

1 **Title**

2 BLOod Test Trend for cancEr Detection (BLOTTED): protocol for an observational and prediction
3 model development study using English primary care electronic health records data.

5 **Authors**

6 Pradeep S. Virdee¹, Clare Bankhead¹, Constantinos Koshari¹, Cynthia Wright Drakesmith¹, Jason
7 Oke¹, Diana Withrow¹, Subhashisa Swain¹, Kiana Collins¹, Lara Chammas², Andres Tamm², Tingting
8 Zhu¹, Eva Morris³, Tim Holt¹, Jacqueline Birks³, Rafael Perera¹, FD Richard Hobbs¹, Brian D. Nicholson¹

10 ¹Nuffield Department of Primary Care Health Sciences, University of Oxford, Oxford, UK

11 ²Big Data Institute, University of Oxford, Oxford, UK

12 ³Centre for Statistics in Medicine, NDORMS, University of Oxford, Oxford, UK

14 Pradeep S. Virdee (PSV): pradeep.virdee@phc.ox.ac.uk

15 Clare Bankhead (CB): clare.bankhead@phc.ox.ac.uk

16 Constantinos Koshari (CK): constantinos.koshari@phc.ox.ac.uk

17 Cynthia Wright Drakesmith (CWD): cynthia.wright@phc.ox.ac.uk

18 Jason Oke (JO): jason.oke@phc.ox.ac.uk

19 Diana Withrow (DW): diana.withrow@phc.ox.ac.uk

20 Subhashisa Swain (SS): subhashisa.swain@phc.ox.ac.uk

21 Kiana Collins (KC): kiana.collins@st-hughs.ox.ac.uk

22 Lara Chammas (LC): lara.chammas@st-annes.ox.ac.uk

23 Andres Tamm (AT): andres.tamm@wolfson.ox.ac.uk

24 Tingting Zhu (TZ): tingting.zhu@eng.ox.ac.uk

25 Eva Morris (EM): eva.morris@ndph.ox.ac.uk

26 Tim Holt (TH): tim.holt@phc.ox.ac.uk
27 Jacqueline Birks (JB): jacqueline.birks@csm.ox.ac.uk
28 Rafael Perera (RP): rafael.perera@phc.ox.ac.uk
29 Richard Hobbs (RH): richard.hobbs@phc.ox.ac.uk
30 Brian D. Nicholson (BDN): brian.nicholson@phc.ox.ac.uk

31

32 **Corresponding author**

33 Author: Brian D. Nicholson

34 Address: Radcliffe Primary Care Building, Radcliffe Observatory Quarter, Woodstock Road, University
35 of Oxford, Oxford, OX2 6GG, UK

36 Email: brian.nicholson@phc.ox.ac.uk

37 ORCID: 0000-0003-0661-7362

38

39 **Abstract**

40 **Background**

41 Simple blood tests can play an important role in identifying patients for cancer investigation. The
42 current evidence base is limited almost entirely to tests used in isolation. However, recent evidence
43 suggests combining multiple types of blood tests and investigating trends in blood test results over
44 time could be more useful to select patients for further cancer investigation. Such trends could
45 increase cancer yield and reduce unnecessary referrals. We aim to explore whether trends in blood
46 test results are more useful than symptoms or single blood test results in selecting primary care
47 patients for cancer investigation. We aim to develop clinical prediction models that incorporate
48 trends in blood tests to identify risk of cancer.

49

50 **Methods**

51 Primary care electronic health records data from the English Clinical Practice Research Datalink
52 Aurum primary care database will be accessed and linked to cancer registrations and secondary care
53 datasets. Using a cohort study design, we will describe patterns in blood testing (Aim 1) and explore
54 associations between covariates and trends in blood tests with cancer using mixed-effects, Cox, and
55 joint models (Aim 2). To build the predictive models for risk of cancer, we will use multivariate joint
56 modelling and machine-learning, incorporating simultaneous trends in multiple blood tests, together
57 with other covariates (Aim 3). Model performance will be assessed using various performance
58 measures, including c-statistic and calibration plots.

59

60 **Discussion**

61 These models will form decision rules to help general practitioners find patients who need referral
62 for further investigation of cancer. This could increase cancer yield, reduce unnecessary referrals,
63 and give more patients the opportunity for treatment and improved outcomes.

