

- 1 BLINDED TITLE PAGE:
- 2 **Do Radiological Research Articles Apply the Term “Pilot Study” Correctly?**
- 3 **Systematic Review.**

4 INTRODUCTION

5 It has been suggested that radiological research exhibits less methodological expertise
6 than other medical disciplines^{1,2}. Methodological expertise is reflected partly by a clear
7 understanding of different study designs and their specific purpose. To statisticians and
8 methodologists, the term “pilot study” has precise implications, namely a “small-scale
9 test of the methods and procedures to be used on a larger scale” subsequently³. Pilot
10 studies therefore aim to evaluate the viability and applicability of critical research
11 components and assumptions needed for larger studies to be successful. They play a
12 pivotal role in mitigating unwanted surprises, saving cost and time when later
13 undertaking large trials. They may prevent studies that are doomed to fail, for example
14 due to inadequate recruitment. For these reasons, the UK Medical Research Council
15 guidance on designing and evaluating complex interventions recommends that pilot
16 studies be conducted before any definitive large-scale evaluation⁴.

17 By virtue of their design, radiological pilot studies are not powered to determine
18 diagnostic accuracy (or a “pilot study” would be unnecessary). Rather, their endpoints
19 should revolve around metrics necessary to determine if a larger study/trial ought to
20 proceed. The most obvious endpoint is recruitment rate, which tends to be over-
21 optimistic when planning research. Test method feasibility in a trial setting and data
22 capture and retention are examples of other common endpoints for pilot studies.
23 However, we have noted anecdotally that the description “pilot study” is used by
24 radiological studies that instead report measures of diagnostic accuracy while often
25 appearing underpowered. If so, this would imply that the researchers do not appreciate
26 the precise implications of this study design and are using the term “pilot” to excuse

27 underpowering. Therefore, by systematic review we aimed to determine the proportion
28 of radiological “pilot” studies that use this description correctly.

29

30 **MATERIAL AND METHODS**

31

32 **Ethical approval**

33 Our institution does not require ethical approval for systematic review of indexed
34 literature.

35

36 **Review question**

37 We aimed to determine what proportion of radiological “pilot” studies genuinely used
38 this study design.

39

40 **Search strategy and study eligibility**

41 Following investigator discussion, we identified studies from four well-recognised,
42 representative, indexed radiological journals: Radiology (RADIOL), European Radiology
43 (ER), American Journal of Roentgenology (AJR), and British Journal of Radiology
44 (BJR). We decided to limit our review to four well-established, internationally visible
45 journals believing that these would provide data representative of pilot study reporting in
46 the radiological domain, and that little additional information would be gained by
47 extending the search to all potential journals. Limiting the search allowed us to study
48 each individual journal in greater depth. To be eligible for inclusion the study title had to
49 include the word “pilot” and investigate human subjects (of any age). We elected to

50 exclude narrative and systematic reviews if encountered. We decided *a priori* that a
51 maximum of 20 studies per journal (for a maximum sample size of 80) would achieve
52 sufficient saturation for information to be representative, i.e. we anticipated that by that
53 stage our findings would be clear and that a larger sample would add no further, useful
54 information.

55

56 **Search strategy and string**

57 The senior author (BLINDED) searched The National Library of Science via PubMed
58 using the following terms: “Radiology [ta] AND pilot”, “European Radiology [ta] AND
59 pilot”, “American Journal of Roentgenology [ta] AND pilot”, “British Journal of Radiology
60 [ta] AND pilot”. Articles were selected from present day (search performed 20th
61 September 2018) retrospectively, without date restriction.

62

63 **Search process and citation management**

64 Potentially eligible citations were passed to a junior researcher (BLINDED) who then
65 screened the electronic abstracts, discarding “clearly unsuitable” articles (e.g. “pilot”
66 absent from the study title, radiotherapy studies, animal studies, interventions without
67 imaging). An exclusion log was kept for all excluded papers (Figure 1). The junior
68 researcher had no prior experience in performing systematic reviews and so was
69 supervised closely by more senior members of the research team. Duplicate extraction
70 of the first 20 articles was performed by another member of the team (BLINDED) to
71 check consistency, which was acceptable. The remaining extraction was performed by
72 the junior researcher thereafter. Any uncertainty was discussed face-to-face with the

73 senior author. Potentially eligible citations were entered into a datasheet designed
74 specifically for the review (Microsoft Excel for Mac 2011 v. 14.5.9, Microsoft
75 Corporation, Washington).

