

1 **Title:** Assessing the impact of specialist home visiting upon maltreatment in England:
2 a feasibility study of data linkage from a public health trial to routine health and
3 social care data

4

5 **Authors:**

6

- 7 • Fiona Lugg-Widger¹ LuggFV@cardiff.ac.uk
- 8 • Rebecca Cannings-John¹ CanningsRL@cardiff.ac.uk
- 9 • Lianna Angel¹ AngelL@cardiff.ac.uk
- 10 • Gwenllian Moody¹ MoodyG@cardiff.ac.uk
- 11 • Jeremy Segrott^{1 2} SegrottJ@cardiff.ac.uk
- 12 • Joyce Kenkre³ Joyce.Kenkre@southwales.ac.uk
- 13 • Michael Robling^{1 2} RoblingMR@cardiff.ac.uk

14

15 Author affiliations

- 16 1 Centre for Trials Research, Cardiff University, Neuadd Meirionnydd, Heath Park,
17 Cardiff. CF14 4YS
- 18 2 DECIPHer Centre, Cardiff University, Cardiff, UK.
- 19 3 Faculty of Life Sciences and Education, University of South Wales

20

21

22

23 Corresponding author: Dr Fiona Lugg-Widger

24

25 **Abstract (350/350)**

26

27 **Background**

28 Follow-up for public health trials may benefit from greater use of routine data. Our
29 trial of a home visiting intervention for first-time teenage mothers assessed
30 outcomes to the child's second birthday. To examine its medium-term impact,
31 particularly upon maltreatment outcomes, we designed a study using routine
32 records.

33

34 **Methods**

35 We aimed to establish the feasibility of our study design, which combines trial data
36 with routine health, social care and education data using a dissent-based linkage
37 model. Trial participant identifiers were linked to routine health, social care and
38 education data if women did not dissent. Data were forwarded to a safe haven and
39 further linked to de-identified trial outcome data. The feasibility study aimed first to
40 establish the acceptability of data linkage through a discussion group of young
41 mothers and by levels of dissent received by the research team. Second, we assessed
42 levels of accurate linkage to both health (via NHS Digital) and education and social
43 care (both via National Pupil Database, NPD). Third, we assessed the availability of
44 data and levels of missingness for key outcomes received for a sample of target
45 study years.

46

47 **Results**

48 Of 1545 mother-child dyads contacted, eight women opted out. The engagement
49 exercise with stakeholders found support for the principle of data linkage, including
50 in the context of maltreatment. Some contributors preferred opt-in consent. Most
51 (99.9%) health records were matched on either three or all four identifiers. Fifty
52 participants were not matched to any health data. Primary outcome data from NPD
53 are derived from any one of three fields, all of which were satisfactorily returned and
54 provided an indication of cases for analysis. Missing data for secondary outcomes
55 varied from 0%(Child looked after status) to 70%(Anatomical Area A&E diagnosis)
56 however when combined with other variables the levels of missingness for outcomes
57 decrease.

58

59 ***Conclusions***

60 Through study set-up and in this pilot, we provide evidence that the main study is
61 feasible, satisfies governance requirements and is likely to generate data of sufficient
62 quality to address our main research questions. Observed levels of missingness or
63 low event rates are likely to affect some secondary analysis (eg, state transition
64 modelling) although overall were satisfactory.

65

66 ***Trial registration:*** NA

67

68 ***Keywords:*** routine data, feasibility study, progression criteria, teenage mothers,
69 home visiting, child maltreatment

70

71 **Background:**

72 Achieving a successful start in life can be particularly challenging for children born to
73 teenage mothers who themselves may struggle to achieve longer-term socio-
74 economic stability [1,2]. We have previously reported on the Building Blocks (BB:0-2)
75 trial of a nurse-led home-visiting intervention, the Family Nurse Partnership (FNP)
76 being delivered to teenage first-time mothers living in eighteen sites in England [3–
77 6]. Over 1600 women participated in the trial, which randomly allocated women to
78 either usually available supportive health and social care alone (usual care) or visits
79 provided in addition to usual care from specially trained FNP nurses from the end of
80 the first trimester until their first child was aged two years.

81

82 Large-scale evaluations of community-based home-visiting and similar public health
83 interventions present a number of challenges. Adverse socio-economic
84 circumstances facing families may create barriers to identifying and recruiting
85 women in the first place, and retaining engagement over an extended period of time
86 can impact outcome assessment. For example, while our self-report follow-up rate
87 at 24 months was 70%, in two contemporary trials of the same or similar
88 intervention in the Netherlands and in Germany the rates were 48% and 46%
89 respectively [7,8]. The BB:0-2 trial also made use of routinely provided healthcare
90 data which was used solely or in combination with other data for both primary and
91 secondary outcomes, enabling more data to be available for analysis compared with
92 just self-report data.

93 In BB:0-2, families were followed-up at 24 months post-partum but the programme
94 was expected to impart beneficial effects on child health and development and
95 maternal life course that would accrue many years after visiting ended. These
96 benefits would be expected to extend into multiple sectors such as education and
97 social care, and so beyond the original primarily healthcare setting of the
98 intervention. Key outcomes would include domains that are sensitive in nature such
99 as maltreatment, which may be subject to reporting bias and non-response bias if
100 solely assessed by maternal self-report.

