

Rodriguez-Flores et al.

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

Version 33

20 March 2024

1

## Title

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

## Authors

Juan L. Rodriguez-Flores<sup>1</sup>, Shareef Khalid<sup>2</sup>, Neelroop Parikshak<sup>1</sup>, Asif Rasheed<sup>3</sup>, Bin Ye<sup>1</sup>, Manav Kapoor<sup>1</sup>, Joshua Backman<sup>1</sup>, Farshid Sepehrband<sup>1</sup>, Silvio Alessandro DiGioia<sup>4</sup>, Sahar Gelfman<sup>1</sup>, Tanima De<sup>1</sup>, Nilanjana Banerjee<sup>1</sup>, Deepika Sharma<sup>1</sup>, Hector Martinez<sup>4</sup>, Sofia Castaneda<sup>5</sup>, David D'Ambrosio<sup>4</sup>, Xingmin A. Zhang<sup>1</sup>, Pengcheng Xun<sup>4</sup>, Ellen Tsai<sup>6</sup>, I-Chun Tsai<sup>4</sup>, Regeneron Genetics Center<sup>1</sup>, Maleeha Zaman Khan<sup>3</sup>, Muhammad Jahanzaib<sup>3</sup>, Muhammad Rehan Mian<sup>3</sup>, Muhammad Bilal Liaqat<sup>3</sup>, Khalid Mahmood<sup>7</sup>, Tanvir Us Salam<sup>8</sup>, Muhammad Hussain<sup>8</sup>, Javed Iqbal<sup>10</sup>, Faizan Aslam<sup>11</sup>, Michael N. Cantor<sup>1</sup>, Gannie Tzoneva<sup>1</sup>, John Overton<sup>1</sup>, Jonathan Marchini<sup>1</sup>, Jeff Reid<sup>1</sup>, Aris Baras<sup>1</sup>, Niek Verweij<sup>1</sup>, Luca A. Lotta<sup>1</sup>, Giovanni Coppola<sup>1</sup>, Katia Karalis<sup>1</sup>, Aris Economides<sup>1</sup>, Sergio Fazio<sup>4</sup>, Wolfgang Liedtke<sup>4</sup>, John Danesh<sup>12</sup>, Ayeesha Kamal<sup>9</sup>, Philippe Frossard<sup>3</sup>, Thomas Coleman<sup>1</sup>, Alan R. Shuldiner<sup>1</sup>, Danish Saleheen<sup>2,3</sup>

## Affiliations

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

## Correspondence

Danish Saleheen, danish.saleheen@cncdpc.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 33

20 March 2024

2

## Abstract

1  
2  
3 The genetic factors of stroke in South Asians are largely unexplored. Exome-wide sequencing  
4 and association analysis (ExWAS) in 75 K Pakistanis identified  
5 NM\_000435.3(*NOTCH3*):c.3691C>T, encoding the missense amino acid substitution  
6 p.Arg1231Cys, enriched in South Asians (alternate allele frequency = 0.58% compared to  
7 0.019% in Western Europeans), and associated with subcortical hemorrhagic stroke [odds ratio  
8 (OR) = 3.39, 95% confidence interval (CI) = [2.26, 5.10], p value =  $3.87 \times 10^{-9}$ ], and all strokes  
9 (OR [CI] = 2.30 [3.01, 1.77], p value =  $7.79 \times 10^{-10}$ ). *NOTCH3* p.Arg231Cys was strongly  
10 associated with white matter hyperintensity on MRI in United Kingdom Biobank (UKB)  
11 participants (effect [95% CI] in SD units = 1.1 [0.61, 1.5], p value =  $3.0 \times 10^{-6}$ ). The variant is  
12 attributable for approximately 5.5% of hemorrhagic strokes and 1% of all strokes in South  
13 Asians. These findings highlight the value of diversity in genetic studies and have major  
14 implications for genomic medicine and therapeutic development in South Asian populations.  
15

## Introduction

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% of strokes and 5% 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.

## Results

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

Characteristics of the 5,135 stroke cases and 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 =  $3.87 \times 10^{-9}$ ; 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 =  $4.44 \times 10^{-8}$ ) 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 31<sup>st</sup> 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

Rodriguez-Flores et al.

5

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

Version 33

20 March 2024

1 An additional 44 K Pakistani were recruited and whole-exome sequenced by the Broad Institute,  
2 approximately 30 K of whom self-reported stroke case:control status was obtained (160 cases  
3 and 30,239 controls. Replication of the association was observed in this independent cohort (OR  
4 [95% CI] = 3.49 [1.56, 7.83], p value =  $5.00 \times 10^{-3}$ . Meta-analysis for all strokes in the combined  
5 61 K cohort achieved a genome-wide significant p value (OR [95% CI]) = 2.30 [1.76, 2.99], p  
6 value =  $7.08 \times 10^{-10}$ ) (**Figure 3**).

### 7 8 **Recall by Genotype**

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

15  
16 Eight out of twelve (66%) p.Arg1231Cys homozygotes had a history of hypertension. Among  
17 the 128 callback participants, both systolic and diastolic blood pressure were trending higher  
18 (**Supplementary Table 8**). While there was no association with hypertension in the discovery  
19 cohort, p.Arg1231Cys homozygotes in the PGR had nominally higher diastolic blood pressure  
20 than heterozygotes or homozygous reference individuals (median = 95 mmHg, interquartile  
21 range 86 to 100; heterozygotes (median = 80 mmHg, interquartile range = 80 to 90), p value =  
22 0.016 (**Supplementary Table 9**).

### 23 24 **Allele Frequency and Population Attributable Risk**

25  
26 The allele frequency of p.Arg1231Cys was 1.1% across the PGR 75 K, equivalent to a  
27 population prevalence of 1 in 46. After removing cases recruited for cardiovascular diseases  
28 (individuals enrolled at time of acute stroke, MI, and heart failure), the allele frequency of the  
29 variant was 0.51%, equivalent to a population prevalence of 1 in 98. This frequency was orders  
30 of magnitude higher relative to exomes of European ancestry from UK Biobank (AAF =  
31 0.019%), corresponding to a population prevalence of 1 in 2,614). The variant was enriched  
32 (AAF > 0.1%) in other South Asian and West Asian populations [15] both within and outside of  
33 Pakistan (**Supplementary Data, Supplementary Tables 10 and 11, Supplementary Figure 5**).

34  
35 We estimate that 5.5% [bootstrap 95% CI based on 10,000 resamples: 4.2% to 6.7%] of  
36 hemorrhagic strokes and 1.0% [bootstrap 95% CI based on 10,000 resamples: 0.6% to 1.4%] of  
37 all strokes in the Pakistani population are attributable to p.Arg1231Cys. Thus, this variant is a  
38 common cause of strokes in SAS and WAS populations, a finding that has implications for  
39 medical care as well as global health in these populations.

### 40 41 **Suggestive Associations at Other Loci**

42  
43 Although *NOTCH3* p.Arg1231Cys was the only variant associated with stroke at a genome-wide  
44 significant p value below  $5.0 \times 10^{-8}$ , there were a total of 9 associations (5 loci) with p values  
45 below  $1.0 \times 10^{-6}$  and at least 10 variant carriers (**Supplementary Table 12**). In addition to

Rodriguez-Flores et al.

