

30 **Abstract**

31 **Background:** Stroke in UK Biobank (UKB) is ascertained via linkages to coded administrative
32 datasets and self-report. We studied the accuracy of these codes using genetic validation.

33 **Methods:** We compiled stroke-specific and broad cerebrovascular disease (CVD) code lists (Read
34 V2/V3, ICD-9/-10) for medical settings (hospital, death record, primary care) and self-report.

35 Among 408,210 UKB participants we identified all with a relevant code, creating 12 stroke
36 definitions based on the code type and source. We performed genome-wide association studies
37 (GWASs) for each definition, comparing summary results against the largest published stroke
38 GWAS (MEGASTROKE), assessing genetic correlations, and replicating 32 stroke-associated loci.

39 **Results:** Stroke case numbers identified varied widely from 3,976 (primary care stroke-specific
40 codes) to 19,449 (all codes, all sources). All 12 UKB stroke definitions were significantly correlated
41 with the MEGASTROKE summary GWAS results (r_g 0.81-1) and each other (r_g 0.4-1). However,
42 Bonferroni-corrected confidence intervals were wide, suggesting limited precision of some results.
43 Six previously reported stroke-associated loci were replicated using ≥ 1 UKB stroke definitions.

44 **Conclusions:** Stroke case numbers in UKB depend on the code source and type used, with a 5-fold
45 difference in the maximum case-sample size. All stroke definitions are significantly genetically
46 correlated with the largest stroke GWAS to date.

47

48

49

50

51

52

53 **Introduction**

54 UK Biobank (UKB) is a prospective population-based cohort study with extensive phenotype and
55 genotype information on >500,000 participants from England, Scotland and Wales
56 (www.ukbiobank.ac.uk). It is an open access resource, established to facilitate research into the
57 determinants of a wide range of health outcomes, particularly those relevant in middle and older
58 age (Sudlow, 2015). An example of such a disease is stroke, the second most common cause of
59 death worldwide and a major global cause of disability (Lozano, 2012).

60 Disease outcomes in UKB are ascertained chiefly via linkages to routinely collected, coded,
61 national administrative health datasets. In addition, data on self-reported medical conditions was
62 collected at recruitment. However, to use these data appropriately, researchers need to select
63 which particular disease codes to use for their study and have an understanding of their accuracy.
64 For example, to identify stroke cases, existing codes can be divided into those that are stroke-
65 specific and those that fall under the broad cerebrovascular disease (CVD) category. Stroke-
66 specific codes are used to code acute stroke events where the clinician is confident about the
67 diagnosis and can usually assign a subtype. In contrast, broad CVD codes also capture cases with:
68 (i) phenotypes which pose a high risk for a subsequent stroke (e.g. a code for transient ischaemic
69 attack, an unruptured aneurysm, or carotid artery stenosis); (ii) a past history of stroke with
70 residual symptoms (e.g. a code for sequelae of cerebral infarction); (iii) events where there may be
71 some diagnostic uncertainty (e.g. a code for unspecified cerebrovascular disease); and (iv)
72 intracranial haemorrhages other than intracerebral or subarachnoid haemorrhage (e.g., extradural
73 or subdural haemorrhages, which most clinicians consider different from stroke). Including codes
74 from the broad CVD category will therefore significantly increase the overall number of cases
75 identified, but while this is likely to include at least some misclassified true acute stroke cases,
76 non-stroke cases will also be included.

77 In a systematic review of studies validating stroke code accuracy from case-note review, the
78 overall positive predictive value (proportion of true-positive cases among all identified cases) for
79 identifying acute stroke cases was consistently >70% for stroke-specific codes, dropping to <50% in
80 many studies when broad CVD codes were included (Woodfield, 2015 & Rannikmäe, 2020). For
81 self-reported stroke events, the positive predictive value ranged from 22% to 87% across different
82 studies, making it hard to draw firm conclusions (Woodfield, 2015). While case-note review for
83 code validation is often considered a gold-standard, this method also has its limitations. It is time-
84 consuming and labour-intensive, so can only be achieved in relatively small numbers of cases, with
85 limited precision of the results. In addition, the results rely on: (i) accessing the complete relevant
86 medical record; (ii) the detail and quality of the medical record; (iii) the qualification of the person
87 reviewing the notes; (iv) the inter-adjudicator agreement, which we know is not perfect even
88 between highly specialised clinicians; and (v) the consistency of results across different healthcare
89 settings/providers (Lieberman, 2018).

90 We set out to supplement current knowledge about the accuracy of stroke codes with a method
91 making use of large-scale genetic data, which we refer to as 'genetic validation'. The fundamental
92 idea is to use existing knowledge of genetic associations with a disease (in this case acute stroke),
93 to assess how well various potential code lists capture people who truly have this disease, which in
94 turn could be used to harmonise disease definitions across cohorts and health systems (Manolio,
95 2020). If the code list captures true-positive cases, we would expect the genetic associations that
96 result from stroke cases identified through coded data to closely mirror the genetic association
97 results from previous studies of stroke.

98

99

100

101 **Methods**

102 Study setting

103 We included all 408,210 UKB white British ancestry participants in this study. We restricted our
104 analyses to this ancestry subgroup because it covers 94% of the UKB participants and allowed us
105 to achieve a good balance between attaining sufficient case numbers while reducing population
106 stratification and analytic complexity. As part of the UK Biobank recruitment process, informed
107 consent was obtained from all individual participants included in the study. At the time of the
108 study, UKB had linked hospital admissions and death registry administrative coded data available
109 for all participants, and primary care administrative coded data for 47% of the cohort (191,146),
110 covering the time period up to March and September 2019, respectively (S1 Table). In addition, all
111 participants self-reported pre-existing health conditions during an interview at recruitment. The
112 subset of the cohort with primary care data available was similar to the whole cohort with respect
113 to age at recruitment, sex and Townsend deprivation index (S2 Table).

114 Identifying stroke cases in UKB

115 We compiled stroke-specific and broad cerebrovascular disease (CVD) code lists for each medical
116 setting (hospital admission, death record, primary care) and self-report. This process was informed
117 by previously published codes where available (Woodfield, 2015 & Rannikmäe, 2020),
118 supplemented by the selection of additional codes by expert clinicians (authors KR, CLMS, ED, RW)
119 on discussion and mutual agreement. This resulted in a total of 8 code lists, covering the ICD-
120 9/ICD-10, Read Version 2, Clinical Terms Version 3 (Read Version 3) and UKB self-report illness
121 coding systems (S3 Table).

