



King's Research Portal

DOI:

[10.1093/ehjqcco/qcv005](https://doi.org/10.1093/ehjqcco/qcv005)

Document Version

Peer reviewed version

[Link to publication record in King's Research Portal](#)

Citation for published version (APA):

Denaxas, S. C., & Morley, K. I. (2015). Big biomedical data and cardiovascular disease research: opportunities and challenges. *European Heart Journal - Quality of Care and Clinical Outcomes*, 1(1), 9-16.
<https://doi.org/10.1093/ehjqcco/qcv005>

Citing this paper

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

General rights

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

Take down policy

If you believe that this document breaches copyright please contact librarypure@kcl.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.

1 **Big biomedical data and cardiovascular disease research:**
2 **opportunities and challenges**

3

4 **Spiros C. Denaxas (1,2,*), Katherine I. Morley (1,3)**

5

6 1. Farr Institute of Health Informatics Research, University College London, United
7 Kingdom

8 2. Institute of Health Informatics, University College London, United Kingdom

9 3. National Addiction Centre, Institute of Psychiatry, Psychology, and
10 Neuroscience, King's College London, United Kingdom.

11

12 * Corresponding author:

13 Farr Institute of Health Informatics Research, University College London

14 222 Euston Road, London, NW1 2DA

15 Email: s.denaxas@ucl.ac.uk

16

17

18 Over few past years, “*big data*” has become a frequently used catchall phrase
19 for research approaches involving the use of complex, large-scale data sets ^{1,2}. There
20 are many types of data that may fit this description, but within the sphere of
21 clinically-oriented research this term is often considered synonymous to Electronic
22 Health Record (EHR) data, or Electronic Medical Record (EMR) data. The powerful
23 potential of these data for advancing biomedical research has been recognised by
24 many ³⁻⁵. Funders in both the USA and the UK have recently made substantial
25 investments in the area, such as the USD\$215 million Precision Medicine Initiative ⁶
26 announced by the US government and Genomics England, aiming to sequence
27 100,000 whole genomes during routine clinical care⁷. Additionally, funding
28 organizations are actively encouraging research utilizing large-scale biomedical data
29 through specific initiatives. The Big Data To Knowledge (BD2K) programme⁸ was
30 established by the National Institutes of Health (NIH) in 2012 to address the
31 challenges and opportunities presented by big biomedical data through the provision
32 of seed funding for biomedical data-science based research, methods and training
33 material development. In the UK, a consortium of 10 UK government and charity

34 funders, led by the Medical Research Council (MRC) have committed over £90
35 million across several initiatives that are aimed at supporting translational research
36 using big data such as the national Farr Institute of Health Informatics Research⁹, the
37 UK Health Informatics Research Network (UKHIRN) and the Medical
38 Bioinformatics Initiative. The amount of EHR data being digitally generated and
39 collected is vast and rapidly expanding, and presents multiple opportunities that have
40 the potential to transform medical practice and cardiovascular research across all
41 stages of translation.

42

43 However, big data is not a panacea for all research problems, and for many
44 researchers the path from big data to clinical impact for a specific research question is
45 unclear. There are many factors that must be considered when planning to use EHR
46 data for research related not just to ethical and policy issues raised by combining data
47 sources^{10 11} but also the logistical and analytical decisions the process entails. One of
48 the major impediments to the use of EHRs for research is that is the data they contain
49 differs from data collected in a conventional cohort study or randomised controlled
50 trial (RCT) in terms of both why and how it is recorded, and requires substantial
51 processing before it can be statistically analysed. These data are generated and
52 recorded throughout the patient pathway during interactions with primary, secondary
53 and tertiary healthcare providers. Data from specialised disease registries, which
54 were originally set up for auditing clinical standards and benchmarking quality
55 improvement initiatives, may also be incorporated. These different sources also
56 record information in different ways. EHR data can be structured (e.g. diagnosis
57 recorded using medical classification systems such as the International Classification
58 of Diseases 10th revision (ICD-10)¹² or SNOMED-CT¹³) or unstructured (e.g. textual
59 narrative in clinical notes or coronary angiography reports in hospital information
60 systems¹⁴). EHR are also increasingly including cardiovascular imaging data from
61 procedures such as echocardiography, angiography, magnetic resonance imaging or
62 computed tomography¹⁵. For all sources of information, data collection will have been
63 motivated by clinical care, administrative, or other reasons, and will be recorded using
64 a variety of ways. The research-user is faced with substantial missing or incomplete
65 information, data collected at irregular time-points, information that may be
66 temporally inconsistent, and potentially the task of integrating and harmonising
67 information contained in multiple sources.

68 These challenges are not insurmountable and do not mean that EHR data cannot be
69 widely used for research, but do require a clear identification of research areas that
70 can best leverage EHR data, and the development of tools that smooth the path from
71 research question to research result.

72

73 **Research opportunities well-placed to leverage EHR data**

74 *High-resolution observational cohort studies*

75 Linkage of multiple EHR data sources permits the creation of large-scale
76 cohorts of patients for whom extensive follow-up data is already available. This
77 allows researchers to answer questions that reliance upon traditional investigator-led
78 cohort studies would otherwise make impossible due to the scale, diagnostic
79 resolution, timeframe, or cost. In addition, it allows researchers to define and examine
80 the entire patient journey, from early presentations of non-acute manifestations
81 through the various syndrome transitions to cardiac (or non-cardiac) death. This
82 enables them to resolve the time sequence, examine and understand, the aetiological
83 and prognostic differences between different coronary disease phenotypes¹⁸.

