

Title

NOTCH3 p.Arg1231Cys is Markedly Enriched in South Asians and Associated with Stroke

Authors

Juan L. Rodriguez-Flores¹, Shareef Khalid^{2,3}, Neelroop Parikshak¹, Asif Rasheed³, Bin Ye¹, Manav Kapoor¹, Joshua Backman¹, Farshid Sepehrband¹, Silvio Alessandro DiGioia⁴, Sahar Gelfman¹, Tanima De¹, Nilanjana Banerjee¹, Deepika Sharma¹, Hector Martinez⁴, Sofia Castaneda⁵, David D'Ambrosio⁴, Xingmin A. Zhang¹, Pengcheng Xun⁴, Ellen Tsai⁶, I-Chun Tsai⁴, Regeneron Genetics Center¹, Maleeha Zaman Khan³, Muhammad Jahanzaib³, Muhammad Rehan Mian³, Muhammad Bilal Liaqat³, Khalid Mahmood⁷, Tanvir Us Salam⁸, Muhammad Hussain⁸, Javed Iqbal⁹, Faizan Aslam¹⁰, Michael N. Cantor¹, Gannie Tzoneva¹, John Overton¹, Jonathan Marchini¹, Jeff Reid¹, Aris Baras¹, Niek Verweij¹, Luca A. Lotta¹, Giovanni Coppola¹, Katia Karalis¹, Aris Economides¹, Sergio Fazio⁴, Wolfgang Liedtke⁴, John Danesh¹¹, Ayesha Kamal¹², Philippe Frossard³, Thomas Coleman¹, Alan R. Shuldiner¹, Danish Saleheen^{2,3}

Affiliations

- ¹. Regeneron Genetics Center, Regeneron Pharmaceuticals Inc, Tarrytown, NY, USA
². Columbia University, New York, NY, USA
³. Center for Non-Communicable Diseases, Karachi, Pakistan
⁴. Regeneron Pharmaceuticals Inc, Tarrytown, NY, USA
⁵. Rye Country Day School, Rye, NY, USA
⁶. University of California at Los Angeles, Los Angeles, CA, USA
⁷. Dow University of Health Sciences and Civil Hospital, Karachi, Pakistan
⁸. Lahore General Hospital, Lahore, Pakistan
⁹. Department of Neurology, Allied Hospital, Faisalabad, Pakistan.
¹⁰. Department of Neurology, Aziz Fatima Hospital, Faisalabad, Pakistan.
¹¹. Department of Public Health and Primary Care, University of Cambridge, Cambridge, UK.
¹². Section of Neurology, Department of Medicine, Aga Khan University, Karachi, Pakistan.

Correspondence

Danish Saleheen, danish.saleheen@cncdpk.com
Alan Shuldiner, alan.shuldiner@regeneron.com

Rodriguez-Flores et al.

NOTCH3 p.Arg1231Cys is Markedly Enriched in South Asians and Associated with Stroke

Version 34

17 June 2024

2

Abstract

39
40
41
42
43
44
45
46
47
48
49
50
51
52
53

The genetic factors of stroke in South Asians are largely unexplored. Exome-wide sequencing and association analysis (ExWAS) in 75 K Pakistanis identified NM_000435.3(*NOTCH3*):c.3691C>T, encoding the missense amino acid substitution p.Arg1231Cys, enriched in South Asians (alternate allele frequency = 0.58% compared to 0.019% in Western Europeans), and associated with subcortical hemorrhagic stroke [odds ratio (OR) = 3.39, 95% confidence interval (CI) = [2.26, 5.10], p value = 3.87×10^{-9}], and all strokes (OR [CI] = 2.30 [1.77, 3.01], p value = 7.79×10^{-10}). *NOTCH3* p.Arg231Cys was strongly associated with white matter hyperintensity on MRI in United Kingdom Biobank (UKB) participants (effect [95% CI] in SD units = 1.1 [0.61, 1.5], p value = 3.0×10^{-6}). The variant is attributable for approximately 2.0% of hemorrhagic strokes and 1.1% of all strokes in South Asians. These findings highlight the value of diversity in genetic studies and have major implications for genomic medicine and therapeutic development in South Asian populations.

Introduction

54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99

Pakistan, a country in South Asia, comprises over 231 million inhabitants. It is the fifth most populous country in the world with diverse ancestral backgrounds from South and Central Asia, West Asia, and Africa. Pakistan, and in general South Asia, represents an understudied region in large-scale genetic studies [1], thus providing an opportunity for novel discoveries of the genetic basis of diseases.

Stroke is a leading cause of death globally [2], and epidemiological studies suggest an elevated incidence and prevalence of stroke in Pakistan [2, 3] relative to Europe. The disparities in incidence and prevalence between Pakistan and Europe could be due to many factors, including difference in access to healthcare facilities with high-quality diagnostic capabilities and public health awareness and education. These disparities also may reflect differences in prevalence of risk factors such as hypertension and diabetes, lifestyle factors such as diet, physical activity and smoking, and genetic predispositions [4-7]. Studies of the genetic underpinnings of stroke in Pakistani populations have been limited, making this understudied population an opportune venue for stroke exome-wide sequencing and association studies (ExWAS).

At least 9 rare monogenic disorders are characterized by increased stroke risk, such as cerebral autosomal dominant arteriopathy with subcortical infarcts and leukoencephalopathy (CADASIL) due to mutations in *NOTCH3* [8]. CADASIL is distinct from other hereditary stroke diseases because it is characterized by vascular smooth muscle cell (VSMC) degeneration in small arteries and accumulation of protein aggregates known as granular osmophilic material (GOM) that contain aggregates of misfolded *NOTCH3* extracellular domain (ECD). The more common forms of stroke are likely polygenic with substantial contributions from behavioral and environmental factors as well as age. Major risk factors include high systolic blood pressure, high body mass index, hyperlipidemia, elevated glucose, and smoking [9]. Recent genome-wide association studies (GWAS) identified single nucleotide variants (SNVs) in more than 28 loci associated with stroke [10]. These variants are common non-coding variants with small effect sizes and were identified in predominantly European ancestry populations.

The aim of this study was to identify protein coding deleterious missense or loss-of-function (LoF) variants associated with stroke phenotypes in the Pakistani population. We performed exome sequencing in a 31 K discovery cohort consisting of 5,135 stroke cases and 26,602 controls of Pakistani origin. ExWAS identified NM_000435.3(*NOTCH3*):c.3691C>T, encoding the missense amino acid substitution p.Arg1231Cys, with an approximately three-fold increased risk of hemorrhagic stroke in heterozygotes. Follow-up meta-analysis of 61 K Pakistani (including an additional 160 cases and 30,239 controls) provided further support for association of *NOTCH3* p.Arg1231Cys with stroke (combined ischemic and hemorrhagic). This variant was present in approximately 1% of Pakistani and was markedly enriched with respect to Europeans in multiple South Asian (SAS) and West Asian (WAS) (also referred to as Greater Middle Eastern [11]) populations ranging from Turkey to India. The variant was estimated to explain up to 1.1% of strokes and 2.0% of hemorrhagic strokes in South Asia, a region having a population of > 2 billion people, thus having significant medical implications in these very large yet understudied populations and their global diaspora.

100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145

Results

ExWAS in the Pakistan Genomics Resource (PGR) discovery cohort identifies a markedly enriched missense variant in *NOTCH3* associated with stroke

Characteristics of the $n = 5,135$ stroke cases and $n = 26,602$ controls in the discovery cohort are summarized in Table 1. Compared to controls, stroke cases were modestly older and had a higher prevalence of known risk factors for vascular disease including hypertension, diabetes, myocardial infarction, and tobacco use (all $p < 0.01$). As expected, in this cohort ascertained for stroke, most cases were ischemic strokes and most hemorrhagic strokes were subcortical (Supplementary Figures 1 and 2, Supplementary Table 1).

Case:control ExWAS for all stroke cases and 4 stroke subtypes with sufficient case counts to provide statistical power (Supplementary Table 1) identified a genome-wide significant (p value $< 5.0 \times 10^{-8}$) association for NM_000435.3(*NOTCH3*):c.3691C>T (rs201680145), encoding the missense amino acid substitution p.Arg1231Cys, with subcortical hemorrhagic stroke (OR [95% CI] = 3.39 [2.26, 5.1], p value = 3.87×10^{-9} ; AAF = 0.58%) (Figure 1B, Supplementary Table 1 and Supplementary Figures 3 and 4). The p.Arg1231Cys variant also showed evidence for association with all strokes combined (OR [95% CI] = 2.18 [1.65, 2.89], p value = 4.44×10^{-8}) and other sub-categories of stroke (Supplementary Table 2). No other variants in the locus were associated with stroke (Figure 1C). We did not observe an association between p.Arg1231Cys and history of hypertension, elevated systolic or diastolic blood pressure, or smoking, known major risk factors for stroke (Supplementary Tables 3 and 4), and inclusion of these risk factors in regression analysis did not appreciably alter the effect of p.Arg1231Cys on stroke risk (Supplementary Table 5).

NOTCH3 encodes Notch Receptor 3, a transmembrane signaling protein and part of an evolutionarily conserved family that plays a pleiotropic role in cell-cell interaction and neural development [12]. The extra-cellular domain (ECD) of *NOTCH3* consists of 34 Epidermal Growth Factor-like repeat (EGFr) domains [13], each containing 6 Cysteine (Cys) residues that form three disulfide bonds (Figure 2). Adding or removing Cys residues in the first 6 EGFr domains cause classical CADASIL, a highly penetrant rare autosomal dominant disease clinically characterized by migraine with aura, early-onset recurrent strokes, dementia, and behavioral changes [8]. The Cys-altering variant associated with stroke in this study, p.Arg1231Cys, occurs in the 31st EGFr [14] and is predicted deleterious (Supplementary Table 6). Stroke cases heterozygous for p.Arg1231Cys and stroke cases without the variant were similar with respect to age, age of stroke onset, type of stroke, and stroke risk factors, suggesting a milder form of CADASIL not obviously clinically distinguishable from common forms of stroke in this population, although a detailed history for migraine or other manifestations of CADASIL were not available (Table 1).

Replication and Meta Analysis in PGR

146 An additional 30 K of whom self-reported stroke case:control status was obtained (160 cases and
147 30,239 controls) were sequenced by CNCD. Replication of the association was observed in this
148 independent cohort (OR [95% CI] = 3.49 [1.56, 7.83], p value = 5.00×10^{-3} . Meta-analysis for all
149 strokes in the combined 61 K cohort achieved a genome-wide significant p value (OR [95% CI])
150 = 2.30 [1.76, 2.99], p value = 7.08×10^{-10}) (Figure 3).

151

152 **Recall by Genotype**

153

154 A total of 12 p.Arg1231Cys homozygotes from 9 nuclear families were identified in the PGR,
155 including 9 discovery cohort probands and 3 follow-up cohort relatives identified through a
156 callback of 128 call-back participants. Baseline characteristics of the 12 homozygotes are shown
157 in Supplementary Table 7. Three of twelve (25%) homozygotes had a history of stroke; of note,
158 all with stroke were >65 years of age while all without a history of stroke were <55 years of age.
159

159

160 Eight out of twelve (66%) p.Arg1231Cys homozygotes had a history of hypertension. Among
161 the 128 callback participants, both systolic and diastolic blood pressure were trending higher
162 (Supplementary Table 8). While there was no association with hypertension in the discovery
163 cohort, p.Arg1231Cys homozygotes in the PGR had nominally higher diastolic blood pressure
164 than heterozygotes or homozygous reference individuals (median = 95 mmHg, interquartile
165 range 86 to 100; heterozygotes (median = 80 mmHg, interquartile range = 80 to 90), p value =
166 0.016 (Supplementary Table 9).

