

Supplemental information

Table of Contents

SUPPLEMENT METHODS	ERROR! BOOKMARK NOT DEFINED.
Theoretical Heterozygous SNP Sensitivity (THS)	2
Variant Pathogenicity Interpretation	2
Reference	3
SUPPLEMENT TABLES	ERROR! BOOKMARK NOT DEFINED.
Supplement Table S1. Genes sequenced	4
Supplement Table S2. Sequence ontology terms and functional impact definition	6
Supplement Table S3. Association analysis for common variants	7
Supplement Table S4. FH-FHRs fusion proteins identified in 400 aHUS patients	8
Supplement Table S5. Rare variants with $0.1\% < \text{MAF (NFE)} < 1\%$ identified in <i>CFH</i> , <i>CD46</i> , <i>C3</i> , <i>CFI</i> and <i>CFB</i> in 400 aHUS patients	9
SUPPLEMENT FIGURES	ERROR! BOOKMARK NOT DEFINED.
Supplement Figure S1. Rare variants in the <i>CFH</i> gene are significantly more abundant in the non-Finnish European (NFE) subpopulation as compared to the Finnish (FIN) subpopulation from gnomAD due to population stratification.	10
Supplement Figure S2. Cluster analysis shows that UI controls and aHUS cases are mostly similar to NFE	11
Supplement Figure S3. Population stratification within aHUS cases and UI controls was used to remove outliers.	12
Supplement Figure S4. Mean coverage is correlated with theoretical heterozygous SNP sensitivity (HET SNP sensitivity) in patients and controls	13
Supplement Figure S5. Six samples with a shifted ratio of ref/alt reads were excluded from the study cohort.	14
Supplement Figure S6. Prediction score distribution of neutral and pathogenic variants on <i>CFH</i> gene.	15
Supplement Figure S7. Median age of female patients is significantly lower than that of male patients	16
Supplement Figure S8: Enrichment of ultra-rare variants “contaminates” the result of association analysis with higher MAF cut-off	17

SUPPLEMENTAL METHODS

Theoretical Heterozygous SNP Sensitivity (THS)

THS is a quality metric that estimates the theoretical sensitivity to detect heterozygous variants based on coverage distribution and base quality distribution from massively parallel sequencing data.^{1,2} Under the assumptions: 1) DNA is diploid; 2) at a HET site with genotype AB, the only possible calls are A and B; 3) there is no reference bias; and 4) coverage distribution $P(n)$ and base quality distribution $P(q)$ are known and statistically independent, the model of HET detection is based on Bernoulli distribution as:

$$\frac{\binom{n}{m} \left(\frac{1}{2}\right)^n}{\binom{n}{m} \prod_{j=1}^m e_j \prod_{j=m+1}^n (1 - e_j)} > T$$

where, n is the depth from the coverage distribution $P(n)$; m is the number of true alternate alleles from $m \sim \text{binomial}(n, 0.5)$ covering the HET site; $e_j = 10^{-q_j/10}$ is the probability of error, and q_j is from the base quality distribution $P(q)$.

Variant Pathogenicity Interpretation

Pathogenicity of variant was based on absence in large health populations, presence/enrichment in aHUS patients, and functional association.

1. Large health populations refer to the gnomAD database (138,632 subjects).
2. Presence in aHUS patients is defined as: 1) reported in the literature, 2) reported in an aHUS disease mutation database, or 3) observed in our patient cohort.
3. Enrichment in aHUS patients is determined by association analysis in patients and controls with adjustment for population stratification.
4. Functional association is defined as: 1) well-studied functional changes that contribute to aHUS development, 2) truncating protein where loss of function is a known disease mechanism, 3) known disruption of protein structure (e.g. cysteine-related missense variants in SCRs of *CFH* and *CD46*), or 4) localization in well-defined aHUS-related domains.

Pathogenic is defined as: 1) absent in gnomAD *and* reported at least once in the literature or an aHUS database *and* observed at least once in our patient cohort; OR 2) absent in gnomAD *and* observed at least twice in our patient cohort *and* with a functional impact.