122 Next, we identified all participants with a relevant code from any of the code lists and created 12
123 different ways of defining stroke cases in UKB based on the code type (stroke-specific, broad CVD)
124 and source (hospital admission, death record, primary care, self-report). This resulted in 12

125 partially overlapping case-control groups, where cases were all the individuals with a stroke code
126 for the particular stroke definition, and all the remaining participants acted as controls. A specific
127 UKB participant could therefore be a stroke case for one definition and a control for another
128 definition.

129 Genome-Wide Association Studies

130 We performed 12 genome-wide association studies (GWASs), one for each case-control set (i.e. for
131 each definition of stroke cases and their controls). We applied a linear mixed model method using
132 the BoltLMM software package (v2.3.4) software (Loh, 2015). We included the following as
133 covariates: genotyping array, UKB assessment centre, sex, age at recruitment, and principal
134 components one to ten. We filtered the results for single nucleotide polymorphisms (SNPs) with
135 an imputation quality INFO score ≥ 0.9 and minor allele frequency $\geq 1\%$. After filtering the results
136 for SNP imputation quality and minor allele frequency, we included 9,524,428 SNPs. We converted
137 the linear mixed model effects to odds ratios using R code provided in
138 <https://shiny.cnsgenomics.com/LMOR/> (Lloyd-Jones, 2018).

139 Analyses of GWAS results

140 We compared summary results from our 12 GWASs against the largest published stroke GWAS
141 meta-analysis project - the MEGASTROKE study. The MEGASTROKE study is a meta-analysis of 29
142 stroke GWASs (17 including individuals of European ancestry) and does not include UKB data.
143 Almost all studies included in MEGASTROKE (covering >95% included cases) required the stroke
144 diagnosis to be confirmed by a medical professional or required evidence of stroke from >1
145 source, even if the initial case ascertainment included using administrative codes (Malik, 2018). All
146 analyses were done using R software version 3.6.2.

147 *Genetic correlation with the MEGASTROKE study results*

148 We applied a high-definition likelihood method using the HDL software (Ning, 2020) to assess the
149 genetic correlation between our GWAS results using the 12 stroke definitions, and the
150 MEGASTROKE study GWAS summary results for any stroke subtype in European samples. Genetic
151 correlation (r_g) is the proportion of variance that two stroke definitions share due to genetic
152 causes. A genetic correlation of 0 implies that the genetic effects on one definition are
153 independent of the other, while a correlation of 1 implies that all of the genetic influences on the
154 two definitions are identical. We assessed if the correlation was significantly different from 0 and
155 from 1, setting the p-value significance threshold to 0.0042 after a Bonferroni correction for the 12
156 tests. We displayed the results (correlation measured as r_g) on a heatmap. We also display
157 Bonferroni corrected confidence intervals to aid interpretation.

158 *Genetic correlation within our study definitions*

159 We then used the HDL software to assess genetic correlations within our study across the 12
160 definitions. We set the significance threshold to 0.0024 after a Bonferroni correction for 7
161 independent non-overlapping case-control definitions (definitions not in bold in Table 1), resulting
162 in 21 correlation tests. We also display Bonferroni corrected confidence intervals to aid
163 interpretation.

164 *Replicating the MEGASTROKE study stroke-significant loci*

165 The MEGASTROKE study identified 32 genetic loci significantly associated with stroke. We
166 identified these loci (the lead SNP for each locus) in our GWAS summary results and considered a
167 locus to be replicated (i.e. also significantly associated with the respective stroke definition in our
168 data) if the p-value of association in our GWAS was <0.00156 (Bonferroni corrected for 32 loci).
169 We compared the number of replicated loci across our summary definitions. We compared the
170 effect sizes of the associations between MEGASTROKE trans-ethnic and European ancestry GWASs
171 and our GWAS summary results. Where the lead SNP was not available in our data, we identified

172 SNPs in moderate LD ($r^2 > 0.7$ in the 1000 Genomes GBR population using the Ensembl LD
173 calculator https://www.ensembl.org/Homo_sapiens/Tools/LD) with the lead SNP, and if any SNPs
174 in LD available in our data were identified, we examined their associations instead. We displayed
175 results for five of our summary definitions of stroke cases and their controls: stroke-specific code
176 from any medical setting; broad CVD code from any medical setting; stroke-specific or broad CVD
177 code from any medical setting; specific or non-specific self-reported stroke event; any code or self-
178 reported event. We highlighted significantly associated (i.e. replicated) loci.
179 We also calculated our expected power to replicate the 32 loci using the Genetic Association Study
180 (GAS) Power Calculator (http://csg.sph.umich.edu/abecasis/gas_power_calculator/index.html),
181 assuming a stroke prevalence of 2.26% and inputting the disease allele frequency and genotype
182 relative risk estimates from the MEGASTROKE publication Table 1 (Malik, 2018).

183 **Results**

184 Stroke cases in UKB

185 The number of relevant codes for identifying stroke cases varied widely depending on the coding
186 system (ICD versus Read versus self-report) and code type (stroke-specific versus broad CVD code)
187 – from less than five codes for a specific self-reported stroke event, to >500 codes when including
188 all possible codes across all coding systems. The stroke-specific and broad cerebrovascular disease
189 (CVD) code lists for each medical setting and self-report are shown in S3 Table.

190 The number of stroke cases identified among the 408,210 participants also varied widely
191 depending on the code type and source used – from 3,976 cases in primary care when using
192 stroke-specific codes, to 19,449 cases when including all possible code combinations (stroke-
193 specific and broad CVD) across all sources (hospital admission, death record, primary care, self-
194 report) (Table1).

195 The code source for cases with a stroke-specific code was: self-report only in 27%, primary care
196 only for 9%, hospital/death record code only for 29%, and >1 source for 35% (S1 Figure). The code
197 source for cases with either a stroke-specific or a broad CVD code was: self-report only in 14%,
198 primary care only for 15%, hospital/death record code only for 34%, and >1 source for 37% (S1
199 Figure). These proportions are calculated based on the primary care data being currently available
200 only for ~50% of the participants, and so will change when primary care data for the whole cohort
201 become available.

202 The overall proportion of prevalent codes (i.e. first code predates participant's recruitment to
203 UKB) versus incident codes (i.e. first code date occurs after participant's recruitment to UKB) was
204 the same for stroke-specific and broad CVD categories: 38% prevalent versus 62% incident codes.
205 These proportions are dependent on the updates to different linked health datasets and the
206 proportion of incident codes will continue to increase with increasing duration of follow up. (S1
207 Table).