84

85 Chung *et al.*¹⁹ were able to take advantage of available EHR data in this way
86 to conduct a comparative effectiveness study of acute coronary care on an
87 international scale. Currently, Sweden and the UK are the only two countries in the
88 world with ongoing, national registries for acute coronary syndrome events that cover
89 all hospital care. Using these data, the authors showed that 30-day mortality
90 following acute myocardial infarction was substantially higher in the UK, and that
91 uptake of effective treatment was slower in the UK. The richness of the data meant
92 that a substantial amount of clinical information could be incorporated into the
93 casemix, including demography, risk factor comorbidity, and pre-hospital treatment.
94 The researchers were also able to determine that diagnoses made in the two countries
95 were comparable by examining troponin values and propensity to make a diagnosis.
96 The results from this study are thus more robust than those based on a simple
97 comparison of mortality rates, or focused on data from bespoke studies undertaken in
98 hospitals that may not be representative of the broader healthcare system.

99

100 EHR cohorts can also be used to make timely contributions to debates of
101 clinical importance, such as the controversy over the relationship between varenicline

102 and adverse cardiovascular events. In 2011 a meta-analysis of 14 RCTs raised
103 concerns that use of varenicline for smoking cessation may increase risk for adverse
104 cardiovascular events (ischemia, arrhythmia, congestive heart failure, sudden death or
105 cardiovascular-related death)²⁰. Three subsequent meta-analyses of RCTs did not find
106 a significant association²¹⁻²³. However, the question remains controversial, partly due
107 to disagreements over analytical methods used in these studies, but also because meta-
108 analyses are limited to the analysis of existing studies^{22 24 25}. Svanström and
109 colleagues were able to rapidly contribute new data to the debate by investigating the
110 question in a cohort made up of the EHR data of over 35,000 Danish individuals who
111 used either varenicline or bupropion for smoking cessation ²⁶. In this observational
112 study, published in 2012, there was no evidence for a higher number of adverse
113 events in patients using varenicline (acute coronary syndrome, ischaemic stroke, and
114 cardiovascular death). It would not have been feasible to take a comparable
115 traditional cohort study from study design to publication within a similar timeframe,
116 especially as very large number of patients would be required to ensure sufficient
117 outcome numbers (only 117 were observed amongst the 35,000 patients in the EHR
118 study).

119

120 The capacity to investigate novel research questions has generally been
121 limited by available data and funding for obtaining new data, but EHR data can
122 potentially be used to address this problem. The relationship between auto-immune
123 inflammatory conditions and atrial fibrillation is one example where EHR data have
124 been able to fill a research niche. Although there is substantial research interest in
125 this area ²⁷, many of the large cardiovascular cohort studies (e.g. Framingham²⁸) have
126 limited data available on inflammatory conditions as this was not part of the original
127 study design. However, researchers in the UK, USA, and Denmark have been able to
128 use EHR resources to explore this research area using very large samples, finding
129 associations with increased risk of AF and a range of conditions including rheumatoid
130 arthritis and psoriasis²⁹⁻³². Other researchers have taken an even broader, non-
131 hypothesis-driven approach, using advanced computation techniques that consider
132 any and all disease information available in EHR data to identify novel associations
133 between diseases ³³. The costs associated with using EHR data for these studies
134 would have been much lower than comparable data-collection, making them a cost-
135 effective entry point into new areas of cardiovascular research.

136

137 *Enhanced clinical trials*

138 There is growing concern that current model of discovering new interventions,
139 evaluating them through clinical trials and implementing the findings as part of
140 clinical care is significantly inefficient. The translation process itself is taking too
141 long, with an average figure of 17 years reported in some cases³⁴. Additionally, the
142 number of new drugs introduced to the market per year has been broadly flat since the
143 1950s yet the costs have steadily grown³⁵ and the cost of bringing a new licensed drug
144 to the market has been estimated between 5 and 11 billion USD³⁶.

145

146 In cardiovascular diseases, the problem is more acutely manifested through
147 problems observed in the current clinical trials pipeline. There is a lack of
148 contemporary and representative population data that can be utilized to draw accurate
149 estimates of events and inform the selection of appropriate primary and secondary
150 endpoints for clinical trial. Clinical trials are often conducted in highly selected
151 populations that are not necessarily representative of the populations presented in
152 routine clinical care and as such, results obtained have limited generalizability and
153 external validity³⁷. For example, the clinical characteristics, treatments and inpatient
154 outcomes of patients enrolled in a large trial of acute heart failure (Acute Study of
155 Nesiritide in Decompensated Heart Failure) were found to be significantly different
156 that those found in a contemporary disease registry³⁸. Furthermore, despite their
157 growing importance in CVD research, non-drug interventions such as interventions
158 based on clinical algorithms and decision support tools are not systematically
159 evaluated through clinical trials since the process of randomization and outcome
160 ascertainment is not seamlessly integrated into the clinical care pathway.

161

162 This has had a significant negative impact on clinical trial conduct and
163 findings. For example, recently there have been several late drug failures occurring
164 within phase III clinical trials of therapeutic agents each costing several hundred
165 million USD\$. High-Density Lipoprotein Cholesterol (HDL-C) raising agents such as
166 niacin, fibrates and cholesteryl ester transfer protein (CETP) failed to reduce all cause
167 mortality, coronary heart disease mortality and myocardial infarction event rates in
168 patients treated with statins³⁹. Likewise, heart rate lowering agents such as ivabradine
169 when introduced to patients with stable coronary artery disease without clinical heart

170 failure failed to improve cardiovascular mortality and non-fatal myocardial infarctions
171 rates ⁴⁰.

