH)* programme
858 (EP/F063822/1 and EP/G012393/1). The research continued and was completed independent of any
859 support. The funding body did not have any role in the study design, the collection, analysis and

860 interpretation of the data, in the writing of the paper, and in the decision to submit the manuscript
861 for publication. The views expressed are those of the authors alone.

862 **Authors' contributions:** JB and TY conceived the study; KV, JB, and TY designed the study; KV
863 identified the articles and extracted the data; KV and JB assessed eligibility of articles; KV, JB, ST, and
864 TY contributed to the analysis of the data, discussed the findings and early drafts of the paper; KV
865 developed the final manuscript; KV, JB, ST, and TY read and approved the manuscript.

866 **Acknowledgments:** We would like to thank Dr Paula Smith and Katharine Lee for their comments on
867 a previous draft of this paper as well as Natalie Ann Mitchell and Meron Teferra for assisting us with
868 data extraction.

869

870 **References**

- 871 1. Spencer L, Ritchie J, Lewis J, Dillon L. Quality in qualitative evaluation: a framework for
872 assessing research evidence. National Centre for Social Research. 2003
873 https://www.heacademy.ac.uk/system/files/166_policy_hub_a_quality_framework.pdf
874 Accessed 11 May 2018.
- 875 2. Fusch PI, Ness LR. Are we there yet? Data saturation in qualitative research. Qual Rep.
876 2015;20(9):1408-16.
- 877 3. Robinson OC. Sampling in interview-based qualitative research: a theoretical and practical
878 guide. Qual Res Psychol. 2014;11(1):25-41.
- 879 4. Sandelowski M. Sample size in qualitative research. Res Nurs Health. 1995;18(2): 179-83.
- 880 5. Sandelowski M. One is the liveliest number: the case orientation of qualitative research. Res
881 Nurs Health. 1996;19(6):525-9.
- 882 6. Luborsky MR, Rubinstein RL. Sampling in qualitative research: rationale, issues, and
883 methods. Res Aging. 1995;17(1): 89-113.

- 884 7. Marshall MN. Sampling for qualitative research. *Fam Pract*. 1996;13(6):522-526.
- 885 8. Patton MQ. *Qualitative evaluation and research methods*. 2nd ed. Newbury Park, CA: Sage;
- 886 1990.
- 887 9. van Rijnsvoever FJ. (I Can't Get No) Saturation: a simulation and guidelines for sample sizes in
- 888 qualitative research. *PLoS One*. 2017;12(7):e0181689.
- 889 10. Morse JM. The significance of saturation. *Qual Health Res*. 1995;5(2):147-9.
- 890 11. Morse JM. Determining sample size. *Qual Health Res*. 2000;10(1):3-5.
- 891 12. Gergen KJ, Josselson R, Freeman M. The promises of qualitative inquiry. *Am Psychol*.
- 892 2015;70(1):1-9.
- 893 13. Borsci S, Macredie RD, Barnett J, Martin J, Kuljis J, Young T. Reviewing and extending the
- 894 five-user assumption: A grounded procedure for interaction evaluation. *ACM Trans Comput*
- 895 *Hum Interact*. 2013;20(5):29.
- 896 14. Borsci S, Macredie RD, Martin JL, Young T. How many testers are needed to assure the
- 897 usability of medical devices?. *Expert Rev Med Devices*. 2014;11(5):513-25.
- 898 15. Glaser BG, Strauss, AL. *The discovery of grounded theory: strategies for qualitative research*.
- 899 Chicago, IL: Aldine; 1967.
- 900 16. Kerr C, Nixon A, Wild D. Assessing and demonstrating data saturation in qualitative inquiry
- 901 supporting patient-reported outcomes research. *Expert Rev Pharmacoecon Outcomes Res*.
- 902 2010;10(3):269-81.
- 903 17. Lincoln YS, Guba EG. *Naturalistic inquiry*. London: Sage; 1985.
- 904 18. Malterud K, Siersma VD, Guassora AD. Sample size in qualitative interview studies: guided by
- 905 information power. *Qual Health Res*. 2015;26:1753-60.
- 906 19. Nelson J. Using conceptual depth criteria: addressing the challenge of reaching saturation in
- 907 qualitative research. *Qual Res*. 2017;17(5):554-70.