101

102 Therefore we designed a study which used routine data to evaluate longer term
103 programme impact [9]. This seeks to create a pseudonymised (i.e. replaces key
104 characteristics in the dataset so that individuals can't be directly identified) research
105 database comprising the original trial dataset with a further four years of data from
106 health, social care and education records. Unlike the original BB:0-2 trial, which
107 involved prospectively recorded participant consent, this study would require a
108 dissent process and no additional recourse to self-reported data. As the BB:0-2 trial
109 made use of routinely collected data we had some reassurance about the feasibility
110 of using some of the expected data for longer-term evaluation but not all data
111 sources and not with the dissent model.

112

113 Key remaining questions about the viability of the research design were addressed
114 through a two-stage pilot. The first stage summarised the integrity of programme
115 delivery, potential effect, and ability to access the routinely collected data. This was

116 in response to funder review comments requiring these elements to be addressed
117 early on in the study and treated as start-stop criteria for the continuation of the
118 study. The second stage addressed a range of feasibility parameters related to
119 participant identification, matching and record linkage and data quality, including
120 missingness and numbers available for analysis. Criteria thresholds for progression
121 were not set; rather the project team reported findings to their management group
122 and independent steering committee for information. The rationale for this second
123 stage was to ensure the final datasets could answer the research questions robustly
124 and timely. Records supplied via the Department for Education's National Pupil
125 Database include those from several linked datasets including safeguarding data
126 from local authority departments of children's social care, a primary outcome for the
127 study. The data providers for our study use different unique identifiers to match (e.g.
128 NHS Digital primarily use NHS Number, National Pupil Database use name, postcode,
129 date of birth and gender). When we collected baseline trial data, we were intending
130 to solely link to healthcare records. How well these identifiers from hard to reach
131 young families could be used to allow us to develop a research database of sufficient
132 coverage and data quality needed to be verified. Finally, our research plan involves
133 unconsented access to identifiable records for families who had previously
134 consented to trial participation. We considered that we needed to explore the views
135 of similar members of the public (acting as stakeholders) about such activity and its
136 acceptability (To note, not directly with participants). The aim of this paper is to
137 describe our study to establish the feasibility of using trial data linked to records

138 from multiple sources to study the longer-term impact of a specialist home visiting
139 programme to support teenage first-time mothers in England.

140

141 **Methods:**

142 The study design of this follow-on study using routinely collected data has been
143 described previously in a published protocol paper [9]. For convenience we briefly
144 summarise the essential design here before describing the methods specific to the
145 pilot phase evaluation.

146 *Overview of study design*

147 Building Blocks:2-6 (BB:2-6) aims to extend the duration of follow-up for participants
148 exiting the BB:0-2 trial of the Family Nurse Partnership intervention [3]. It will do so
149 by identifying and linking to routine data from three principal sources, NHS Digital
150 (health data), the National Pupil Database (NPD; education and select social care
151 data) and the Office of National Statistics (ONS; mortality data). These data will be
152 matched to participant identifiers held by the trial team at Cardiff University. The
153 two primary data centres, NHS Digital and the NPD use a different combination of
154 identifiers for matching. NHS Digital use NHS numbers, date of birth, postcode and
155 gender. The NPD uses a combination of forename, surname, date of birth and
156 postcode. Both data centres use exact matching, with NPD fuzzy matching by
157 forename if required. Unique Pupil Number – a unique matching variable used by
158 NPD, was not collected during the trial however will be assigned following matching
159 to enable linkage across NPD datasets.

160 Retrieved data minus personal identifiers will be provided direct to a trusted third
161 party data safe haven which will use the project specific identifiers to link these data
162 to trial outcome data sent from the trial team. Project identifiers are removed and
163 replaced with an encrypted anonymised linking field (ALF-E). Data are accessible to
164 named researchers via a secured remote portal. Data are processed legally under
165 section 251 approval provided via the Confidentiality Advisory Group, Health
166 Research Authority [10]. Trial participants were offered the opportunity to dissent
167 from the study following contacts made by post, email and text messages.

168 *Two-stage pilot*

169 Evidence to support progression of the study (pilot stage one) was gathered via a
170 number of sources. Some of these related to conduct of the original trial (specifically
171 adequacy of intervention delivery and of short-term effect) and are briefly
172 summarised in additional file 1. The key feasibility elements addressed in the second
173 pilot stage were (i) developing an adequate participant dissent model, (ii)
174 establishing acceptable levels of record linkage and (iii) establishing adequate data
175 quality. The governance model outlining required approvals has been described in
176 the BB:2-6 protocol paper.[9] The full follow-up period for the study will include
177 records to 31st March 2017, representing the end of the six-year follow-up. For the
178 pilot stage 2 we requested data from centres to enable a sufficiently informative
179 assessment of data linkage and quality. For NHS Digital (and ONS mortality data
180 which is accessed via the same provider) this included data from study entry of the
181 first mother (June 2009) to 31st March 2015. Local authority safeguarding data

182 accessed via the modular NPD datasets were requested to 31st March 2014 and

183 education data requested from NPD to other end-points in 2014 (table 1).