76

77 **Data extraction**

78 The junior researcher extracted data from component primary studies into the datasheet
79 using the following categories: Citation; Prospective/retrospective design; Single/multi-
80 centre; Total subjects recruited; Sample size rationale and/or power calculation stated
81 (yes/no); Primary endpoint stated explicitly (yes [description]/no); Secondary endpoints
82 stated explicitly (yes [description]/no); If no primary endpoint stated explicitly, what
83 appeared to be the primary endpoint; Did the endpoint suggest a pilot study (yes/no);
84 Were diagnostic test accuracy data presented (yes [description, e.g.
85 sensitivity/specificity]/no); Ultimately, in the Reviewer's opinion, was the study a genuine
86 "pilot" (yes/no [description]). We judged studies to be genuine pilots if it was a, "small-
87 scale test of the methods and procedures to be used on a larger scale" subsequently³.
88 For example, the design could potentially investigate recruitment rate for a subsequent
89 trial, intervention allocation (e.g. by randomisation) and acceptability (i.e. willingness to
90 be randomised), assessment procedures, data retention, etc. If believed to be another
91 study design, the nature of this was recorded. Any uncertainty was resolved via face-to-
92 face discussion with other members of the team.

93

94 **Statistical analysis**

95 We reported our review according to PRISMA guidelines⁵. No risk of bias assessment
96 was performed because methodological quality of component studies was not our
97 primary concern and meta-analysis (where risk of bias is used to assess precision of the
98 point estimate) was not anticipated. Data were expressed as simple frequencies and
99 proportions.

100

101 **RESULTS**

102

103 The PRISMA flow diagram for studies identified by the systematic review is shown in
104 Figure 1. The search string identified a total of 658 records over the four targeted
105 journals as follows: 236 RADIOL, 158 ER, 146 AJR, 118 BJR. Our target of 20
106 consecutive articles was achieved for all journals with the exception of BJR, where only
107 18 articles were eligible from the 118 identified, giving a total sample size of 78. The 78
108 component articles selected for the systematic review are listed in Appendix 1.

109

110 55 (70.5%) studies were prospective, 21 (26.9%) retrospective and 2 (2.6%) mixed. The
111 overwhelming majority of studies were conducted in single centres; 76 (94.7%). 62
112 studies (79.5%) investigated a single imaging modality (including 26 MRI; 18 CT; 12
113 ultrasound; 2 PET-CT; 2 PET-MRI), 13 studies (16.7%) investigated multiple imaging
114 modalities, two studies investigated radio-frequency ablation and a single study
115 investigated a biliary stent.

116

117 5,572 patients were reported across 77 studies; the remaining study reported individual
118 lesions only and omitted the number of individual patients⁶. Median sample size per
119 study was 20 patients, range 7⁷ to 1666⁸. The majority of studies (55, 70.5%) stated no
120 rationale for their sample size. Furthermore, no study presented a power calculation to
121 justify sample size. A primary endpoint was stated explicitly by 70 (89.7%) studies and
122 could be inferred in all the remaining 8 studies. Secondary endpoints were stated
123 explicitly by 20 (25.6%) studies.

124
125 Ultimately, we judged that no individual study qualified as a genuine pilot study when
126 assessed against our *a priori* criteria. The individual study types encountered are
127 described in Table 1. Notably, 66 (84.6%) studies presented measures of test accuracy
128 and were framed as studies of diagnostic test accuracy. 12 studies presented elements
129 of feasibility, and 8 elements of technology assessment (Table 1).