167

168 **Allele Frequency and Population Attributable Risk**

169

170 The allele frequency of p.Arg1231Cys was 1.1% across the PGR 75 K, equivalent to a
171 population prevalence of 1 in 46. After removing cases recruited for cardiovascular diseases
172 (individuals enrolled at time of acute stroke, MI, and heart failure), the allele frequency of the
173 variant was 0.51%, equivalent to a population prevalence of 1 in 98. This frequency was orders
174 of magnitude higher relative to exomes of European ancestry from UK Biobank (AAF =
175 0.019%), corresponding to a population prevalence of 1 in 2,614). The variant was enriched
176 (AAF > 0.1%) in other South Asian and West Asian populations [15] both within and outside of
177 Pakistan (Supplementary Data, Supplementary Tables 10 and 11, Supplementary Figure 5).
178

178

179 We estimate that 2.0% [bootstrap 95% CI based on 10,000 resamples: 1.0% to 2.9%] of
180 hemorrhagic strokes and 1.1% [bootstrap 95% CI based on 10,000 resamples: 0.6% to 1.6%] of
181 all strokes in the Pakistani population are attributable to p.Arg1231Cys. Thus, this variant is a
182 common cause of strokes in SAS and WAS populations, a finding that has implications for
183 medical care as well as global health in these populations.
184

184

185 **Suggestive Associations at Other Loci**

186

187 Although *NOTCH3* p.Arg1231Cys was the only variant associated with stroke at a genome-wide
188 significant p value below 5.0×10^{-8} , there were a total of 9 associations (5 loci) with p values
189 below 1.0×10^{-6} and at least 10 variant carriers (Supplementary Table 12). In addition to
190 *NOTCH3*, these included one known locus previously associated with stroke in a recent GWAS,
191 lymphocyte specific protein *LSPI* [10].

192
193 The *LSP1* locus variant with suggestive association with intracerebral hemorrhagic stroke in this
194 study was a common intronic variant (rs661348, OR [95% CI] = 1.3 [1.2, 1.4], p value = 8.0×10^{-8} ,
195 AAF = 0.27). While rs661348 was not previously associated with stroke, two common non-
196 coding variants in the *LSP1* locus were previously reported to be associated with stroke
197 (rs569550 and rs1973765) [10]. Both variants were in linkage disequilibrium with rs661348 ($r^2 >$
198 0.4 in 10 K unrelated PGR participants). A test of association for rs661348 conditional on these
199 variants reduced the strength of the association for rs661348 to nominal (8.83×10^{-4})
200 (Supplementary Table 13), suggesting that rs661348 represents the same known stroke risk
201 locus. The *LSP1* locus was previously reported to be associated with hypertension [16]. In
202 PGR there was a nominal association with hypertension (rs661348, OR [95% CI] = 1.0 [1.0,
203 1.1], p value = 2.61×10^{-2}) (Supplementary Tables 14 and 15) also observed in UKB (OR [95%
204 CI] = 1.0 [1.0, 1.1], p value = 8.05×10^{-25}) (Supplementary Tables 15 and 16).

205 206 ***NOTCH3* p.Arg1231Cys is associated with stroke and CADASIL-like phenotypes in UK** 207 **Biobank**

208
209 To investigate stroke and CADASIL-related phenotypes in an independent cohort, UK Biobank
210 data was reviewed for associations with *NOTCH3* p.Arg1231Cys. A total of 255 heterozygotes
211 for *NOTCH3* p.Arg1231Cys were observed in 450 K exome-sequenced individuals from this
212 predominantly European cohort, with a markedly lower allele frequency (AAF = 0.019%)
213 (Supplementary Table 18). Phenome-wide association (PheWAS) of 10,168 phenotypes revealed
214 nominally significant association of p.Arg1231Cys with ischemic stroke (OR [95% CI] = 4.0
215 [1.9, 8.6]), p value = 4.1×10^{-4}), all strokes combined (OR [95% CI] = 1.9 [1.1, 3.5], p value =
216 0.031), hypertension (ICD 10 code I10) (OR [95% CI] = 1.5 [1.1, 2.2], p value = 0.019), and
217 recurrent major depression (OR [95% CI] = 3.2 [1.5, 6.8], p value = 0.0031) (Supplementary
218 Table 19). No association was observed for hemorrhagic stroke, migraine, dementia, mood
219 changes, Alzheimer's disease, or urinary incontinence. The lack of association (p value = 0.066)
220 with hemorrhagic stroke in UKB Europeans was likely due to low statistical power, given the
221 lower variant allele frequency and lower hemorrhagic stroke prevalence in UKB compared to
222 PGR. Nonetheless, the odds ratio (OR [95% CI] = 5.8 [0.88, 39.1]) was high.

223
224 In addition to recurrent strokes, brain white matter loss is a major and early phenotype
225 characteristic of CADASIL that is focused on particular brain regions [8, 14]. *NOTCH3*
226 p.Arg1231Cys was strongly associated with a cluster of brain MRI quantitative phenotypes, e.g.,
227 total volume of white matter hyperintensities (WMH) from T1 and T2 FLAIR images (effect
228 [95% CI] in SD units = 1.1 [0.61, 1.5], p value = 3.0×10^{-6}) with carriers having 7.4 cm³ more
229 WMH volume than controls (Supplementary Figure 6 and Supplementary Table 20). The most
230 prominent alterations in WMH in p.Arg1231Cys carriers were observed in the centrum
231 semiovale and periventricular white matter (Supplementary Figure 7). Taken together, these
232 results demonstrate *NOTCH3* p.Arg1231Cys carriers have increased risk of established markers
233 of small vessel disease and clinical phenotypes observed in CADASIL [8].

234 235 **Pathogenic burden of all Cys-altering variants within *NOTCH3* EGFr domains specifically** 236 **associated with CADASIL phenotypes in UK Biobank**

237

238 Burden test analysis allows for increased statistical power to detect association by combining
239 signal across multiple rare variants. Prior studies have shown that pathogenic variants in
240 CADASIL are limited to variants that add or remove a Cysteine (Cys-altering) in *NOTCH3* EGFr
241 domains normally containing 6 Cysteines. Furthermore, patients with Cys-altering variants in the
242 first 6 EGFr domains have more severe symptoms than in EGFr domains 7 to 34 [14, 17, 18],
243 including larger regions of brain white matter loss [19], more granular osmophilic material
244 (GOM) aggregates in blood vessels [19], and worse prognosis [14].
245

246 In order to test these hypotheses, a set of custom gene burden tests were designed and compared
247 to single variant test results for *NOTCH3* p.Arg1231Cys. In UKB, 758 individuals carried one of
248 98 unique Cys-altering variants across the 34 EGFr domains in *NOTCH3* (Supplementary Tables
249 21 and 22). A burden test aggregating all UKB EGFr domain Cys-altering variants into a single
250 statistical test was strongly associated with stroke (OR [95% CI] = 2.86 [2.14, 3.82], p value =
251 6.29×10^{-10} ; AAF = 0.01%) (Table 2). In contrast to Cys-altering variants within EGFr domains,
252 Cys-altering variants outside of EGFr domains were not associated with stroke (OR [95% CI] =
253 0.97 [0.46, 2.03], p value = 9.3×10^{-1} ; AAF = 0.039%) (Table 2). In order to rule out the
254 possibility that any missense variants in EGFr domains are associated with stroke, a test limited
255 to the most commonly altered (added or removed) amino acid in *NOTCH3*, serine (Ser), was
256 tested and did not show any evidence of association with stroke (OR [95% CI] = 0.98 [(0.8, 1.2],
257 p value = 0.084; AAF = 0.54%) (Table 2, Supplementary Table 23). Interestingly, a burden test
258 limited only to predicted loss of function (pLoF) variants (frameshift, splice variant, stop gain)
259 did not show significant evidence for association with stroke (OR [95% CI] = 1.38 [0.50, 3.85], p
260 value = 0.54; AAF = 0.019%) (Table 2, Supplementary Table 24). These results provide
261 evidence to support the hypothesis that EGFr domain Cys-altering variants within *NOTCH3* are
262 associated with stroke, in contrast to other protein-altering variants.
263

264 While hemorrhagic stroke represents a small proportion of the strokes reported in the UKB, the
265 set of Cys-altering variants were also tested for association with hemorrhagic stroke. A nominal
266 association with hemorrhagic stroke (OR [95% CI] = 3.61 [1.39, 9.34], p value = 8.31×10^{-3} ; AAF
267 0.025%) was observed, despite low statistical power.
268

269 Consistent with stroke risk, in MRI data of 35,344 UKB individuals, Cys-altering variants in
270 *NOTCH3* EGFr domains were strongly associated with WMH volume (p value = 3.7×10^{-13} ; with
271 carriers having 5.4 cm³ greater WMH volume than controls). These WMH differences were
272 strongest in the centrum semiovale and periventricular white matter (Supplementary Figure 7).
273 Additionally, we found strong WMH signal in the external capsule, which is known to be
274 involved in CADASIL. We found weaker evidence for association of *NOTCH3* LoF variants
275 with WMH (effect size [95% CI] = 6.8 cm³ [4.2 cm³, 10.9 cm³], p-value = 1.68×10^{-4}).
276

277 Prior studies have binned *NOTCH3* EGFr domain Cys-altering variants in up to three distinct
278 severity or risk groups based on EGFr domain number [10, 14, 17]. Indeed, we observed a much
279 larger effect size for Cys-altering variants in high-risk EGFr domains 1-6 (OR [95% CI] = 29.5
280 [10.4, 83.8], p value = 1.37×10^{-7} ; AAF = 0.002%) compared to Cys-altering variants in EGFr
281 domains 7-34 (OR [95% CI] = 2.55 [1.87, 3.46], p value = 1.59×10^{-7} ; AAF = 0.098%) (Table 2,
282 Supplementary Table 23). These results are consistent with prior reports of differences in stroke

283 risk between EGFr domain risk groups not correlated with differences in signaling activity
284 between EGFr risk groups [17].

285

286

Discussion

287

288 This report describes the largest ExWAS of stroke conducted thus far in a South Asian
289 population and highlights a Cys-altering missense variant in the 31st EGFr domain of *NOTCH3*
290 associated with stroke at a genome-wide level of statistical significance. This is the first study to
291 report a genome-wide-significant association between *NOTCH3* and stroke, a discovery enabled
292 because *NOTCH3* p.Arg1231Cys is markedly enriched in Pakistanis compared to Western
293 European and non-Eurasian populations. Harbored in ~1 percent of Pakistani, p.Arg1231Cys is
294 associated with a ~3-fold increased risk of hemorrhagic stroke. While some regional variability
295 in the allele frequency is observed, p.Arg1231Cys is enriched in populations ranging from
296 Turkey in West Asia to India in South Asia, suggesting a substantial contribution to stroke risk in
297 millions of individuals across South Asia and West Asia as well as their global diaspora.

298

299 *NOTCH3* was not previously associated with stroke in the largest GWAS predominantly
300 consisting of European-derived participants [10]. In contrast to prior studies, both the discovery
301 and replication cohorts in this study were South Asian, hence avoiding the bias encountered in
302 studies with a European discovery cohort. Given the much lower allele frequency of
303 p.Arg1231Cys in European populations, we observed a nominal association between
304 p.Arg1231Cys and stroke in the UK Biobank study, showing a similar effect size as in South
305 Asians. Nominal associations were also observed for phenotypes related to CADASIL, such as
306 hypertension and depression. While brain images were not available for the Pakistani cohort, a
307 strong association was observed between p.Arg1231Cys and quantitative brain MRI phenotypes
308 in UKB data, such as white matter hyperintensity.

309

310 Cys-altering mutations in proximal EGFr domains of *NOTCH3* are known to cause autosomal
311 dominant CADASIL, a rare highly penetrant distinct syndrome that includes early onset
312 recurrent subcortical strokes. In contrast to classical CADASIL pathogenic variants,
313 p.Arg1231Cys is in the 31st of 34 EGFr domains, appears to have more moderate penetrance, and
314 is not obviously clinically distinguishable from more common multi-factorial forms of stroke in
315 South Asians. The p.Arg1231Cys *NOTCH3* variant is currently classified in ClinVar and recent
316 reviews [14] as a variant of uncertain significance [20] or “low risk” [17]; however, based on our
317 current findings, there is strong genetic, computational, and imaging evidence of pathogenicity
318 for this variant despite reduced penetrance and severity compared to “classic” Cys-altering
319 CADASIL pathogenic variants in EGFr domains 1 to 6 [14, 19].