Likely pathogenic is defined as: 1) absent in gnomAD *and* observed at least twice in our patient cohort; OR 2) absent in gnomAD *and* observed at least once in our patient cohort *and* with a functional impact; OR 3) significantly enriched in our patient cohort compare to the corresponding gnomAD population *and* with a functional impact.

Likely benign is defined as: frequency > 0.1% in any gnomAD population *and* not enriched in patients *and* lacking a functional effect.

Benign is defined as: frequency > 1% in any gnomAD population *and* not enriched in patients.

Reference

1. Yossi Farjoun Jon Bloom: Theoretical HET Sensitivity.

https://www.broadinstitute.org/files/shared/mia/theoretical_HET_sensitivity.pdf

2. Kylee Degatano, David Benjamin, Jonathan M. Bloom, Maura Costello, Jason Rose, Kathleen TibbeFs, CharloFe Tolonen, Yossi Farjoun: Optimizing Delivered Sequencing Data with a Theoretical Sensitivity to Heterozygous SNPs. ASHG, 2016.

http://www.genomics.broadinstitute.org/data-sheets/POS_OptimizingDeliveredSequencingDataTheoreticalSensitivityHeteroSNPs_ASHG_2016.pdf

SUPPLEMENTAL TABLES

Supplement Table S1. Genes sequenced

	Gene	Full Name	RefSeq ID
1	<i>A2M</i>	Alpha-2-Macroglobulin	NM_000014
2	<i>ABCD4</i>	ATP Binding Cassette Subfamily D Member 4	NM_005050
3	<i>ADAMTS13</i>	ADAM Metallopeptidase With Thrombospondin Type 1 Motif 13	NM_139025
4	<i>ADM</i>	Adrenomedullin	NM_001124
5	<i>ADM2</i>	Adrenomedullin 2	NM_001253845
6	<i>APCS</i>	Amyloid P Component, Serum	NM_001639
7	<i>C1QA</i>	Complement C1q A Chain	NM_015991
8	<i>C1QB</i>	Complement C1q B Chain	NM_000491
9	<i>C1QC</i>	Complement C1q C Chain	NM_172369
10	<i>C1R</i>	Complement Component 1, R Subcomponent	NM_001733
11	<i>C1S</i>	Complement Component 1, S Subcomponent	NM_201442
12	<i>*C2</i>	Complement C2	NM_000063
13	<i>C3</i>	Complement C3	NM_000064
14	<i>C3AR1</i>	Complement C3a Receptor 1	NM_004054
15	<i>*C4A</i>	Complement C4A	NM_007293
16	<i>*C4B</i>	Complement C4B	NM_000715
17	<i>C4BPA</i>	Complement Component 4 Binding Protein Alpha	NM_000715
18	<i>C4BPB</i>	Complement Component 4 Binding Protein Beta	NM_000716
19	<i>C5</i>	Complement C5	NM_001735
20	<i>C5AR1</i>	Complement C5a Receptor 1	NM_001736
21	<i>C5AR2</i>	Complement C5a Receptor 2	NM_018485
22	<i>C6</i>	Complement C6	NM_000065
23	<i>C7</i>	Complement C7	NM_000587
24	<i>C8A</i>	Complement C8 Alpha Chain	NM_000562
25	<i>C8B</i>	Complement C8 Beta Chain	NM_000066
26	<i>C8G</i>	Complement C8 Gamma Chain	NM_000606
27	<i>C9</i>	Complement C9	NM_001737
28	<i>CD46</i>	Membrane Cofactor Protein	NM_002389
29	<i>CD55</i>	Decay Accelerating Factor For Complement	NM_000574
30	<i>CD59</i>	Membrane Attack Complex Inhibition Factor	NM_000611
31	<i>CFB</i>	Complement Factor B	NM_001710
32	<i>CFD</i>	Complement Factor D	NM_001928
33	<i>CFH</i>	Complement Factor H	NM_000186
34	<i>*CFHR1</i>	Complement Factor H Related 1	NM_002113
35	<i>CFHR2</i>	Complement Factor H Related 2	NM_005666
36	<i>*CFHR3</i>	Complement Factor H Related 3	NM_021023
37	<i>CFHR4</i>	Complement Factor H Related 4	NM_006684
38	<i>CFHR5</i>	Complement Factor H Related 5	NM_030787
39	<i>CFI</i>	Complement Factor I	NM_000204
40	<i>CFP</i>	Complement Factor Properdin	NM_002621
41	<i>CLU</i>	Clusterin	NM_001831
42	<i>COLEC11</i>	Collectin Subfamily Member 11	NM_199235
43	<i>CPN1</i>	Anaphylatoxin Inactivator	NM_001308
44	<i>*CR1</i>	Complement C3b/C4b Receptor 1	NM_000651
45	<i>CR2</i>	Complement C3b/C4b Receptor 2	NM_001006658
46	<i>CRP</i>	C-Reactive Protein	NM_000567
47	<i>DGKE</i>	Diacylglycerol Kinase Epsilon	NM_003647
48	<i>F10</i>	Coagulation Factor X	NM_000504
49	<i>F11</i>	Coagulation Factor XI	NM_000128
50	<i>F12</i>	Coagulation Factor XII	NM_000505
51	<i>F2</i>	Coagulation Factor II, Thrombin	NM_000506
52	<i>F2RL2</i>	Coagulation Factor II Thrombin Receptor Like 2	NM_004101
53	<i>F3</i>	Coagulation Factor III, Tissue Factor	NM_001993
54	<i>F5</i>	Coagulation Factor V	NM_000130
55	<i>F7</i>	Coagulation Factor VII	NM_000131
56	<i>F8</i>	Coagulation Factor VIII, Procoagulant Component	NM_000132