130

131 **DISCUSSION**

132

133 This systematic review investigated our anecdotal observation that most radiological
134 studies describing themselves as “pilot” studies are actually small, underpowered
135 studies reporting diagnostic test accuracy. Our hypothesis was correct; not one of the
136 78 studies ultimately included in our systematic review satisfied our *a priori* criteria for a
137 genuine “pilot” study, i.e. a “small-scale test of the methods and procedures to be used
138 on a larger scale” subsequently³. As predicted, we found that the greatest proportion of
139 studies described as pilots actually attempted to investigate diagnostic accuracy,

140 frequently presenting inferential statistics and performing hypothesis testing.
141 Furthermore, while good studies of test accuracy should be powered sufficiently to
142 estimate diagnostic performance (for example sensitivity and specificity, and/or area
143 under the receiver operating characteristic curve), no study presented a power
144 calculation.

145
146 While the majority of study designs we encountered were focussed on diagnostic test
147 accuracy, we did encounter other designs and it is worth discussing the distinction
148 between these and true “pilot” studies. We found that “feasibility studies” either alone or
149 in combination with other study types were the second largest group identified.

150 However, while the terms “feasibility” and “pilot” study are often used interchangeably,
151 this is erroneous since they describe different study designs. The UK National Institute
152 for Health Research defines a pilot study as, “a smaller version of the main study used
153 to test whether the components of the main study can work together. It is focussed on
154 processes of the main study, for example to ensure that recruitment, randomisation,
155 treatment, and follow-up assessments all run smoothly”⁹. A pilot study, therefore, does
156 not aim to test whether the primary intervention is effective. Where uncertainty exists
157 regarding whether a radiological test “works” or not, i.e. does it measure (or do) what is
158 intended, then the appropriate design is a “feasibility” study.

159
160 Feasibility studies examine, “Characteristics of the proposed outcome measure and in
161 some cases feasibility studies might involve designing a suitable outcome measure”⁹. A
162 radiological example might be a study to determine whether diffusion weighted MR

163 identifies cancer cells. However, the sensitivity and specificity with which that is
164 achieved can only be determined via a properly powered diagnostic test accuracy study,
165 which might require a “pilot” beforehand to see if adequate numbers of representative
166 patients can be accrued. Feasibility studies ask, “can this study be done?”⁹, and have a
167 narrower focus than pilots. Notably, while “Feasibility studies do not evaluate the
168 outcome of interest”⁹, pilot studies do so but not in numbers sufficient for a properly
169 powered study of diagnostic accuracy. However, pilots may contribute individual data to
170 a subsequent main trial, in which case they are termed “internal” pilots (vs, “external”
171 pilots, which exclude their data from any main trial).

172

173 We also encountered designs that aimed to investigate whether a specified technology
174 worked for its intended purpose or not, which can be termed, "technology assessment".
175 In its broadest sense, "health technology assessment" can incorporate evaluation of
176 intended and unintended effects of an intervention, often including clinical and cost
177 effectiveness. "Efficacy assessment" study designs attempt to determine the relative
178 balance of benefits versus harms and are often concerned about safety. We identified
179 one “incidence study”, described as a “pilot”¹⁰. Incidence studies are epidemiological
180 designs that measure incidence of an outcome, myocardial infarction in this case, in a
181 population following their exposure¹¹.

182

183 An initial objective of our research was to investigate the sample size of radiological
184 pilot studies in order to determine how closely they matched general recommendations
185 for pilot sample sizes (hypothesising that they would be too small generally) but our

186 failure to identify genuine pilots precluded this. Notwithstanding this, we found a median
187 sample size of just 20 individuals recruited and 75% of studies described samples
188 below 40. As noted previously, no article presented a power calculation to justify their
189 sample. Several recommendations to estimate pilot sample size are available, for
190 example using a confidence interval around the anticipated standard deviation¹², or 3%
191 the anticipated size of the definitive trial¹³. Several other alternatives have been
192 proposed¹⁴⁻¹⁶. The sample sizes we observed were largely below these
193 recommendations. Furthermore, the few authors presenting a rationale for sample size,
194 usually based this on practicality and convenience rather than factors important to plan
195 a future definitive study. Even those researchers who argue that sample size calculation
196 may not always be necessary, agree that some justification should always be given by
197 authors¹⁷. It has been argued frequently that reporting underpowered research is
198 unethical because such studies can encourage clinical practice based on invalid
199 results¹⁸⁻²¹. While it is self-evident that all genuine pilot studies must be prospective, we
200 nevertheless identified 21 purely retrospective studies. It is also self-evident that
201 planning for a subsequent large trial might occur over several centres but the
202 overwhelming majority of studies in our review were single centre.