320

321 Prior studies have debated if the mechanism whereby Cys-altering variants contribute to
322 CADASIL-related pathology is through toxic aggregate gain of function (GoF) or a loss of
323 normal signaling function (LoF). One study demonstrated excess risk of CADASIL-related
324 phenotypes for Cys-altering variants in EGFr domains 1 to 34 [18], while other studies showed
325 greater risk in EGFr domains 1 to 6 relative to EGFr domains 7 to 34 [14, 19]. A recent study
326 showed evidence for expanding the high-risk tier of EGFr domains to include domains 8, 11, and
327 26 [17]. The current study provides three additional refinements. First, this study is the first to
328 assess risk for LoF variants, and did not observe significant association signal (**Table 2**),

329 although the number of LoF carriers was small and thus power is limited to detect such
330 associations. These findings suggest that pathological mechanisms driven by dysfunctional
331 disulfide bridge formation and subsequent protein misfolding and aggregation, as is commonly
332 observed in CADASIL, may be more pathologic than simple LoF (haploinsufficiency) [19].
333 Second, we demonstrate association between p.Arg1231Cys with stroke, thus demonstrating that
334 CADASIL-related stroke is not uncommon as was previously thought. While prior studies have
335 shown enrichment of p.Arg1231Cys in South Asians [21], and have used this information as
336 evidence to classify p.Arg1231Cys variant as “low-risk” [17], the current study provides
337 evidence contrary to that verdict. Furthermore, we have demonstrated a broader enrichment of
338 the variant across the region, including multiple West Asian and South Asian populations. Third,
339 the prior studies demonstrated a brain-wide association with WMH, while the current study
340 identifies pathology focused in the external capsule and other brain regions known for
341 CADASIL pathology.

342
343 CADASIL is characterized by both ischemic and hemorrhagic strokes, although the factors that
344 contribute to the manifestation of one versus the other stroke type awaits further clarification [8].
345 Hemorrhagic stroke appears to represent a larger proportion of strokes in South Asia than in
346 Europe [2]. In this study, the p.Arg1231Cys association signal was stronger in PGR for
347 hemorrhagic strokes than for ischemic strokes, despite nearly two-fold larger ischemic stroke
348 case counts. In contrast, the UKB association signal appeared stronger for ischemic stroke,
349 possibly due to low hemorrhagic stroke case count and thus statistical power in this cohort.
350 Statistical power issues aside, differences in manifestation of p.Arg1231Cys in South Asians
351 compared to Europeans may be attributable to differences in risk factors such as age, blood
352 pressure, diabetes, air pollution, smoking, medications such as anti-platelet and anti-coagulants
353 used to manage atherosclerotic disease, genetic background, or study-specific differences in
354 criteria to categorize stroke sub-types. Further research is needed to better ascertain the
355 mechanism behind cerebral arterial wall pathobiology and clinical presentation of ischemic
356 versus hemorrhagic stroke in p.Arg1231Cys carriers.

357
358 Currently there are no known effective preventive or therapeutic interventions for CADASIL or
359 less penetrant forms of *NOTCH3* related stroke. However, our analyses provide clues toward
360 their development. First, in contrast to our analyses of EGFr domain Cys-altering missense
361 variants in *NOTCH3* that were significantly associated with stroke, predicted loss of function
362 variants that would be expected to not produce a functional protein were not significantly
363 associated with stroke. These observations suggest that targeting therapeutic interventions that
364 decrease expression of mutant protein (such as siRNA, antisense oligonucleotides, and CRISPR),
365 induce exon skipping of altered EGFr domains [22, 23], or accelerate removal of GOM may
366 prove beneficial for prevention and/or treatment [24].

367
368 Our analyses also suggest that p.Arg1231Cys is modestly associated with hypertension, although
369 p.Arg1231Cys association with stroke risk appears independent of hypertension or other stroke
370 risk factors such as smoking, age and sex. Animal models of CADASIL show decreased vascular
371 tone and contractility, most likely driven by loss of physiologic function and subsequent
372 degeneration of vascular smooth muscle cells (VSMCs) [25]. These observations suggest that
373 while management of hypertension and smoking cessation are effective modalities for primary
374 and secondary prevention of stroke, those with *NOTCH3* mutation related strokes will need

375 additional therapeutic interventions, as existing hypertensive medications cannot restore VSMC
376 function.

377
378 A limitation of this study is the lack of brain imaging analysis for the Pakistani carriers, such that
379 specific brain regions affected by the ischemic and hemorrhagic strokes could be ascertained and
380 compared. In addition, we lacked more detailed clinical data such as presence of migraines and
381 longitudinal data of disease course including stroke recurrence, dementia, and depression.
382 Further characterization of p.Arg1231Cys carriers will be necessary to obtain better estimates of
383 penetrance as well as to identify distinguishing clinical or biomarker characteristics that may
384 have utility in early diagnosis, prevention and treatment, and for recommendations for cascade
385 screening in family members. Migraine symptoms typically precede stroke by 10+ years in
386 CADASIL patients [8], thus the combination of migraine with aura, depression and family
387 history of stroke could be sufficient evidence to prescribe *NOTCH3* genetic testing. Finally, the
388 effect of LoFs on stroke risk will require larger sample sizes for more definitive comparison to
389 stroke risk of Cys-altering variants.

390
391 In conclusion, we identified a highly enriched Cys-altering variant in *NOTCH3* in South Asians
392 that expands the phenotypic spectrum of CADASIL from rare and highly penetrant to common
393 and moderately penetrant. Based on our estimates, this single variant may be responsible for
394 ~1.1% of all strokes combined and ~2.0% of hemorrhagic strokes in South Asians. Among 1.9
395 billion South Asians there could be over 26 million carriers for the variant. Thus, this work has
396 important implications for genetic screening and early identification of at-risk individuals, and
397 the future opportunity for rationally targeted therapeutic interventions.

399 **Methods**

400 **Summary**

401
402
403 Details of methods below. Briefly, 75 K individuals were recruited and consented in Pakistan for
404 whole exome sequencing, including a stroke case:control discovery cohort of 31 K (including n
405 = 5,135 cases and n = 26,602 controls) sequenced by the Regeneron Genetics Center and a
406 follow-up cohort of 44 K (n=44,082), including 30 K with self-report stroke case:control status
407 used for replication and meta-analysis (n = 160 cases and n = 30,239 controls). ExWAS was
408 conducted for stroke and 4 overlapping stroke subtypes (intracerebral hemorrhage, subcortical
409 intracerebral hemorrhage, ischemic stroke, and partial anterior circulating infarct) in the
410 discovery cohort and combined in a meta-analysis of stroke with the replication cohort using
411 both single-variant and gene burden test models [26]. Population attributable fraction of stroke
412 for associated variants was calculated as based on the prevalence of mutation among cases and
413 odds ratio (OR) for risk of stroke in the discovery cohort using the standard definition [27].
414 Consented callbacks were conducted in n = 128 individuals within families of homozygotes for
415 associated variants. For comparison and validation, analyses were conducted in UK Biobank data
416 using publicly available datasets and methodologies [28-31], including association analysis of
417 p.Arg1231Cys *NOTCH3* with stroke phenotypes in 380 K participants, and brain imaging
418 phenotypes in 35 K participants (Figure 1A).

419 **Study Populations**

421
422 This study focused on two distinct cohorts, including 75 thousand (K) individuals from the
423 Pakistan Genomic Resource (PGR) and 370 K individuals from the United Kingdom Biobank
424 (UKB) The PGR 75 K individuals were recruited and consented in Pakistan for whole exome
425 sequencing (WES) (n = 75,819), including a stroke case:control discovery cohort of 31 K (n =
426 31,737, including n = 5,135 cases and n = 26,602 controls) sequenced by the Regeneron Genetics
427 Center and a follow-up cohort of 44 K (n = 44,082), including 30 K with self-report stroke
428 case:control status used for replication and meta-analysis (n = 30,399, including n = 160 cases
429 and n = 30,239 controls). The remaining n = 13,684 in the follow-up cohort had sequence data
430 but not stroke case:control status known, including n = 6,067 produced by WES and n = 7,616
431 produced by whole genome sequencing (WGS). (Figure 1A).

432

433 **Pakistan Genomic Resource (PGR)**

434

435 PGR is a growing biobank that aims to enroll 1 million participants across Pakistan and as of
436 September 2023, ~250,000 participants across 48 clinical sites in 17 cities from all over Pakistan
437 have been enrolled. Following the success of a case-control study design in genetic studies
438 adopted by several international (e.g., Wellcome Trust case-control consortium) and regional
439 studies (e.g., PROMIS), PGR is a national consortium of several case-control studies focused on
440 50 distinct phenotypes, including: stroke, myocardial infarction, angiographically confirmed
441 coronary artery disease, heart failure, age-related macular degeneration, keratoconus, diabetic
442 retinopathy, glaucoma, asthma, chronic obstructive pulmonary disease (COPD), non-alcoholic
443 fatty liver disease (NAFLD), type-2 diabetes, chronic kidney disease, Alzheimer's disease,
444 Parkinson's disease, dementia, progressive multiple sclerosis, autism, Huntington's disease,
445 hematological cancers, breast cancer, ovarian cancer, cancers of head and neck, esophageal
446 cancer, lung cancer, gastric cancer, colorectal cancer, melanoma, cancers of the urinary tract,
447 cervical cancer, prostate cancer, rheumatoid arthritis, systemic lupus erythematosus (SLE),
448 psoriatic arthritis, ankylosing spondylitis, osteoarthritis, scleroderma, juvenile arthritis, systemic
449 sclerosis, inflammatory myositis, alopecia areata, acne rosacea, primary Sjogren's syndrome,
450 sarcoidosis, idiopathic pulmonary cholangitis, idiopathic pulmonary fibrosis, vitiligo, longevity /
451 healthy aging, and previously uncharacterized Mendelian disorder. For each of these phenotypes,
452 screening is carried out at specialized clinical sites across Pakistan by trained research medical
453 officers who review inclusion and exclusion criteria and approach eligible participants for
454 recruitment into PGR. In a similar manner, for each of the phenotypes, controls are frequency-
455 matched to cases on sex and age (in 5-year bands). Controls are being recruited in the following
456 order of priority: (1) visitors of patients attending the out-patient department; (2) patients
457 attending the out-patient department for routine non-phenotype related complaints, or (3) non-
458 blood related visitors. Following informed consent, both cases and controls are enrolled.
459 Research medical officers administer pre-piloted epidemiological questionnaires to participants
460 that seek a total of >200 items of information in relation to: ethnicity (e.g, personal and paternal
461 ethnicity, spoken language, place of birth and any known consanguinity); demographic
462 characteristics; lifestyle factors (e.g., tobacco and alcohol consumption, dietary intake and
463 physical activity); and personal and family history of disease; and medication usage. The Center
464 for Non-Communicable Diseases (CNCD), Pakistan serves as the sponsor and the coordinating
465 center of PGR.

466

467 Using standardized procedures and equipment, research officers obtain measurements of height,
468 weight, waist and hip circumference, systolic and diastolic blood pressure, and heart rate. Waist
469 circumference is assessed over the abdomen at the widest diameter between the costal margin
470 and the iliac crest, and hip circumference is assessed at the level of the greater trochanters.
471 Information extracted from questionnaires and physical measurements is entered by two different
472 operators into the central database, which is securely held at CNCND, Karachi, Pakistan. Non-
473 fasting blood samples (with the time since last meal recorded) are drawn by phlebotomists from
474 each participant and centrifuged within 45 min of venipuncture. A total of 29 ml of blood is
475 drawn from each participant in 2×6 ml serum tubes and 3×5 ml EDTA tubes. Hence, a total of
476 five blood tubes are collected per participant, including serum, EDTA plasma and whole blood
477 which are all stored in cryogenic vials. All samples are stored temporarily at each recruitment
478 center at -20°C . Serum, plasma and whole blood samples are transported daily to the central
479 laboratory at CNCND where they are stored at -80°C . The long-term -80°C sample repository.
480 Measurements of total cholesterol, high-density lipoprotein-cholesterol, triglycerides, AST,
481 ALT, glucose, creatinine, and HbA1c (in a subset) are performed centrally using (Roche
482 Diagnostics GmbH, USA) in all study participants.

483
484 Research technicians trained in accordance with standard operating procedures in laboratories at
485 CNCND extracted DNA from leukocytes using a reference phenol-chloroform protocol. DNA
486 concentrations are determined. The yield of DNA per participant is typically between 600 and
487 800 ng/ μl in a total volume of about 500 μl . To minimize any systematic biases arising from
488 plate- or batch-specific genotyping error and/or nonrandom missingness, stock plates are used to
489 generate genotyping plates which contain a mixture of cases and controls along with negative
490 and positive controls designed to address genotyping quality control (QC), plate identification
491 and orientation.