64

65 **Keywords**

66 Cancer; early detection; blood test; trend; primary care; CPRD

67

68 **Introduction**

69 A recent clinical review concluded that simple blood tests have an important role in identifying
70 patients for cancer investigation¹. However, analysis of National Cancer Diagnosis Audit in Primary
71 Care data suggests that primary care investigations may delay referral². Smarter use of blood tests to
72 select patients for further cancer investigation could increase cancer yield and reduce unnecessary
73 referrals, minimising the psychological and physical harm to patients and economic costs of
74 unnecessary testing.

75

76 The current evidence base for using blood tests to identify patients at risk of cancer in primary care
77 is almost entirely limited to single blood tests¹. Except for anaemia and jaundice, the risk of
78 individual cancers associated with blood test abnormalities is too low to warrant urgent cancer
79 investigation³. For all cancers combined, the following blood test abnormalities increase the risk of
80 cancer above the three percent threshold recommended by the National Institute for Health and
81 Care Excellence (NICE) for urgent investigation: low albumin, raised platelets, raised calcium, and
82 raised inflammatory markers¹. Whilst these test abnormalities give general practitioners (GPs) an
83 indication that cancer may be present, they leave uncertainty over which cancer(s) should be
84 investigated. In practice, GPs may interpret test results in combination with other contemporaneous
85 tests or previous results of the same test, particularly if the current result is abnormal.

86

87 Methodological innovation is required to understand whether incorporating blood test change over
88 time may provide more accurate cancer prediction for cancer overall and for individual cancer sites.
89 For example, a patient with a low-normal haemoglobin may not be regarded as high risk if the
90 haemoglobin result is interpreted in isolation and an abnormal/normal binary threshold is used.
91 However, a low-normal result following years of high-normal results may represent an opportunity
92 for cancer investigation. Methods for incorporating repeated measures data into clinical decision
93 rules are well known⁴. Our group has recently completed analyses investigating trends in multiple
94 components of the full blood count blood test in the ten years prior to colorectal cancer diagnosis⁵.
95 We also developed prediction models (using joint modelling) designed for early detection of
96 colorectal cancer, using only data earlier than two years before diagnosis, which incorporated trends
97 to identify two-year risk of diagnosis⁶. The sex-stratified trends indicated that a simultaneous
98 patient-level decline in haemoglobin and mean corpuscular volume and rise in platelet count from a
99 steady trend increased the risk of colorectal cancer diagnosis in two years from the current blood

100 test, with good discrimination and calibration⁶. Serial testing may hold greater potential to rule-in
101 and rule-out further cancer investigation.

102

103 In addition to serial testing, test combinations may hold greater potential to rule-in and rule-out
104 patients for further cancer investigation⁷. A recent case-control study illustrated that a normal
105 erythrocyte sedimentation rate plus a normal haemoglobin reduced the risk of myeloma sufficiently
106 to rule-out the need for further investigation⁸. A cohort study of patients referred to a Danish
107 Multidisciplinary Diagnostic Centre (MDC) showed combinations of abnormal tests markedly
108 increased the probability of cancer being diagnosed⁹. However, multivariable cancer-risk prediction
109 models widely accessible to National Health Service (NHS) GPs incorporate multiple symptoms and
110 risk factors but not blood test results^{10, 11}. Simple clinical scores including age-group, sex, and seven
111 simple primary care blood tests (albumin, alkaline phosphatase, C-reactive protein, haemoglobin,
112 liver enzymes, platelets, and total white cell count) could be used to select patients with unexpected
113 weight loss who do and do not warrant further cancer investigation¹². Internal validation of these
114 risk scores have shown good discrimination between patients with and without cancer, were well
115 calibrated at the levels of risk that decisions to investigate are made in primary care, and have
116 shown superior clinical utility compared to models including only age, sex, and symptoms¹². It
117 remains unclear whether these findings are valid in external NHS and primary care datasets from
118 abroad and whether similar scores could be developed for patients with other cancer symptoms.

119

120 There is therefore a pressing need for studies to establish the optimal and most rational use of
121 trends in simple and commonly available blood tests to select primary care patients for urgent
122 cancer investigation.

123

124 **Aims and objectives**

125 The overall aim of BLOTTED is to understand whether incorporating changes in blood tests over time
126 could optimise the selection of primary care patients for cancer investigation compared to models
127 including symptoms, signs, and single blood test results.