- 908 20. Saunders B, Sim J, Kingstone T, Baker S, Waterfield J, Bartlam B, et al. Saturation in
909 qualitative research: exploring its conceptualization and operationalization. *Qual Quant*.
910 2017; doi:10.1007/s11135-017-0574-8.
- 911 21. Caine K. Local standards for sample size at CHI. In *Proceedings of the 2016 CHI Conference*
912 *on Human Factors in Computing Systems*. 2016;981-992. ACM.
- 913 22. Carlsen B, Glenton C. What about N? A methodological study of sample-size reporting in
914 focus group studies. *BMC Med Res Methodol*. 2011;11(1):26.
- 915 23. Constantinou CS, Georgiou M, Perdikogianni M. A comparative method for themes
916 saturation (CoMeTS) in qualitative interviews. *Qual Res*. 2017;17(5):571-88.
- 917 24. Dai NT, Free C, Gendron Y. Interview-based research in accounting 2000-2014: a review.
918 November 2016. <https://ssrn.com/abstract=2711022> or
919 <http://dx.doi.org/10.2139/ssrn.2711022>. Accessed 17 May 2018
- 920 25. Francis JJ, Johnston M, Robertson C, Glidewell L, Entwistle V, Eccles MP, et al. What is an
921 adequate sample size? Operationalising data saturation for theory-based interview studies.
922 *Psychol Health*. 2010;25(10):1229-45.
- 923 26. Guest G, Bunce A, Johnson L. How many interviews are enough? An experiment with data
924 saturation and variability. *Field Methods*. 2006;18(1):59-82.
- 925 27. Guetterman TC. Descriptions of sampling practices within five approaches to qualitative
926 research in education and the health sciences. *Forum Qual Soc Res*. 2015;16(2):25.
927 <http://nbn-resolving.de/urn:nbn:de:0114-fqs1502256>. Accessed 17 May 2018
- 928 28. Hagaman AK, Wutich A. How many interviews are enough to identify metathemes in
929 multisited and cross-cultural research? Another perspective on Guest, Bunce, and Johnson's
930 (2006) landmark study. *Field Methods*. 2017;29(1):23-41.
- 931 29. Hennink MM, Kaiser BN, Marconi VC. Code saturation versus meaning saturation: how many
932 interviews are enough?. *Qual Health Res*. 2017;27(4):591-608.

- 933 30. Marshall B, Cardon P, Poddar A, Fontenot R. Does sample size matter in qualitative
934 research?: a review of qualitative interviews in IS research. *J Comput Inform Syst.*
935 2013;54(1):11-22.
- 936 31. Mason M. Sample size and saturation in PhD studies using qualitative interviews. *Forum*
937 *Qual Soc Res.* 2010;11(3):8. <http://nbn-resolving.de/urn:nbn:de:0114-fqs100387>. Accessed
938 17 May 2018
- 939 32. Safman RM, Sobal J. Qualitative sample extensiveness in health education research. *Health*
940 *Educ Behav.* 2004;31(1):9-21.
- 941 33. Saunders MN, Townsend K. Reporting and justifying the number of interview participants in
942 organization and workplace research. *Br J Manag.* 2016;27(4):836-52.
- 943 34. Sobal J. 2001. Sample extensiveness in qualitative nutrition education research. *J Nutr Educ.*
944 2001;33(4):184-92.
- 945 35. Thomson SB. 2010. Sample size and grounded theory. *JOAAG.* 2010;5(1).
946 http://www.joaag.com/uploads/5_1_Research_Note_1_Thomson.pdf. Accessed 17 May
947 2018
- 948 36. Baker SE, Edwards R. How many qualitative interviews is enough?: expert voices and early
949 career reflections on sampling and cases in qualitative research. National Centre for
950 Research Methods Review Paper. 2012.
951 http://eprints.ncrm.ac.uk/2273/4/how_many_interviews.pdf. Accessed 17 May 2018.
- 952 37. Ogden J, Cornwell D. The role of topic, interviewee, and question in predicting rich interview
953 data in the field of health research. *Sociol Health Illn.* 2010;32(7):1059-71.
- 954 38. Green J, Thorogood N. *Qualitative methods for health research.* London: Sage; 2004.
- 955 39. Ritchie J, Lewis J, Elam G. Designing and selecting samples. In: Ritchie J, Lewis J, Editors.
956 *Qualitative research practice: a guide for social science students and researchers.* London:
957 Sage; 2003. p. 77–108.