492
493 PGR has received approval by the relevant research ethics committee of each of the institutions
494 involved in participant recruitment, as well as centrally by the IRB board of the Center for Non-
495 Communicable Diseases which is registered with the National Institutes of Health, USA. PGR
496 has also been approved by the National Bioethics Committee, Islamabad Health Research
497 Institute, National Institutes of Health of Pakistan.

498
499 Eligibility criteria was defined as described below. Ischemic stroke sub-types in the PGR cohort
500 were defined using TOAST [32] and Oxfordshire [33] clinical criteria.

501 502 **UK Biobank**

503
504 The UK Biobank (UKB) cohort had detailed medical records and lifestyle data as described
505 online in the UKB Showcase (<https://biobank.ndph.ox.ac.uk/showcase/>) [28]. Stroke case:control
506 status was available in 380 K UK Biobank participants, of which 280 K also had smoking and
507 hypertension status (referred to as UKB replication cohort). A sub-cohort of $n = 35,977$ UKB
508 participants had brain MRI data which was produced and analyzed as described online
509 (https://biobank.ctsu.ox.ac.uk/crystal/crystal/docs/brain_mri.pdf) [29-31]. MRI images for $n = 19$
510 p.Arg1231Cys carriers were re-analyzed in order to identify brain regions affected (referred to as
511 UKB 35 K brain imaging cohort). The description of phenotypes and methods for normalizing
512 the data, including rank-inverse normal transformation (RINT) are described online [28].

513 Exome Sample Preparation, Sequencing, and QC

514 Genomic DNA samples were transferred to the Regeneron Genetics Center from the CNCD and
515 stored in an automated sample biobank at -80°C before sample preparation. DNA libraries were
516 created by enzymatically shearing DNA to a mean fragment size of 200 bp, and a common Y-
517 shaped adapter was ligated to all DNA libraries. Unique, asymmetric 10 bp barcodes were added
518 to the DNA fragment during library amplification to facilitate multiplexed exome capture and
519 sequencing. Equal amounts of sample were pooled before overnight exome capture, with a
520 slightly modified version of IDT's xGen v1 probe library; all samples were captured on the same
521 lot of oligonucleotides. The captured DNA was PCR amplified and quantified by quantitative
522 PCR. The multiplexed samples were pooled and then sequenced using 75 bp paired-end reads
523 with two 10 bp index reads on an Illumina NovaSeq 6000 platform on S4 flow cells. A total of n
524 = 42,695 samples were made available for processing. We were unable to process $n = 1,948$
525 samples, most of which failed QC during processing owing to low or no DNA being present.

526 A total of $n = 40,747$ samples were sequenced, of which $n = 2,943$ (7%) did not pass one or more
527 of our QC metrics and were subsequently excluded. Criteria for exclusion were as follows:
528 disagreement between genetically determined and reported sex ($n = 900$); high rates of
529 heterozygosity or contamination ($\text{VBID} > 5\%$) ($n = 709$); low sequence coverage (less than 80%
530 of targeted bases achieving $20\times$ coverage) ($n = 115$); genetically identified sample duplicates
531 ($n = 1,662$ total samples); WES variants discordant with the genotyping chip ($n = 43$). The
532 remaining $n = 37,804$ (37 K) samples were then used to compile a project-level VCF (PVCF) for
533 downstream analysis using the GLnexus joint genotyping tool. This final dataset contained $n =$
534 7,655,430 variants. Within this dataset of 37 K exomes, stroke case:control status was known for
535 $n = 31,737$, referred to as the 31 K discovery cohort. The remaining $n = 6,067$ were part of the 41
536 K follow-up cohort.

537 Exome sequencing in the replication 30 K cohort ($n = 39,399$) was conducted by the CNCD
538 using a publicly available protocol [34]. Briefly, blood derived DNA samples, with 10 to 100 ng
539 concentration of initial genomic DNA, underwent hybridization and capture using Illumina
540 Rapid Capture Exome Kit or Agilent's SureSelect Human Exon v2. Samples were denatured and
541 amplified HiSeq v3 cluster chemistry and HiSeq 2000 or 2500 flowcells based on the
542 manufacturers protocol. Reads were aligned to the GRCh38 genome reference and variants were
543 called using GATK v.30 followed by variant recalibration to remove false positive variants.
544

545 The remaining $n = 7,616$ exome samples of the 75 K PGR consisted of whole genome sequence
546 (WGS) data that produced a VCF subsequently filtered to include only variants in protein coding
547 sequence. WGS samples were sequenced and process as described in a publicly available
548 protocol (<https://www.nature.com/articles/s41586-021-03205-y>). Briefly, 30x whole genome
549 sequencing was performed using Illumina HiSeqX instruments. Reads were aligned to the
550 GRCh38 reference using BWA-align and variants were called using the publicly available
551 GotCloud pipeline (<https://genome.sph.umich.edu/wiki/GotCloud>), which includes QCing
552 variants based on a support vector machine trained on specific site quality metrics.

553 Variant calling

554 The PGR discovery cohort WES data was reference-aligned using the OQFE protocol [35],
555 which uses BWA MEM to map all reads to the GRCh38 reference in an alt-aware manner, marks
556 read duplicates and adds additional per-read tags. The OQFE protocol retains all reads and
557 original quality scores such that the original FASTQ is completely recoverable from the resulting
558 CRAM file. Single-sample variants were called using DeepVariant with custom exome
559 parameters [35], generating a gVCF for each input OQFE CRAM file. These gVCFs were
560 aggregated and joint-genotyped using GLnexus (v.1.3.1). All constituent steps of this protocol
561 were executed using open-source software. The PGR replication and follow-up cohort were
562 analyzed using the publicly available GotCloud workflow
563 (<https://genome.sph.umich.edu/wiki/GotCloud>).

564 **Identification of low-quality variants from sequencing using machine learning**

565 Similar to other recent large-scale sequencing efforts, we implemented a supervised machine-
566 learning algorithm to discriminate between probable low-quality and high-quality variants [36,
567 37]. In brief, we defined a set of positive control and negative control variants based on the
568 following criteria: (1) concordance in genotype calls between array and exome-sequencing data;
569 (2) transmitted singletons; (3) an external set of likely ‘high quality’ sites; and (4) an external set
570 of likely ‘low quality’ sites. To define the external high-quality set, we first generated the
571 intersection of variants that passed QC in both TOPMed Freeze 8 and GnomAD v.3.1 genomes.
572 This set was additionally restricted to 1000 Genomes Phase 1 high-confidence SNPs from the
573 1000 Genomes Project [38] and gold-standard insertions and deletions from the 1000 Genomes
574 Project and a previous study [39], both available through the GATK resource bundle
575 (<https://gatk.broadinstitute.org/hc/en-us/articles/360035890811-Resource-bundle>). To define the
576 external low-quality set, we intersected GnomAD v3.1 fail variants with TOPMed Freeze 8
577 Mendelian or duplicate discordant variants. Before model training, the control set of variants
578 were binned by allele frequency and then randomly sampled such that an equal number of
579 variants were retained in the positive and negative labels across each frequency bin. A support
580 vector machine using a radial basis function kernel was then trained on up to 33 available site
581 quality metrics, including, for example, the median value for allele balance in heterozygote calls
582 and whether a variant was split from a multi-allelic site. We split the data into training (80%) and
583 test (20%) sets. We performed a grid search with fivefold cross-validation on the training set to
584 identify the hyperparameters that returned the highest accuracy during cross-validation, which
585 were then applied to the test set to confirm accuracy. This approach identified a total of $n =$
586 931,823 WES variants as low-quality, resulting in a dataset of $n = 6,723,607$ variants.

587 **Variant annotation**

588 Variants were annotated as described in a publicly available pipeline 38. In brief, variants were
589 annotated using Ensembl variant effect predictor, with the most severe consequence for each
590 variant chosen across all protein-coding transcripts. In addition, we derived canonical transcript
591 annotations based on a combination of MANE, APPRIS and Ensembl canonical tags. MANE
592 annotation was given the highest priority followed by APPRIS. When neither MANE nor
593 APPRIS annotation tags were available for a gene, the canonical transcript definition of Ensembl
594 was used. Gene regions were defined using Ensembl release 100. Variants annotated as stop
595 gained, start lost, splice donor, splice acceptor, stop lost or frameshift, for which the allele of

596 interest was not the ancestral allele, were considered predicted loss-of-function variants. Five
597 annotation resources were utilized to assign deleteriousness to missense variants: SIFT,
598 Polyphen2_HDIV, Polyphen2_HVAR, LRT, MutationTaster [40-43], and LRT, obtained using
599 dbNSFP [44]. Missense variants were considered ‘likely deleterious’ if predicted deleterious by
600 all five algorithms, ‘possibly deleterious’ if predicted deleterious by at least one algorithm and
601 ‘likely benign’ if not predicted deleterious by any algorithm.

602 **Pakistan Genome Resource Statistical Analysis**

603
604 ExWAS of SNPs with minor allele count > 5 in the 31 K PGR discovery cohort was conducted
605 with 5 binary stroke phenotypes using REGENIE (v 3.1.1) [26] with age, age², sex, age*sex,
606 exome batch, 10 genotyping array principal components (PCs), 10 common variant exome PCs
607 and 10 rare variant exome PCs as covariates. The minimum of 1,000 cases was selected based on
608 a power calculation [45] (1,000 cases; 25,000 controls; significance threshold 0.005; prevalence
609 0.0012; disease allele frequency 0.005; genotype relative risk 3.0; > 80% power). Follow-up
610 analyses were conducted with the added covariates of hypertension and tobacco use, or as
611 environmental factors in a gene-by-environment interaction test using REGENIE [26]. Gene
612 burden analysis was conducted using REGENIE with separate masks for pLoF, pLoF +
613 missense, pLoF + deleterious missense (as predicted by at least 1 of 5 algorithms), and pLoF +
614 deleterious missense (as predicted by 5 of 5 algorithms). Analysis in the 61 K PGR meta-analysis
615 cohort, including 31 K discovery and 30 K replication cohorts, was conducted using REGENIE.

617 **PGR Population Genetic Analysis**

618
619 Using principal components (PCs) [46] and Uniform Manifold Approximation and Projection
620 (UMAP) [47] based analyses PGR and UKB South Asian sub-populations were mapped to
621 distinct groups or clusters. Specifically, we used the imputed genotypes to merge the PGR
622 dataset with UKB and 1000 genome datasets. Imputed data was used to maximize the number of
623 common variants between all three datasets. The Plink [48] “--bmerge” option was used to
624 merge datasets. A minimal QC was applied to the merged genotypes to exclude variants with
625 MAF less than 5%, missing genotype rate greater than 10%, and Hardy Weinberg equilibrium P
626 value less than 5×10^{-5} . Variants mapping within the HLA region were excluded. Merged
627 datasets were pruned for linkage disequilibrium ($r^2 > 0.25$). A total of 20 PCAs were calculated
628 in the merged data using the Plink -pca option. Calculated PCAs were imported to R and merged
629 with reported ethnicities or country of birth information. The first 6 PCs calculated on the
630 merged data were reduced to two dimensions using the UMAP package in R. The two
631 eigenvectors of UMAP were calculated using an alpha value of 1.1 and beta value of 0.8. Two
632 eigenvectors were plotted along with ethnicity and country of birth labels using the Plotly
633 package in R. UKB self-reported ethnicities or country of birth was confirmed to be highly
634 correlated with data obtained from UMAP.

636 **Population Attributable Fraction**

637
638 Estimation of the proportion of all strokes combined or hemorrhagic stroke in Pakistan
639 population attributable to p.Arg1231Cys was calculated using the formula [27],

640

641

$$AF_p = P_c \times AF_e = P_c \left(1 - \frac{1}{OR}\right)$$

642

643 where P_c is the prevalence of mutation among cases, AF_e is the attributable fraction in the
644 exposure, and OR was for risk of stroke (i.e., all stroke combined or hemorrhagic stroke)
645 comparing mutant vs wild type of p.Arg1231Cys in the discovery cohort.

646

647 OR was obtained directly from Supplementary Table 2. The related 95% confidence intervals
648 were constructed using bootstrap method with 10,000 resamples [49], which was implemented
649 with the command “Bootstrap” using Stata (College Station, Texas 77845 USA).