128

129 **Aim 1.1**

130 To describe results and trends of blood tests in individuals attending NHS primary care between
131 2000 and 2019 overall, by age, sex, ethnicity, deprivation, and comorbidity. This aim will evaluate
132 the underlying epidemiology of blood test change over time in a primary care population.

133

134 **Aim 1.2**

135 To describe whether changes in blood test trends occur prior to cancer symptoms, single
136 measurement blood test abnormalities, screening findings, and referrals for suspected cancer overall
137 and by type of blood test, symptom, referral pathway, and cancer characteristics (diagnosis route,
138 site, grade, histology, stage). This aim will assess opportunities for earlier cancer diagnosis, such as
139 referral for urgent cancer investigation based on changes in the trend of blood tests instead of
140 symptoms or abnormal single blood test measurements.

141

142 **Aim 2.1**

143 To test the association between blood test trends and subsequent cancer overall and by age, sex,
144 and cancer characteristics (diagnosis route, site, grade, histology, stage). This aim will establish the
145 association between cancer diagnosis and trends in blood tests in the years leading up to cancer
146 diagnosis and identify when blood test trends could be used to prompt cancer investigation.

147

148 **Aim 2.2**

149 To explore whether blood test trends are independent of trends in other blood tests and clinical
150 features of cancer (symptoms, signs, and individual blood test results). This aim will examine
151 whether blood test trends should be considered in isolation or if their predictive value increases
152 when combined with other blood test trends and/or clinical features.

153

154 **Aim 3**

155 To develop and validate decision rules (prediction models) to select individuals attending NHS
156 primary care for cancer investigation. This aim will pull together the learnings from the prior aims
157 into prediction models to derive and test a clinical strategy for GPs to guide patient selection for
158 cancer investigation, comparing the performance of these models with performance of existing
159 models.

160

161 **Methods**

162 Aims 1 and 2 will be reported following the STROBE guidelines for observational research. Aim 3 will
163 be reported following the TRIPOD guidelines for development and/or validation of prediction
164 models.

165

166 **Data**

167 Primary care electronic health records data will be obtained from the English Clinical Practice
168 Research Datalink (CPRD) Aurum¹³. Cancer diagnoses will be obtained from the National Cancer
169 Registration and Analysis Service (NCRAS), Hospital Episode Statistics (HES) database, and Office of
170 National Statistics (ONS) (if related to death). Our blood test data preparation and reporting will
171 follow the steps outlined in our previous work¹⁴.

172

173 **Study design**

174 A cohort study will be used for each aim, allowing for an appropriate assessment of absolute risk of
175 cancer diagnosis.

176

177 **Participants**

178 We will include patients registered in CPRD between January 1st 2000 and December 31st 2019.

179 Patients will be eligible for linkage with NDRS, HES, and ONS databases. We will exclude patients
180 with less than 12 months registered with the general practice or less than 2 years of follow-up data
181 following study entry.

182

183 **Outcome**

184 An incident diagnosis of any cancer made during the study period recorded. Diagnoses will be
185 obtained primarily from the NDRS database, with additional diagnoses obtained from the CPRD, HES,
186 and ONS databases. Patients without a diagnosis will be censored at the earliest of date of leaving
187 the practice, death, or 31st December 2019 (data cut).

188

189 A validated library of Read/SNOMED CT/ICD-10 codes developed previously by this group will be
190 used to identify all incident cancer diagnoses throughout the study period. Data will also be
191 extracted on cancer site, stage, grade, and histology at diagnosis.

192

193 **Predictors**

194 We will explore trends in the blood test results including the full blood count, liver function, and
195 renal function tests (a full list is in Box 1):

196

197 **Box 1: blood tests under investigation in BLOTTED**

Full Blood Count: red blood cell count, white blood cell count, haemoglobin, haematocrit, mean cell volume, mean cell haemoglobin, mean cell haemoglobin concentration, red blood cell distribution width, platelet count, mean platelet volume, lymphocyte count, eosinophil count, neutrophil count, basophil count, monocyte count, lymphocyte %, eosinophil %, neutrophil %, basophil %, monocyte %

Liver Function Tests: alanine aminotransaminase, albumin, alkaline phosphatase, aspartate transaminase, bilirubin, alpha fetoprotein