172

173 There is growing optimism that EHR can enrich RCT design, delivery and
174 follow up.. EHR data can offer real-world phenotype-rich data that can directly
175 inform trial design, enable the identification of optimal target populations and offer
176 accurate event rate estimates similar to those encountered in clinical care. The entire
177 trial conduct pipeline, from recruitment at the point of care to randomization and
178 adverse event capture can be integrated with routine clinical care enabling the cost-
179 effective and efficient trialling of non-drug interventions. Additionally, EHR can
180 provide richer contemporary data on trial participants at a fraction of the cost thus
181 enabling the generalization of trial results to external populations ⁴¹. For example, the
182 Thrombus Aspiration during ST-Segment Elevation Myocardial Infarction (TASTE)
183 trial⁴² for assessing the clinical effect of routine intracoronary thrombus aspiration
184 before primary percutaneous coronary intervention in patients with ST-segment
185 elevation myocardial infarction recruited patients by enrolling patients through the
186 Swedish Coronary Angiography and Angioplasty Registry and utilizing national EHR
187 and registry data for defining trial endpoints. Finally, EHR data provide valid,
188 complete, long-term follow-up of phase III trials that would otherwise be too costly
189 and complex to establish and too narrow in focus ⁴³. While EHR offer a rich data-
190 scaffolding for designing and implementing clinical trials, significant challenges still
191 exist, mainly around information governance and recruitment of clinicians as outlined
192 in the evaluation by van Staa and colleagues^{44 45}.

193

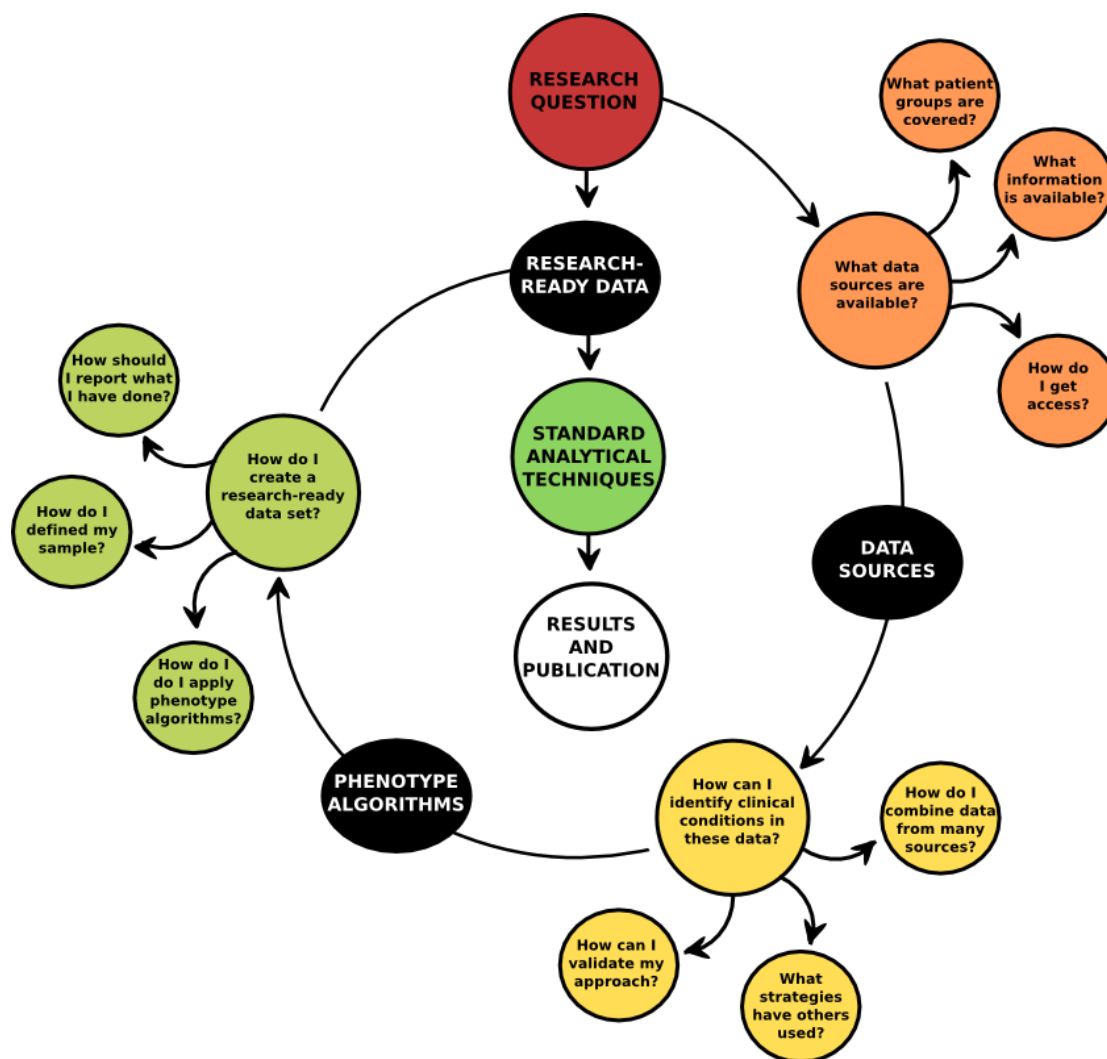
194 **Challenges in the pathway from EHR data to research results**

195

196 Although the benefits of using EHR data for research are potentially large, the
197 widespread use of EHR data is hampered by the fact that there are currently a number
198 of additional steps, and many associated queries, in the pathway from research
199 question to results and publication. As an example, consider a research project using
200 existing data to investigate whether there is a relationship between gender and onset
201 of atrial fibrillation (AF). Most projects would involve applying standard analytical
202 techniques to a bespoke investigator-led cohort of healthy individuals followed-up for
203 cardiovascular conditions including AF (e.g. The Framingham Heart Study). For an

204 existing data set, only relatively minimal data preparation would be required before
 205 analyses could be conducted and data are often provided with detailed documentation.
 206 However, using existing EHR data to answer the same question would require a
 207 number of additional preparatory steps before statistical analyses could be conducted.
 208 Broadly speaking, these relate to: (i) identifying the EHR source(s) that contain the
 209 data needed for the research question; (ii) developing strategies for extracting the
 210 required information from the data source(s), and combining it where necessary; (iii)
 211 creating a data set that is ready for analysis using standard statistical techniques (see
 212 Figure 1).