650

651 **Recall by Genotype**

652

653 A subset of carriers of *NOTCH3* p.Arg1231Cys were contacted by the Center of Non-
654 Communicable Diseases in Karachi Pakistan under protocols approved by the IRB committee of
655 the Center for Non-Communicable Diseases (NIH registered IRB 00007048). After obtaining
656 consent from the proband and from the family members, questionnaires regarding past medical
657 and family history were administered by trained research staff, in the local language. Physical
658 measurements such as height, weight, hip and waist circumference were obtained in the standing
659 position by using height and weight scales. Blood pressure and heart rate were recorded sitting
660 by using OMRON healthcare M2 blood pressure monitors. Non-fasting blood samples were
661 collected from each participant in EDTA and Gel Tubes. Serum and plasma were separated
662 within 45 minutes of venipuncture. A random urine sample was also collected from each
663 participant. The samples were stored temporarily in dry ice in the field and transported to the
664 central laboratory based at CNCD and stored at -80 degrees Celsius. Measurements for total-
665 cholesterol, HDL cholesterol, LDL cholesterol, triglycerides, VLDL, AST, ALT and creatinine
666 were obtained from serum samples using enzymatic assays. HbA1c was measured using a
667 turbidimetric assay in whole-blood samples (Roche Diagnostics). All measurements were done
668 at a central laboratory at CNCD. Statistical analysis comparing across genotypes was conducted
669 using the numpy library in Python 3.11.4.

670

671 **PGR Stroke Case Control Definitions**

672

- 673 • Controls
 - 674 • Inclusion
 - 675 • No medical history of stroke, myocardial infarction (MI), coronary artery
676 disease, heart failure (HF), valvular disease, or pacemaker
 - 677 • Cases
 - 678 • Inclusion
 - 679 • General Criteria
 - 680 ○ Stroke: Diagnosis of ‘Stroke’
 - 681 ○ Ischemic: Diagnosis of 'Ischemic stroke'
 - 682 ○ Hemorrhagic: Diagnosis of 'Hemorrhagic stroke'
 - 683 ■ Subcortical: Type of intracerebral hemorrhage = 'Subcortical'

- 684 ▪ Parenchymal: Type of intracerebral hemorrhage =
685 ‘Parenchymal’
- 686 • Oxfordshire Criteria
 - 687 ○ Partial anterior circulation infarct (PACI): Partial anterior circulation
 - 688 infarcts (PACI) stroke sub-type
 - 689 ○ Posterior circulation infarction (POCI): Posterior circulation infarcts
 - 690 POCI stroke sub-type
 - 691 ○ Total anterior circulation infarct (TACI): Total anterior circulation
 - 692 infarcts TACI' stroke sub-type
 - 693 ○ Lacunar infarct (LACI): Lacunar infarcts stroke sub-type
 - 694 • TOAST Criteria
 - 695 ○ Cardioembolism (CE): Cardioembolism ischemic stroke subtype
 - 696 ○ CE probable: CE criteria, in addition diagnosis is made if the clinical
 - 697 findings, neuroimaging data, and results of diagnostic studies are
 - 698 consistent with one subtype and other etiologies have been excluded
 - 699 ○ Large artery atherosclerosis (LAA): Large artery atherosclerosis
 - 700 ischemic stroke subtype based on TOAST classification
 - 701 ○ LAA probable: LAA criteria, in addition in addition diagnosis is made
 - 702 if the clinical findings, neuroimaging data, and results of diagnostic
 - 703 studies are consistent with one subtype and other etiologies have been
 - 704 excluded
 - 705 ○ Small artery atherosclerosis (SAA): Small artery atherosclerosis
 - 706 ischemic stroke subtype
 - 707 ○ SAA probable: SAA criteria, in addition in addition diagnosis is made
 - 708 if the clinical findings, neuroimaging data, and results of diagnostic
 - 709 studies are consistent with one subtype and other etiologies have been
 - 710 excluded

712 UK Biobank Stroke Case Control Definition

713
714 UK Biobank stroke case control definitions were based on ICD10 codes as follows.

- 715
- 716 • Cases
 - 717 ○ Inclusion
 - 718 ▪ Phe10_I63, ICD10 3D: Cerebral infarction
 - 719 ▪ Phe10_I630, ICD10 4D: Cerebral infarction due to thrombosis of
 - 720 precerebral arteries
 - 721 ▪ Phe10_I631, ICD10 4D: Cerebral infarction due to embolism of
 - 722 precerebral arteries
 - 723 ▪ Phe10_I632, ICD10 4D: Cerebral infarction due to unspecified occlusion
 - 724 or stenosis of precerebral arteries
 - 725 ▪ Phe10_I633, ICD10 4D: Cerebral infarction due to thrombosis of cerebral
 - 726 arteries
 - 727 ▪ Phe10_I634, ICD10 4D: Cerebral infarction due to embolism of cerebral
 - 728 arteries

- 729 ▪ Phe10_I635, ICD10 4D: Cerebral infarction due to unspecified occlusion
730 or stenosis of cerebral arteries
731 ▪ Phe10_I638, ICD10 4D: Other cerebral infarction
732 ▪ Phe10_I639, ICD10 4D: Cerebral infarction, unspecified
733 ▪ Self-reported: SR_1583_ischaemic_stroke
734 ▪ Primary and secondary cause of death using above ICD codes.
735 ○ Exclusion
736 ▪ Phe10_I636, ICD10 4D: Cerebral infarction due to cerebral venous
737 thrombosis, nonpyogenic
738 • Controls
739 • Inclusion
740 • Negative for the above codes
741 • Negative for Phe10_Z823, ICD10 4D: Family history of stroke
742 2. Exclusion:
743 • Phe10_G45, ICD10 3D: Transient cerebral ischemic attacks and related
744 syndromes
745 • Phe10_G458, ICD10 4D: Other transient cerebral ischemic attacks and related
746 syndromes
747 • Phe10_G459, ICD10 4D: Transient cerebral ischemic attack, unspecified
748
749

750 Custom Burden Tests in UKB 450 K

751
752 Burden tests aim to boost statistical power by aggregating association signal across multiple rare
753 variants. Prior studies in human and animal models have debated the role of various variant
754 classes on *NOTCH3* function, CADASIL pathology and patient prognosis, including experiments
755 designed to determine if the pathogenicity of CADASIL variants follows a loss of function (LoF)
756 mechanism [8, 25, 50]. Using data from hundreds of missense and LoF variants in *NOTCH3*
757 observed in 450 K UKB participant exomes, burden tests were conducted to assess the impact of
758 LoF and missense variants.
759

760 Ten distinct gene burden tests of association with stroke were conducted using REGENIE [26],
761 divided into three distinct groups. The Group I tests assessed the impact of Cys-altering variants,
762 Group II tests assessed the impact of Ser-altering variants, and Group III tests assessed LoF and
763 all missense variants. Group I and II consisted of four distinct tests, including (1) a test of all
764 group variants in EGFr domains 1 to 34, (2) a test of all group variants in EGFr domains 1 to 6,
765 (3) a test of all group variants in EGFr domains 7 to 34, and (4) a test of all group variants
766 outside of EGFr domains. The difference between Group I and Group II was Group I variants are
767 missense variants that either add or remove a Cysteine (Cys-altering), while Group II variants are
768 missense variants that either add or remove a Serine (Ser-altering). While the role of Cys-altering
769 variants in CADASIL is well known [15, 19, 20], Ser-altering variants were chosen based on
770 being the most common class of variants among *NOTCH3* variants in UKB 450 K exomes.
771 Group III consisted of two tests, including (1) a test limited to LoF variants and (2) a test limited
772 to LoF and missense variants.
773
774

Rodriguez-Flores et al.

19

NOTCH3 p.Arg1231Cys is Markedly Enriched in South Asians and Associated with Stroke

Version 34

17 June 2024

775

776

777

Data Availability

778 The primary data in this study not already presented in the manuscript and supplement consists

779 of ExWAS summary statistics for 5 stroke phenotypes analyzed in the PGR cohort. This data is

780 publicly available in the GWAS Catalog under accessions GCST90432122, GCST90432123,

781 GCST90432124, GCST90432125, and GCST90432126.

782

783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810

Code Availability

Data collection and analysis software, tools, algorithms, and packages used in this manuscript are publicly available in software repositories. Below is a list of the softwares used, their versions, and contemporary links to public repositories.

Data Production

Read alignment was conducted using BWA MEM v.0.7.17. Variant calling was conducted using Deep Variant v.0.10.0. Joint genotyping was conducted using GLnexus v1.3.1. Variant deleteriousness was calculated using SIFT v.2011, Polyphen 2 v.2011, LRT v.2013, MutationTaster v.4.3, and dbNSFP v.3.2. Single variant ExWAS and gene burden test of PGR and UKB data was conducted using REGENIE v.3.1.3.

Data Analysis

Figures were plotted using R v.4.4.1. Data management was conducted using Python v.3.11.4. Statistical tests in Table 1 were conducted using R stats library v.4.3.0. Supplementary Figure 1 and 2 were produced using R UpSetR library v.1.4.0. Supplementary Figures 3 and 4 were produced using R QQMAN library v.0.1.8. Supplementary Figure 5 was produced using R Plotly library v.4.10.1. Supplementary Figure 6 was produced using R v.4.4.1. Supplementary Figure 7 was produced using ITK-SNAP v.3.8.0. Principal Components Analysis was conducted using PLINK v.1.9. Uniform Manifold Approximation and Projection was conducted using UMAP v.0.2.10.0. Power calculation was conducted using GAS Power Calculator v.2017. Meta Analysis in Figure 3 was conducted using METAL v.2011-03-25. MRI Image Analysis was conducted using FMRIB Software Library (FSL) v.5.0.10 and FreeSurfer v.6.0. Alignment in Figure 2 was conducted using BLAST v.2.14.

811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856

References

1. Mills, M.C. and C. Rahal, *The GWAS Diversity Monitor tracks diversity by disease in real time*. Nat Genet, 2020. **52**(3): p. 242-243.
2. Collaborators, G.B.D.S., *Global, regional, and national burden of stroke and its risk factors, 1990-2019: a systematic analysis for the Global Burden of Disease Study 2019*. Lancet Neurol, 2021. **20**(10): p. 795-820.
3. Sherin, A., et al., *Prevalence of stroke in Pakistan: Findings from Khyber Pakhtunkhwa integrated population health survey (KP-IPHS) 2016-17*. Pak J Med Sci, 2020. **36**(7): p. 1435-1440.
4. Valcarcel-Nazco, C., et al., *Variability in the use of neuroimaging techniques for diagnosis and follow-up of stroke patients*. Neurologia (Engl Ed), 2019. **34**(6): p. 360-366.
5. Farooq, A., N. Venketasubramanian, and M. Wasay, *Stroke Care in Pakistan*. Cerebrovasc Dis Extra, 2021. **11**(3): p. 118-121.
6. Farooq, M.U., et al., *The epidemiology of stroke in Pakistan: past, present, and future*. Int J Stroke, 2009. **4**(5): p. 381-9.
7. Mullen, M.T., et al., *Hospital-Level Variability in Reporting of Ischemic Stroke Subtypes and Supporting Diagnostic Evaluation in GWTG-Stroke Registry*. J Am Heart Assoc, 2023. **12**(24): p. e031303.
8. Chabriat, H., et al., *Cadasil*. Lancet Neurol, 2009. **8**(7): p. 643-53.
9. Markidan, J., et al., *Smoking and Risk of Ischemic Stroke in Young Men*. Stroke, 2018. **49**(5): p. 1276-1278.
10. Mishra, A., et al., *Stroke genetics informs drug discovery and risk prediction across ancestries*. Nature, 2022. **611**(7934): p. 115-123.
11. Scott, E.M., et al., *Characterization of Greater Middle Eastern genetic variation for enhanced disease gene discovery*. Nat Genet, 2016. **48**(9): p. 1071-6.
12. Wang, T., M. Baron, and D. Trump, *An overview of Notch3 function in vascular smooth muscle cells*. Prog Biophys Mol Biol, 2008. **96**(1-3): p. 499-509.
13. Duvaud, S., et al., *Expasy, the Swiss Bioinformatics Resource Portal, as designed by its users*. Nucleic Acids Res, 2021. **49**(W1): p. W216-w227.
14. Rutten, J.W., et al., *Broad phenotype of cysteine-altering NOTCH3 variants in UK Biobank: CADASIL to nonpenetrance*. Neurology, 2020. **95**(13): p. e1835-e1843.
15. Rodriguez-Flores, J.L., et al., *The QChip1 knowledgebase and microarray for precision medicine in Qatar*. NPJ Genom Med, 2022. **7**(1): p. 3.
16. Hoffmann, T.J., et al., *Genome-wide association analyses using electronic health records identify new loci influencing blood pressure variation*. Nat Genet, 2017. **49**(1): p. 54-64.
17. Hack, R.J., et al., *Three-tiered EGFR domain risk stratification for individualized NOTCH3-small vessel disease prediction*. Brain, 2023. **146**(7): p. 2913-2927.
18. Cho, B.P.H., et al., *Association of Vascular Risk Factors and Genetic Factors With Penetrance of Variants Causing Monogenic Stroke*. JAMA Neurol, 2022. **79**(12): p. 1303-1311.
19. Rutten, J.W., et al., *The effect of NOTCH3 pathogenic variant position on CADASIL disease severity: NOTCH3 EGFR 1-6 pathogenic variant are associated with a more severe phenotype and lower survival compared with EGFR 7-34 pathogenic variant*. Genet Med, 2019. **21**(3): p. 676-682.