208 Mean and median age at recruitment were higher among stroke cases (for all stroke definitions)
209 than for the whole cohort of UKB participants (mean age 61 to 62 years vs 57 years; median age
210 62 to 63 years vs 58 years). Mean and median age at the time of stroke (in case of multiple events,
211 age at the earliest event was taken) was higher for coded diagnoses from the medical setting
212 compared to self-reported events (mean age 62 years vs 53 years, median age 63 years vs 55
213 years, respectively). This is to be expected, considering that all self-reported events were recorded
214 at the time of recruitment, whereas medical codes also capture diagnoses after recruitment
215 during follow up. The proportion of women was lower among stroke cases than across all UKB
216 participants (43% for those with any medical setting or self-reported code vs 54% for all UKB). This
217 is to be expected as age-specific incidence rates are substantially lower in women than men in
218 younger and middle-age groups, but these differences narrow down so that in the oldest age

219 groups, incidence rates in women are approximately equal to or even higher than in men (Virani,
220 2020). (Table 1).

221 Analyses of GWAS results

222 *Genetic correlation with the MEGASTROKE study results*

223 All 12 UKB stroke definitions were significantly correlated with the MEGASTROKE summary GWAS
224 results, with genetic correlations (r_g) ranging from 0.81 to 1, and confidence intervals overlapping.
225 The p-values for difference from 1 were not significant, compatible with perfect correlation.
226 However, the Bonferroni corrected CIs were wide, especially for 5 of the 12 tests, where the lower
227 confidence limit was <0.7 , limiting the precision of some of these results (Figure 1, Figure 2, S4
228 Table).

229 *Genetic correlation within our study definitions*

230 The UKB summary definitions in our study were all significantly correlated with each other (all p-
231 values significantly different from 0), with r_g ranging from 0.4 to 1. Again, the Bonferroni corrected
232 CIs were wide, with the lower confidence limit even suggesting the possibility of a negative
233 correlation for two comparisons. (Figure 1, Figure 2, S5 Table).

234 *Replicating the MEGASTROKE study stroke-significant loci*

235 Within our GWASs, 6 of the 32 previously reported stroke-associated loci were replicated by one
236 or more definitions. Analyses using stroke-specific codes and analyses using any code or self-
237 reported event both replicated the biggest number of known stroke loci (5 of 32). The power from
238 additional cases for the latter category did not result in replicating more loci than stroke-specific
239 codes alone. However, for three of the five replicated loci, the p-values were smaller in the larger
240 dataset (analyses using any code or self-reported event) suggesting a more robust replication
241 when using the broadest definition of stroke in UKB. Within our data, effect sizes (expressed as

242 odds ratios) were similar across the stroke definitions, with overlapping Bonferroni corrected
243 confidence intervals.

244 For two of the six replicated loci (*PITX2* and *HDAC9–TWIST1*), the effect size of the association
245 (odds ratio) was bigger in the MEGASTROKE dataset than in our data (across all five summary
246 stroke definitions). These two loci are known to be associated with particular stroke subtypes –
247 *PITX2* with cardioembolic and *HDAC9–TWIST1* with large artery stroke [10]. (Figure 3, Table 2, S6
248 Table).

249 Power calculations suggested we had $\geq 80\%$ power to replicate all 32 loci for the any code or self-
250 reported event definition, while having $\geq 80\%$ power for only 11/32 loci for the definition
251 including a stroke-specific code from any medical setting. (S7 Table).

252 **Discussion**

253 Our analyses show, that depending on the code source and type used for identifying stroke cases
254 in the UKB, the currently achieved maximum case-sample size can range from ~4,000 to ~20,000 –
255 a remarkable 5-fold difference. We go on to demonstrate, that regardless of the code source and
256 type used, the resulting GWAS summary results are significantly genetically correlated with the
257 largest stroke GWAS to date. Finally, when we try to replicate known stroke-significant loci in our
258 data, both stroke-specific codes from any medical setting as well as a broad definition including
259 any code or self-reported event, replicate 5 of the 32 loci. Replication generated broadly similar
260 effect sizes for all but 2 stroke subtype specific loci, which is likely explained by our dataset
261 including a mix of stroke subtypes. Another possible explanation is the “winner's curse”
262 phenomenon (i.e. the estimated effect of a marker allele from the initial study reporting the
263 marker-allele association is often exaggerated relative to the estimated effect in follow-up
264 studies).

265 The correlation of all definitions with the MEGASTROKE study results suggests one or more of the
266 following: (i) all definitions retrieve true-positive acute stroke cases, meaning that broad CVD
267 codes include additional true-positive cases not identified by stroke-specific codes; (ii) cases coded
268 with a broad CVD code have not necessarily suffered an acute stroke, but represent a range of
269 phenotypes with a similar genetic architecture to acute stroke (e.g. previous research has shown
270 at least one overlapping locus for carotid artery disease and acute stroke (Malik, 2018)); (iii) the
271 MEGASTROKE study includes some misclassified broad CVD cases as false-positive acute stroke
272 cases. It is most likely that a combination of these factors is contributing to our findings, but we
273 are unable to dissect their separate contributions in the current study.

274 Previous case-note validation studies suggested that broad CVD codes are better at identifying the
275 broad conditions they signify as opposed to ascertaining acute stroke cases (McCormick, 2015),

276 supporting a role for option two above. An example of this would be a case-note review of
277 patients with a code for an unruptured intracranial aneurysm or carotid artery stenosis confirming
278 that the diagnosis was also most likely an unruptured intracranial aneurysm or carotid artery
279 stenosis, rather than rather than the reviewing clinician deciding it was an acute stroke that had
280 been miscoded as an unruptured intracranial aneurysm or carotid artery stenosis.

281 Despite the definition using any code or self-reported event increasing the sample size by more
282 than two-fold compared to the definition using only stroke-specific codes from a medical setting, it
283 did not replicate a higher number of known stroke-associated loci. This could suggest that there is
284 still insufficient power to replicate additional loci using any of our definitions despite power
285 calculations suggesting $\geq 80\%$ power for all loci. Also, associations for nine of the 32 loci in
286 MEGASTROKE were only found for specific stroke subtypes and 11 of the 32 loci were significant in
287 analyses including only ischaemic stroke cases, the proportions of which are unlikely to be
288 identical between the two datasets. For example, the MEGASTROKE study sample included 90%
289 confirmed ischaemic stroke cases. The stroke subtype breakdown among the UKB participants is
290 available only for stroke-specific codes from hospital, death record and self-reported data (UK
291 Biobank datafields '42009', '42011' and '42013') and shows a proportion of confirmed ischaemic
292 stroke cases of 47%, with 10% cases being intracerebral haemorrhage and 11% subarachnoid
293 haemorrhage and the remainder of unspecified stroke subtype. Alternatively, it could also suggest
294 that the additional cases identified by using any code or self-report are not true-positive stroke
295 cases or that some of these known stroke-associated loci are false-positive findings.