6

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

Version 33

20 March 2024

1 *NOTCH3*, these included one known locus previously associated with stroke in a recent GWAS,  
2 lymphocyte specific protein *LSP1* [10].

3  
4 The *LSP1* locus variant with suggestive association with intracerebral hemorrhagic stroke in this  
5 study was a common intronic variant (rs661348, OR [95% CI] = 1.3 [1.2, 1.4], p value =  $8.0 \times 10^{-8}$ ,  
6 AAF = 0.27). While rs661348 was not previously associated with stroke, two common non-  
7 coding variants in the *LSP1* locus were previously reported to be associated with stroke  
8 (rs569550 and rs1973765) [10]. Both variants were in linkage disequilibrium with rs661348 ( $r^2 >$   
9 0.4 in 10 K unrelated PGR participants). A test of association for rs661348 conditional on these  
10 variants reduced the strength of the association for rs661348 to nominal ( $8.83 \times 10^{-4}$ )  
11 (**Supplementary Table 17**), suggesting that rs661348 represents the same known stroke risk  
12 locus. The *LSP1* locus was previously reported to be associated with hypertension [16]. In  
13 PGR there was a weak association with hypertension (p value =  $2.61 \times 10^{-2}$ ) (**Supplementary**  
14 **Tables 13, 14, 15, and 16**).

15  
16 ***NOTCH3* p.Arg1231Cys is associated with stroke and CADASIL-like phenotypes in UK**  
17 **Biobank**

18  
19 To investigate stroke and CADASIL-related phenotypes in an independent cohort, UK Biobank  
20 data was reviewed for associations with *NOTCH3* p.Arg1231Cys. A total of 255 heterozygotes  
21 for *NOTCH3* p.Arg1231Cys were observed in 450 K exome-sequenced individuals from this  
22 predominantly European cohort, with a markedly lower allele frequency (AAF = 0.019%)  
23 (**Supplementary Table 18**). Phenome-wide association (PheWAS) of 10,168 phenotypes  
24 revealed nominally significant association of p.Arg1231Cys with ischemic stroke (OR [95% CI]  
25 = 4.0 [1.9, 8.6]), p value =  $4.1 \times 10^{-4}$ ), all strokes combined (OR [95% CI] = 1.9 [1.1, 3.5], p value  
26 = 0.031), hypertension (ICD 10 code I10) (OR [95% CI] = 1.5 [1.1, 2.2], p value = 0.019), and  
27 recurrent major depression (OR [95% CI] = 3.2 [1.5, 6.8], p value = 0.0031) (**Supplementary**  
28 **Table 19**). No association was observed for hemorrhagic stroke, migraine, dementia, mood  
29 changes, Alzheimer's disease, or urinary incontinence. The lack of association (p value = 0.066)  
30 with hemorrhagic stroke in UKB Europeans was likely due to low statistical power, given the  
31 lower variant allele frequency and lower hemorrhagic stroke prevalence in UKB compared to  
32 PGR. Nonetheless, the odds ratio (OR [95% CI] = 5.8 [0.88, 39.1]) was high.

33  
34 In addition to recurrent strokes, brain white matter loss is a major and early phenotype  
35 characteristic of CADASIL that is focused on particular brain regions [8, 14]. *NOTCH3*  
36 p.Arg1231Cys was strongly associated with a cluster of brain MRI quantitative phenotypes, e.g.,  
37 total volume of white matter hyperintensities (WMH) from T1 and T2 FLAIR images (effect  
38 [95% CI] in SD units = 1.1 [0.61, 1.5], p value =  $3.0 \times 10^{-6}$ ) with carriers having 7.4 cm<sup>3</sup> more  
39 WMH volume than controls (**Supplementary Figure 6** and **Supplementary Table 20**). The  
40 most prominent alterations in WMH in p.Arg1231Cys carriers were observed in the centrum  
41 semiovale and periventricular white matter (**Supplementary Figure 7**). Taken together, these  
42 results demonstrate *NOTCH3* p.Arg1231Cys carriers have increased risk of established markers  
43 of small vessel disease and clinical phenotypes observed in CADASIL [8].

44  
45 **Pathogenic burden of all Cys-altering variants within *NOTCH3* EGFr domains specifically**  
46 **associated with CADASIL phenotypes in UK Biobank**

Rodriguez-Flores et al.

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

Version 33

20 March 2024

7

1  
2 Burden test analysis allows for increased statistical power to detect association by combining  
3 signal across multiple rare variants. Prior studies have shown that pathogenic variants in  
4 CADASIL are limited to variants that add or remove a Cysteine (Cys-altering) in *NOTCH3* EGFr  
5 domains normally containing 6 Cysteines. Furthermore, patients with Cys-altering variants in the  
6 first 6 EGFr domains have more severe symptoms than in EGFr domains 7 to 34 [14, 17, 18],  
7 including larger regions of brain white matter loss [19], more granular osmophilic material  
8 (GOM) aggregates in blood vessels [19], and worse prognosis [14].

9  
10 In order to test these hypotheses, a set of custom gene burden tests were designed and compared  
11 to single variant test results for *NOTCH3* p.Arg1231Cys. In UKB, 758 individuals carried one of  
12 98 unique Cys-altering variants across the 34 EGFr domains in *NOTCH3* (**Supplementary**  
13 **Tables 21 and 22**). A burden test aggregating all UKB EGFr domain Cys-altering variants into  
14 a single statistical test was strongly associated with stroke (OR [95% CI] = 2.86 [2.14, 3.82], p  
15 value =  $6.29 \times 10^{-10}$ ; AAF = 0.01%) (**Table 2**). In contrast to Cys-altering variants within EGFr  
16 domains, Cys-altering variants outside of EGFr domains were not associated with stroke (OR  
17 [95% CI] = 0.97 [0.46, 2.03], p value =  $9.3 \times 10^{-1}$ ; AAF = 0.039%) (**Table 2**). In order to rule out  
18 the possibility that any missense variants in EGFr domains are associated with stroke, a test  
19 limited to the most commonly altered (added or removed) amino acid in *NOTCH3*, serine (Ser),  
20 was tested and did not show any evidence of association with stroke (OR [95% CI] = 0.98 [(0.8,  
21 1.2], p value = 0.084; AAF = 0.54%) (**Table 2, Supplementary Table 23**). Interestingly, a  
22 burden test limited only to predicted loss of function (pLoF) variants (frameshift, splice variant,  
23 stop gain) did not show significant evidence for association with stroke (OR [95% CI] = 1.38  
24 [0.50, 3.85], p value = 0.54; AAF = 0.019%) (**Table 2, Supplementary Table 24**). These results  
25 provide evidence to support the hypothesis that EGFr domain Cys-altering variants within  
26 *NOTCH3* are associated with stroke, in contrast to other protein-altering variants.

27  
28 While hemorrhagic stroke represents a small proportion of the strokes reported in the UKB, the  
29 set of Cys-altering variants were also tested for association with hemorrhagic stroke. A nominal  
30 association with hemorrhagic stroke (OR [95% CI] = 3.61 [1.39, 9.34], p value =  $8.31 \times 10^{-3}$ ; AAF  
31 0.025%) was observed, despite low statistical power.