203
204 While we believe ours is the first study to investigate pilot studies of radiological
205 interventions, other authors have found this term misused commonly by other medical
206 disciplines. For example, Shanyinde and co-workers randomly sampled 50 medical
207 studies described as pilots and found that the large majority actually focussed on
208 efficacy²². Kannan and Gowri investigated 93 Indian “pilot” studies, finding, like us, that

209 none satisfied this description²³. Similarly, Arain and co-workers concluded that, “Pilot
210 studies are still poorly reported, with inappropriate emphasis on hypothesis testing”²⁴.
211 While it could be argued that Reviewers (and certainly Journal Editors) will have more
212 knowledge of study design and terminology than the average researcher, our (and
213 others’) data suggest that deficiencies in study description are not being rectified after
214 acceptance for publication.

215
216 Our study does have limitations. Most obviously, we considered only those articles
217 using “pilot” in their title. It is possible that a study may be declared a “pilot” elsewhere,
218 for example in the methods section. For example, reporting of a large scale multi-centre
219 study may conceivably incorporate an internal pilot. However, it has been suggested
220 that the large majority (87%) of studies describing themselves as pilots do so in their
221 title²⁵. It may have been informative to identify how many “pilot” studies were followed
222 by a subsequent larger trial but our failure to identify genuine pilots diminished our
223 enthusiasm to do this. *A priori, we decided to restrict our search to four prominent*
224 *radiological journals, hypothesising that if the issue of misinterpretation was a problem*
225 *for major journals (who we might assume enjoy greatest reviewer and Editorial skills),*
226 *then it would be plausible to assume the problem would be even greater for “lesser”*
227 *journals. Restricting a systematic search to specific journals is a well-recognised*
228 *methodological tactic where this facilitates the review while being unlikely to impact the*
229 *outcome. Lancaster and co-workers²⁶ restricted their search of pilot studies preceding*
230 *randomised controlled trials to seven journals, a tactic also adopted by Arain and co-*
231 *workers²⁴.*

232 In summary, by systematic review we observed that, in reality, radiological “pilot”
233 studies mostly report underpowered studies of diagnostic test accuracy. In order to have
234 scientific credibility, we encourage authors, reviewers, and Editors of radiologic journals
235 to familiarise themselves with different methodological study designs and their precise
236 implications.
237

238 **REFERENCES**

239

240 1. Dwyer AJ, Doppman JL, Black WC. The poor quality of early evaluations of MRI.

241 JAMA 1988;260:2661-4.

242 2. Smidt N, Rutjes AW, van der Windt DA, et al. Quality of reporting of diagnostic

243 accuracy studies. Radiology 2005;235:347-53.

244 3. Thabane L, Ma J, Chu R, et al. A tutorial on pilot studies: the what, why and how.

245 BMC Med Res Methodol 2010;10:1.

246 4. Craig P, Dieppe P, Macintyre S, et al. Developing and evaluating complex

247 interventions: the new Medical Research Council guidance. BMJ 2008;337:a1655.

248 5. Moher D, Liberati A, Tetzlaff J, Altman DG, Group P. Preferred reporting items

249 for systematic reviews and meta-analyses: the PRISMA statement. Ann Intern Med

250 2009;151:264-9, W64.

251 6. Aubry S, Dubut J, Nueffer JP, Chaigneau L, Vidal C, Kastler B. Prospective 1-

252 year follow-up pilot study of CT-guided microwave ablation in the treatment of bone and

253 soft-tissue malignant tumours. Eur Radiol 2017; 27:1477-85.

254 7. Elfatairy KK, Filson CP, Sanda MG, Osunkoya AO, Geller RL, Nour SG. In-bore

255 MRI-guided biopsy: can it optimize the need for periodic biopsies in prostate cancer

256 patients undergoing active surveillance? A pilot test-retest reliability study. Br J Radiol

257 2018;91:20170603.