Renal Function: sodium, potassium, creatinine, urea

Inflammatory Markers: C-reactive protein, erythrocyte sedimentation rate, plasma viscosity

Other tests: amylase, HBA1c, calcium, calcium adjusted, total protein, blood glucose, fasting glucose, thyroid stimulating hormone

198

199

200 Data will also be extracted using SNOMED CT codes to explore the effect of the following covariates,
201 which could independently affect the predictive value of blood test trends and likelihood of cancer:

202

203 1) Personal characteristics – age, sex, ethnicity, smoking history, alcohol intake, family/personal
204 history of cancer, and patient-level Index of Multiple Deprivation (IMD) score.

205 2) Cancer symptoms and signs – symptoms shown to have an independent association with cancer
206 as described by NICE³ and the primary care literature. All occurrences in the study period will be
207 identified, allowing analyses of the first occurrence and repeat occurrences.

208 3) Results of other basic investigations available in primary care: B12 level, folate level,
209 Carbohydrate Antigen 19-9 (CA19-9), Carcinoembryonic Antigen (CEA), Cancer Antigen 125
210 (CA125), Chest X-ray, Lactate Dehydrogenase (LDH), Prostate Specific Antigen (PSA), Iron Studies
211 (Ferritin, Serum Iron, Total Iron Binding Capacity, Transferrin), Cholesterol.

212 4) Co-morbidity – recorded or implied from the prescribing record.

213 5) Prescribed medications.

214 6) Referrals for the urgent investigation of cancer, participation in cancer screening, and diagnosis
215 route.

216

217 **Sample size**

218 A sample size calculation was performed using the '*pmsampsize*' package in Stata software, recently
219 developed for prediction models by Riley *et al*¹⁵. The package uses the number of proposed
220 predictor parameters in the model and the expected mean follow-up time, event rate, Cox-Snell R²,
221 and amount of shrinkage (to adjust for overfitting). We based the calculation on our recent
222 prediction model for 6-month risk of cancer following unexpected weight loss¹². However, the
223 appropriate risk window will be explored (Aim 3) so we performed the calculation twice: for 6-
224 month risk and 1-year risk, to show the range of patient numbers required.

225

226 The number of proposed predictor variables in this study is 70, expected mean follow-up is 11 years
227 (obtained from a second study using blood tests from CPRD¹⁶), 0.046 Cox-Snell R² (from our recent
228 model¹²) and expected 0.9 shrinkage factor. The 6-month event rate was 0.014 in our recent model
229 (908 of 63,973 diagnosed¹²). We found no study that reported 1-year risk of any cancer in patients
230 visiting primary care, so we used a conservative approach to obtain the 1-year event rate: as the
231 outcome window is doubled from 6 months, we also doubled the event rate to give ~1,800 events.
232 However, patients are less likely to be diagnosed around the 1-year time-point so we decreased the
233 number of events to 1,500. The 1-year event rate is then assumed $1500/63973 = 0.023$. These inputs
234 give 13,315 patients required at minimum. From these, we expect 146,465 person-years of follow-
235 up, with 2,079 diagnosed in 6 months (~30 events-per-variable) and 3,435 in 1 year (~50 events-per-
236 variable) (note: we request all patients, not 13,315 alone).

237

238 We expect to exceed these numbers, with another study reporting 69,942 incident cases of cancer
239 within two years of reporting symptoms among 3,850,712 primary care patients^{10, 11}. Another study

240 using CPRD GOLD included all primary care patients available, reporting 10,875,556 full blood count
241 blood tests among 1,893,641 primary care patients between 01/01/2000 and 14/01/2014¹⁷,
242 providing reassurance that there will be sufficient data for analysis.

243

244 **Data/statistical analysis**

245 Aim 1 (descriptive statistics):

246

247 1. Descriptive statistics will be used to describe patterns in blood testing, such as the frequency
248 and time between serial tests, and summarise the results of each blood test overall and
249 stratified by personal characteristics and comorbidity status. For the test results, descriptive
250 statistics will include means (standard deviation) and medians (inter-quartile range).
251 Associations between patient characteristics and patterns in blood testing will be examined.

252

253 2. The feasibility and comparison of established monitoring methods to summarise and identify
254 blood test change over time (e.g. absolute change, percentage change, and rate of change).