57	<i>F9</i>	Coagulation Factor IX	NM_000133
58	<i>FCN1</i>	Ficolin 1	NM_002003
59	<i>FCN2</i>	Ficolin 2	NM_004108
60	<i>FCN3</i>	Ficolin 3	NM_003665
61	<i>FGL2</i>	Fibrinogen Like 2	NM_006682
62	<i>IFNG</i>	Interferon Gamma	NM_000619
63	<i>INF2</i>	Inverted Formin, FH2 And WH2 Domain Containing	NM_022489
64	<i>ITGAM</i>	Integrin Subunit Alpha M	NM_000632
65	<i>KLKB1</i>	Kallikrein B1	NM_000892
66	<i>LMBRD1</i>	LMBR1 Domain Containing 1	NM_018368
67	<i>MAP3K5</i>	Mitogen-Activated Protein Kinase Kinase Kinase 5	NM_005923
68	<i>MASP1</i>	Mannan Binding Lectin Serine Peptidase 1	NM_001879
69	<i>MASP2</i>	Mannan Binding Lectin Serine Peptidase 2	NM_006610
70	<i>MBL2</i>	Mannose Binding Lectin 2	NM_000242
71	<i>MBTPS1</i>	Membrane Bound Transcription Factor Peptidase, Site 1	NM_003791
72	<i>MMACHC</i>	Methylmalonic Aciduria CblC Type, With Homocystinuria	NM_015506
73	<i>MMADHC</i>	Methylmalonic Aciduria CblD Type, With Homocystinuria	NM_015702
74	<i>MTR</i>	5-Methyltetrahydrofolate-Homocysteine Methyltransferase	NM_000254
75	<i>MTRR</i>	5-Methyltetrahydrofolate-Homocysteine Methyltransferase Reductase	NM_002454
76	<i>PHB</i>	Prohibitin	NM_002634
77	<i>PLAT</i>	Plasminogen Activator, Tissue	NM_000930
78	<i>PLAU</i>	Plasminogen Activator, Urokinase	NM_002658
79	<i>PLG</i>	Plasminogen	NM_000301
80	<i>PROC</i>	Protein C, Inactivator Of Coagulation Factors Va And VIIIa	NM_000312
81	<i>PROS1</i>	Protein S alpha	NM_000313
82	<i>PTX3</i>	Pentraxin 3	NM_002852
83	<i>SERPINA1</i>	Serpin Family A Member 1	NM_000295
84	<i>SERPINA5</i>	Serpin Family A Member 5	NM_000624
85	<i>SERPINC1</i>	Serpin Family C Member 1	NM_000488
86	<i>SERPIND1</i>	Serpin Family D Member 1	NM_000185
87	<i>SERPINE1</i>	Serpin Family E Member 1	NM_000602
88	<i>SERPINF2</i>	Serpin Family F Member 2	NM_000934
89	<i>SERPING1</i>	Serpin Family G Member 1	NM_000062
90	<i>THBD</i>	Thrombomodulin	NM_000361
91	<i>VSIG4</i>	V-Set And Immunoglobulin Domain Containing 4	NM_007268
92	<i>VTN</i>	Vitronectin	NM_000638
93	<i>VWF</i>	von Willebrand Factor	NM_000552