- 857 20. ClinVar, N.C.f.B.I.
858 21. Rutten, J.W., et al., *Archetypal NOTCH3 mutations frequent in public exome: implications for CADASIL*. Ann Clin Transl Neurol, 2016. **3**(11): p. 844-853.
859
860 22. Rutten, J.W., et al., *Therapeutic NOTCH3 cysteine correction in CADASIL using exon skipping: in vitro proof of concept*. Brain, 2016. **139**(Pt 4): p. 1123-35.
861
862 23. Gravesteijn, G., et al., *Naturally occurring NOTCH3 exon skipping attenuates NOTCH3 protein aggregation and disease severity in CADASIL patients*. Hum Mol Genet, 2020. **29**(11): p. 1853-1863.
863
864
865 24. Ghezali, L., et al., *Notch3(ECD) immunotherapy improves cerebrovascular responses in CADASIL mice*. Ann Neurol, 2018. **84**(2): p. 246-259.
866
867 25. Belin de Chantemele, E.J., et al., *Notch3 is a major regulator of vascular tone in cerebral and tail resistance arteries*. Arterioscler Thromb Vasc Biol, 2008. **28**(12): p. 2216-24.
868
869 26. Mbatchou, J., et al., *Computationally efficient whole-genome regression for quantitative and binary traits*. Nat Genet, 2021. **53**(7): p. 1097-1103.
870
871 27. Greenland, S., *Applications of Stratified Analysis Methods*, in *Modern epidemiology: Third edition*. 2008, Lippincott Williams & Wilkins: Philadelphia. p. 295-297.
872
873 28. Backman, J.D., et al., *Exome sequencing and analysis of 454,787 UK Biobank participants*. Nature, 2021. **599**(7886): p. 628-634.
874
875 29. Benjamini, Y. and Y. Hochberg, *Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing*. Journal of the Royal Statistical Society. Series B (Methodological) 1995. **57**(1): p. 289-300.
876
877
878 30. Elliott, L.T., et al., *Genome-wide association studies of brain imaging phenotypes in UK Biobank*. Nature, 2018. **562**(7726): p. 210-216.
879
880 31. Griffanti, L., et al., *BIANCA (Brain Intensity AbNormality Classification Algorithm): A new tool for automated segmentation of white matter hyperintensities*. Neuroimage, 2016. **141**: p. 191-205.
881
882
883 32. Adams, H.P., Jr., et al., *Classification of subtype of acute ischemic stroke. Definitions for use in a multicenter clinical trial. TOAST. Trial of Org 10172 in Acute Stroke Treatment*. Stroke, 1993. **24**(1): p. 35-41.
884
885
886 33. Bamford, J., et al., *Classification and natural history of clinically identifiable subtypes of cerebral infarction*. Lancet, 1991. **337**(8756): p. 1521-6.
887
888 34. Saleheen, D., et al., *Human knockouts and phenotypic analysis in a cohort with a high rate of consanguinity*. Nature, 2017. **544**(7649): p. 235-239.
889
890 35. Olga Krasheninina, Y.-C.H., Xiaodong Bai, Aleksandra Zalcman, Evan Maxwell, Jeffrey G. Reid, William J. Salerno Jr., *Open-source mapping and variant calling for large-scale NGS data from original base-quality scores*. Biorxiv, 2020.
891
892
893 36. Lin, M., et al., *Admixed Populations Improve Power for Variant Discovery and Portability in Genome-Wide Association Studies*. Front Genet, 2021. **12**: p. 673167.
894
895 37. Karczewski, K.J., et al., *The mutational constraint spectrum quantified from variation in 141,456 humans*. Nature, 2020. **581**(7809): p. 434-443.
896
897 38. Genomes Project, C., et al., *A global reference for human genetic variation*. Nature, 2015. **526**(7571): p. 68-74.
898
899 39. Mills, R.E., et al., *An initial map of insertion and deletion (INDEL) variation in the human genome*. Genome Res, 2006. **16**(9): p. 1182-90.
900
901 40. Sim, N.L., et al., *SIFT web server: predicting effects of amino acid substitutions on proteins*. Nucleic Acids Res, 2012. **40**(Web Server issue): p. W452-7.
902

- 903 41. Adzhubei, I., D.M. Jordan, and S.R. Sunyaev, *Predicting functional effect of human*
904 *missense mutations using PolyPhen-2*. *Curr Protoc Hum Genet*, 2013. **Chapter 7**: p.
905 Unit7 20.
- 906 42. Chun, S. and J.C. Fay, *Identification of deleterious mutations within three human*
907 *genomes*. *Genome Res*, 2009. **19**(9): p. 1553-61.
- 908 43. Steinhaus, R., et al., *MutationTaster2021*. *Nucleic Acids Res*, 2021. **49**(W1): p. W446-
909 W451.
- 910 44. Liu, X., et al., *dbNSFP v4: a comprehensive database of transcript-specific functional*
911 *predictions and annotations for human nonsynonymous and splice-site SNVs*. *Genome*
912 *Med*, 2020. **12**(1): p. 103.
- 913 45. Skol, A.D., et al., *Joint analysis is more efficient than replication-based analysis for two-*
914 *stage genome-wide association studies*. *Nat Genet*, 2006. **38**(2): p. 209-13.
- 915 46. Price, A.L., et al., *Principal components analysis corrects for stratification in genome-*
916 *wide association studies*. *Nat Genet*, 2006. **38**(8): p. 904-9.
- 917 47. McInnes, L.A., John%ASaul, Nathaniel%AGroßberger, Lukas%BJournal Name: Journal
918 of Open Source Software, J.V. 3, and J.I. 29, *UMAP: Uniform Manifold Approximation*
919 *and Projection*. Journal Name: Journal of Open Source Software; Journal Volume: 3;
920 Journal Issue: 29, 2018: p. Medium: X; Size: 861.
- 921 48. Chang, C.C., et al., *Second-generation PLINK: rising to the challenge of larger and*
922 *richer datasets*. *Gigascience*, 2015. **4**: p. 7.
- 923 49. Efron, B. and R. Tibshirani, *Bootstrap Methods for Standard Errors, Confidence*
924 *Intervals, and Other Measures of Statistical Accuracy*. *Statistical Science*, 1986. **1**.
- 925 50. Joutel, A., et al., *Cerebrovascular dysfunction and microcirculation rarefaction precede*
926 *white matter lesions in a mouse genetic model of cerebral ischemic small vessel disease*.
927 *J Clin Invest*, 2010. **120**(2): p. 433-45.
928
929
930

Acknowledgements

931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948

Supported by Regeneron Pharmaceuticals.

This research has been conducted using the UK Biobank Resource (project 26041). The authors thank everyone who made this work possible, particularly the UK Biobank team, their funders, the professionals from the member institutions who contributed to and supported this work, and most especially the UK Biobank participants, without whom this research would not be possible. The exome sequencing was funded by the UK Biobank Exome Sequencing Consortium (Bristol Myers Squibb, Regeneron, Biogen, Takeda, Abbvie, Alnylam, AstraZeneca and Pfizer). Ethical approval for the UK Biobank was previously obtained from the North West Centre for Research Ethics Committee (11/NW/0382).

Disclosure forms provided by the authors are available with the full text of this article.

Drs. Rodriguez-Flores, Khalid, Shuldiner and Saleheen contributed equally to this article.

Rodriguez-Flores et al.

25

NOTCH3 p.Arg1231Cys is Markedly Enriched in South Asians and Associated with Stroke

Version 34

17 June 2024

949

Inclusion and Ethics

950

951 The Institutional Review Board (IRB) at the Center for Non-Communicable Diseases (IRB:

952 00007048, IORG0005843, FWAS00014490) approved the study. All participants gave

953 written informed consent.

954

955

Author Contributions Statement

956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993

Conceptualization by Juan L. Rodriguez-Flores (JLRF), Shareef Khalid (SK), Alan R. Shuldiner (ARS), and Danish Saleheen (DS).

Data curation by Nilanjana Banerjee (NB), Deepika Sharma (DeS), Michael Cantor (MC), John Overton (JO), and Jeff Reid (JR).

Formal analysis by Bin Ye (BY), Manav Kapoor (MK), Joshua Backman (JB), Gannie Tzoneva (GT), Ellen Tsai (ET), Sahar Gelfman (SG), Tanima De (TD), Niek Verweij (NV), Luca A. Lotta (LAL), Aaron Zhang (AZ), Neelroop Parikshak (NP), Farshid Sepehrband (FS), Jonathan Marchini (JM), Giovanni Coppola (GC), Sofia Castaneda (SC), Pengcheng Xun (PX), Ellen Tsai (ET), and Regeneron Genetics Center (RGC).

Funding acquisition by DS, ARS, Aris Baras (AB), and RGC.

Investigation by Silvio Alessandro DiGioia (SADG), Hector Martinez (HM), I-Chun Tsai (IT), Katia Karalis (KK), Aris Economides (AE), David D'Ambrosio (DDA), and Asif Rasheed (AR).
Methodology by JLRF, SK, and RGC.

Project administration by Thomas Coleman (TC), RGC, and AR.

Resources by Muhammad Jahanzaib (MJ), Maleeha Zaman (MZ), Muhammad Rehan Mian (MRM), Muhammad Bilal Liaqat (MBL), Khalid Mahmood (KM), Tanvir-us-Salam (TUS), Muhammad Hussain (MH), Ayeesha Kamal (AK), Javed Iqbal (JI), and Faizan Aslam (FA).

Software by RGC, JLRF, and SK. Supervision by DS and ARS. Validation by SK and JLRF.

Visualization by JLRF, NP, FS, JM, and RGC.

Writing of original draft by JLRF, SK, ARS, DS, PX, Sergio Fazio (SF), Wolfgang Liedtke (WL), John Danesh (JD), Ayeesha Kamal (AK), Philippe Frossard (PF), and RGC.

Writing review and editing by JLRF, PX, SK, ARS, DS, and RGC.

994

Competing Interests Statement

995

996 Funding

997

998 Fieldwork for this study was funded by the Center for Non-Communicable Diseases, Pakistan.

999 DNA sequencing was funded by Regeneron Pharmaceuticals Inc.

1000

1001 Employment

1002

1003 JLRF, ARS, NB, DeS, MC, JO, JR, BY, MK, JB, GT, SG, TD, NV, LAL, AZ, NP, FS, JM, GC,

1004 PX, AB, SADG, HM, IT, KK, AE, DDA, SF, WL, and TC are or were employees of Regeneron

1005 Genetics Center LLC or Regeneron Pharmaceuticals Inc. and contributed to this manuscript as

1006 part of their regular duties as salaried employees.

1007

1008 ET and SC are or were student interns of Regeneron Genetics Center LLC or Regeneron

1009 Pharmaceuticals Inc. and contributed to this manuscript as part of their internship activities.

1010

1011 AR, MJ, MZK, MRM, MBL, PF, and DS and SK are or were employees of the Center for Non-

1012 Communicable Disease and received salaried compensation for their contribution to this

1013 manuscript.

1014

1015 Personal Financial Interests

1016

1017 JLRF, ARS, NB, DeS, MC, JO, JR, BY, MK, JB, GT, SG, TD, NV, LAL, AZ, NP, FS, JM, GC,

1018 PX, AB, SADG, HM, IT, KK, AE, DDA, SF, WL, TC are or were employees of Regeneron

1019 Genetics Center LLC or Regeneron Pharmaceuticals Inc. and received stock and stock options as

1020 part of their compensation as employees.