32  
33 Consistent with stroke risk, in MRI data of 35,344 UKB individuals, Cys-altering variants in  
34 *NOTCH3* EGFr domains were strongly associated with WMH volume (p value =  $3.7 \times 10^{-13}$ ; with  
35 carriers having 5.4 cm<sup>3</sup> greater WMH volume than controls). These WMH differences were  
36 strongest in the centrum semiovale and periventricular white matter (**Supplementary Figure 7**).  
37 Additionally, we found strong WMH signal in the external capsule, which is known to be  
38 involved in CADASIL. We found weaker evidence for association of *NOTCH3* LoF variants  
39 with WMH (effect size [95% CI] = 6.8 cm<sup>3</sup> [4.2 cm<sup>3</sup>, 10.9 cm<sup>3</sup>], p-value =  $1.68 \times 10^{-4}$ ).

40  
41 Prior studies have binned *NOTCH3* EGFr domain Cys-altering variants in up to three distinct  
42 severity or risk groups based on EGFr domain number [10, 14, 17]. Indeed, we observed a much  
43 larger effect size for Cys-altering variants in high-risk EGFr domains 1-6 (OR [95% CI] = 29.5  
44 [10.4, 83.8], p value =  $1.37 \times 10^{-7}$ ; AAF = 0.002%) compared to Cys-altering variants in EGFr  
45 domains 7-34 (OR [95% CI] = 2.55 [1.87, 3.46], p value =  $1.59 \times 10^{-7}$ ; AAF = 0.098%) (**Table 2,**  
46 **Supplementary Table 23**). These results are consistent with prior reports of differences in

Rodriguez-Flores et al.

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

Version 33

20 March 2024

8

1 stroke risk between EGFr domain risk groups not correlated with differences in signaling activity  
2 between EGFr risk groups [17].

## 3 4 Discussion

5  
6 This report describes the largest ExWAS of stroke conducted thus far in a South Asian  
7 population and highlights a Cys-altering missense variant in the 31<sup>st</sup> EGFr domain of *NOTCH3*  
8 associated with stroke at a genome-wide level of statistical significance. This is the first study to  
9 report a genome-wide-significant association between *NOTCH3* and stroke, a discovery enabled  
10 because *NOTCH3* p.Arg1231Cys is markedly enriched in Pakistanis compared to Western  
11 European and non-Eurasian populations. Harbored in ~1 percent of Pakistani, p.Arg1231Cys is  
12 associated with a ~3-fold increased risk of hemorrhagic stroke. While some regional variability  
13 in the allele frequency is observed, p.Arg1231Cys is enriched in populations ranging from  
14 Turkey in West Asia to India in South Asia, suggesting a substantial contribution to stroke risk in  
15 millions of individuals across South Asia and West Asia as well as their global diaspora.

16  
17 *NOTCH3* was not previously associated with stroke in the largest GWAS predominantly  
18 consisting of European-derived participants [10]. In contrast to prior studies, both the discovery  
19 and replication cohorts in this study were South Asian, hence avoiding the bias encountered in  
20 studies with a European discovery cohort. Given the much lower allele frequency of  
21 p.Arg1231Cys in European populations, we observed a nominal association between  
22 p.Arg1231Cys and stroke in the UK Biobank study, showing a similar effect size as in South  
23 Asians. Nominal associations were also observed for phenotypes related to CADASIL, such as  
24 hypertension and depression. While brain images were not available for the Pakistani cohort, a  
25 strong association was observed between p.Arg1231Cys and quantitative brain MRI phenotypes  
26 in UKB data, such as white matter hyperintensity.

27  
28 Cys-altering mutations in proximal EGFr domains of *NOTCH3* are known to cause autosomal  
29 dominant CADASIL, a rare highly penetrant distinct syndrome that includes early onset  
30 recurrent subcortical strokes. In contrast to classical CADASIL pathogenic variants,  
31 p.Arg1231Cys is in the 31<sup>st</sup> of 34 EGFr domains, appears to have more moderate penetrance, and  
32 is not obviously clinically distinguishable from more common multi-factorial forms of stroke in  
33 South Asians. The p.Arg1231Cys *NOTCH3* variant is currently classified in ClinVar and recent  
34 reviews [14] as a variant of uncertain significance [20] or “low risk” [17]; however, based on our  
35 current findings, there is strong genetic, computational, and imaging evidence of pathogenicity  
36 for this variant despite reduced penetrance and severity compared to “classic” Cys-altering  
37 CADASIL pathogenic variants in EGFr domains 1 to 6 [14, 19].

38  
39 Prior studies have debated if the mechanism whereby Cys-altering variants contribute to  
40 CADASIL-related pathology is through toxic aggregate gain of function (GoF) or a loss of  
41 normal signaling function (LoF). One study demonstrated excess risk of CADASIL-related  
42 phenotypes for Cys-altering variants in EGFr domains 1 to 34 [18], while other studies showed  
43 greater risk in EGFr domains 1 to 6 relative to EGFr domains 7 to 34 [14, 19]. A recent study  
44 showed evidence for expanding the high-risk tier of EGFr domains to include domains 8, 11, and  
45 26 [17]. The current study provides three additional refinements. First, this study is the first to  
46 assess risk for LoF variants, and did not observe significant association signal (**Table 2**),



Rodriguez-Flores et al.

9

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

Version 33

20 March 2024

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

14  
15 CADASIL is characterized by both ischemic and hemorrhagic strokes, although the factors that  
16 contribute to the manifestation of one versus the other stroke type awaits further clarification [8].  
17 Hemorrhagic stroke appears to represent a larger proportion of strokes in South Asia than in  
18 Europe [2]. In this study, the p.Arg1231Cys association signal was stronger in PGR for  
19 hemorrhagic strokes than for ischemic strokes, despite nearly two-fold larger ischemic stroke  
20 case counts. In contrast, the UKB association signal appeared stronger for ischemic stroke,  
21 possibly due to low hemorrhagic stroke case count and thus statistical power in this cohort.  
22 Statistical power issues aside, differences in manifestation of p.Arg1231Cys in South Asians  
23 compared to Europeans may be attributable to differences in risk factors such as age, blood  
24 pressure, diabetes, air pollution, smoking, medications such as anti-platelet and anti-coagulants  
25 used to manage atherosclerotic disease, genetic background, or study-specific differences in  
26 criteria to categorize stroke sub-types. Further research is needed to better ascertain the  
27 mechanism behind cerebral arterial wall pathobiology and clinical presentation of ischemic  
28 versus hemorrhagic stroke in p.Arg1231Cys carriers.

29  
30 Currently there are no known effective preventive or therapeutic interventions for CADASIL or  
31 less penetrant forms of *NOTCH3* related stroke. However, our analyses provide clues toward  
32 their development. First, in contrast to our analyses of EGFr domain Cys-altering missense  
33 variants in *NOTCH3* that were significantly associated with stroke, predicted loss of function  
34 variants that would be expected to not produce a functional protein were not significantly  
35 associated with stroke. These observations suggest that targeting therapeutic interventions that  
36 decrease expression of mutant protein (such as siRNA, antisense oligonucleotides, and CRISPR),  
37 induce exon skipping of altered EGFr domains [22, 23], or accelerate removal of GOM may  
38 prove beneficial for prevention and/or treatment [24].