255

256 3. Locally Weighted Scatterplot Smoothing (LOWESS) will be used to graphically describe trends in
257 each blood test, summarised overall and by personal characteristics, cancer characteristics
258 (diagnosis route, site, grade, histology, stage), and comorbidity status.

259

260 4. The mean (standard deviation) and median (inter-quartile range) number of days will be used to
261 summarise the intervals between changes in blood test trends, cancer symptoms, single blood
262 test abnormalities based on existing thresholds, referrals for cancer investigation or specialist
263 review, and cancer diagnosis.

264

265 Aim 2 (association with cancer):

266

267 5. Mixed-effects models: The trend in each blood level will be analysed using mixed-effects models.
268 Multivariable models (one for each blood level) will include time (to model trends), personal,
269 and clinical characteristics as fixed effects. Non-linearity will be assessed, such as non-linear
270 relationships in blood levels over time. A random intercept for patient and random slope for
271 time will be used to allow for differences in trends between patients. An unstructured
272 correlation matrix will be used to account for repeated measures. Estimates will be presented
273 with corresponding 95% confidence intervals (CIs). Lessons learned from this analysis, such as
274 the relevant blood tests, confounders, and methods for handling non-linear trends, will inform
275 the development of the prediction models in Aim 3.

276

277 6. Time-to-event model: Cox modelling will be used to establish temporal associations between
278 personal and clinical characteristics and subsequent cancer diagnoses. Models will be repeated
279 by cancer characteristics (diagnosis route, site, grade, histology, stage) to assess, for example,
280 whether covariates are associated with individual or grouped cancer sites or can distinguish
281 between early-stage and late-stage cancer. Estimates will be presented with corresponding 95%
282 CIs. Lessons learned from this analysis will inform the development of the prediction models in
283 Aim 3.

284

285 7. Joint models: Joint modelling of longitudinal blood test data and time-to-event data will be used
286 for each blood level separately to assess the association between blood level trend and
287 subsequent cancer diagnosis. In joint modelling, the trend is identified using a mixed-effects
288 model and included as a covariate in the Cox model to assess association with subsequent
289 cancer diagnosis. Initially, univariable joint models will be developed, which include only the
290 trend (adjusted for personal and clinical characteristics, as per the mixed-effects model above).
291 Multivariable joint models will then include the trend and personal and clinical characteristics

292 related to cancer diagnosis (as per the Cox model above). Estimates will be presented with
293 corresponding 95% CI. Lessons learnt from this analysis will inform the development of the
294 prediction models in Aim 3.

295

296 8. Machine learning models, such as supervised clustering of trajectories using neural networks,
297 will be considered to benchmark against the joint models. Instead of a traditional clustering
298 method, which is unsupervised, the proposed method will cluster trends based on the outcome
299 labels. As a result, clusters of different cancer characteristics (diagnosis route, site, grade,
300 histology, stage) can be obtained, which allow us to provide patient-specific temporal
301 association between outcome vs. exposure and covariates (such as symptoms and blood tests).

302

303 Aim 3 (prediction):

304

305 9. Multivariate joint model: Multivariate joint modelling of longitudinal and time-to-event data will
306 be developed to incorporate simultaneous trends in multiple blood tests to identify risk of
307 cancer. The most promising blood test trends (as identified in the associations analysis described
308 above) will be combined into a multivariate joint model, which will also include personal and
309 clinical characteristics as covariates for subsequent cancer diagnosis. We will explore the
310 predictive value of various clinical features and blood test trends, both as single covariates and
311 in combinations. Pearson/Spearman correlation will be used to assess the correlation between
312 blood components. Model coefficients will be presented with corresponding 95% CI.

313

314 10. Machine learning models, such as supervised clustering of trajectories using neural networks,
315 will be considered to benchmark against multivariate joint modelling.

316

317 11. Model validation: Internal k-fold cross validation (where k = NHS geographical region) is planned
318 to internally validate the model. However, this approach may not be possible due to the
319 computational intensiveness of joint models. Alternatively, a split sample approach may be
320 considered, whilst ensuring sample size requirements are met for model development.
321 Performance statistics will be derived, such as sensitivity, specificity, predictive values, area
322 under the receiver operating characteristic curve (or c-statistic), D-statistic for discrimination,
323 Brier score, R^2 statistic for explained variation, calibration slope, and calibration plots.