184 Table 1 Data requested for the second pilot stage

185

Provided by	Dataset	Eligibility / Coverage		Mother	Child	Requested for the Pilot
Dept. of Health	Abortions	England and Wales All abortions performed in the NHS or an approved independent sector		✓	✗	✗
ONS	Mortality records	UK		✓	✓	entry – 31 st March 2015
NHS Digital	Inpatient	Any NHS hospital in England		✓	✓	entry – 31 st March 2015
	Outpatient					entry – 31 st March 2015
	Accident & Emergency					entry – 31 st March 2015
Dept. for Education	Child In Need	< 18 years Registered with social services in England		✓	✓	entry - 31 st March 2014
	Child Looked After					entry - 31 st March 2014
	Early Years Foundation Stage Profile	Public Schools in England	4 yrs	✗	✓	Assessment day July 2013 & July 2014
	Early Years Census		3-4 yrs	✗	✓	Census day Jan 2013 & Jan 2014
	Alternative Provision		2-19 yrs	✓	✓	Census day Jan 2013 & Jan 2014
	Pupil Referral Unit		2-19 yrs	✓	✓	Census day Jan 2014
	School census		2-19 yrs	✓	✓	Winter term 2012 – Summer term 2014
	Key stage One		5-7 yrs	✗	✓	✗

186

187 *(i) Developing and assessing adequacy of dissent model*

188 There were two components to this assessment. First, we implemented the process

189 that provided trial participants an opportunity to register their dissent. Dissent could

190 be registered through a variety of channels (online, email, text message, phone,

191 post). Early on in the study set-up we worked with a group of care-experienced

192 young people (i.e. have spent time in the care of the local authority) to develop our

193 written letter to trial participants [11,12]. Numbers of trial participants approached

194 and numbers of dissenting responses received were recorded.

195 Through a public engagement / involvement process we explored key factors which
196 influenced the acceptability of the planned data linkage and the importance of
197 anonymity to a contact group of young mothers, and how we could develop
198 materials to support dissemination of study findings (and the research methods
199 used) to interested lay parties. Two researchers (JS, JK) met with an on-going young
200 mothers group ('Our Place') who had previously provided lay input to the Building
201 Blocks trial [13] as external stakeholders (i.e. these were not trial participants). A
202 plan for the meeting was jointly developed within the research team, including
203 audio-recording this single session with the approximately 20 mothers, who were
204 expected to attend the group's own regular meeting place in South Wales. Verbal
205 agreement from stakeholders was obtained prior to their participation in the session
206 with the researchers. This was following previous communication between the
207 research team and group coordinators including the provision of information to
208 mothers in advance of the meeting. An initial discussion was held with the group as a
209 whole and then the researchers worked with two smaller (self-selected) groups to
210 gather their views of the use of linked datasets. The discussion was supported by the
211 use of visual aids, which provided further information about the topic (e.g. A4 cards
212 describing datasets being linked). These mothers were also not participants of
213 research, instead external stakeholders. Although the output from the discussions
214 are presented descriptively in line with topic headings, this public involvement was
215 not undertaken as qualitative research and no formal qualitative methodology was
216 applied.
217

218 *(ii) Establishing acceptable record linkage*

219 The number and proportion of participant identifiers matched to records by each
220 data centre was assessed. For NHS Digital data this also included an assessment of
221 the match rate by each step in the matching algorithm (table 2). For both data
222 centres, matches would include both mothers and children. We also assessed
223 descriptively the process for receipt, de-identification and linkage of datasets by the
224 data safe haven.

225 **Table 2 NHS Digital match algorithm**

226

Step (match rate)¹	NHS number	DoB	Sex	Postcode
1	Exact	Exact	Exact	Exact
2	Exact	Exact	Exact	-
3	Exact	Partial	Exact	Exact
4	Exact	Partial	Exact	-
5	Exact	-	-	Exact
6 ²	-	Exact	Exact	Exact
7 ³	-	Exact	Exact	Exact
8	Exact	-	-	-

227 *1: Matching at step 1 or 2 would provide greatest*
228 *reassurance of valid match*

229 *2: Where NHS number does not contradict the match and*
230 *DOB is not 1 January and the POSTCODE is not in the*
231 *'ignore' list*

232 *3: Where NHS number does not contradict the match and*
233 *DOB is not 1 January*

234

235 *(iii) Establishing adequacy of data quality*

236 We assessed data availability and completeness for all variables supplied from both
237 data centres required for primary and secondary analysis. Priority was placed on
238 primary and key secondary outcomes. Numbers of available records, reasons for
239 missingness and narrative assessment of potential impact was undertaken to
240 indicate potential feasibility of the main study.

241

242 **Results**

243 *(i) Adequacy of dissent model*

244 *Retaining eligible participants:* One mother and child dyad was removed due to a
245 child death, leaving 1545 mother-child dyads to contact (Figure 1). Of these, 93 had
246 electively withdrawn during the original trial. In October 2014 letters were sent to all
247 1452 women who had not electively withdrawn and additionally, SMS text messages
248 (n=653, 45%) and emails (n=386, 27%) to those women where contact details were
249 available. Following additional approval to contact women who had electively
250 withdrawn from the trial, we contacted all 93 women by letter in September 2016,
251 and of these we also sent text messages (n=60, 65%) and emails (n=16, 17%) where
252 possible. Of the 1545 mothers contacted, eight (0.5%) dissented and were excluded
253 from the research database. This was made up of seven and one from the 2014 and
254 2016 letters respectively. Additional approval was required from ethics and the
255 confidentiality advisory group for the withdrawn population to ensure the letter sent
256 to them reflected that they had previously withdrawn from the trial.