213



214

215 **Figure 1: Diagram of steps from research question to results and publication. The four central circles show**
 216 **the path from research question to results for a conventional study using existing data. Circles on outside of**
 217 **the spiral indicate the additional steps needed to conduct a research project using electronic health record**
 218 **data.**

219

220 *What EHR data sources are available?*

221 The availability of diverse data sources, including EHR, is rapidly expanding,
222 making the identification of relevant sources for a single project overwhelming.
223 Selecting appropriate data sources for research is dependent upon knowledge of the
224 patients included (e.g. in-patients, ambulatory care, specialist treatment), the types of
225 data recorded (e.g. diagnosis, prescriptions, test results, procedures), and the format of
226 those data (e.g. diagnostic codes, imaging, free text), but often much of this can be
227 difficult to determine in detail until data access has been granted. A recent Wellcome
228 Trust report on the discoverability of EHR and other biomedical datasets for
229 research⁴⁶ found that for the vast majority of sources, no systematic method is used to
230 capture, curate, and display information about the data contained in each source, or to
231 provide guidance on the information governance restrictions attached to them which
232 determine how they can be accessed and used for research. The limited use of
233 standardised methods (e.g. metadata) for describing such information hinders
234 recognition of the limitations and opportunities these data sources present, and
235 potentially results in under-utilisation of data sources due to lack of knowledge about
236 what they contain.

237

238 However, overcoming this challenge is worth the additional effort, as
239 combining data from multiple sources strengthens EHR-based cardiovascular
240 research. For example, Herrett *et al.* explored the completeness of recording for acute
241 myocardial infarction (AMI) in four EHR sources: primary care (Clinical Practice
242 Research Datalink; CPRD), hospital admissions (Hospital Episode Statistics; HES); a
243 MI disease registry (Myocardial Ischaemia National Audit Project; MINAP), and
244 national mortality data (Office of National Statistics; ONS). Compared to the disease
245 registry, which was treated as the gold standard data source, none of the other data
246 sources captured all MI events and consequently incidence rates based on data from a
247 single source were underestimated by 25-50%⁴⁷. This finding is not limited to AMI; a
248 similar investigation of AF diagnoses found that only about 40% of the 72,793 AF
249 patients identified had a diagnosis recorded in both primary and secondary care⁴⁸.

250

251 Thus, for our example research question regarding gender and AF, we would
252 likely decide to combine multiple EHR data sources, such as CPRD, HES, and ONS.
253 This would enable us to use a sample of individuals broadly representative of the UK
254 general population, and would include a more representative set of AF cases as

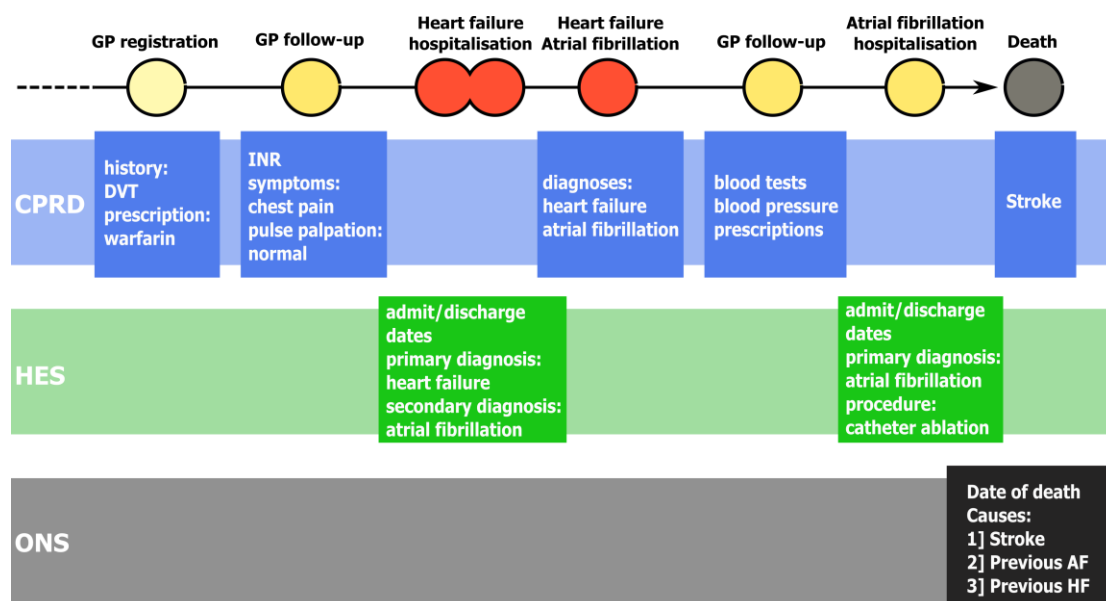
255 diagnoses made in both primary and secondary care would be identified. However,
 256 individual access applications would need to be made for each EHR source prior to
 257 linkage of the different data sources, and information about what is contained within
 258 each would currently be limited to knowledge of the clinical coding systems used.