1021

1022 JLRF, ARS, DS, AB, and SK are named inventors on patent pending US 20230000897A1 that

1023 discloses methods of treating subjects having a cerebrovascular disease by administering

1024 Neurogenic Locus Notch Homolog Protein 3 (*NOTCH3*) agents, and methods of identifying

1025 subjects having an increased risk of developing a cerebrovascular disease.

1026

Tables

Table 1. Baseline Characteristics of PGR Stroke Case-control Discovery Cohort¹

	Case (n = 5,135)		Control (n = 26,602)		P value	Case p.Arg1231Cys Carrier (n=103) ³		Case Non- Carrier (n=4,998)		P value
	mean / n	SD / %	mean / n	SD / %		mean / n	SD / %	mean / n	SD / %	
Female, n (%)	2,277	44.3	9,330	35.1	< 0.01	50	48.5	2,211	44.2	0.44
Age at enrolment, years	58.9	13.1	52.8	11.4	< 0.01	58.7	12.1	58.9	13.2	0.87
BMI, kg/m ²	25	3.9	27.5	4.4	< 0.01	24.6	3.4	24.9	3.9	0.35
Cholesterol mg/dL	175.3	55.5	172.1	48.2	< 0.01	174.8	49.1	175.3	55.7	0.94
LDL-C mg/dL	111.1	45.8	101.2	38.2	< 0.01	110.8	41.1	111.1	45.9	0.94
HDL-C mg/dL	37.8	12.7	35	10.8	< 0.01	38.8	12.1	37.8	12.7	0.43
Triglyceride mg/dL	140.3	81.8	185.7	120.1	< 0.01	126.3	62.7	140.5	82.3	0.04
Glucose, mg/dL	148.7	75.5	143.4	84.9	< 0.01	145.2	74.7	148.6	75.3	0.68
HbA1c %	6.9	1.8	6.6	2	< 0.01	7.3	1.9	6.9	1.9	0.15
Creatinine mg/dl	1.2	0.8	0.9	0.5	< 0.01	1.1	0.7	1.2	0.8	0.39
Tobacco or other stimulant user ²	1,850	36	9,242	34.7	< 0.01	32	31.1	0	0	0.34
Comorbidities										
Hypertension , n (%)	2,953	57.5	9,385	35.3	< 0.01	57	55.3	2,879	57.6	1
Diabetes, n (%)	1,195	23.3	5,766	21.7	< 0.01	29	28.2	1,157	23.1	0.29
Myocardial infarction, n (%)	371	6.2	622	2.3	< 0.01	5	4.9	215	4.3	0.98
Family history of										
Stroke, n (%)	324	6.3	0	0	< 0.01	14	13.6	306	6.1	< 0.01
Hypertension, n (%)	573	11.2	4,299	16.2	< 0.01	15	14.6	556	11.1	0.35
Diabetes, n (%)	349	6.8	5,199	19.5	< 0.01	11	10.7	334	6.7	0.16
Sudden death, n (%)	150	2.9	569	2.1	< 0.01	8	7.8	142	2.8	< 0.01

1031
1032 ^{1.} Shown are mean or total n and standard deviation or percentage for n = 5,135 cases versus n
1033 = 26,602 controls on the left and n = 103 case p.Arg1231Cys *NOTCH3* carriers versus n =
1034 4,998 case non-carriers on the right. Comparison was conducted using R and p values are
1035 shown from chi-square test for categorical variables (Fisher's exact test if cell size was <5)
1036 and from T-test for continuous variables.

1037 ^{2.} Tobacco or other stimulants include cigarettes, paan (chewed betel leaf and areca nut),
1038 naswar (snuff), gutka (chewing tobacco), huqqa (water pipe), chillum (hasish pipe).

1039 ^{3.} Genotypes for the p.Arg1231Cys variant were not available for n=34.

1040

1041
1042
1043

Table 2. UKB Ischemic Stroke Association Across *NOTCH3* Variant Classes and Domains¹

<i>NOTCH3</i> Variant Class	Effect [OR (95%CI)]	P value	AAF (%)	Cases (RR RA AA)	Controls (RR RA AA)
p.Arg1231Cys	3.38 (1.65,6.94)	8.8x10 ⁻⁴	0.02	9124 11 0	370986 139 0
EGFr 1-34 Cys-altering	2.86 (2.14,3.82)	6.3x10 ⁻¹⁰	0.01	9094 49 0	370693 709 1
EGFr 1-6 Cys-altering	29.51 (10.39,83.82)	1.4x10 ⁻⁷	0.002	9137 6 0	371394 9 0
EGFr 7-34 Cys-altering	2.55 (1.87,3.46)	1.6x10 ⁻⁷	0.098	9100 43 0	370702 700 1
Non-EGFr Cys-altering	0.97 (0.46,2.03)	9.3x10 ⁻¹	0.039	9136 7 0	371111 292 0
EGFr 1-34 Ser-altering	0.98 (0.8,1.2)	8.4x10 ⁻¹	0.54	9046 97 0	367355 4042 6
EGFr 1-6 Ser-altering	0.76 (0.23,2.56)	6.6x10 ⁻¹	0.018	9141 2 0	371271 132 0
EGFr 7-34 Ser-altering	0.99 (0.8,1.21)	8.9x10 ⁻¹	0.53	9048 95 0	367486 3911 6
Non-EGFr Ser-altering	0.84 (0.52,1.34)	4.5x10 ⁻¹	0.10	9128 15 0	370656 744 3
LoF variants	1.38 (0.50,3.85)	5.4x10 ⁻¹	0.019	9138 5 0	371261 142 0
LoF + any missense	1.09 (1.01,1.19)	3.3x10 ⁻²	3.30	8490 652 1	346648 24710 45

1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058

1. Burden test with age, age², gender, 10 PCs as covariates was conducted using REGENIE to compare the association signal and effect size across variant classes and domains for ischemic stroke in n = 9,143 cases and n = 371,403 controls from UK Biobank. *NOTCH3* variant classes include the single variant association (p.Arg1231Cys) for reference at the top and burden tests limited to: Cys-altering variants in EGFr domains 1 to 34; Cys-altering variants in EGFr domains 1 to 6; Cys-altering variants in EGFr domains 7 to 34; Cys-altering variants outside of EGFr domains; Ser-altering variants in EGFr domains 1 to 6; Ser-altering variants in EGFr domains 7 to 34; Ser-altering variants outside of EGFr domains; Loss of function (LoF) variants (defined as stop-gain, stop-loss, frameshift and splice-site variants with AAF <1%); and LoF and any missense variants (AAF < 1%). Abbreviations: EGFr, epidermal growth factor-like repeat domain; AAF, alternate allele frequency; RR, reference allele homozygote; RA, reference/alternate allele heterozygote; AA, alternate allele homozygote; 95% CI, 95% confidence interval

Figure Legends

1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082

Figure 1. ExWAS Identifies *NOTCH3* p.Arg1231Cys Associated with Subcortical Hemorrhagic Stroke in Pakistan Genome Resource 31 K Discovery Cohort. **A.** Flow chart of the study described in this report. The discovery cohort consisted of a 31 K stroke case-control cohort (n = 31,737, including n = 5,135 stroke cases and n = 26,602 controls) from the Pakistan Genome Resource (PGR) (green boxes). A second PGR follow-up cohort of 44 K (n = 44,082) included 30 K participants with self-reported stroke case:control status for replication (n = 30,399, including n = 160 cases and n = 30,239 controls). UK Biobank data from 450 K sequenced participants was used for further analysis in a predominantly European ancestry population (blue boxes), 380 K of whom had stroke case:control status known (n = 9,143 cases and n = 371,403 controls), and 35 K of whom had brain MRI data (n = 35,344). **B.** Manhattan plot of subcortical hemorrhagic stroke ExWAS in 31 K PGR discovery cohort participants (n = 1,388 cases and n = 26,602 controls) with likelihood ratio test $-\log_{10}$ p-values of calculated using REGENIE (y-axis) across chromosomes (alternating grey and black dots) and variants (x-axis). A single variant (NC_000019.10:g.15179052G>A) on chromosome 19 predicting a missense variant p.Arg1231Cys in *NOTCH3* (pink diamond) exceeded the genome-wide significance threshold of 5×10^{-8} (red line). **C.** *NOTCH3* locus zoom plot of subcortical stroke ExWAS. The likelihood ratio test $-\log_{10}$ p values for variants tested are shown on the y-axis. The p.Arg1231Cys variant is labeled as a diamond. Other variants (circles) are colored based on linkage disequilibrium with the reference variant in 1000 Genomes [38]. Gene exon (thick line) and intron (thin line) model shown below the graph.

1083 **Figure 2. *NOTCH3* EGFr Domain Disruption by p.Arg1231Cys.** Shown is *NOTCH3*
1084 p.Arg1231 in context of human *NOTCH3* protein domains and cross-species alignment of
1085 *NOTCH3* amino acid sequences. **A.** Human *NOTCH3* Protein Domains. Shown is the position of
1086 the associated variant in context of transcript exons (top, alternating blue and purple with
1087 numbering) and protein domains (bottom, color coded). *NOTCH3* can be divided into four major
1088 regions, from left-to-right the signal peptide (light blue), the extra-cellular domain (ECD,
1089 brown), the transmembrane domain (orange), and the intra-cellular domain (ICD, blue). The
1090 majority of the ECD is composed of 34 Epidermal Growth Factor-like repeat (EGFr) domains (in
1091 purple with white numbers). Domains involved in signaling are highlighted, including EGFr
1092 domains 10 to 11 involved in ligand binding (light purple, numbered), three cleavage domains
1093 (S1 in yellow, S2 in yellow with diagonal black stripes, S3 in yellow with black horizontal
1094 stripes), and three *Lin12/NOTCH* repeats (light green, numbered). The ICD contains the
1095 Recombination signal-binding protein for Ig of κ region (RAM) domain for transcription factor
1096 interaction (green and white checkers), the Nuclear Localization Sequences (NLS, orange with
1097 black stripes), five Ankyrin repeats involved in signal transduction (green with white numbers),
1098 and the Proline, glutamic acid, serine, and threonine-rich (PEST) domain essential for
1099 degradation (green with white stripes). The p.Arg1231Cys variant (red line top to bottom)
1100 removes a disulfide-bridge-forming cysteine in the 31st EGFr domain of the ECD, coded by the
1101 22nd exon. **B.** Cross-species Alignment of *NOTCH3* (EGFr) Domain # 31 Amino Acid
1102 Sequences calculated using BLAST. Shown is an amino acid alignment of 31st EGFr domain of
1103 *NOTCH3* (human sequence amino acids 1205 to 1244), including (top-to-bottom) human
1104 reference, human p.Arg1231Cys mutant, chimpanzee (*Pan troglodytes*), mouse (*Mus musculus*),
1105 zebrafish (*Danio rerio*), western clawed frog (*Xenopus tropicalis*), and green sea turtle (*Chelonia*
1106 *mvdas*), indicating conservation of the arginine (R) at position 1231 in mammals. Highly-
1107 conserved cysteine (C) residues (normally 6 per EGFr) are highlighted in yellow.
1108
1109

1110 **Figure 3. Forest Plot Showing Replication of *NOTCH3* p.Arg1231Cys Association with**
1111 **Stroke Across 61 K Pakistan Genome Resource Meta-Analysis.** Shown is the cohort name,
1112 trait, odds ratio with 95% confidence interval, likelihood ratio test p value calculated using
1113 REGENIE, alternate allele frequency, case count, and control count for five stroke phenotypes
1114 in the PGR 31 K discovery cohort (n = 5,135 stroke cases and n = 26,602 controls) (top), and
1115 Inverse Variance Weighted (IVW) Meta-Analysis using METAL of stroke in 61 K PGR Cohort,
1116 including 31 K PGR discovery cohort and 30 K PGR replication cohort (n = 160 cases and n =
1117 30,239 controls) subset of the 44 K PGR follow-up cohort (bottom).
1118
1119

Rodriguez-Flores et al.