257

258 **Figure 1 Participant flowchart: families recruited in BB:0-2 trial and followed up in BB:2-6**
259 **feasibility study**

260

261 *Stakeholder views on linkage:* Audio-recording of the discussion with Our Place
262 mothers proved impractical due to background noise and contemporaneous notes
263 were taken instead. 20 mothers were in attendance at the meeting and their
264 children so two groups with six and five mothers each spent time separate to the

265 main group with the researchers to discuss data linkage and its use in more detail.
266 Representing the data linkage process using A3 sheets (for organisations) and A4
267 sheets (for datasets) and how anonymity was preserved when data were accessed
268 by the research team appeared to be informative for participants. The group
269 expressed preferences for a greater use of visual methods (for example, using
270 computers, pictures to represent organisations). The ease by which individuals could
271 be identified through combining data across datasets arose as a question from the
272 group.

273 Although focus group stakeholders were content with the data linkage procedure
274 described and with reassurances about anonymity, there was nevertheless concern
275 expressed about data security against hacking. The nature of data being held (e.g.
276 more sensitive data on maltreatment) did not affect the perceived acceptability of
277 the linkage approach. One participant asked about the possibility of individuals
278 requesting their own data, which may suggest that there remained some lack of
279 clarity about the non-reversibility of anonymisation. One important area where
280 some disagreement within the group arose was the use the of dissent model. The
281 group appeared to be mostly supportive of this approach given the original consent
282 provided in the preceding trial, the efforts made by the researchers to contact
283 women and the pseudonymisation of data involved. However, some participants
284 preferred an opt-in approach as a general principle.

285 *(ii) Adequacy of data linkage*

286 Match rates to NHS Digital and NPD datasets are shown in tables 3 and 4. For NHS
287 Digital, 2851 unique records were sent and 2804 (98.4%) participants were matched,

288 of these, 2801 participants (99.9%) were matched at either step 1 or 2 (see table 2
289 for definition) indicating a greater reassurance of matching to correct individual.
290 There were 64 participants (31 mothers, 33 children) missing from the Inpatients
291 dataset where they would have been expected (i.e. as there should be at least a
292 birth record). However, 15 of these were present in other NHS Digital datasets,
293 indicating a successful match but missing an inpatient record. 49 participants did not
294 appear in any dataset, which is likely due to matching failure or National opt-out
295 (whereby NHS patients in England electively opt-out of their clinical data being used
296 for purposes other than their own direct care).
297 For NPD data, mothers would appear only if aged under 19 years and a child in need
298 or looked after, or in school (table 5). The denominator for planned study primary
299 outcome analysis would be the 90% of children adequately matched. The number of
300 records matched per NPD dataset reflected age of children and duration of coverage
301 of each requested dataset.

302 **Table 3 Match¹ rates for NHS Digital and NPD**
303

	Participants sent (n)	Participants matched (n)	Proportion matched
NHS Digital			
Mother	1434	1407	98.1%
Child	1419	1397	98.4%
NPD			
Mother	1428	99	6.9%
Child	1412	1272	90%

304 *1 Any type of match*

305 **Table 4 Data received from NHS Digital**

Dataset name	N participants in dataset	N Mothers in dataset	N Children in dataset	N Match step 1 & 2 (% based on N participants in dataset)	N records in dataset (multiple records per participant)
--------------	---------------------------	----------------------	-----------------------	---	---

c:\users\siskm1\appdata\local\microsoft\windows\temporary internet

files\content.ie5\1yp9vcd7\clean manuscript of bb2 pilot resubmit 25apr2018.docx

A&E	2451	1205	1246	2446 (99.8%)	13,211
Outpatients	2338	1398	940	2336 (99.9%)	39,067
Inpatients	2789	1403 ¹	1386 ²	2786 (99.9%)	11,882

306 *1: 31 missing, 27 unmatched and 4 of these present in A&E dataset; 2: 33 missing, 22*

307 *unmatched and 11 of these present in A&E & Outpatients datasets*

308

309 **Table 5** **Data received from National Pupil Database (NPD)**

NPD Dataset name	Years provided	Records in dataset (n)	Participants in dataset (n)	Mothers in dataset (n)	Children in dataset (n)
Pupil Level and School Census (PLASC)	2012/13; 2013/14	760	760	4	756 ¹
Pupil Referral Unit (PRU) Census	2013/14	2	2	2	0
Alternative Provision	2012/13; 2013/14	1	1	0	1
Early Years Census (EYC)	2012/13; 2013/14	581	565	0	565 ¹
Child Looked After (CLA)	2008/9 -2013/14	23	23	10	13
Child in Need	2008/9 -2013/14	331*	169*	98*	71

310 *1: 54% of 1412 children were identified in PLASC; 40% in EYC. Summer 2014 was the*
311 *last school census dataset requested for the pilot thus not all children would have*
312 *been expected to be in school (i.e. only by March 2014, would all children have*
313 *turned 3 years of age). *1 record received does not contain any data and therefore*
314 *following further data cleaning may be removed.*

315

316 *(iii) Adequacy of data quality*

317 Assessment included establishing that key outcomes could be adequately derived
318 from supplied data. The primary study outcome is Child in Need (CIN) status to be
319 derived from a combination of three NPD CIN dataset fields (Referral date, Referral
320 but no further action, Reason for closure). For these and all fields retrieved we
321 undertook an impact assessment to clarify the field's role in analysis, number of
322 records retrieved, explanatory notes regarding missingness, and impact on planned
323 analysis. A field's purpose would include acting as primary or secondary outcome
324 (either in combination or with other fields), for cross-checking / validation of other
325 data, and for planned exploratory analysis. Impact was assessed as either No, Low,
326 Medium, High or Not required, with explanations where justified.