296 Self-reported cases (both stroke-specific and broad CVD) also showed a close genetic correlation
297 with the MEGASTROKE study, supporting the use of self-report as a means of identifying additional
298 stroke cases in the UKB. This was so despite the highly variable results from previous case-note
299 based validation studies of self-report for ascertaining stroke cases. Studying this by case-note

300 validation in UKB itself would be challenging, given the difficulties accessing NHS records which
301 predate recruitment by many years and the fact that participants may have moved between UK
302 regions during their life-course. Other studies have also reported a close genetic correlation
303 between a wide range of self-reported diseases and medical setting diagnoses. Examples include
304 both acute and chronic conditions (e.g. depression, myocardial infarction, rheumatoid arthritis)
305 (Wray, 2018; Howard, 2019; DeBoever, 2020).

306 For some of our definitions, the r_g was >1 . The estimated r_g is a combination of the true r_g and
307 variation. When the true r_g is close to the boundary (-1 or 1) and/or variation is large, the
308 estimated r_g can go beyond the boundary (Ning, 2020). In r_g estimation, some common reasons
309 for generating large variation are: (i) at least one of the h^2 estimates is very low; (ii) small sample
310 size; (iii) many SNPs in the reference panel are absent in one of the two GWASs; (iv) there is a
311 severe mismatch between the GWAS population and the population for computing reference
312 panel. We can exclude the last two options, and the small sample size is therefore the likely
313 explanation.

314 In our analyses, the case-control groups were partially overlapping and a specific UKB participant
315 could therefore be a stroke case for one definition and a control for another definition. We used
316 this study design to mimic the 'real world' situation, creating binary case-control definitions based
317 on each code list. In theory this could reduce the power of some of the analyses, since it means
318 controls can end up including some true-positive stroke cases. However, in reality it is unlikely to
319 have a significant effect given the overall large number of controls. For example, for analyses using
320 stroke-specific codes from any medical setting, just over 2% of controls have a broad CVD code
321 and/or have self-reported a stroke event.

322 We used BoltLMM for running the GWAS (Loh, 2015). Our case-fraction ranged from 1% to 5%
323 depending on the case-definition and we limited our analyses to SNPs with a minor allele

324 frequency of at least 1%. Based on simulations done using BoltLMM, the authors of the software
325 suggest that with these case-fraction and minor allele frequency parameters, they did not find a
326 statistically significant inflation of type I error rates (Supplementary table 8 in Loh, 2015).

327 The strengths of our study are: (i) we included - and have made available to re-use - a clinically
328 informed, comprehensive set of codes across all relevant coding systems; (ii) we compared our
329 results against the largest stroke GWAS to date; (iii) we used multiple methods for comparison
330 accounting for both GWAS significant loci but also SNPs across the whole genome – i.e. correlation
331 and replication; (iv) we have added novel data to what is already known from case-note validation.

332 Our study also has some limitations: (i) some of our definitions included relatively small case
333 numbers compared to the MEGASTROKE study, reducing our power to replicate known loci; (ii)
334 uneven numbers across definitions not allowing direct comparisons, but rather reflecting the real-
335 world situation; (iii) our definitions included the subarachnoid haemorrhage stroke subtype codes,
336 whereas the MEGASTROKE study did not, resulting in a slightly different mix of stroke cases; (iv)
337 the UKB participants' demographic characteristics differ from those of the UK general population
338 with evidence of a healthy-volunteer selection bias, which needs to be considered when
339 extrapolating these results to other settings (Fry, 2017).

340 We have shown that the selection of codes and code sources used to ascertain stroke cases has a
341 major impact on the overall stroke case numbers in the UKB. Given the close genetic correlation
342 between stroke cases identified using broad CVD codes, self-report, and physician-confirmed
343 stroke cases, we suggest that for studies accepting of more crude stroke and cerebrovascular
344 disease outcomes, researchers may wish to include all codes and self-reported events for
345 increased power. Alternatively, this information is also helpful in informing the selection of
346 controls for various studies. Including a large number of broad CVD coded cases among controls
347 might weaken any association seen for certain study designs. However, since we cannot exclude

348 the effects of a shared genetic control of broad CVD phenotypes and acute stroke, this evidence is
349 not sufficient to support using broad CVD codes in studies that need to define acute stroke
350 outcomes very accurately (e.g. clinical trials).

351 Further research is needed: to better understand the underlying reasons for the close genetic
352 correlation between stroke-specific and broad CVD codes; to dissect the underlying explanation
353 for our results with targeted case-note validation; and to replicate our results in other datasets. In
354 addition, more data is needed on the accuracy of different coding systems for identifying specific
355 pathological stroke subtypes (ischaemic stroke versus intracerebral haemorrhage versus
356 subarachnoid haemorrhage) and aetiological stroke subtypes (e.g., small vessel disease versus
357 large artery disease versus cardioembolic stroke versus other / unknown cause).

358 **Acknowledgements:** This work was undertaken under a UKB project 2532 “UK Biobank Stroke
359 Study (UKBISS): developing an in-depth understanding of the determinants of stroke and its
360 subtypes.” Authors acknowledge Dr Spiros Denaxas, UCL, for his valuable comments on the
361 manuscript. The MEGASTROKE project received funding from sources specified at
362 <http://www.megastroke.org/acknowledgments.html> and the author list for the MEGASTROKE
363 consortium is added in Appendix 1.

364 **Funding:** K Rannikmae is funded by Health Data Research UK Rutherford fellowship
365 MR/S004130/1. AF is funded by BHF award RE/18/5/34216. AT is funded by HDR-UK awards HDR-
366 9004 and HDR-9003. The funders had no role in study design, data collection and analysis, decision
367 to publish, or preparation of the manuscript

368 **Disclosures:** Authors report no competing interests.

369 **Author contributions:** All authors made substantial contributions to the conception or design of
370 the work; or the acquisition, analysis, or interpretation of data for the work; AND drafting the
371 work or revising it critically for important intellectual content; AND final approval of the version to
372 be published; AND agreement to be accountable for all aspects of the work in ensuring that
373 questions related to the accuracy or integrity of any part of the work are appropriately
374 investigated and resolved.

375 **Data sharing:** The data that supports the findings of this study are available in the supplementary
376 material of this article.