- 958 40. Britten N. Qualitative research: qualitative interviews in medical research. *BMJ*.
959 1995;311(6999):251-3.
- 960 41. Creswell JW. *Qualitative inquiry and research design: choosing among five approaches*. 2nd
961 ed. London: Sage; 2007.
- 962 42. Fugard AJ, Potts HW. Supporting thinking on sample sizes for thematic analyses: a
963 quantitative tool. *Int J Soc Res Methodol*. 2015;18(6):669-84.
- 964 43. Emmel N. Themes, variables, and the limits to calculating sample size in qualitative research:
965 a response to Fugard and Potts. *Int J Soc Res Methodol*. 2015;18(6):685-6.
- 966 44. Braun V, Clarke V. (Mis) conceptualising themes, thematic analysis, and other problems with
967 Fugard and Potts' (2015) sample-size tool for thematic analysis. *Int J Soc Res Methodol*.
968 2016;19(6):739-43.
- 969 45. Hammersley M. Sampling and thematic analysis: a response to Fugard and Potts. *Int J Soc*
970 *Res Methodol*. 2015;18(6):687-8.
- 971 46. Charmaz K. *Constructing grounded theory: a practical guide through qualitative analysis*.
972 London: Sage; 2006.
- 973 47. Bowen GA. Naturalistic inquiry and the saturation concept: a research note. *Qual Res*.
974 2008;8(1):137-52.
- 975 48. O'Reilly M, Parker N. 'Unsatisfactory Saturation': a critical exploration of the notion of
976 saturated sample sizes in qualitative research. *Qual Res*. 2013;13(2):190-7.
- 977 49. Morse JM. "Data were saturated...". *Qual Health Res*. 2015;25(5):587-8.
- 978 50. Manen M, Higgins I, Riet P. A conversation with Max van Manen on phenomenology in its
979 original sense. *Nurs Health Sci*. 2016;18(1):4-7.
- 980 51. Dey I. *Grounding grounded theory*. San Francisco, CA: Academic Press; 1999.
- 981 52. Hays DG, Wood C, Dahl H, Kirk-Jenkins A. Methodological rigor in *Journal of Counseling &*
982 *Development* qualitative research articles: a 15-year review. *J Couns Dev*. 2016;94(2):172-
983 83.

- 984 53. Moher D, Liberati A, Tetzlaff J, Altman DG, Prisma Group. Preferred reporting items for
985 systematic reviews and meta-analyses: the PRISMA statement. *PLoS Med.* 2009; 6(7):
986 e1000097.
- 987 54. Hsieh HF, Shannon SE. Three approaches to qualitative content analysis. *Qual Health Res.*
988 2005;15(9):1277-88.
- 989 55. Boyatzis RE. Transforming qualitative information: thematic analysis and code development.
990 Thousand Oaks, CA: Sage; 1998.
- 991 56. Levitt HM, Motulsky SL, Wertz FJ, Morrow SL, Ponterotto JG. Recommendations for
992 designing and reviewing qualitative research in psychology: promoting methodological
993 integrity. *Qual Psychol.* 2017;4(1):2-22.
- 994 57. Morrow SL. Quality and trustworthiness in qualitative research in counseling psychology. *J*
995 *Couns Psychol.* 2005;52(2):250-60.
- 996 58. Barroso J, Sandelowski M. Sample reporting in qualitative studies of women with HIV
997 infection. *Field Methods.* 2003;15(4):386-404.
- 998 59. Glenton C, Carlsen B, Lewin S, Munthe-Kaas H, Colvin CJ, Tunçalp Ö, et al. Applying GRADE-
999 CERQual to qualitative evidence synthesis findings—paper 5: how to assess adequacy of
1000 data. *Implement Sci.* 2018;13(Suppl 1):14.
- 1001 60. Onwuegbuzie AJ, Leech NL. A call for qualitative power analyses. *Qual. Quant.*
1002 2007;41(1):105-121.
- 1003 61. Sandelowski M. Real qualitative researchers do not count: the use of numbers in qualitative
1004 research. *Res Nurs Health.* 2001;24(3):230-40.
- 1005 62. Erickson F. Qualitative methods in research on teaching. In: Wittrock M, editor. *Handbook of*
1006 *research on teaching.* 3rd ed. New York: Macmillan; 1986. p. 119–61.
- 1007 63. Bradbury-Jones C, Taylor J, Herber O. How theory is used and articulated in qualitative
1008 research: development of a new typology. *Soc Sci Med.* 2014;120:135-41.

1009 64. Greenhalgh T, Annandale E, Ashcroft R, Barlow J, Black N, Bleakley A, et al. An open letter to
1010 The BMJ editors on qualitative research. *BMJ*. 2016;352:i563.

1011

1012 **Figure Legends**

1013 *Figure 1.* PRISMA flow diagram.

1014 *Figure 2.* Number of eligible articles published each year per journal

1015

1016 **Additional Files**

1017 *Additional File 1.* Editorial positions on qualitative research and sample considerations (where
1018 available)

1019 *Additional File 2.* List of eligible articles included in the review ($N = 214$)

1020 *Additional File 3.* Data Extraction Form

1021 *Additional File 4.* Citations used by articles to support their position on saturation

1022

ⁱ A non-parametric test of difference for independent samples was performed since the variable *number of interviews* violated assumptions of normality according to the standardized scores of skewness and kurtosis (BMJ: z skewness = 3.23, z kurtosis = 1.52; BJHP: z skewness = 4.73, z kurtosis = 4.85; SHI: z skewness = 12.04, z kurtosis = 21.72) and the Shapiro-Wilk test of normality ($p < .001$).