324

325 12. The chosen model(s) will be used to generate individual patient predictions. Predictions will be
326 generated at multiple time-points along the blood test trajectory and performance measures
327 derived for each time-point to identify when the optimal trend/threshold for referral for further
328 investigations is reached.

329

330 13. Existing models incorporating blood test data to predict cancer diagnosis in primary care will be
331 externally validated (by running the equations generated in the original derivation study on our
332 CPRD data). Predictive performance will be compared to our new model(s). Performance of
333 models incorporating blood test trend for cancer-risk will be compared to existing models
334 including risk factors, symptoms, signs, static (or single) blood test data, and screening
335 assessments. Performance measures described above will be used to assess diagnostic ability.

336

337 14. Decision curve analysis will be used to determine whether models incorporating blood test trend
338 have superior clinical utility to models including static (or single) blood test data.

339

340 **Missing data**

341 The amount of missing data and reasons, such as informative missingness and dropout, will be
342 assessed for each blood test and other covariates. Levels of missing data will guide subsequent
343 approaches to address missing data:

344

- 345 • Cox modelling - missing values for patient characteristics and blood test data will be replaced
346 using a multiple imputation model including all predictors, the outcome, and auxiliary variables.
- 347 • Mixed-effects and joint modelling – due to the computationally intensive nature of mixed-
348 effects and joint modelling and sporadic nature of blood testing over time, multiple imputation
349 may not be possible. Therefore, we will initially model the data as-is and explore methods to
350 account for missing longitudinal data.
- 351 • Machine learning – missing value imputation using a generative adversarial network will be
352 explored and compared with the standard imputation methods in Cox.

353

354 **Patient and Public Involvement**

355 We have set up a patient and public involvement (PPI) group, consisting of eight PPI advisors. The
356 group will convene routinely throughout the project, sharing their experience with screening,
357 symptoms, diagnostic pathways, treatment, outcome, and more. They will input into study
358 dissemination.

359

360 **Discussion**

361 In this paper, we describe the approved CPRD protocol designed to investigate blood test trends in
362 patients attending primary care and the association between blood test trend(s) and diagnosis of
363 cancer. The overall aim of BLOTTED is to explore whether incorporating blood test trends into the
364 assessment of cancer risk in primary care may offer superior predictive performance to existing
365 approaches to risk stratification. It is hypothesised that historical blood test trend may be used to

366 trigger an urgent referral or direct access cancer investigation before symptoms develop or before
367 individual blood test values reach an “actionable” threshold, as defined in current clinical guidance.
368 We are particularly interested to understand whether the incorporation of blood test trajectory may
369 increase the estimated risk of cancer sufficiently to trigger cancer investigation for patients with
370 symptoms and risk factors that would not currently lead to a recommendation for cancer
371 investigation. If patterns of blood test trend predict cancer in the absence of currently recognised
372 risk factors and symptoms, an additional at-risk cohort may be identified for cancer investigation.
373 The identification and referral of these additional patients may expedite their cancer diagnosis,
374 reducing diagnostic delay. The individual and health system consequences of investigating these
375 patients would merit further prospective and health economic evaluation.

376

377 **Limitations**

378 A major limiting factor for BLOTTED is likely to be computational capacity. Employing advanced
379 statistical modelling and machine learning to analyse longitudinal multivariable health records data
380 from around 40 million patients in CPRD Aurum will lead to significant computational burden. We
381 will work to develop and test our analytical approach on representative limited datasets before
382 deploying the code on the main cohort. Additional technical resources and a secure centralised
383 computing facility are also available to the research team at their host institution, if required to
384 ensure the project is feasible.

385

386 A major consideration when developing an appropriate analytical approach will be to take into
387 account the clinical indication for the available blood test data. Blood tests are ordered in routine
388 clinical practice for many reasons besides cancer investigation. Patients without cancer may have
389 another acute or chronic disease, or undergo treatment, that influence blood levels in a similar way
390 to a malignant process. The frequency of testing is additionally dependent on a patient’s underlying
391 consulting behaviour, their GPs clinical practice, and test access within the local health system.