377

378

379

380

381 **References:**

- 382 1. Sudlow, C. L. M., Gallacher, J., Allen, N., Beral, V., Burton, B., Danesh, J. et al. (2015). UK
383 Biobank: an open access resource for identifying the causes of a wide range of complex
384 diseases of middle and old age. *PLoS Medicine*. 31;12(3):e1001779.
- 385 2. Lozano, R., Naghavi, M., Foreman, K., Lim, S., Shibuya, K., Aboyans, V. et al. (2010). Global
386 and regional mortality from 235 causes of death for 20 age groups in 1990 and 2010: a
387 systematic analysis for the Global Burden of Disease Study 2010. *Lancet*. 380(9859):2095-
388 2128.
- 389 3. Woodfield, R., Grant, I., UK Biobank Stroke Outcomes group, UK Biobank Follow-Up and
390 Outcomes Working Group, Sudlow, C. L. M. (2015). Accuracy of Electronic Health Record
391 Data for Identifying Stroke Cases in Large-Scale Epidemiological Studies: A Systematic
392 Review from the UK Biobank Stroke Outcomes Group. *PLoS One*. 10(10):e0140533.
- 393 4. Rannikmäe, K., Ngoh, K., Bush, K., Al-Shahi Salman, R., Doubal, F., Flaig, R. et al. (2020).
394 Accuracy of identifying incident stroke cases from linked health care data in UK Biobank.
395 *Neurology*. 95(6):e697-e707.
- 396 5. Woodfield, R., UK Biobank Stroke Outcomes Group, UK Biobank Follow-up and Outcomes
397 Working Group, Sudlow, C. L. M. (2015). Accuracy of Patient Self-Report of Stroke: A
398 Systematic Review from the UK Biobank Stroke Outcomes Group. *Plos One*.
399 10(9):e0137538.
- 400 6. Liberman, A. L., Rostanski, S. K., Ruff, I. M., Meyer, A. N. D., Maas, M. B., Prabhakaran, S.
401 (2018). Inter-rater Agreement for the Diagnosis of Stroke Versus Stroke Mimic. *Neurologist*.
402 23(4):118-121.

- 403 7. Manolio, T.A., Goodhand, P., Ginsburg, G. (2020). The International Hundred Thousand Plus
404 Cohort Consortium: integrating large-scale cohorts to address global scientific challenges.
405 *Lancet Digital Health*. 2(11):e567-e568.
- 406 8. Loh, P-R., Tucker, G., Bulik-Sullivan, B. K., Vilhjálmsson, B. J., Finucane, H. K., Salem, R. M. et
407 al. (2015). Efficient Bayesian mixed-model analysis increases association power in large
408 cohorts. *Nature Genetics*. 47(3):284-290.
- 409 9. Lloyd-Jones, L. R., Robinson, M.R., Yang, J., Visscher, P. M. (2018). Transformation of
410 Summary Statistics from Linear Mixed Model Association on All-or-None Traits to Odds
411 Ratio. *Genetics*. 208(4):1397-1408.
- 412 10. Malik, R., Chauhan, G., Traylor, M., Sargurupremraj, M., Okada, Y., Mishra, A. et al. (2018).
413 Multiancestry genome-wide association study of 520,000 subjects identifies 32 loci
414 associated with stroke and stroke subtypes. *Nature Genetics*. 50(4):524-537.
- 415 11. Ning, Z., Pawitan, Y., Shen, X. (2020). High-definition likelihood inference of genetic
416 correlations across human complex traits. *Nature Genetics*. 52:859–864.
- 417 12. Virani, S. S., Alonso, A., Benjamin, E. J., Bittencourt, M. S., Callaway, C. W., Carson, A.P. et
418 al. (2020). Heart Disease and Stroke Statistics - 2020 Update: A Report From the American
419 Heart Association. *Circulation*. 141(9):e139.
- 420 13. McCormick, N., Bhole, V., Lacaille, D., Avina-Zubieta, J. A. (2015). Validity of Diagnostic
421 Codes for Acute Stroke in Administrative Databases: A Systematic Review. *PLoS One*.
422 10(8):e0135834.
- 423 14. Wray, N. R., Ripke, S., Mattheien, M., Trzaskowski, M., Byrne, E. M., Abdellaoui, A. et al.
424 (2018). Genome-wide association analyses identify 44 risk variants and refine the genetic
425 architecture of major depression. *Nature Genetics*. 50:668–681.

- 426 15. Howard, D. M., Adams, M. J., Clarke, T. K., Hafferty, J. D., Gibson, J., Shiralil, M. et al. (2019).
427 Genome-wide meta-analysis of depression identifies 102 independent variants and
428 highlights the importance of the prefrontal brain regions. *Nature Neuroscience*. 22:343–
429 352.
- 430 16. DeBoever, C., Tanigawa, Y., Aguirre, M., McInnes, G., Lavertu, A., Rivas, M. A. (2020).
431 Assessing Digital Phenotyping to Enhance Genetic Studies of Human Diseases. *American*
432 *Journal of Human Genetics*. 106(5):611-622.
- 433 17. Fry, A., Littlejohns, T. J., Sudlow, C. L. M., Doherty, N., Adamska, L., Sprosen, T. et al. (2017).
434 Comparison of Sociodemographic and Health-Related Characteristics of UK Biobank
435 Participants With Those of the General Population. *American Journal of Epidemiology*.
436 186(9):1026-1034.
- 437
438
439
440
441
442
443
444
445
446
447

Table 1. Number and demographic characteristics of stroke cases identified in UKB

STROKE DEFINITION	NUMBER OF CASES	NUMBER OF CONTROLS	Mean (median) age at recruitment (years)	Mean (median) age; age range at stroke (years)	Sex (% female)
Stroke-specific code from hospital/death records	6,887	401,323	61 (63)	63 (64); 31 to 79	40
Stroke-specific code from primary care	3,976	404,234	61 (63)	59 (61); 1 to 79	41
Stroke-specific code from any medical setting	8,665	399,545	61 (63)	61 (63); 1 to 79	41
Broad CVD code from hospital/death records	5,725	402,485	62 (63)	64 (65); 31 to 79	43
Broad CVD code from primary care	4003	404,207	62 (63)	62 (63); 1 to 79	41
Broad CVD code from any medical setting	8,085	400,125	62 (63)	63 (64); 1 to 79	44
Stroke-specific or broad CVD code from hospital/death records	12,612	395,598	61 (63)	63 (64); 31 to 79	41
Stroke-specific or broad CVD code from primary care	7,979	400,231	62 (63)	60 (62); 1 to 79	41
Stroke-specific or broad CVD code from any medical setting	16,750	391,460	62 (63)	62 (63); 1 to 79	42
Specific self-reported stroke event	5,915	402,295	61 (62)	53 (55); 0 to 70	41
Specific or non-specific self-reported stroke event	7,536	400,674	61 (63)	53 (55); 0 to 70	42
Any code or self-reported event	19,449	388,761	61 (63)	60 (61); 0 to 79	43
Across all UKB participants	408,210		57 (58)	Not applicable	54