259

260 ***How can I define clinical conditions in EHR data?***

261 Once relevant the data source(s) have been identified, researchers face another
 262 challenge: how to determine which patients have been diagnosed with a particular
 263 condition. Extracting phenotypic information (i.e. disease status), a process known as
 264 *phenotyping*, is a time-consuming and challenging task even in relation to a single
 265 data source, as multiple diagnosis codes may be used to describe similar or related
 266 conditions and their data. This challenge is amplified when data from multiple
 267 sources, recorded using different coding systems, are combined. Figure 2 illustrates
 268 this, using as an example data for one individual from the three EHR sources in our
 269 hypothetical research question. In this example, an AF diagnosis is recorded at three
 270 different time-points: as a secondary diagnosis during a hospital admission, in the
 271 primary care record after hospital admission information is transferred to their GP,
 272 and as a primary diagnosis when the patient is admitted to hospital for an AF-related
 273 surgical procedure. This information needs to be reconciled in order to determine not
 274 only if, but also when, a diagnosis occurred.

275



276

277 **Figure 2: Illustration of linked primary care data (Clinical Practice Research Datalink; CPRD), secondary**
 278 **care data (Hospital Episode Statistics; HES), and mortality data (Office of National Statistics; ONS) for a**
 279 **single patient. Circles on the top line show events recorded in one or more sources; red circles indicate a**

280 diagnosis. DVT indicates deep vein thrombosis; INR indicates International Normalisation Ratio; AF
281 indicates atrial fibrillation; HF indicates heart failure.

282

283 Reconciling coded information from multiple sources is made more
284 challenging by the different medical classification systems that are used by each
285 source. For example, in the UK, primary care sources use Read codes, a subset of the
286 Systematic Nomenclature of Medicine – Clinical Terms (SNOMED-CT) clinical
287 terminology, whereas secondary care and mortality sources use the International
288 Classification of Disease – 10th Revision (ICD-10). Combining data recorded using
289 these systems for a single condition, such as AF, is not straightforward as the clinical
290 resolution they offer can vary substantially; there are 23 Read codes relating to AF,
291 including disease subtype classification, but only one ICD-10 code. Data-driven
292 computational methodologies, such as support vector machines (SVM), can be
293 applied on unstructured data (e.g. clinical text, electrocardiographic (ECG)
294 monitoring data) to further enhance and fine-tune the accuracy of algorithms utilizing
295 coded data^{4 49 50}. For example, Mohebbi *et al.* created an algorithm which consisted
296 of a linear discriminant analysis based feature reduction scheme and a SVM-based
297 classifier and were able to accurately (sensitivity 99.07%, specificity 100%, positive
298 predictive value 100%) detect AF cases using RR intervals extracted from ECG
299 signals⁵¹.

300

301 No standardised methodologies and mechanisms exist to help research-users
302 define, share and evaluate EHR-derived phenotypes in a consistent way, or to apply
303 algorithms for creating these phenotypes to their own data, although development of
304 tools for this is very active⁵²⁻⁵⁴. The USA-based eMERGE Consortium have
305 developed an AF phenotype algorithm⁵⁵ which focuses on clinical notes and
306 electrocardiogram impression data. These data are not available in CPRD, HES, or
307 ONS, although there is a UK-oriented EHR phenotype resource called CALIBER that
308 does contain an AF phenotype based on coded data from primary and secondary care
309⁴⁸, which could be applied in this situation. However, if no phenotype algorithm
310 existed, we would need to go through the process of developing a new phenotype
311 algorithm for AF, and we would need to repeat this process for every other variable
312 we wanted to include in our final data set such as gender and any covariates such as
313 other cardiovascular diseases, smoking status, or hypertension.

314

315 Validation, preferably against a gold standard, is a key step of defining disease
316 phenotyping algorithms⁵⁶. The goal of the validation exercise is to evaluate the
317 accuracy of the algorithm: is the phenotyping algorithm including all patients that are
318 eligible and excluding all patients that are ineligible, thus accurately allocating them
319 in the case and control groups. Some phenotypes, such as type 2 diabetes⁵⁷, are
320 inherently complex as they make use of multiple data elements (e.g. diagnostic codes,
321 medication information, laboratory measurements, clinical text) and should ideally be
322 validated through manual review of case notes in primary or secondary healthcare
323 providers in order to understand the information the physician had available at the
324 time of diagnosis. Clinical notes however are not available at scale due to information
325 governance restrictions and scaling this process for large cohorts of patients is
326 challenging and time-consuming. An alternative approach is to validate the developed
327 phenotyping algorithms by conducting epidemiological analyses of the association of
328 known risk factors and the phenotype in question and compare associations found in
329 other studies. Other phenotypes, such as white blood cell count, the goal of the
330 validation exercise is to ensure that the algorithm included all eligible patients and
331 discarded outliers and incorrect values.

332

333 ***How do I create a research-ready EHR data set?***

334 The process of applying phenotype algorithms to raw EHR data and creating a
335 data set that is ready to be statistically analysed requires several data transformations
336 that are challenging due to data heterogeneity and complexity. Description of the
337 process is rarely provided as part of academic outputs, and there is increasing
338 recognition of the weaknesses that pervade the current landscape of EHR research in
339 relation to sharing and standardisation of data transformation methods⁵⁸. The
340 prevalent scientific culture does not promote or reward sharing of standardised and re-
341 usable data transformation libraries, which leads to substantial duplication of effort
342 and increases the potential for a lack of reproducible results from EHR-based studies.

343

344 As for a conventional study, an EHR-based study requires a clear definition
345 including the population from which individuals are sampled, inclusion and exclusion
346 criteria, follow-up, and handling of missing data. For our example question, we may
347 need to specify the age range of our patients, whether we are including individuals

348 with prior cardiovascular conditions such as heart failure, and how missing data were
349 handled, but there is additional information that should be reported for EHR data
350 including: the data sources included, the end date of our follow-up data, whether there
351 are exclusion/inclusion criteria based on data quality or other administrative
352 information, details of new phenotype algorithms, and how data were multiply
353 imputed if applicable. While this information can be described to some extent in the
354 Methods section of a scientific paper, the associated computational manipulation and
355 analyses are not standardised for EHR data, and there is currently no provision in
356 scientific papers for detailed explanations of these methods or distribution of
357 associated phenotype algorithms, computer software, or scripts.