392 Overall, the trends observed in the study population will represent a combination of testing driven
393 by attendance, comorbidity, intra- and inter-person physiological variation over time, medication
394 use, consultation and clinical behaviour. We plan to take into account for these factors during
395 multivariable modelling so that false positives (patients determined to be high risk who are not at an
396 increased risk of cancer) and false negatives (no discernible cancer signal despite high cancer risk)
397 may be minimised.

398

399 **Implications for clinical practice**

400 There are currently no diagnostic prediction models that utilise available data by incorporating blood
401 test trends integrated into the electronic health record. Therefore, a vast amount of historical
402 information is missing from approaches to risk stratification in primary care. Alongside the
403 identification of promising discriminative approaches to risk stratification that incorporate blood test
404 trend, we will conduct future research to understand the implementation challenges of integrating
405 them in clinical practice and explore whether there are existing integrated approaches that could be
406 modified to incorporate our findings. If methods require significant processing power, significant
407 technological development will be required to enable timely actionable risk stratification in the
408 clinic. Once available in clinical practice, new research will be required to understand the
409 implementation challenges, not least GP uptake and patient acceptance of this method of patient
410 selection.

411

412 **Abbreviations**

413 Cancer Antigen 125 (CA125)

414 Carbohydrate Antigen 19-9 (CA19-9)

415 Carcinoembryonic Antigen (CEA)

416 Clinical Practice Research Datalink (CPRD)

417 Confidence interval (CI)

- 418 General practitioner (GP)
- 419 Hospital Episode Statistics (HES)
- 420 Index of Multiple Deprivation (IMD)
- 421 Lactate Dehydrogenase (LDH)
- 422 Locally Weighted Scatterplot Smoothing (LOWESS)
- 423 Multidisciplinary Diagnostic Centre (MDC)
- 424 National Cancer Registration and Analysis Service (NCRAS)
- 425 National Health Service (NHS)
- 426 National Institute for Health and Care Excellence (NICE)
- 427 Office of National Statistics (ONS)
- 428 Patient and public involvement (PPI)
- 429 Prostate Specific Antigen (PSA)

430

431 **Additional information**

432 **Ethical approval and consent**

433 CPRD has ethical approval from the Health Research Authority to hold anonymised patient data and
434 to support research using that data. CPRD's approval of data access for individual research projects
435 includes ethics approval and consent for those projects. Ethical approval will therefore be covered
436 for this study by the CPRD.

437

438 **Consent for publication**

439 Not applicable

440

441 **Availability of data and materials**

442 Not applicable

443

444 **Competing interests**

445 The authors declare that they have no competing interests.

446

447 **Funding**

448 BN is the Principal Investigator for this study, funded by a Cancer Research UK Population Research
449 Committee Postdoctoral Fellowship (RCCPDF\100005).

450

451 **Author contributions**

452 PSV and BDN wrote the first draft of the protocol. All authors reviewed and updated draft versions
453 and contributed to revisions. All authors approved the final manuscript.

454

455 **Acknowledgements**

456 The authors would like to thank PPI representatives Bernard Gudgin, Julian Ashton, Clara Martins de
457 Barros, Shannon Draisey, Susan Lynne, Emily Lam, Ian Belloch, and Margaret Ogden for input into
458 protocol development.

459

460 **References**

461 [1] Watson J, Mounce L, Bailey SE, Cooper SL, Hamilton W. Blood markers for cancer. *BMJ* (Clinical
462 researched). 2019;367:l5774.

463

464 [2] Rubin GP, Saunders CL, Abel GA, McPhail S, Lyratzopoulos G, Neal RD. Impact of investigations in
465 general practice on timeliness of referral for patients subsequently diagnosed with cancer: analysis
466 of national primary care audit data. *Br J Cancer*. 2015;112(4):676-87.

467

468 [3] NICE. Suspected cancer: recognition and referral (NG12). National Institute for Health and Care
469 Excellence 2015. Available online at <<https://www.nice.org.uk/guidance/ng12>>. Last accessed 4th
470 March 2022.

471

472 [4] Bull LM, Lunt M, Martin GP, Hyrich K, Sergeant JC, Harnessing repeated measurements of
473 predictor variables for clinical risk prediction: a review of existing methods. *Diagn Progn Res*, 2020.
474 4: p. 9.