39  
40 Our analyses also suggest that p.Arg1231Cys is modestly associated with hypertension, although  
41 p.Arg1231Cys association with stroke risk appears independent of hypertension or other stroke  
42 risk factors such as smoking, age and sex. Animal models of CADASIL show decreased vascular  
43 tone and contractility, most likely driven by loss of physiologic function and subsequent  
44 degeneration of vascular smooth muscle cells (VSMCs) [25]. These observations suggest that  
45 while management of hypertension and smoking cessation are effective modalities for primary  
46 and secondary prevention of stroke, those with *NOTCH3* mutation related strokes will need

1 additional therapeutic interventions, as existing hypertensive medications cannot restore VSMC  
2 function.

3  
4 A limitation of this study is the lack of brain imaging analysis for the Pakistani carriers, such that  
5 specific brain regions affected by the ischemic and hemorrhagic strokes could be ascertained and  
6 compared. In addition, we lacked more detailed clinical data such as presence of migraines and  
7 longitudinal data of disease course including stroke recurrence, dementia, and depression.

8 Further characterization of p.Arg1231Cys carriers will be necessary to obtain better estimates of  
9 penetrance as well as to identify distinguishing clinical or biomarker characteristics that may  
10 have utility in early diagnosis, prevention and treatment, and for recommendations for cascade  
11 screening in family members. Migraine symptoms typically precede stroke by 10+ years in  
12 CADASIL patients [8], thus the combination of migraine with aura, depression and family  
13 history of stroke could be sufficient evidence to prescribe *NOTCH3* genetic testing. Finally, the  
14 effect of LoFs on stroke risk will require larger sample sizes for more definitive comparison to  
15 stroke risk of Cys-altering variants.

16  
17 In conclusion, we identified a highly enriched Cys-altering variant in *NOTCH3* in South Asians  
18 that expands the phenotypic spectrum of CADASIL from rare and highly penetrant to common  
19 and moderately penetrant. Based on our estimates, this single variant may be responsible for  
20 ~1.0% of all strokes combined and ~5.5% of hemorrhagic strokes in South Asians. Among 1.9  
21 billion South Asians there could be over 26 million carriers for the variant. Thus, this work has  
22 important implications for genetic screening and early identification of at-risk individuals, and  
23 the future opportunity for rationally targeted therapeutic interventions.

## 24 25 **Online Methods**

26  
27 Details of methods are available in the Supplementary text. Briefly, 75 K individuals were  
28 recruited and consented in Pakistan for exome sequencing, including a stroke case:control  
29 discovery cohort of 31 K (5,135 cases and 26,602 controls) sequenced by the Regeneron  
30 Genetics Center and a follow-up cohort of 44 K, including 30 K with self-report stroke  
31 case:control status used for replication and meta-analysis (160 cases and 30,239 controls).  
32 ExWAS was conducted for stroke and 4 overlapping stroke subtypes (intracerebral hemorrhage,  
33 subcortical intracerebral hemorrhage, ischemic stroke, and partial anterior circulating infarct) in  
34 the discovery cohort and combined in a meta-analysis of stroke with the replication cohort using  
35 both single-variant and gene burden test models [26]. Population attributable fraction of stroke  
36 for associated variants was calculated as described previously [27]. Consented callbacks were  
37 conducted in families of homozygotes for associated variants. For comparison and validation,  
38 analyses were conducted in UK Biobank data as described previously [28-31], including  
39 association analysis of p.Arg1231Cys *NOTCH3* with stroke phenotypes in 380 K participants,  
40 and brain imaging phenotypes in 35 K participants (**Supplementary Methods; Figure 1A**).

## 41 42 **Custom Burden Tests in UKB 450 K**

43  
44 Burden tests aim to boost statistical power by aggregating association signal across multiple rare  
45 variants. Prior studies in human and animal models have debated the role of various variant  
46 classes on *NOTCH3* function, CADASIL pathology and patient prognosis, including experiments

1 designed to determine if the pathogenicity of CADASIL variants follows a loss of function (LoF)  
2 mechanism [8, 25, 32]. Using data from hundreds of missense and LoF variants in *NOTCH3*  
3 observed in 450 K UKB participant exomes, burden tests were conducted to assess the impact of  
4 LoF and missense variants. We first assessed the impact of Cys-altering variants in EGFr  
5 domains 1 to 34, 1 to 6, and 7 to 34 on stroke risk. To compare the effect of other missense  
6 variants in these same domains, we further tested the effect of Ser-altering variants. Finally, we  
7 tested across the full protein-coding region for association with stroke of (i) LoFs and (ii) LoFs  
8 and missense variants.

### Figure Legends

9  
10

11 **Figure 1. ExWAS Identifies *NOTCH3* p.Arg1231Cys Associated with Subcortical**  
12 **Hemorrhagic Stroke in Pakistan Genome Resource 31 K Discovery Cohort. A.** Flow chart  
13 of the study described in this report. The discovery cohort consisted of a 31 K stroke case-control  
14 cohort from the Pakistan Genome Resource (PGR). A second PGR follow-up cohort of 44 K  
15 included 30 K participants with self-reported stroke case:control status for replication. UK  
16 Biobank data from 450 K sequenced participants was used for further analysis in a  
17 predominantly European ancestry population, 380 K of whom had stroke case:control status  
18 known, and 35 K of whom had brain MRI data (see Methods and Supplementary Methods). **B.**  
19 Manhattan plot of subcortical hemorrhagic stroke ExWAS in 31 K PGR discovery cohort  
20 participants with  $-\log_{10}$  p-values (y-axis) across chromosomes and variants (x-axis). A single  
21 variant (NC\_000019.10:g.15179052G>A) on chromosome 19 predicting a missense variant  
22 p.Arg1231Cys in *NOTCH3* exceeded the genome-wide significance threshold of  $5 \times 10^{-8}$  (red  
23 line). **C.** *NOTCH3* locus zoom plot of subcortical stroke ExWAS. The  $-\log_{10}$  p-values for  
24 variants tested are shown on the y-axis. The p.Arg1231Cys variant is labeled as a diamond.  
25 Other variants (circles) are colored based on linkage disequilibrium with the reference variant in  
26 1000 Genomes [33]. Gene exon (thick line) and intron (thin line) model shown below the graph.

27

28

1 **Figure 2. *NOTCH3* EGFR Domain Disruption by p.Arg1231Cys.** Shown is *NOTCH3*  
2 p.Arg1231 in context of human *NOTCH3* protein domains and cross-species alignment of  
3 *NOTCH3* amino acid sequences. **A.** Human *NOTCH3* Protein Domains. Shown is the position of  
4 the associated variant in context of transcript exons (top, alternating blue and purple with  
5 numbering) and protein domains (bottom, color coded). *NOTCH3* can be divided into four major  
6 regions, from left-to-right the signal peptide (light blue), the extra-cellular domain (ECD,  
7 brown), the transmembrane domain (orange), and the intra-cellular domain (ICD, blue). The  
8 majority of the ECD is composed of 34 Epidermal Growth Factor-like repeat (EGFr) domains (in  
9 purple with white numbers). Domains involved in signaling are highlighted, including EGFr  
10 domains 10 to 11 involved in ligand binding (light purple, numbered), three cleavage domains  
11 (S1 in yellow, S2 in yellow with diagonal black stripes, S3 in yellow with black horizontal  
12 stripes), and three *Lin12/NOTCH* repeats (light green, numbered). The ICD contains the  
13 Recombination signal-binding protein for Ig of  $\kappa$  region (RAM) domain for transcription factor  
14 interaction (green and white checkers), the Nuclear Localization Sequences (NLS, orange with  
15 black stripes), five Ankyrin repeats involved in signal transduction (green with white numbers),  
16 and the Proline, glutamic acid, serine, and threonine-rich (PEST) domain essential for  
17 degradation (green with white stripes). The p.Arg1231Cys variant (red line top to bottom)  
18 removes a disulfide-bridge-forming cysteine in the 31<sup>st</sup> EGFr domain of the ECD, coded by the  
19 22<sup>nd</sup> exon. **B.** Cross-species Alignment of *NOTCH3* (EGFr) Domain # 31 Amino Acid  
20 Sequences. Shown is an amino acid alignment of 31<sup>st</sup> EGFr domain of *NOTCH3* (human  
21 sequence amino acids 1205 to 1244), including (top-to-bottom) human reference, human  
22 p.Arg1231Cys mutant, chimpanzee (*Pan troglodytes*), mouse (*Mus musculus*), zebrafish (*Danio*  
23 *rerio*), western clawed frog (*Xenopus tropicalis*), and green sea turtle (*Chelonia mydas*),  
24 indicating conservation of the arginine (R) at position 1231 in mammals. Highly-conserved  
25 cysteine (C) residues (normally 6 per EGFR) are highlighted in yellow.