358

359 **Recommendations for advancing EHR research**

360

361 Many countries in Europe, and internationally, have EHR systems that could
362 be utilised for research; national, centralised resources that facilitate the steps from
363 research question to research data set would substantially enhance the research
364 potential of these data sources. Initiatives are already underway to achieve this in
365 some countries, but few tackle all aspects of this process.

366

367 The UK-based CALIBER platform ⁵⁹ combines a repository of EHR
368 phenotypes with curated record linkages combining primary care (Clinical Practice
369 Research Datalink), hospital discharge (Hospital Episode Statistics), disease registry
370 (Myocardial Ischaemia National Audit Project⁶⁰) and death registry (Office of
371 National Statistics) data in over 2 million adults with 10 million person years of
372 follow-up. However, this resource does not provide any tools for bidirectional
373 interactions with EHR data sources. In contrast, the Clinical Record Interactive
374 Search (CRIS) system (based at the NIHR Mental Health Biomedical Research Centre
375 and Dementia Unit at the South London and Maudsley NHS Foundation Trust) allows
376 researchers to investigate anonymised secondary care data, including clinical notes
377 and other text, via novel user-friendly tools that facilitate identification of patients
378 meeting certain criteria and development of text-mining algorithms ⁶¹. The Electronic
379 Medical Records and Genomics (eMERGE) Network ⁵², a US National Human
380 Genome Research Institute-funded consortium, combines a phenotype repository with

381 EHR data from multiple secondary healthcare providers, including imaging and text,
382 linked to genotypic data for all participants.

383

384 National EHR portals could combine the strengths of all these projects by
385 including: (i) a national catalogue of contemporary EHR sources curated using
386 metadata standards; (ii) an interactive thesaurus of EHR-derived phenotype
387 algorithms; (iii) standards-driven tools that will enable researchers to visually create
388 observational and interventional research studies (population, inclusion/exclusion
389 criteria, sources, phenotypes). The national catalogue should support the harvesting
390 and integration of metadata from external sources, and manual curation by researchers
391 within a standardised and reproducible framework, as well as providing guidance on
392 data access and data content. This will allow users to identify data sources that can
393 provide information both within and across disease areas. The EHR phenotype
394 algorithms and data set creation tools need to be implemented in a fashion that
395 supports reuse and modification by other users, as well as appropriate academic credit
396 and/or citation. Creating this type of resource will help to foster an "open source"
397 approach to EHR research in which researchers can collaborate and learn from each
398 other, and this will ultimately produce a greater advance in EHR research than could
399 be achieved by any research group in isolation.

400

401

402

403

404

405 **Electronic Health Records:** Electronic Health Records (EHR) are data generated and
406 recorded during routine clinical care. EHRs are diverse and encompass nationally and
407 regionally available structured and unstructured data from primary care, hospitals,
408 administrative data, and disease, procedure and death registries; increasingly
409 including genomic, imaging and patient sensor data.

410 **Medical ontology:** a structured controlled vocabulary of medical concepts and their
411 semantic relations used to record, store and transmit medical knowledge and patient-
412 related clinical information efficiently

413 **Metadata:** is data that describes aspects around a particular data element. For an EHR
414 source metadata can include information about the manner in which the data get
415 generated and recorded, the medical ontologies used to record information and the
416 methods by which researchers can access the data for research.

417 **Phenotyping:** In the context of EHR, phenotyping is defined as the process of
418 creating algorithms that define an observable trait (physical or biochemical) such as a
419 clinical condition within EHR data.

420 **Box 1.** Definitions

421

422 **References**

- 423 1. Community cleverness required. *Nature* 2008;**455**(7209):1-1.
- 424 2. Challenges and Opportunities. 02/11/ 2011.
425 <http://dx.doi.org/10.1126/science.331.6018.692>.
- 426 3. Weber G, Mandl K, Kohane I. Finding the Missing Link for Big Biomedical Data.
427 *JAMA* 2014.
- 428 4. Jensen P, Jensen L, Brunak S. Mining electronic health records: towards better
429 research applications and clinical care. *Nat Rev Genet* 2012;**13**(6):395-
430 405.
- 431 5. Khoury M, Lam TK, Ioannidis J, et al. Transforming epidemiology for 21st
432 century medicine and public health. *Cancer epidemiology, biomarkers &*
433 *prevention* 2013;**22**(4):508-16.
- 434 6. Collins F, Varmus H. A New Initiative on Precision Medicine. *N Engl J Med*
435 2015;**372**(9):793-95.
- 436 7. 100,000 Genomes Project. <http://www.genomicsengland.co.uk/>.
- 437 8. Margolis R, Derr L, Dunn M, et al. The National Institutes of Health's Big Data
438 to Knowledge (BD2K) initiative: capitalizing on biomedical big data.
439 *Journal of the American Medical Informatics Association* 2014;**21**(6):957-
440 58.
- 441 9. Farr Institute for Health Informatics Research. <http://www.farrinstitute.org>.
- 442 10. Richards N, King J. Big Data Ethics. Social Science Research Network Working
443 Paper Series 2014.
- 444 11. Boyd D, Crawford K. CRITICAL QUESTIONS FOR BIG DATA. *Information,*
445 *Communication & Society* 2012;**15**(5):662-79.
- 446 12. World Health Organization, International Classification of Diseases (ICD).
447 <http://apps.who.int/classifications/icd10/browse/2015/en>.
- 448 13. Stearns MQ, Price C, Spackman KA, et al. SNOMED clinical terms: overview of
449 the development process and project status. *Proceedings of the American*
450 *Medical Informatics Association Symposium* 2001:662-66.
- 451 14. Wang Z, Shah A, Tate R, et al. Extracting Diagnoses and Investigation Results
452 from Unstructured Text in Electronic Health Records by Semi-Supervised
453 Machine Learning. *PLoS ONE* 2012;**7**(1):e30412.
- 454 15. Petersen S, Selvanayagam J, Wiesmann F, et al. Left ventricular non-
455 compaction: insights from cardiovascular magnetic resonance imaging.
456 *Journal of the American College of Cardiology* 2005;**46**(1):101-05.