33

NOTCH3 p.Arg1231Cys is Markedly Enriched in South Asians and Associated with Stroke

Version 34

17 June 2024

1120

1121

1122

1123

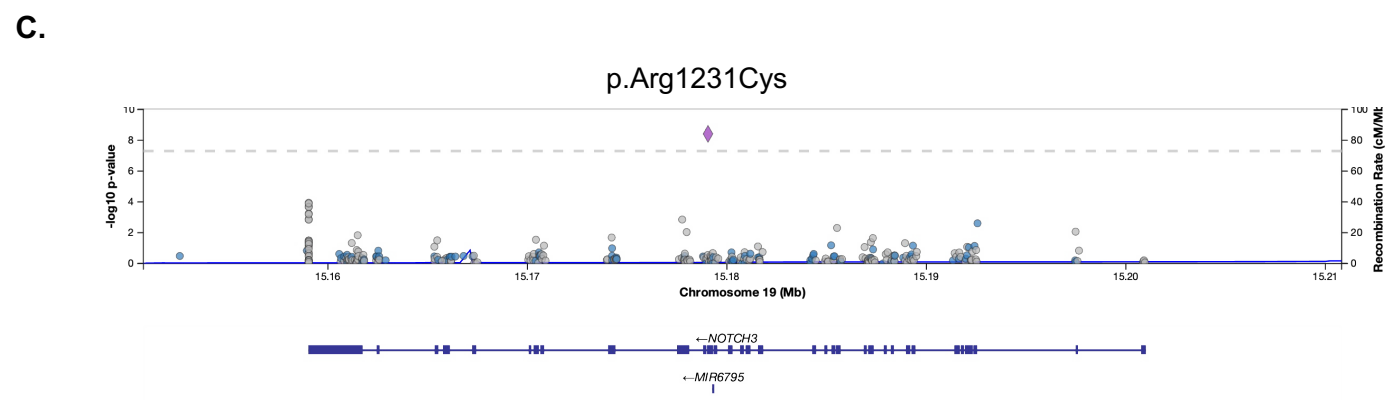
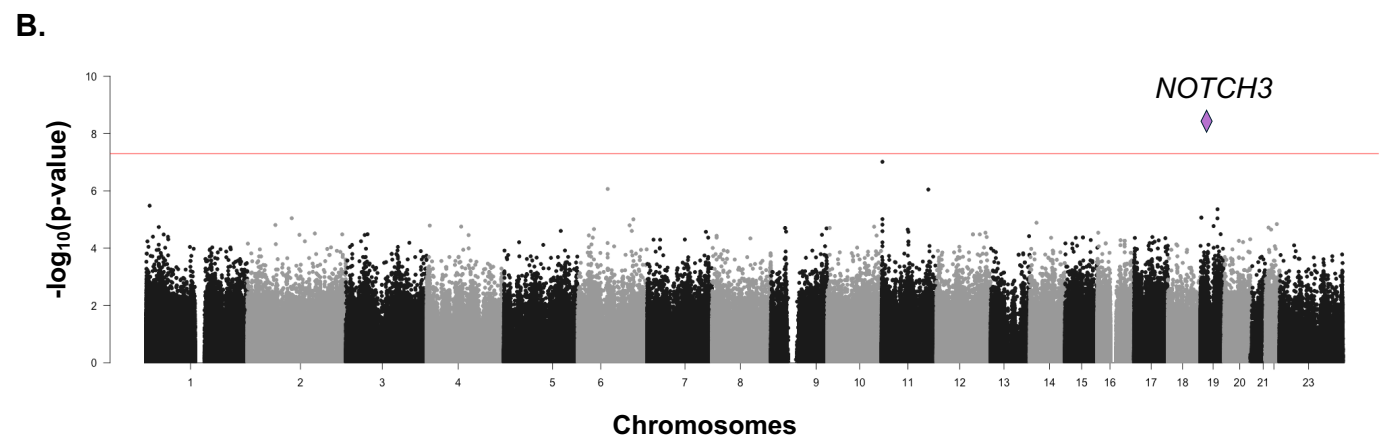
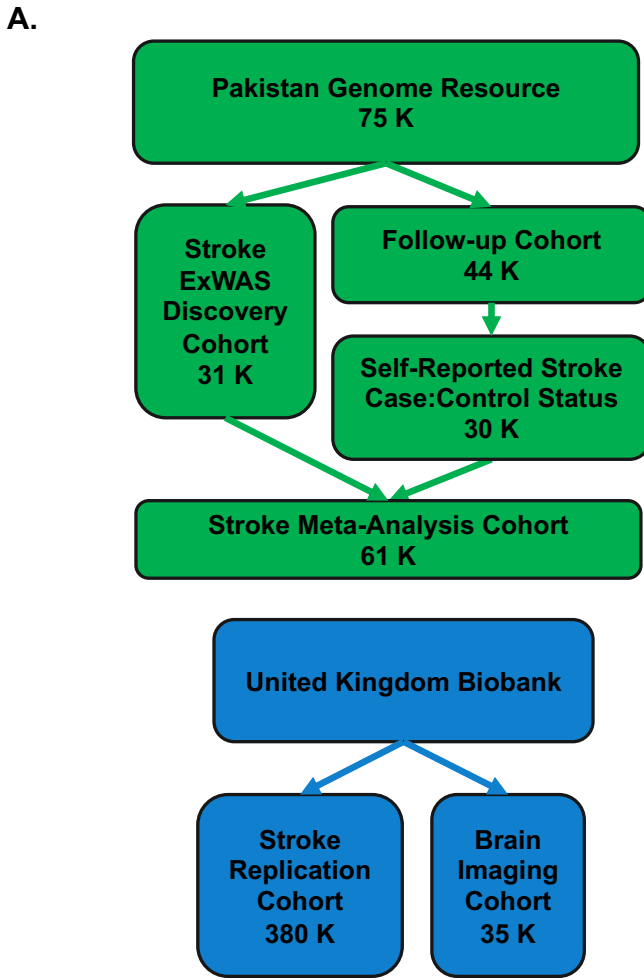
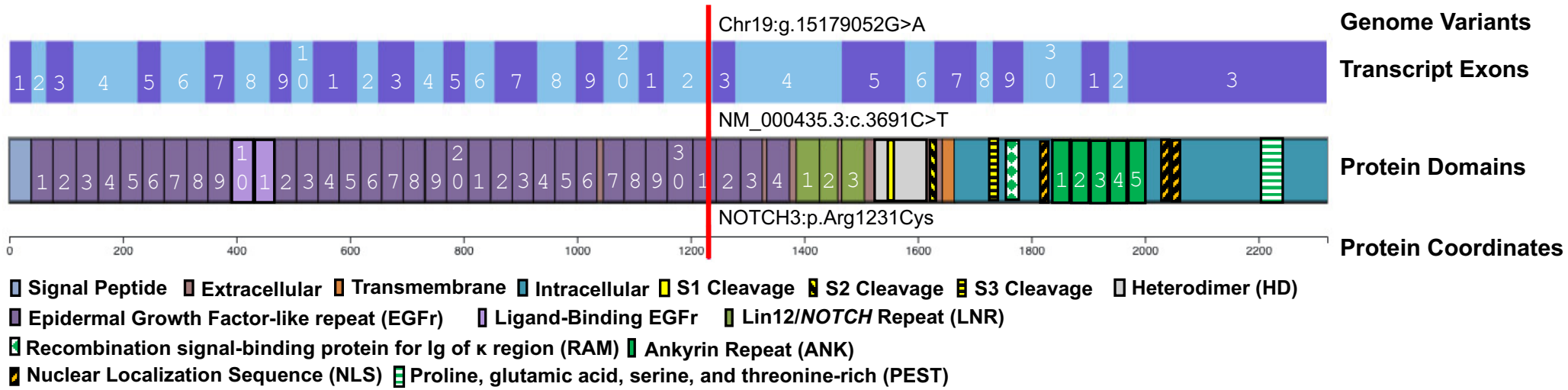


Figure 1

A.



B.

1205	INECRSGA	CHAAHTRD	CLQDPGGGFR	CL	CHAGFSGPR	CQT	Human 1231 Arg (reference)
	INECRSGA	CHAAHTRD	CLQDPGGGF	CCL	CHAGFSGPR	CQT	Human 1231 Cys (mutant)
	INECRSGA	CHAAHTRD	CLQDPGGGFR	CL	CHAGFSGPR	CQT	Chimpanzee
	INECRPGA	CHAAHTRD	CLQDPGGHFR	CV	CHPGFTGPR	CQI	Mouse
	INECLSNP	CNPSNSLD	CIQLPND-YQ	CV	CKPGFTGRG	CQS	Zebrafish
	INECLSGP	CHAQNTRH	CVQLAND-YQ	CV	CKSGYTGRRC	CQS	Western clawed frog
	INECLAKP	CLPQRTLD	CVQGAND-FQ	CL	CKPGYTGRRC	QN	Green sea turtle

Figure 2

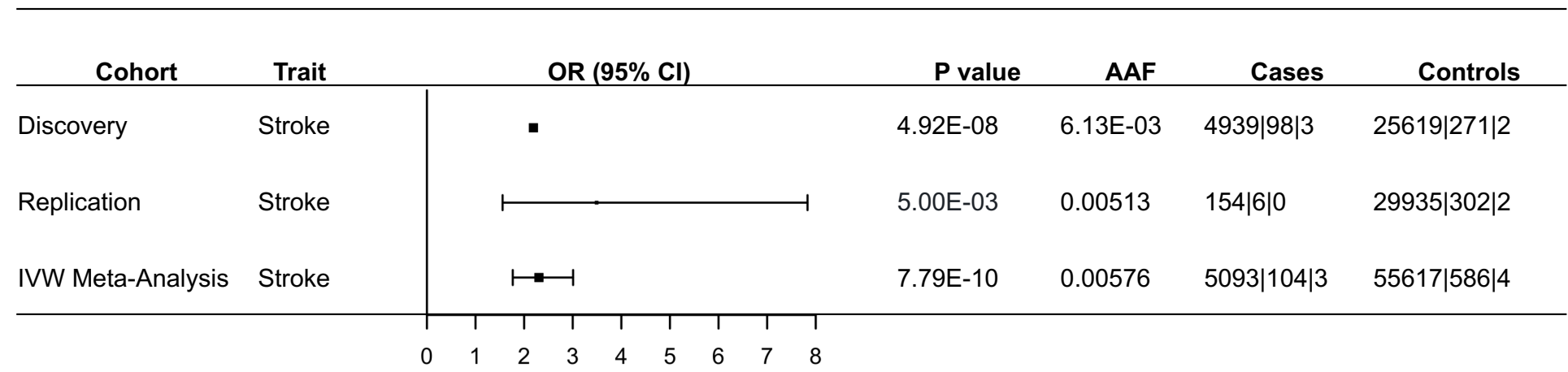
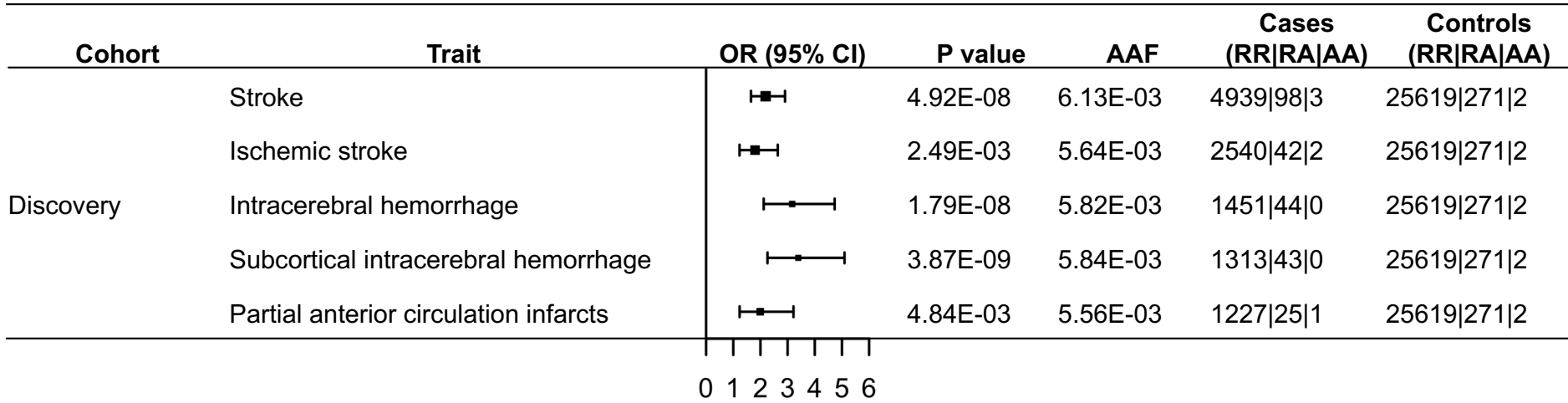


Figure 3