We considered code and self-reported event age missing if it predated the date of birth, contained a negative value, or was recorded by UKB as unreliable or missing – this affected 0.3% to 2.1% cases depending on the stroke definition. CVD: cerebrovascular disease;

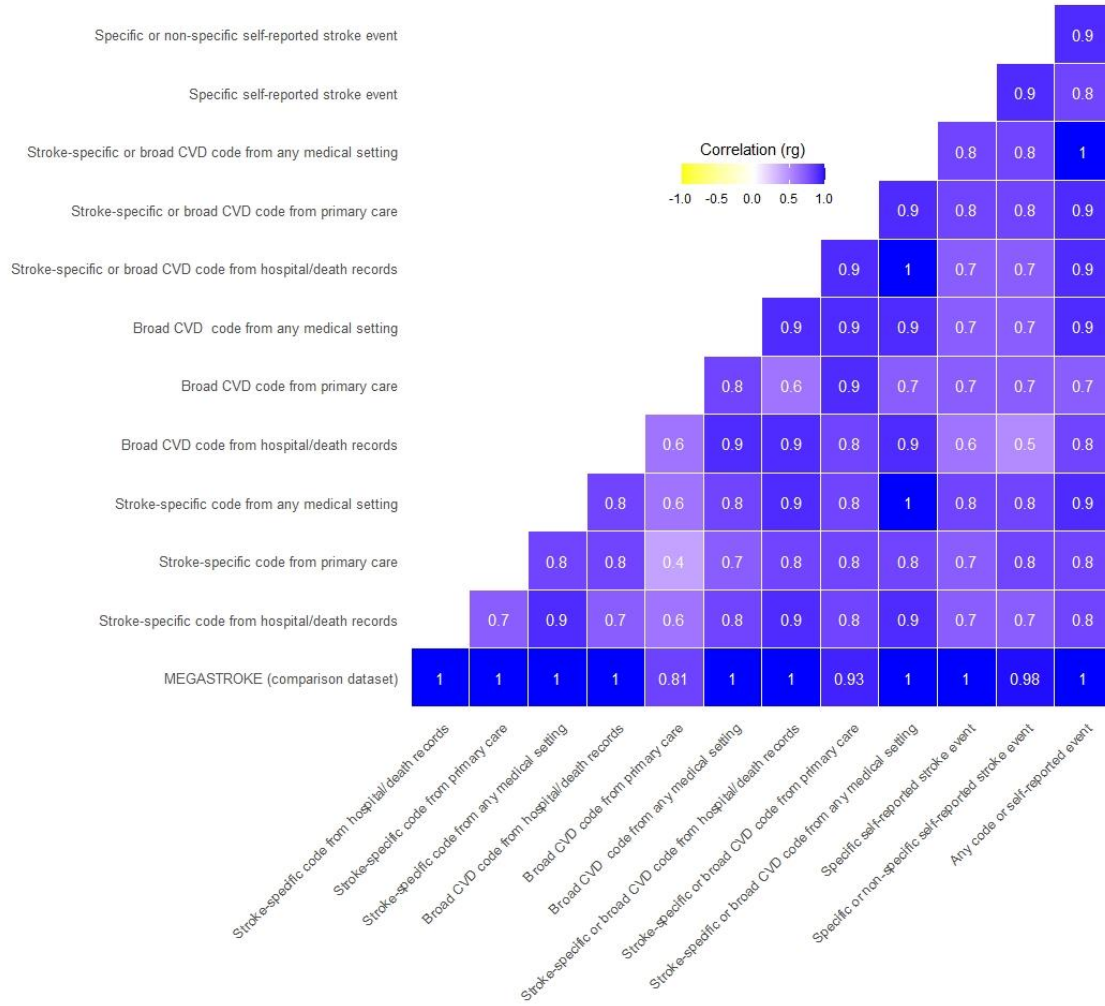
51 **Table 2. MEGASTROKE stroke-significant loci replicated using UKB stroke definitions**

Known stroke genetic locus	Significantly associated stroke type in MEGASTROKE	MEGASTROKE European meta-analyses N = 40,585*	MEGASTROKE transethnic meta-analyses N = 67,162 †	Stroke-specific code from any medical setting N = 8,665	Broad CVD code from any medical setting N = 8,085	Stroke-specific or Broad CVD code from any medical setting N = 16,750	Specific or non-specific self-reported stroke event N = 7,536	Any code or self-reported event N = 19,449
CHR4: rs13143308 (PITX2)	Cardio-embolic stroke	1.34 (1.28-1.40) 5.2x10-41	1.32 (1.27-1.37) 1.9x10-47	1.09 (1.05-1.13) 7.5x10-6	1.06 (1.02-1.1) 0.0036	1.08 (1.05-1.11) 1x10-7	1.04 (0.995-1.08) 0.08	1.07 (1.04-1.1) 3x10-7
CHR7: rs2107595 (HDAC9-TWIST1)	Large-vessel stroke	1.27 (1.19-1.35) 1.4x10-13	1.21 (1.15-1.26) 3.7x10-15	1.07 (1.03-1.12) 1x10-3	1.03 (0.99-1.08) 0.12	1.05 (1.02-1.09) 4.7x10-4	1.08 (1.03-1.13) 8.9x10-4	1.06 (1.03-1.09) 1.1x10-4
CHR10: rs2295786 (SH3PXD2A)	All stroke	1.05 (1.03-1.07) 1.4x10-7	1.05 (1.04-1.07) 1.8x10-10	1.07 (1.04-1.11) 1.5x10-5	1.01 (0.98-1.04) 0.59	1.04 (1.02-1.07) 4x10-4	1.04 (1.01-1.08) 0.021	1.05 (1.03-1.07) 5.1x10-6
CHR12: rs3184504 (SH2B3)	All ischaemic stroke	1.08 (1.06-1.10) 1.2x10-14	1.08 (1.06-1.10) 2.2x10-14	1.04 (1.01-1.07) 0.0096	1 (0.97-1.03) 0.83	1.02 (1-1.05) 0.04	1.06 (1.03-1.1) 1.9x10-4	1.03 (1.01-1.05) 0.0068
CHR9: rs635634 (ABO)	All ischaemic stroke	1.08 (1.05-1.11) 9.2x10-9	1.07 (1.04-1.10) 6.9x10-3	1.08 (1.03-1.12) 1.8x10-4	1 (0.96-1.04) 0.86	1.04 (1.01-1.07) 0.0043	1.1 (1.05-1.14) 5.8x10-6	1.05 (1.02-1.07) 4.3x10-4
CHR9: rs7859727 (Chr9p21)	All stroke	1.05 (1.03-1.07) 7.2x10-8	1.05 (1.03-1.07) 4.2x10-10	1.06 (1.03-1.09) 8.1x10-5	1.04 (1.01-1.07) 0.018	1.05 (1.03-1.07) 5.3x10-6	1.02 (0.98-1.05) 0.32	1.04 (1.02-1.06) 6.3x10-5
Summary: number replicated loci				5/32	0/32	4/32	3/32	5/32

52 Significant associations are in bold in shaded boxes. CVD: cerebrovascular disease; CHR: chromosome; N: number;

53 † MEGASTROKE transethnic meta-analyses included 60,341 ischaemic stroke cases; 9,006 cardio-embolic stroke cases; 6,688 large-vessel stroke cases. MEGASTROKE European meta-analyses
54 included 34,217 ischaemic stroke cases; 7,193 cardio-embolic stroke cases; 4,373 large-vessel stroke cases.