475

476 [5] Virdee PS., Patnick P, Watkinson P, Birks J, Holt T. Trends in the full blood count blood test and
477 colorectal cancer detection: a longitudinal, case-control study of UK primary care patient data. *NIHR*
478 *Open Research*, 2022, 2, 32:1-53. DOI: 10.3310/nihropenres.13266.1.

479

480 [6] Virdee PS, Birks J, Holt T. A dynamic prediction model for early detection of colorectal cancer
481 using routine blood test results from primary care. SAPC ASM 2021 - virtual conference. Available
482 online at <[https://sapc.ac.uk/conference/2021/abstract/dynamic-prediction-model-early-detection-](https://sapc.ac.uk/conference/2021/abstract/dynamic-prediction-model-early-detection-of-colorectal-cancer-using-routine)
483 [of-colorectal-cancer-using-routine](https://sapc.ac.uk/conference/2021/abstract/dynamic-prediction-model-early-detection-of-colorectal-cancer-using-routine)>. Last accessed 4th March 2022.

484

485 [7] Nicholson BD, Perera R, Thompson MJ. The elusive diagnosis of cancer: testing times. *Br J Gen*
486 *Pract.* 2018;68(676):510-1.

487

488 [8] Koshiaris C, Van den Bruel A, Oke JL, Nicholson BD, Shephard E, Braddick M, et al. Early detection
489 of multiple myeloma in primary care using blood tests: a case-control study in primary care. *Br J Gen*
490 *Pract.* 2018;68(674):e586-e93.

491

492 [9] Naeser E, Moller H, Fredberg U, Frystyk J, Vedsted P. Routine blood tests and probability of
493 cancer in patients referred with non-specific serious symptoms: a cohort study. *BMC Cancer*.
494 2017;17(1):817.

495

496 [10] Hippisley-Cox J, Coupland C. Symptoms and risk factors to identify women with suspected
497 cancer in primary care: derivation and validation of an algorithm. *Br J Gen Pract*. 2013;63(606):e11-
498 21.

499

500 [11] Hippisley-Cox J, Coupland C. Symptoms and risk factors to identify men with suspected cancer in
501 primary care: derivation and validation of an algorithm. *Br J Gen Pract*. 2013 Jan;63(606):e1-10.

502

503 [12] Nicholson BD, Aveyard P, Koshiaris C, Perera R, Hamilton W, et al. Combining simple blood tests
504 to identify primary care patients with unexpected weight loss for cancer investigation: Clinical risk
505 score development, internal validation, and net benefit analysis. *PLOS Medicine*.
506 2021;18(8):e1003728.

507

508 [13] Clinical Practice Research Datalink (CPRD). 2022 [Accessed 25 August 2022]; Available from
509 <https://www.cprd.com/>.

510

511 [14] Virdee PS, Fuller A, Jacobs M, Holt T, Birks J, Assessing data quality from the Clinical Practice
512 Research Datalink: a methodological approach applied to the full blood count blood test. *J big Data*,
513 2020. 7, 95:1-18. DOI: 10.1186/s40537-020-00375-w.

514

515 [15] Riley, R.D., Ensor, J., Snell, K.I.E., Harrell, F.E., Jr., Martin, G.P., Reitsma, J.B., Moons, K.G.M.,
516 Collins, G., and van Smeden, M., Calculating the sample size required for developing a clinical
517 prediction model. *BMJ*. 2020;368:p. m441

518

519 [16] Birks J, Bankhead C, Holt T, Fuller A, Patnick J. Evaluation of a prediction model for colorectal
520 cancer: retrospective analysis of 2.5 million patient records. *Cancer Med*. 2017 Oct;6(10):2453-2460.

521

522 [17] Holt T, Birks J, Bankhead C, Nicholson BD, Fuller A, Patnick J. Do changes in full blood count
523 indices predate symptom reporting in people with undiagnosed bowel cancer? Retrospective
524 analysis using cohort and case control designs. SAPC ASM 2021 - virtual conference. Available online
525 at [https://sapc.ac.uk/conference/2021/abstract/do-changes-full-blood-count-indices-predate-](https://sapc.ac.uk/conference/2021/abstract/do-changes-full-blood-count-indices-predate-symptom-reporting-people)
526 [symptom-reporting-people](https://sapc.ac.uk/conference/2021/abstract/do-changes-full-blood-count-indices-predate-symptom-reporting-people)>. Last accessed 4th March 2022.

527