26  
27

Rodriguez-Flores et al.

13

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

Version 33

20 March 2024

1 **Figure 3. Forest Plot Showing Replication of *NOTCH3* p.Arg1231Cys Association with**  
2 **Stroke Across 61 K Pakistan Genome Resource Meta-Analysis.** Shown is the cohort name,  
3 trait, odds ratio with 95% confidence interval, p value, alternate allele frequency, case count, and  
4 control count for five stroke phenotypes in the PGR 31 K discovery cohort (top), and Inverse  
5 Variance Weighted (IVW) Meta-Analysis of stroke in 61 K PGR Cohort, including 31 K PGR  
6 discovery cohort and 30 K PGR replication cohort subset of the 44 K PGR follow-up cohort  
7 (bottom).  
8  
9

Tables

Table 1. Baseline Characteristics of PGR Stroke Case-control Discovery Cohort<sup>1</sup>

	Case (n = 5,135)		Control (n = 26,602)		P value	Case p.Arg1231Cys Carrier (n=103) <sup>3</sup>		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 <sup>2</sup>	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 <sup>2</sup>	1,850	36	9,242	34.7	< 0.01	32	31.1	0	0	0.34
<b>Comorbidities</b>										
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
<b>Family history of</b>										
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

1. Shown are mean or total n and standard deviation or percentage for cases versus controls on the left and case p.Arg1231Cys *NOTCH3* carriers versus case non-carriers on the right. Comparison was conducted and p values are shown from chi-square test for categorical variables (Fisher's exact test if cell size was <5) and from T-test for continuous variables.
2. Tobacco or other stimulants include cigarettes, paan (chewed betel leaf and areca nut), naswar (snuff), gutka (chewing tobacco), huqqa (water pipe), chillum (hasish pipe).
3. Genotypes for the p.Arg1231Cys variant were not available for n=34.

1  
2 **Table 2. UKB Ischemic Stroke Association Across *NOTCH3* Variant Classes and Domains<sup>1</sup>**  
3

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

4  
5 1. Burden test with age, age<sup>2</sup>, gender, 10 PCs as covariates was conducted using REGENIE to  
6 compare the association signal and effect size across variant classes and domains for  
7 ischemic stroke in 9,143 cases and 371,403 controls from UK Biobank. *NOTCH3* variant  
8 classes include the single variant association (p.Arg1231Cys) for reference at the top and  
9 burden tests limited to: Cys-altering variants in EGFr domains 1 to 34; Cys-altering variants  
10 in EGFr domains 1 to 6; Cys-altering variants in EGFr domains 7 to 34; Cys-altering variants  
11 outside of EGFr domains; Ser-altering variants in EGFr domains 1 to 6; Ser-altering variants  
12 in EGFr domains 7 to 34; Ser-altering variants outside of EGFr domains; Loss of function  
13 (LoF) variants (defined as stop-gain, stop-loss, frameshift and splice-site variants with AAF  
14 <1%); and LoF and any missense variants (AAF < 1%). Abbreviations: EGFr, epidermal  
15 growth factor-like repeat domain; AAF, alternate allele frequency; RR, reference allele  
16 homozygote; RA, reference/alternate allele heterozygote; AA, alternate allele homozygote;  
17 95% CI, 95% confidence interval  
18

## References

- 1
- 2
- 3 1. Mills, M.C. and C. Rahal, *The GWAS Diversity Monitor tracks diversity by disease in*
- 4 *real time*. Nat Genet, 2020. **52**(3): p. 242-243.
- 5 2. Collaborators, G.B.D.S., *Global, regional, and national burden of stroke and its risk*
- 6 *factors, 1990-2019: a systematic analysis for the Global Burden of Disease Study 2019*.
- 7 Lancet Neurol, 2021. **20**(10): p. 795-820.
- 8 3. Sherin, A., et al., *Prevalence of stroke in Pakistan: Findings from Khyber Pakhtunkhwa*
- 9 *integrated population health survey (KP-IPHS) 2016-17*. Pak J Med Sci, 2020. **36**(7): p.
- 10 1435-1440.
- 11 4. Valcarcel-Nazco, C., et al., *Variability in the use of neuroimaging techniques for*
- 12 *diagnosis and follow-up of stroke patients*. Neurologia (Engl Ed), 2019. **34**(6): p. 360-
- 13 366.
- 14 5. Farooq, A., N. Venketasubramanian, and M. Wasay, *Stroke Care in Pakistan*.
- 15 Cerebrovasc Dis Extra, 2021. **11**(3): p. 118-121.
- 16 6. Farooq, M.U., et al., *The epidemiology of stroke in Pakistan: past, present, and future*. Int
- 17 J Stroke, 2009. **4**(5): p. 381-9.
- 18 7. Mullen, M.T., et al., *Hospital-Level Variability in Reporting of Ischemic Stroke Subtypes*
- 19 *and Supporting Diagnostic Evaluation in GWTG-Stroke Registry*. J Am Heart Assoc,
- 20 2023. **12**(24): p. e031303.
- 21 8. Chabriat, H., et al., *Cadasil*. Lancet Neurol, 2009. **8**(7): p. 643-53.
- 22 9. Markidan, J., et al., *Smoking and Risk of Ischemic Stroke in Young Men*. Stroke, 2018.
- 23 **49**(5): p. 1276-1278.
- 24 10. Mishra, A., et al., *Stroke genetics informs drug discovery and risk prediction across*
- 25 *ancestries*. Nature, 2022. **611**(7934): p. 115-123.
- 26 11. Scott, E.M., et al., *Characterization of Greater Middle Eastern genetic variation for*
- 27 *enhanced disease gene discovery*. Nat Genet, 2016. **48**(9): p. 1071-6.
- 28 12. Wang, T., M. Baron, and D. Trump, *An overview of Notch3 function in vascular smooth*
- 29 *muscle cells*. Prog Biophys Mol Biol, 2008. **96**(1-3): p. 499-509.
- 30 13. Duvaud, S., et al., *Expasy, the Swiss Bioinformatics Resource Portal, as designed by its*
- 31 *users*. Nucleic Acids Res, 2021. **49**(W1): p. W216-w227.
- 32 14. Rutten, J.W., et al., *Broad phenotype of cysteine-altering NOTCH3 variants in UK*
- 33 *Biobank: CADASIL to nonpenetrance*. Neurology, 2020. **95**(13): p. e1835-e1843.
- 34 15. Rodriguez-Flores, J.L., et al., *The QChip1 knowledgebase and microarray for precision*
- 35 *medicine in Qatar*. NPJ Genom Med, 2022. **7**(1): p. 3.
- 36 16. Hoffmann, T.J., et al., *Genome-wide association analyses using electronic health records*
- 37 *identify new loci influencing blood pressure variation*. Nat Genet, 2017. **49**(1): p. 54-64.
- 38 17. Hack, R.J., et al., *Three-tiered EGFR domain risk stratification for individualized*
- 39 *NOTCH3-small vessel disease prediction*. Brain, 2023. **146**(7): p. 2913-2927.
- 40 18. Cho, B.P.H., et al., *Association of Vascular Risk Factors and Genetic Factors With*
- 41 *Penetrance of Variants Causing Monogenic Stroke*. JAMA Neurol, 2022. **79**(12): p.
- 42 1303-1311.
- 43 19. Rutten, J.W., et al., *The effect of NOTCH3 pathogenic variant position on CADASIL*
- 44 *disease severity: NOTCH3 EGFR 1-6 pathogenic variant are associated with a more*
- 45 *severe phenotype and lower survival compared with EGFR 7-34 pathogenic variant*.
- 46 Genet Med, 2019. **21**(3): p. 676-682.