455 **Figure 1.** Genetic correlation of UKB stroke definitions with MEGASTROKE and each other.



456

457 Legend: Where the rg was >1, we rounded it to 1.

458

459

460

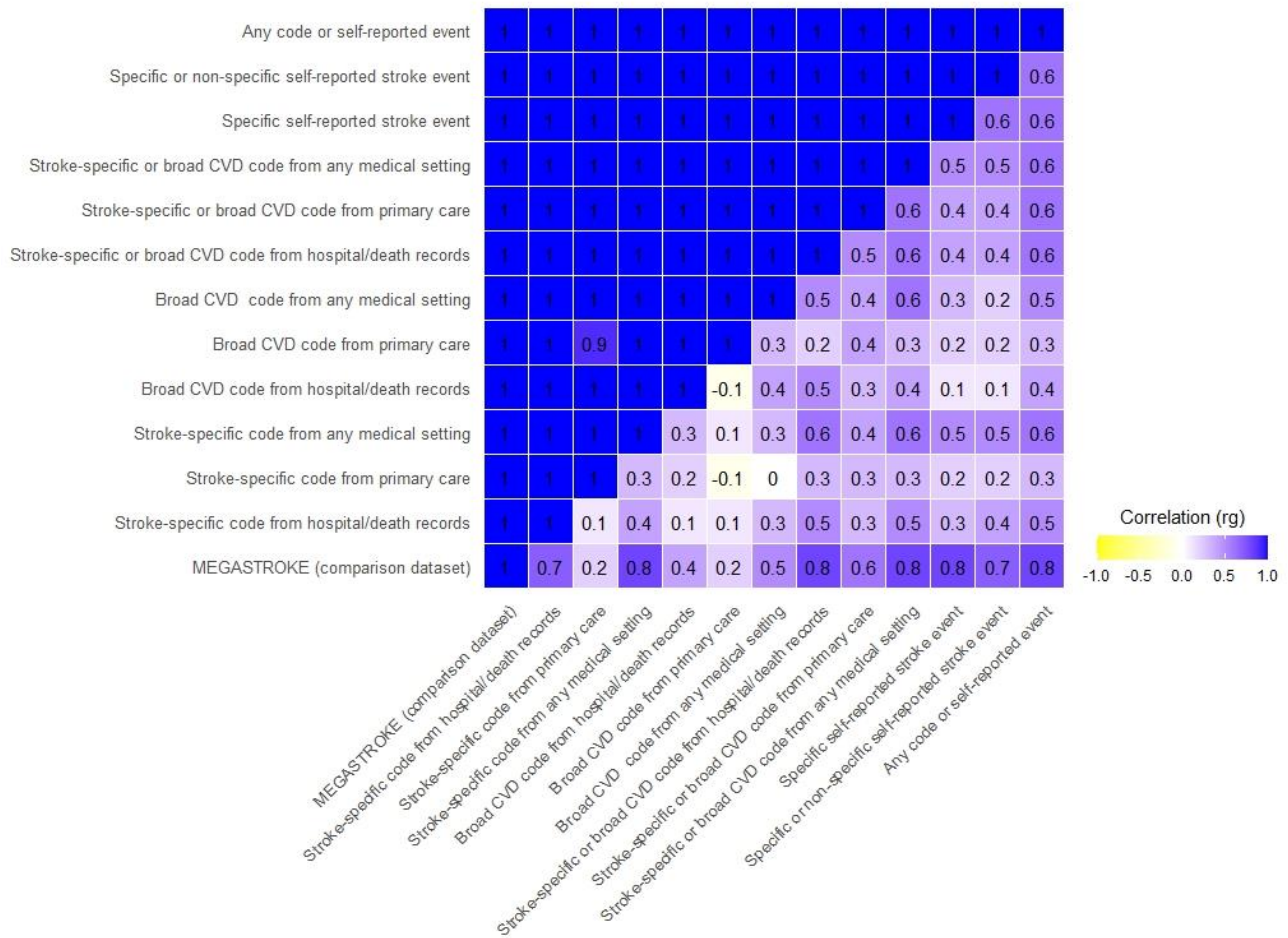
461

462

463

464

465 **Figure 2.** Confidence intervals of genetic correlation of UKB stroke definitions with MEGASTROKE
466 and each other.



467

468 Legend: The upper triangle displays Bonferroni corrected upper confidence intervals, and the
469 lower triangle displays Bonferroni corrected lower confidence intervals. Where the rg was >1 , we
470 rounded it to 1.