258 8. Venturini E, Losio C, Panizza P, et al. Tailored breast cancer screening program

259 with microdose mammography, US, and MR Imaging: short-term results of a pilot study

260 in 40-49-year-old women. Radiology 2013;268:347-55.

- 261 9. National Institute for Health Research: Supporting informatio for applicants
262 applying to the HTA programme. [https://www.nihr.ac.uk/funding-and-](https://www.nihr.ac.uk/funding-and-support/documents/current-funding-opportunities/hta/hta-supporting-information.pdf)
263 [support/documents/current-funding-opportunities/hta/hta-supporting-information.pdf](https://www.nihr.ac.uk/funding-and-support/documents/current-funding-opportunities/hta/hta-supporting-information.pdf)
264 Accessed 18th April 2019.
- 265 10. Paraschin K, Guerra De Andrade A, Rodrigues Parga J. Assessment of
266 myocardial infarction by CT angiography and cardiovascular MRI in patients with
267 cocaine-associated chest pain: a pilot study. *Br J Radiol* 2012;85:e274-8.
- 268 11. Pearce N. Classification of epidemiological study designs. *International Journal of*
269 *Epidemiology* 2012;41:393-397.
- 270 12. Julious SA, Patterson SD. Sample sizes for estimation in clinical research.
271 *Pharm Stat* 2004;3:213-5.
- 272 13. Stallard N. Optimal sample sizes for phase II clinical trials and pilot studies. *Stat*
273 *Med* 2012; 31:1031-42.
- 274 14. Browne RH. On the use of a pilot sample for sample size determination. *Stat*
275 *Med* 1995;14:1933-40.
- 276 15. Julious SA. Sample size of 12 per group rule of thumb for a pilot study. *Pharm*
277 *Stat* 2005;4:287-91.
- 278 16. Teare MD, Dimairo M, Shephard N, Hayman A, Whitehead A, Walters SJ.
279 Sample size requirements to estimate key design parameters from external pilot
280 randomised controlled trials: a simulation study. *Trials* 2014;15:264.
- 281 17. Billingham SA, Whitehead AL, Julious SA. An audit of sample sizes for pilot and
282 feasibility trials being undertaken in the United Kingdom registered in the United
283 Kingdom Clinical Research Network database. *BMC Med Res Methodol* 2013;13:104.

- 284 18. Altman DG. Statistics and ethics in medical research: III How large a sample? Br
285 Med J 1980;281:1336-8.
- 286 19. Altman DG. The scandal of poor medical research. BMJ 1994;308:283-4.
- 287 20. Halpern SD, Karlawish JH, Berlin JA. The continuing unethical conduct of
288 underpowered clinical trials. JAMA 2002;288:358-62.
- 289 21. Ioannidis JP. Why most published research findings are false. PLoS Med
290 2005;2:e124.
- 291 22. Shanyinde M, Pickering RM, Weatherall M. Questions asked and answered in
292 pilot and feasibility randomized controlled trials. BMC Med Res Methodol 2011;11:117.
- 293 23. Kannan S, Gowri S. Pilot studies: Are they appropriately reported? Perspect Clin
294 Res 2015;6:207-10.
- 295 24. Arain M, Campbell MJ, Cooper CL, Lancaster GA. What is a pilot or feasibility
296 study? A review of current practice and editorial policy. BMC Med Res Methodol
297 2010;10:67.
- 298 25. Kaur N, Figueiredo S, Bouchard V, Moriello C, Mayo N. Where have all the pilot
299 studies gone? A follow-up on 30 years of pilot studies in Clinical Rehabilitation. Clin
300 Rehabil 2017; 31:1238-48.
- 301 Lancaster GA, Dodd S, Williamson PR. Design and analysis of pilot studies:
302 recommendations for good practice. J Eval Clin Pract 2004;10:307-12.

303 **Table 1**

304 Study design for 78 radiological studies described as being a “pilot” in their title.

305 Studies can be composed of more than one design type

306

307

Study design type	Number	Percentage (%)
Diagnostic test accuracy	66	84.6
Feasibility	12	15.4
Technology assessment	8	10.3
Efficacy assessment	5	6.4
Incidence study	1	1.3
Pilot	0	0

308

309

310

311