- 457 16. Gymrek M, McGuire A, Golan D, et al. Identifying personal genomes by
458 surname inference. *Science (New York, NY)* 2013;**339**(6117):321-24.
- 459 17. Sheather J, Brannan S. Patient confidentiality in a time of care.data. *BMJ*
460 2013;**347**:f7042.
- 461 18. Timmis A, Feder G, Hemingway H. Prognosis of stable angina pectoris: why
462 we need larger population studies with higher endpoint resolution. *Heart*
463 (British Cardiac Society) 2007;**93**(7):786-91.
- 464 19. Chung S-C, Gedeberg R, Nicholas O, et al. Acute myocardial infarction: a
465 comparison of short-term survival in national outcome registries in
466 Sweden and the UK. *Lancet* 2014;**383**(9925):1305-12.
- 467 20. Singh S, Loke Y, Spangler J, et al. Risk of serious adverse cardiovascular
468 events associated with varenicline: a systematic review and meta-
469 analysis. *Canadian Medical Association journal* 2011;**183**(12):1359-66.
- 470 21. Mills E, Thorlund K, Eapen S, et al. Cardiovascular events associated with
471 smoking cessation pharmacotherapies: a network meta-analysis.
472 *Circulation* 2014;**129**(1):28-41.
- 473 22. Prochaska J, Hilton J. Risk of cardiovascular serious adverse events
474 associated with varenicline use for tobacco cessation: systematic review
475 and meta-analysis. *BMJ* 2012;**344**.
- 476 23. Ware J, Vetrovec G, Miller A, et al. Cardiovascular safety of varenicline:
477 patient-level meta-analysis of randomized, blinded, placebo-controlled
478 trials. *American journal of therapeutics* 2013;**20**(3):235-46.
- 479 24. Prochaska J, Hilton J. Varenicline's adverse events. Choice of summary
480 statistics: relative and absolute measures. *BMJ* 2013;**346**.
- 481 25. Krebs P, Sherman S. ACP Journal Club: review: varenicline for tobacco
482 cessation does not increase CV serious adverse events. *Annals of internal*
483 *medicine* 2012;**157**(4).
- 484 26. Svanström H, Pasternak B, Hviid A. Use of varenicline for smoking cessation
485 and risk of serious cardiovascular events: nationwide cohort study. *BMJ*
486 (Clinical research ed) 2012;**345**.
- 487 27. Hu Y-F, Chen Y-J, Lin Y-J, et al. Inflammation and the pathogenesis of atrial
488 fibrillation. *Nature reviews Cardiology* 2015;**12**(4):230-43.
- 489 28. Kannel WB, McGee DL. Diabetes and cardiovascular disease. The Framingham
490 study. *JAMA* 1979;**241**(19):2035-38.
- 491 29. Kim S, Liu J, Solomon D. The risk of atrial fibrillation in patients with
492 rheumatoid arthritis. *Annals of the rheumatic diseases* 2014;**73**(6):1091-
493 95.
- 494 30. Lindhardsen J, Ahlehoff O, Gislason GH, et al. Risk of atrial fibrillation and
495 stroke in rheumatoid arthritis: Danish nationwide cohort study. *BMJ*
496 (Clinical research ed) 2012;**344**.
- 497 31. Parisi R, Rutter M, Lunt M, et al. Psoriasis and the Risk of Major
498 Cardiovascular Events: Cohort Study Using the Clinical Practice Research
499 Datalink. *The Journal of investigative dermatology* 2015.
- 500 32. Ahlehoff O, Gislason G, Jørgensen C, et al. Psoriasis and risk of atrial
501 fibrillation and ischaemic stroke: a Danish Nationwide Cohort Study.
502 *European heart journal* 2012;**33**(16):2054-64.
- 503 33. Jensen AB, Moseley P, Oprea T, et al. Temporal disease trajectories condensed
504 from population-wide registry data covering 6.2 million patients. *Nature*
505 *communications* 2014;**5**.