471

472

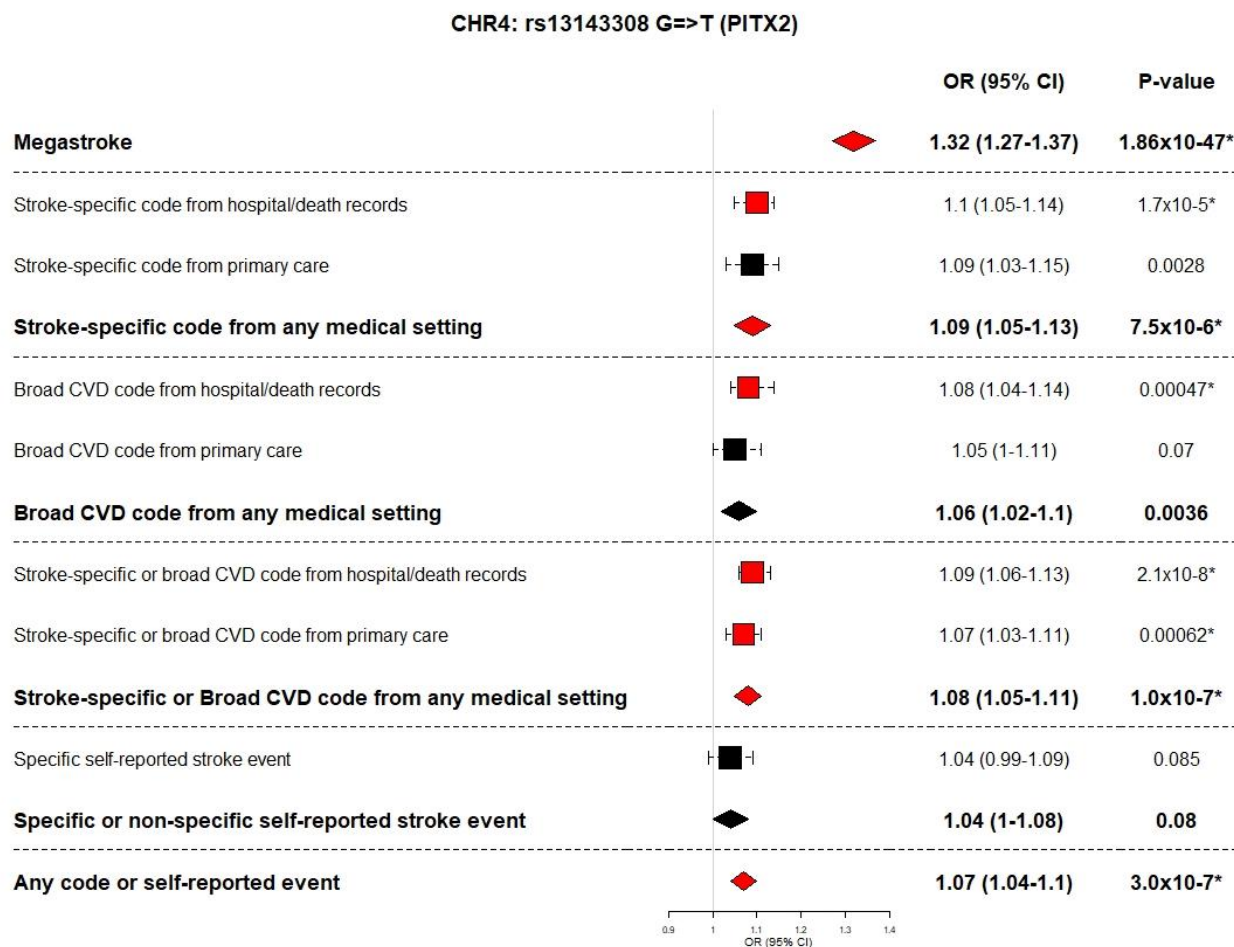
473

474

475 **Figure 3.** Megastroke stroke subtype-significant loci replicated using UK Biobank stroke

476 definitions.

477 **A.**



478

479

480

481

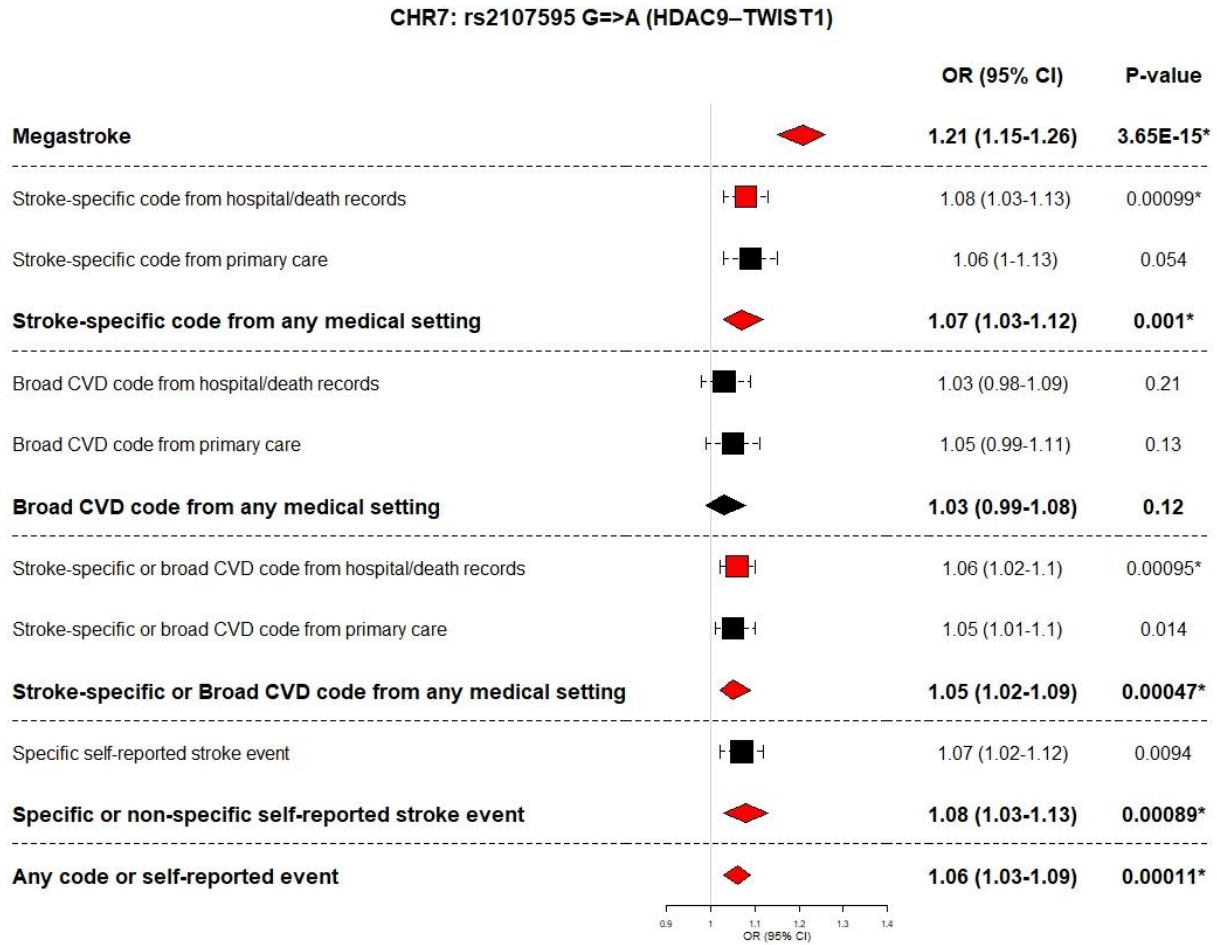
482

483

484

485

486 **B.**



487

488 Legend: MEGASTROKE odds ratio and p-value is shown for the analyses (European or trans-ethnic)

489 showing the lowest p-value