- 1 20. ClinVar, N.C.f.B.I.
- 2 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.
- 3
- 4 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.
- 5
- 6 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.
- 7
- 8
- 9 24. Ghezali, L., et al., *Notch3(ECD) immunotherapy improves cerebrovascular responses in CADASIL mice*. *Ann Neurol*, 2018. **84**(2): p. 246-259.
- 10
- 11 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.
- 12
- 13 26. Mbatchou, J., et al., *Computationally efficient whole-genome regression for quantitative and binary traits*. *Nat Genet*, 2021. **53**(7): p. 1097-1103.
- 14
- 15 27. Greenland, S., *Applications of Stratified Analysis Methods*, in *Modern epidemiology: Third edition*. 2008, Lippincott Williams & Wilkins: Philadelphia. p. 295-297.
- 16
- 17 28. Backman, J.D., et al., *Exome sequencing and analysis of 454,787 UK Biobank participants*. *Nature*, 2021. **599**(7886): p. 628-634.
- 18
- 19 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.
- 20
- 21
- 22 30. Elliott, L.T., et al., *Genome-wide association studies of brain imaging phenotypes in UK Biobank*. *Nature*, 2018. **562**(7726): p. 210-216.
- 23
- 24 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.
- 25
- 26
- 27 32. Joutel, A., et al., *Cerebrovascular dysfunction and microcirculation rarefaction precede white matter lesions in a mouse genetic model of cerebral ischemic small vessel disease*. *J Clin Invest*, 2010. **120**(2): p. 433-45.
- 28
- 29
- 30 33. Genomes Project, C., et al., *A global reference for human genetic variation*. *Nature*, 2015. **526**(7571): p. 68-74.
- 31
- 32
- 33
- 34

## Acknowledgements

1  
2  
3 Supported by Regeneron Pharmaceuticals.

4  
5 The Institutional Review Board (IRB) at the Center for Non-Communicable Diseases (IRB:  
6 00007048, IORG0005843, FWAS00014490) approved the study. All participants gave  
7 written informed consent.

8  
9 This research has been conducted using the UK Biobank Resource (project 26041). The  
10 authors thank everyone who made this work possible, particularly the UK Biobank team,  
11 their funders, the professionals from the member institutions who contributed to and  
12 supported this work, and most especially the UK Biobank participants, without whom this  
13 research would not be possible. The exome sequencing was funded by the UK Biobank  
14 Exome Sequencing Consortium (Bristol Myers Squibb, Regeneron, Biogen, Takeda, Abbvie,  
15 Alnylam, AstraZeneca and Pfizer). Ethical approval for the UK Biobank was previously  
16 obtained from the North West Centre for Research Ethics Committee (11/NW/0382).

17  
18 Disclosure forms provided by the authors are available with the full text of this article.

19  
20 Drs. Rodriguez-Flores, Khalid, Shuldiner and Saleheen contributed equally to this article.  
21

Rodriguez-Flores et al.

19

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

Version 33

20 March 2024

## Data Sharing

1  
2  
3  
4

Data used in this manuscript not already available online can be requested from the authors.

Rodriguez-Flores et al.

20

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

Version 33

20 March 2024

### **Code Sharing**

1  
2  
3  
4  
5  
6  
7  
8  
9  
10

Code used to produce data presented in this manuscript not already available online can be requested from the authors.

1  
2

**Primary Individual Contributions**

<p><b>CONCEPTUALIZATION</b>                      Juan L. Rodriguez-Flores (JLRF)                      Shareef Khalid (SK)                      Alan R. Shuldiner (ARS)                      Danish Saleheen (DS)</p>	<p><b>INVESTIGATION</b>                      Silvio Alessandro DiGioia (SADG)                      Hector Martinez (HM)                      I-Chun Tsai (IT)                      Katia Karalis (KK)                      Aris Economides (AE)                      David D’Ambrosio (DDA)                      Asif Rasheed (AR)</p>	<p><b>VALIDATION</b>                      Regeneron Genetics Center                      Shareef Khalid                      Juan L. Rodriguez-Flores</p>
<p><b>DATA CURATION</b>                      Nilanjana Banerjee (NB)                      Deepika Sharma (DeS)                      Michael Cantor (MC)                      John Overton (JO)                      Jeff Reid (JR)</p>	<p><b>METHODOLOGY</b>                      Juan L. Rodriguez-Flores                      Shareef Khalid                      Regeneron Genetics Center</p>	<p><b>VISUALIZATION</b>                      Juan L. Rodriguez-Flores                      Neelroop Parikshak                      Farshid Sepehrband                      Jonathan Marchini                      Regeneron Genetics Center</p>
<p><b>FORMAL ANALYSIS</b>                      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)                      Regeneron Genetics Center (RGC)</p>	<p><b>PROJECT ADMINISTRATION</b>                      Thomas Coleman (TC)                      Regeneron Genetics Center                      Asif Rasheed</p>	<p><b>WRITING - ORIGINAL DRAFT</b>                      Sergio Fazio (SF)                      Wolfgang Liedtke (WL)                      John Danesh (JD)                      Ayeesha Kamal (AK)                      Philippe Frossard (PF)</p>
<p><b>FUNDING AQUISITION</b>                      Danish Saleheen                      Alan R. Shuldiner                      Aris Baras (AB)                      Regeneron Genetics Center</p>	<p><b>RESOURCES</b>                      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)                      Faizan Aslam (FA)</p>	<p><b>WRITING - REVIEW AND EDITING</b></p>
	<p><b>SOFTWARE</b>                      Regeneron Genetics Center                      Juan L. Rodriguez-Flores                      Shareef Khalid</p>	
	<p><b>SUPERVISION</b>                      Danish Saleheen                      Alan R. Shuldiner</p>	

3  
4

## Competing Interests

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31

### Funding

Fieldwork for this study was funded by the Center for Non-Communicable Diseases, Pakistan.  
DNA sequencing was funded by Regeneron Pharmaceuticals Inc.