- 506 34. Morris Z, Wooding S, Grant J. The answer is 17 years, what is the question:
507 understanding time lags in translational research. *Journal of the Royal*
508 *Society of Medicine* 2011;**104**(12):510-20.
- 509 35. Scannell J, Blanckley A, Boldon H, et al. Diagnosing the decline in
510 pharmaceutical R&D efficiency. *Nat Rev Drug Discov* 2012;**11**(3):191-
511 200.
- 512 36. The Truly Staggering Cost Of Inventing New Drugs - Forbes.
513 [http://www.forbes.com/sites/matthewherper/2012/02/10/the-truly-staggering-](http://www.forbes.com/sites/matthewherper/2012/02/10/the-truly-staggering-cost-of-inventing-new-drugs/)
514 [cost-of-inventing-new-drugs/](http://www.forbes.com/sites/matthewherper/2012/02/10/the-truly-staggering-cost-of-inventing-new-drugs/).
- 515 37. Stuart E, Cole S, Bradshaw C, et al. The use of propensity scores to assess the
516 generalizability of results from randomized trials. *Journal of the Royal*
517 *Statistical Society: Series A (Statistics in Society)* 2011;**174**(2):369-86.
- 518 38. Ezekowitz J, Hu J, Delgado D, et al. Acute Heart Failure. *Circulation: Heart*
519 *Failure* 2012;**5**(6):735-41.
- 520 39. Keene D, Price C, Shun-Shin M, et al. Effect on cardiovascular risk of high
521 density lipoprotein targeted drug treatments niacin, fibrates, and CETP
522 inhibitors: meta-analysis of randomised controlled trials including 117
523 411 patients. *BMJ* 2014;**349**:g4379.
- 524 40. Fox K, Ford I, Steg P, et al. Ivabradine in Stable Coronary Artery Disease
525 without Clinical Heart Failure. *N Engl J Med* 2014;**371**(12):1091-99.
- 526 41. New J, Bakerly N, Leather D, et al. Obtaining real-world evidence: the Salford
527 Lung Study. *Thorax* 2014:thoraxjnl-2014-205259.
- 528 42. Fröbert O, Lagerqvist B, Olivecrona G, et al. Thrombus Aspiration during ST-
529 Segment Elevation Myocardial Infarction. *N Engl J Med*
530 2013;**369**(17):1587-97.
- 531 43. Ford I, Murray H, Packard C, et al. Long-Term Follow-up of the West of
532 Scotland Coronary Prevention Study. *N Engl J Med* 2007;**357**(15):1477-
533 86.
- 534 44. van Staa T-P, Dyson L, McCann G, et al. The opportunities and challenges of
535 pragmatic point-of-care randomised trials using routinely collected
536 electronic records: evaluations of two exemplar trials. *Health technology*
537 *assessment (Winchester, England)* 2014;**18**(43):1-146.
- 538 45. van Staa T-P, Goldacre B, Gulliford M, et al. Pragmatic randomised trials using
539 routine electronic health records: putting them to the test. *BMJ*
540 2012;**344**:e55.
- 541 46. Castillo T, Arofan G, Moore S, et al. Enhancing discoverability of public health
542 and epidemiology, 2014.
- 543 47. Herrett E, Shah A, Boggon R, et al. Completeness and diagnostic validity of
544 recording acute myocardial infarction events in primary care, hospital
545 care, disease registry, and national mortality records: cohort study. *BMJ*
546 2013;**346**.
- 547 48. Morley K, Wallace J, Denaxas S, et al. Defining disease phenotypes using
548 national linked electronic health records: a case study of atrial fibrillation.
549 *PloS one* 2014;**9**(11).
- 550 49. Pathak J, Kho A, Denny J. Electronic health records-driven phenotyping:
551 challenges, recent advances, and perspectives. *Journal of the American*
552 *Medical Informatics Association : JAMIA* 2013;**20**(e2).

- 553 50. Chen Y, Carroll R, Hinz EM, et al. Applying active learning to high-throughput
554 phenotyping algorithms for electronic health records data. *Journal of the*
555 *American Medical Informatics Association* : JAMIA 2013;**20**(e2).
- 556 51. Detection of atrial fibrillation episodes using SVM. *Engineering in Medicine*
557 *and Biology Society, 2008 EMBS 2008 30th Annual International*
558 *Conference of the IEEE*; 2008. IEEE.
- 559 52. Gottesman O, Kuivaniemi H, Tromp G, et al. The Electronic Medical Records
560 and Genomics (eMERGE) Network: past, present, and future. *Genetics in*
561 *medicine* 2013;**15**(10):761-71.
- 562 53. Kho AN, Pacheco Ja, Peissig PL, et al. Electronic medical records for genetic
563 research: results of the eMERGE consortium. *Science translational*
564 *medicine* 2011;**3**(79):79re1-79re1.
- 565 54. Köhler S, Doelken SC, Mungall CJ, et al. The Human Phenotype Ontology
566 project: linking molecular biology and disease through phenotype data.
567 *Nucleic acids research* 2014;**42**(Database issue):D966-74.
- 568 55. Ritchie MD, Denny JC, Crawford DC, et al. Robust replication of genotype-
569 phenotype associations across multiple diseases in an electronic medical
570 record. *American journal of human genetics* 2010;**86**(4):560-72.
- 571 56. Newton K, Peissig P, Kho A, et al. Validation of electronic medical record-
572 based phenotyping algorithms: results and lessons learned from the
573 eMERGE network. *Journal of the American Medical Informatics*
574 *Association* 2013;**20**(e1):e147-e54.
- 575 57. Shah AD, Langenberg C, Rapsomaniki E, et al. Type 2 diabetes and incidence
576 of cardiovascular diseases: a cohort study in 1.9 million people. *The*
577 *lancet Diabetes & endocrinology* 2015;**3**(2):105-13.
- 578 58. Khoury M, Gwinn M, Ioannidis J. The emergence of translational
579 epidemiology: from scientific discovery to population health impact.
580 *American journal of epidemiology* 2010;**172**(5):517-24.
- 581 59. Denaxas S, George J, Herrett E, et al. Data Resource Profile: Cardiovascular
582 disease research using linked bespoke studies and electronic health
583 records (CALIBER). *International Journal of Epidemiology*
584 2012;**41**(6):1625-38.
- 585 60. Herrett E, Smeeth L, Walker L, et al. The Myocardial Ischaemia National Audit
586 Project (MINAP). *Heart* 2010;**96**(16):1264-67.
- 587 61. Stewart R, Soremekun M, Perera G, et al. The South London and Maudsley
588 NHS Foundation Trust Biomedical Research Centre (SLAM BRC) case
589 register: development and descriptive data. *BMC Psychiatry* 2009;**9**:51.
590