327

328 A summary version of the final assessment table is shown to demonstrate these key
329 elements and how that informed the feasibility assessment for each variable (table

330 6). The primary outcome of Child in Need status being recorded by age six years is
331 determined from three fields in the NPD Child In Need data set which shows referral
332 date, further action taken and closure date within the reporting year. A return in any
333 one of these three fields would indicate a positive CIN status. As records would only
334 appear in this dataset following a conditional event (i.e. a referral) it is not possible
335 to assess absent valid cases but does indicate potential number of cases for inclusion
336 in the main analyses. Other secondary outcomes are similarly formed of several
337 fields both for NPD data (e.g. Child protection registration) and NHS Digital data (e.g.
338 Injuries and ingestions) and presence can be inferred by positive entries in one or
339 more of the contributing fields. Levels of missingness in current pilot and original
340 trial data matching are shown where relevant. Some planned analyses were found to
341 be potentially affected by level of missing data (e.g. state transition modelling) or
342 small numbers (Child Looked after status), which would either reduce the scope of
343 analysis or indicate a descriptive approach respectively. Many of the fields in the HES
344 data that show high levels of missing data will be combined (e.g. diagnosis &
345 treatment) and therefore where there is a value in one of these fields it would be
346 assumed that this was an event within the A&E dataset. Missing data may also make
347 some outcomes difficult to derive. In these cases, any assumptions made on the
348 missing information will be stated and if possible varied (worse /best case scenario)
349 and caveats will be made around results to aid interpretation.

350

351 **Table 6. Outcomes and data fields assessed in pilot: records available and**
352 **feasibility assessment**

353

354

355 *Additional work*

356 Data management protocols including de-identification for processing data from
357 project team to data centres and collation at data safe haven were also tested. This
358 included ensuring that the multiple datasets created by each of the two primary data
359 centres could be re-combined while project identifiers known to the project team
360 could be safely removed before data was made available to researchers. Standard
361 data cleaning activities and data re-structuring were enacted but are not otherwise
362 described here.

363

364

365 **Discussion**

366 In this feasibility study we tested a dissent process, which resulted in few trial
367 participants dissenting, and then proceeded to match their identifiers to a high
368 proportion of routine records. The latter include health data matched with a high
369 level of precision using NHS Digital's stepped algorithm process. Fields used in
370 combination will form individual outcomes for the study limiting the impact of some
371 apparent missingness. Some record matching had higher levels of missingness than
372 observed for the same participants in the trial. Nevertheless, the primary outcome
373 analysis appears feasible, as do analyses of many planned secondary outcomes. Low
374 rates of some outcomes may indicate descriptive analysis only and one of the
375 planned analyses of state transition through phases of the child protection process
376 will be limited by the reduced set of fields ultimately available.

377

378 We have established feasibility over two stages. The first required evidence that the
379 evaluation of the nurse-visiting programme had been delivered with sufficient
380 fidelity in the trial phase. A longer-term evaluation also needed to be justified by
381 some indication that the programme was at least not harmful. Progression criteria
382 were developed in discussion with the funder, and the data gathered in the trial's
383 process and outcome evaluation respectively met these criteria. In a second stage,
384 perhaps the most critical set of criteria addressed a range of feasibility parameters
385 many of which could only be determined after the study set-up and through the pilot
386 study presented here. The independent study steering committee has been essential
387 in confirming the scope of, and then progression against these criteria. An inability to
388 meet the criteria at stage one would have probably and correctly led to study
389 closure. It is also possible that serious challenges in the second stage would result in
390 the same decision. However, re-configuring our approach within the same study
391 design was probably the more likely outcome. In practice this is what has happened.
392 Our analysis plan has been adjusted based on what we understand is likely to be
393 available for analysis. The work undertaken to establish the governance
394 infrastructure, mapping and managing the required data linkages and preparing data
395 sets for main analysis (e.g. scripts for data cleaning and re-structuring) provides
396 reassurance for the main study phase. There is greater emphasis being placed on
397 clarity of objectives for feasibility and pilot trials, with detailed criteria and
398 thresholds for progression [14]. Our study adds to that literature with its particular
399 focus on unconsented data linkage from multiple data centres following up from a

400 closed trial sample. Furthermore, we have developed a model for representing data
401 flow relevant to this study type (Figure 1) that will provide the basis for our main
402 results presentation.

403

404 Our study comprises a number of strengths. We have developed and tested the
405 mostly complete model of data linkage required for the main study, with the key
406 data providers and ultimate data safe haven included. This allows us to draw more
407 informed conclusions about how the final model of data linkage will work in concert
408 to produce a viable research database. We have also used actual data from our
409 intended study sample as the basis for the assessment as opposed to simply
410 modelling using dummy data. This therefore provides a more direct test of matching
411 quality and also likely available data, for example, levels of missingness. We haven't
412 presented the data in a way for study results to be interpreted; however, Table 4
413 does describe the number of records found for the cohort. The numbers presented
414 here are consistent with the BB:0-2 trial with regards number of A&E attendance
415 and admissions [3]. In addition, by testing the approach through actual data provider
416 governance systems we have been able to ensure that the final data set can be
417 assembled in a manner that remains acceptable to key stakeholders. Our work with
418 the lay advisory panel has been supportive in this regard too.