### Employment

JLRF, ARS, NB, DeS, MC, JO, JR, BY, MK, JB, GT, SG, TD, NV, LAL, AZ, NP, FS, JM, GC, PX, AB, SADG, HM, IT, KK, AE, DDA, SF, WL, ET, SC, TC are or were employees of Regeneron Genetics Center LLC or Regeneron Pharmaceuticals Inc. and contributed to this manuscript as part of their regular duties as salaried employees.

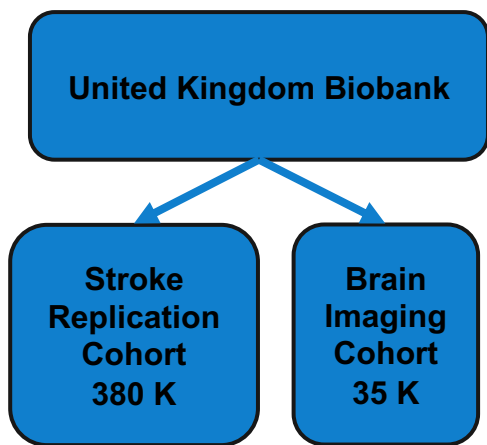
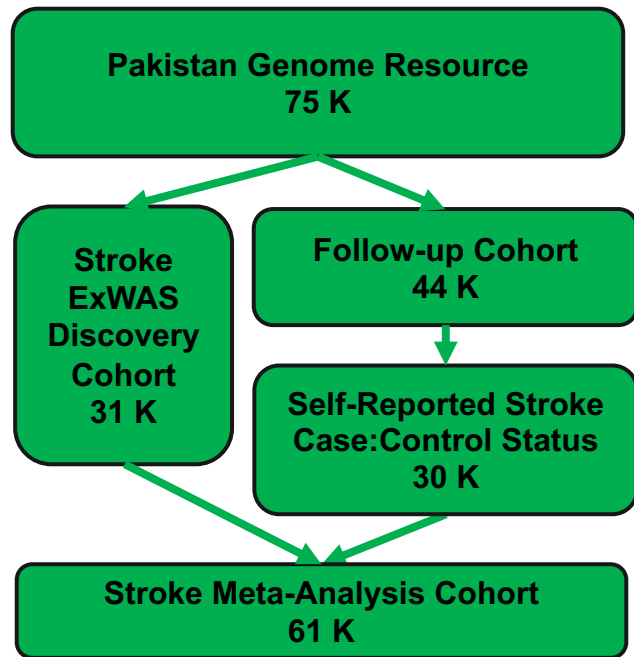
AR, MJ, MZK, MRM, MBL, PF, and DS and SK are employees of the Center for Non-Communicable Disease and received salaried compensation for their contribution to this manuscript.

### Personal Financial Interests

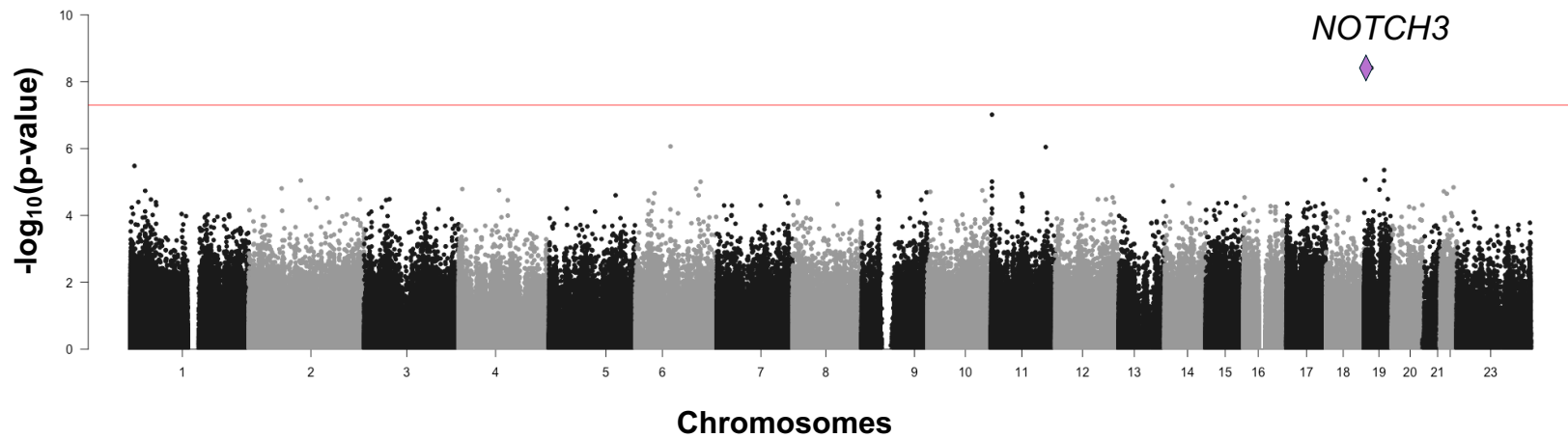
JLRF, ARS, NB, DeS, MC, JO, JR, BY, MK, JB, GT, SG, TD, NV, LAL, AZ, NP, FS, JM, GC, PX, AB, SADG, HM, IT, KK, AE, DDA, SF, WL, TC are or were employees of Regeneron Genetics Center LLC or Regeneron Pharmaceuticals Inc. and received stock and stock options as part of their compensation as employees.

JLRF, ARS, DS, AB, and SK are named inventors on patent pending US 20230000897A1 that discloses methods of treating subjects having a cerebrovascular disease by administering Neurogenic Locus Notch Homolog Protein 3 (*NOTCH3*) agents, and methods of identifying subjects having an increased risk of developing a cerebrovascular disease.

A.



B.



C.