419

420 Nevertheless, some questions remain of either direct or general importance. As
421 some study data are events that may not occur for families (e.g. child protection
422 referrals) the assumption is that the absence of a record from the data set is

423 confirmation that there was no event, which may not be the case. Nevertheless, the
424 presence of other related data for each family can be used to confidently infer an
425 event in some cases and overall rates of data linkage remain high across both health
426 and education data providers. The sample of years used for this pilot provides
427 reassurance of what may be available when all years are subsequently requested.
428 However, assessing education data reliant upon children reaching a certain schooling
429 age means that we are currently less able to determine the quality and availability of
430 data required in the main analysis.

431

432 Our impact assessment placed a greater focus on those variables contributing to
433 primary and key secondary outcomes (e.g. Child in need). We needed to ensure that
434 our main study question could be answered even if some other objectives were at
435 risk. Data may be lost due to a variety of reasons, which also vary by contributing
436 data centre – out of date identifiers (e.g. post code), opt-outs which are general
437 (National opt-out) or specific to the study (dissents) and matching errors. Data may
438 also be lost even before data are provided to the data centre (i.e. invalid returns to
439 NHS Digital and NPD). The cumulative impact can only be fully assessed when the full
440 data set has been retrieved for the main study but our pilot sample provides a good
441 estimate of what is possible.

442

443 An underpinning element of our study design is the extraction of routine data to be
444 held pseudonymously in a third-party data safe-haven and without direct consent.
445 We were unable to obtain ethical approval in our original trial for long-term follow-

446 up using routine data as the parameters for such data collection were not then fixed
447 and the validity of baseline consent for much longer-term consent was also
448 questioned. We collected numerous contact details for all participants (including of
449 key others, such as family members) at trial baseline, which were then periodically
450 refreshed during the trial (during data collection and using a tracing service). While
451 we cannot determine how complete actual notifications about the current study to
452 all trial participants was in practice, this approach has helped to ensure that the
453 process for capturing dissent is as meaningful and valid a process as possible.

454

455 We have explored how our general approach to accessing and using sensitive
456 routine data is understood and judged by members of the public. There is
457 considerable policy interest in routine data in research and some effort to align
458 research practice with public opinion [15–17]. We explored how processes for linking
459 and using data were understood and accepted by lay representatives. While we
460 recognise that only a small number of mothers were involved, they still represented
461 the population who are the subject of the intervention under evaluation.
462 Importantly, we have identified topics for further exploration with the group,
463 particularly dissemination. We will use this as the basis for developing materials to
464 maximise public engagement with likely stakeholders and consumers of study
465 results, including trial participants. We consider that it is incumbent upon
466 researchers to consider the optimum role for public engagement in data linkage
467 studies and proactively support this.

468

469

470 **Conclusions**

471

472 Overall we conclude that the main study objectives are achievable albeit that some
473 secondary outcome analyses may be restricted by data that become available in the
474 main data request phase. The value of public investment in similar trials can be
475 increased through greater use of routine data but questions of feasibility will still
476 need to be answered. We have deployed a two-stage approach for decision-making
477 on progression. The first stage may be characterised by decision options: progress,
478 stop, or substantially adjust, which in this scenario were mostly negotiated directly
479 with the funder. A second stage may be characterised by decisions options: progress,
480 or adjust where possible with the steering committee (on behalf of the funder) and
481 the project team negotiating progress. At this second stage defining exact
482 progression criteria may be less critical than simply understanding how available
483 data have impacted upon study results and their interpretation.

484

485 **List of abbreviations**

486 ALF-E – Anonymised Linking Field(encrypted); BB:0-2 – Building Blocks Trial; BB:2-6 –
487 Building Blocks: 2-6 Study; CAG-Confidentiality Advisory Group; CIN – Child in Need;
488 CLA – Child Looked After; EYC – Early Years Census; FNP – Family Nurse Partnership;
489 HRA – Health Research Authority; NIHR PHR – National Institute for Health Research
490 Public Health Research; NPD – National Pupil Database; ONS – Office for National
491 Statistics; PLASC - Pupil Level and School Census; PRU – Pupil Referral Unit.

492

493 **Declarations**

494 ***Ethics approval and consent to participate***

495 Ethics approval of the study has been given by the Research Ethics Committee for
496 Wales (14/WA10062) and the transfer and use of identifiable data has been
497 approved by the Health Research Authority (HRA) Confidentiality Advisory Group
498 (CAG) (CAG 10-08(b)/2014).

499 ***Consent for publication***

500 The study comprises a pseudonymised data set, which has been developed via the
501 application of a dissenting model in accordance with the Data Protection Act 1998.
502 No individual can be identified in presented data.

503 ***Availability of data and material***

504 Routine data supplied to the study from NHS Digital and from the Department of
505 Education are subject to specific data sharing agreements and the study to the
506 arrangements under Section 251 of the Health and Social Care Act 2006. All study
507 data are maintained in a data safe haven with strict access controls, which restrict
508 further sharing of data.

509 ***Competing interests***

510 The authors declare that they have no competing interests

511 ***Funding***

512 This project was funded by the National Institute for Health Research Public Health
513 Research (NIHR PHR) Programme (project number 11/3002/11). The views and

514 opinions expressed therein are those of the authors and do not necessarily reflect
515 those of the NIHR PHR Programme or the Department of Health.

516 ***Authors' contributions***

517 MR is the chief investigator for the BB:2-6 study, FL is responsible for the study
518 management of the project and data access, RCJ, GM and LA are responsible for the
519 data analysis and management. JS & JK organised and ran the Our Place workshop.
520 All authors contributed to the content of the manuscript and all have read and
521 approved the final version.

522 ***Acknowledgements***

523 We are grateful to all the families who having previously participated in the Building
524 Blocks trial have kindly enabled us to undertake this follow-on study. We are also
525 grateful for the mothers attending the Our Place group who contributed to the focus
526 group discussion. The Centre for Trials Research, Cardiff University receives funding
527 from Health and Care Research Wales. The work was undertaken with the support
528 of The Centre for the Development and Evaluation of Complex Interventions for
529 Public Health Improvement (DECIPHer), a UKCRC Public Health Research Centre of
530 Excellence. Joint funding (MR/KO232331/1) from the British Heart Foundation,
531 Cancer Research UK, Economic and Social Research Council, Medical Research
532 Council, the Welsh Government and the Wellcome Trust, under the auspices of the
533 UK Clinical Research Collaboration, is gratefully acknowledged. Authors are also
534 supported by PRIME centre Wales funding from Health and Care Research Wales.
535 The authors gratefully acknowledge the contribution of NHS Digital (Copyright ©
536 2012, re-used with the permission of the health and social care information centre),

537 all rights reserved) and the Department for Education in providing data to the study.

538 The Department for Education do not accept responsibility for any inferences or

539 conclusions derived from the NPD Data by third parties.

540

541 **References**

542 1. Moffitt TE. Teen-aged mothers in contemporary Britain. *J Child Psychol*

543 *Psychiatry*. 2002;43(6):727–42.

544 2. Botting B, Rosato M, Wood R. Teenage mothers and the health of their

545 children. *Popul Trends*. 1998:19–28.

546 3. Robling M, Bekkers MJ, Bell K, Butler CC, Cannings-John R, Channon S, et al.

547 Effectiveness of a nurse-led intensive home-visitation programme for first-time

548 teenage mothers (Building Blocks): A pragmatic randomised controlled trial. *Lancet*.

549 2016; 387(10014):146-155.

550 4. Olds DL, Henderson CR, Chamberlin R, Tatelbaum R. Preventing Child Abuse

551 and Neglect: A Randomized Trial of Nurse Home Visitation. *Pediatrics*. 1986;78(1).

552 5. Kitzman H, Olds DL, Henderson CR, Hanks C, Cole R, Tatelbaum R, et al. Effect

553 of Prenatal and Infancy Home Visitation by Nurses on Pregnancy Outcomes,

554 Childhood Injuries, and Repeated Childbearing. *JAMA*. American Medical

555 Association; 1997;278(8):644.

556 6. Olds DL, Robinson J, O 'brien R, Luckey DW, Pettitt LM, Henderson CR, et al.

557 Home Visiting by Paraprofessionals and by Nurses: A Randomized, Controlled Trial.

558 *Pediatrics*. 2002;110(3).

- 559 7. Mejdoubi J, Van Den Heijkant SCCM, Van Leerdam FJM, Heymans MW,
560 Crijnen A, Hirasing RA. The effect of VoorZorg, the dutch nurse-family partnership,
561 on child maltreatment and development: A randomized controlled trial. PLoS One.
562 2015;10(4).
- 563 8. Sierau S, Dähne V, Brand T, Kurtz V, von Klitzing K, Jungmann T. Effects of
564 Home Visitation on Maternal Competencies, Family Environment, and Child
565 Development: a Randomized Controlled Trial. Prev Sci. 2016;17(1):40–51.
- 566 9. Lugg-Widger F, Cannings-John R, Channon Sue, Fitzsimmons D, Hood K, Jones
567 K, et al. Assessing the medium-term impact of a home-visiting programme on child
568 maltreatment in England: protocol for a routine data linkage study. BMJ Open.
569 2017;7(e015728).
- 570 10. Health and Social Care Act 2012. Queen’s Printer of Acts of Parliament; 2012
571 Accessed 31 Oct 2017.
- 572 11. Shaw I, Holland S. Doing Qualitative Research in Social Work. SAGE. 2014:91
- 573 12. CASCADE Voices – CASCADE: Children’s Social Care Research and
574 Development Centre - Cardiff University.
575 <http://sites.cardiff.ac.uk/cascade/people/young-peoples-advisory-group/>. Accessed
576 13 Nov 2017.
- 577 13. Robling M, Bekkers MJ, Bell K, Butler CC, Cannings-John R, Channon S, et al.
578 The Building Blocks Trial. 2014.
- 579 14. Eldridge SM, Lancaster GA, Campbell MJ, Thabane L, Hopewell S, Coleman CL,
580 et al. Defining Feasibility and Pilot Studies in Preparation for Randomised Controlled
581 Trials: Development of a Conceptual Framework. PLoS One. 2016;11(3):e0150205.