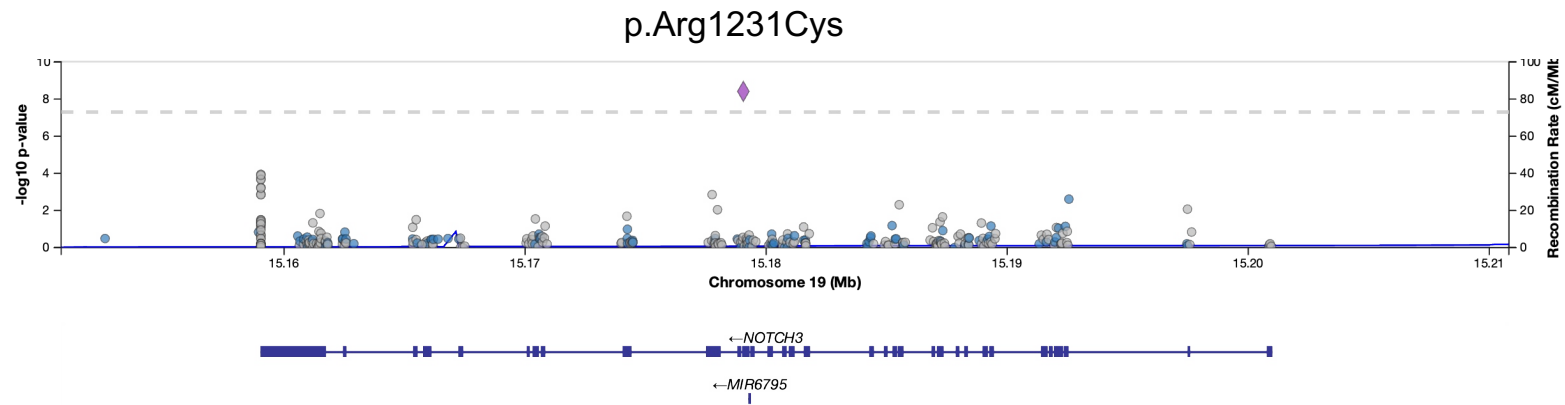
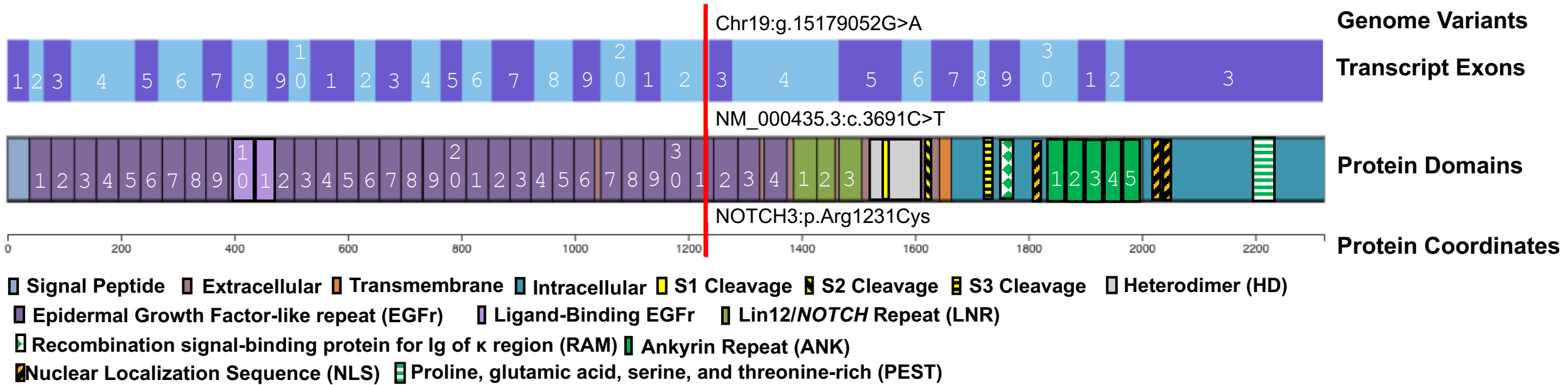


Figure 1

A.



B.

1205	1231	1244	
INECRSGAC	HAAHTRD	CLQDPGGGFR	CLCHAGFSGPRCQT
INECRSGAC	HAAHTRD	CLQDPGGGF	CCCLCHAGFSGPRCQT
INECRSGAC	HAAHTRD	CLQDPGGGFR	CLCHAGFSGPRCQT
INECRPGAC	HAAHTRD	CLQDPGGHFR	CVCHPGFTGPRCQI
INECLSNPC	NPNSLD	CIQLPND-YQ	CVCKPGFTGRGCQS
INECLSGP	CHAQNTRH	CVQLAND-YQ	CVCKSGYTGRRCQS
INECLAKP	CLPQRTL	CVQGAND-FQ	CLCKPGYTGRRCQN
			Human 1231 Arg (reference)
			Human 1231 Cys (mutant)
			Chimpanzee
			Mouse
			Zebrafish
			Western clawed frog
			Green sea turtle

Figure 2



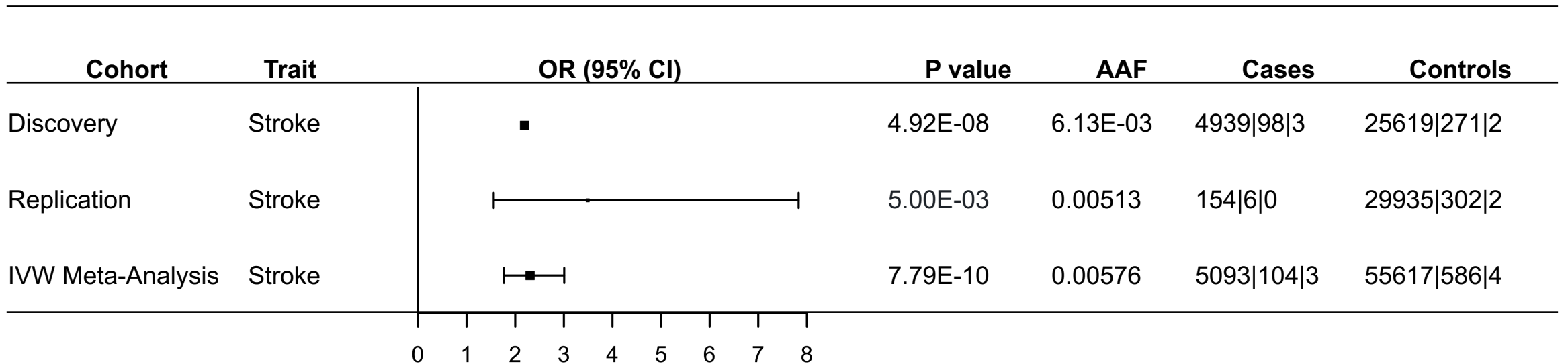
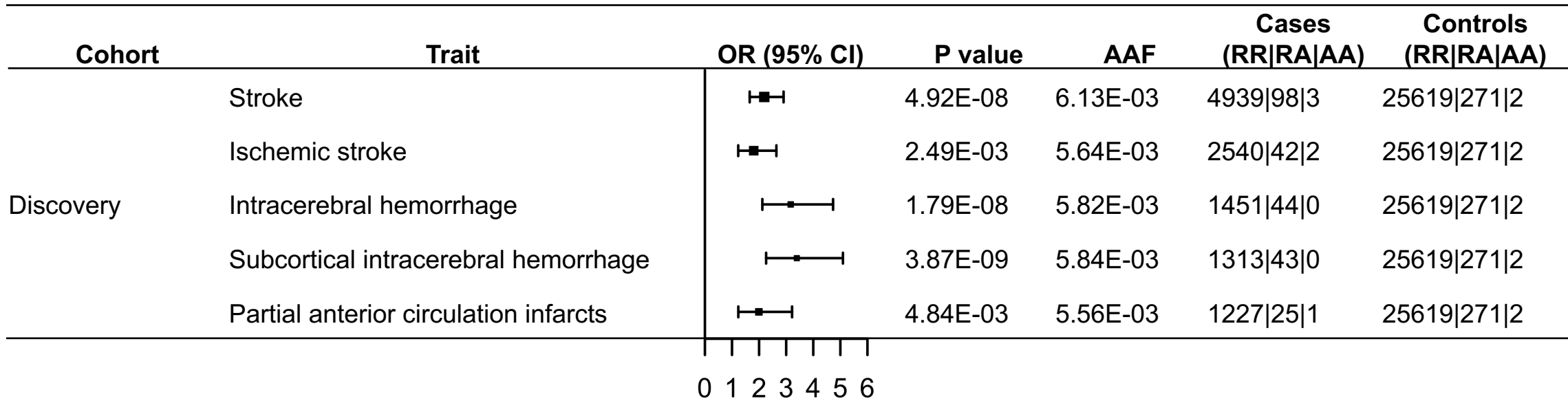


Figure 3