- 582 15. Robling MR, Hood K, Houston H, Pill R, Fay J, Evans HM, et al. Public attitudes
583 towards the use of primary care patient record data in medical research without
584 consent: a qualitative study. *J Med Ethics*. 2004;30:104–9.
- 585 16. Aitken M, De J, Jorre S, Pagliari C, Jepson R, Cunningham-Burley S. Public
586 responses to the sharing and linkage of health data for research purposes: a
587 systematic review and thematic synthesis of qualitative studies. *BMC Med Ethics*.
588 2016;17(1).
- 589 17. Hill EM, Turner EL, Martin RM, Donovan JL. “Let’s get the best quality
590 research we can”: public awareness and acceptance of consent to use existing data
591 in health research: a systematic review and qualitative study. *BMC Med Res*
592 *Methodol*. 2013;13:72.

Table 6 Outcomes and data fields assessed in pilot: records available and feasibility assessment

Outcomes	Data source: Native field name	Missing (n)	Commentary of findings	Impact
Primary				
Child in Need (CIN) status as of 31 March each year	NPD>CIN: <i>Referral date</i>	1	No data appear across one record- record to be excluded (this will apply to all fields below) Note: 34 records with dates prior to time point ranges (1997 – 2007). This is expected.	None
	NPD>CIN: <i>Referral – no further action</i>	42	No data collected in 2008/2009 time point (accounts for 38 records). Some blanks appear in 2009/2010 time point however the referral date on these records is prior to 1st April of that data collection year. *Assumption that time point cycle is Apr-Mar.	Low: assumption missing data indicate further action was required and that the child was in need
	NPD>CIN: <i>Reason for closure</i>	142	No pattern - further investigation required	As above
Secondary^a				
CIN categorisation	<i>NPD>CIN: Category of Abuse</i>	329	Only data from 2008/9 accessed in pilot. For main phase data from 2012/3 will be accessed and also 'NPD> CIN: Latest category of abuse' will be included which may improve data quality.	Still to be determined
Child looked after (CLA) status	<i>NPD>CLA: Category of need; Legal status; Placement; REC</i>	0	All records returned are complete.	None
Child Protection registration (plan)	<i>NPD>CIN: Child Protection Plan (CPP) Indicator</i>	195	Expected - All missings from 2010/2011 time point onwards. Data not collected during these years.	Low: CPP flag can be determined from other fields
	<i>NPD>CIN: No. of previous Child Protection Plans</i>	320	No pattern to missingness. Only 11 records have a value recorded, 9 of these are zero.	As above
	<i>NPD>CIN: Child Protection Plan start date</i>	320	Expected - not all children will have had a CPP. Only 11 records have a date recorded, these correspond with data captured in the 'no. of previous CPPs'	As above

Outcomes	Data source: Native field name	Missing (n)	Commentary of findings	Impact
			above.	
	NPD>CIN: <i>Child Protection Plan end date</i>	327	Expected - only 4 records have an end date recorded. Corresponds with those records where a start date is recorded. Data check done - end dates are after the start date.	As above
Exploratory Markov chain modelling ^b	NPD>CIN: <i>Date of initial child protection conference</i>	327	Expected - not all children would have had a child protection conference. However further checks required to confirm validity of data.	Medium: Low numbers may impact analysis
Injuries and ingestions	NHSD>A&E: <i>A&E diagnosis (diag n)</i>	5981	45% missing (1650/6336 missing in BB Trial – 26% missing)	Medium: All diag / treat / inv fields to be used in combination to define inj / ing ^c
	NHSD>A&E: <i>A&E diagnosis – 2 char (diag2 n D)</i>	3604	27% missing (1849/6336 missing in BB Trial – 29% missing)	As above
	NHSD>A&E: <i>A&E investigation (invest n)</i>	1728	13% missing (1396/6336 missing in BB Trial – 22% missing)	As above
	NHSD>A&E: <i>A&E investigation – 2 char (invest n D)</i>	1712	13% missing (1395/6336 missing in BB Trial – 22% missing)	As above
	NHSD>A&E: <i>A&E treatment (treat n)</i>	2349	18% missing (1411/6336 missing in BB Trial – 22% missing)	As above
	NHSD>A&E: <i>A&E treatment – 2 Char (treat2 n D)</i>	2126	16% missing (1417/6336 missing in BB Trial – 22% missing)	As above
	NHSD>A&E: <i>A&E diagnosis – Anatomical Area (diaga n D)</i>	9281	70% missing (4725/6336 missing in BB Trial – 74% missing)	

(1) small numbers may be an issue - descriptive analysis will be used if necessary; (2) missing a small amount of closure dates;

(a) Additional fields were retrieved for secondary outcomes and assessed solely for presence (Special Educational Needs, Disability, Day care attendance, Early Years assessment, School attendance, Key stage one attainment)

(b) To explore probability of progression through each stage of child protection process

(c) Same fields also contribute to assessment of subsequent pregnancies (via pregnancy-related A&E attendances)

c:\users\siskm1\appdata\local\microsoft\windows\temporary internet
files\content.ie5\1yp9vcd7\clean manuscript of bb2 pilot resubmit 25apr2018.docx

Additional material

File name: Additional file 1 (.docx)

Title of data: Evidence supporting progression derived from trial (BB:0-2) and
feasibility study (BB: 2-6) phases

Description of data: Detailed progression criteria as set out by the study funder at
the start of the project to ensure research objectives could be met.