Genes marked with asterisk were not included in burden analysis due to ambiguous read mapping

Supplement Table S2. Sequence ontology terms and functional impact definition

SO accession	SO term	IMPACT
SO:0001893	transcript_ablation	HIGH
SO:0001574	splice_acceptor_variant	HIGH
SO:0001575	splice_donor_variant	HIGH
SO:0001587	stop_gained	HIGH
SO:0001589	frameshift_variant	HIGH
SO:0001578	stop_lost	HIGH
SO:0002012	start_lost	HIGH
SO:0001889	transcript_amplification	HIGH
SO:0001821	inframe_insertion	MODERATE
SO:0001822	inframe_deletion	MODERATE
SO:0001583	missense_variant	MODERATE
SO:0001818	protein_altering_variant	MODERATE
SO:0001630	splice_region_variant	LOW
SO:0001626	incomplete_terminal_codon_variant	LOW
SO:0001567	stop_retained_variant	LOW
SO:0001819	synonymous_variant	LOW
SO:0001580	coding_sequence_variant	MODIFIER
SO:0001620	mature_miRNA_variant	MODIFIER
SO:0001623	5_prime_UTR_variant	MODIFIER
SO:0001624	3_prime_UTR_variant	MODIFIER
SO:0001792	non_coding_transcript_exon_variant	MODIFIER
SO:0001627	intron_variant	MODIFIER
SO:0001621	NMD_transcript_variant	MODIFIER
SO:0001619	non_coding_transcript_variant	MODIFIER
SO:0001631	upstream_gene_variant	MODIFIER
SO:0001632	downstream_gene_variant	MODIFIER
SO:0001895	TFBS_ablation	MODIFIER
SO:0001892	TFBS_amplification	MODIFIER
SO:0001782	TF_binding_site_variant	MODIFIER
SO:0001894	regulatory_region_ablation	MODERATE
SO:0001891	regulatory_region_amplification	MODIFIER
SO:0001907	feature_elongation	MODIFIER
SO:0001566	regulatory_region_variant	MODIFIER
SO:0001906	feature_truncation	MODIFIER
SO:0001628	intergenic_variant	MODIFIER

High and moderate variants were included in the analysis

Source: http://www.ensembl.org/info/genome/variation/predicted_data.html

Supplement Table S3. Association analysis for common variants

Chr	dbSNP	position	Ref	Alt	Gene	Function	MAF aHUS	MAF UI Control	P ₁	P ₁ adj	OR ₁	MAF gnomAD	P ₂	OR ₂
1	rs9287090	169510380	A	G	<i>F5</i>	p.Leu1316Leu	15.58%	28.00%	1.13E-10	6.41E-08	0.47	21.50%	3.85E-05	0.68
1	rs3753396	196695742	G	A	<i>CFH</i>	p.Gln672Gln	28.88%	15.33%	4.32E-13	6.12E-10	2.24	16.71%	1.86E-17	2.02
1	rs1065489	196709774	T	G	<i>CFH</i>	p.Glu936Asp	28.75%	15.33%	8.82E-13	6.25E-10	2.23	16.84%	8.34E-17	1.99
1	rs3828032	196920178	T	C	<i>CFHR2</i>	intronic	39.38%	28.17%	1.87E-07	5.30E-05	1.66	30.00%	1.85E-08	1.52
1	rs11118580	207959070	C	T	<i>CD46</i>	intronic	30.00%	21.17%	8.57E-06	1.87E-03	1.60	21.37%	1.29E-08	1.58
3	rs3733001	186938956	T	C	<i>MASP1</i>	intronic	30.50%	22.83%	1.42E-04	2.36E-02	1.48	25.18%	6.90E-04	1.30

MAF: minor allele frequency

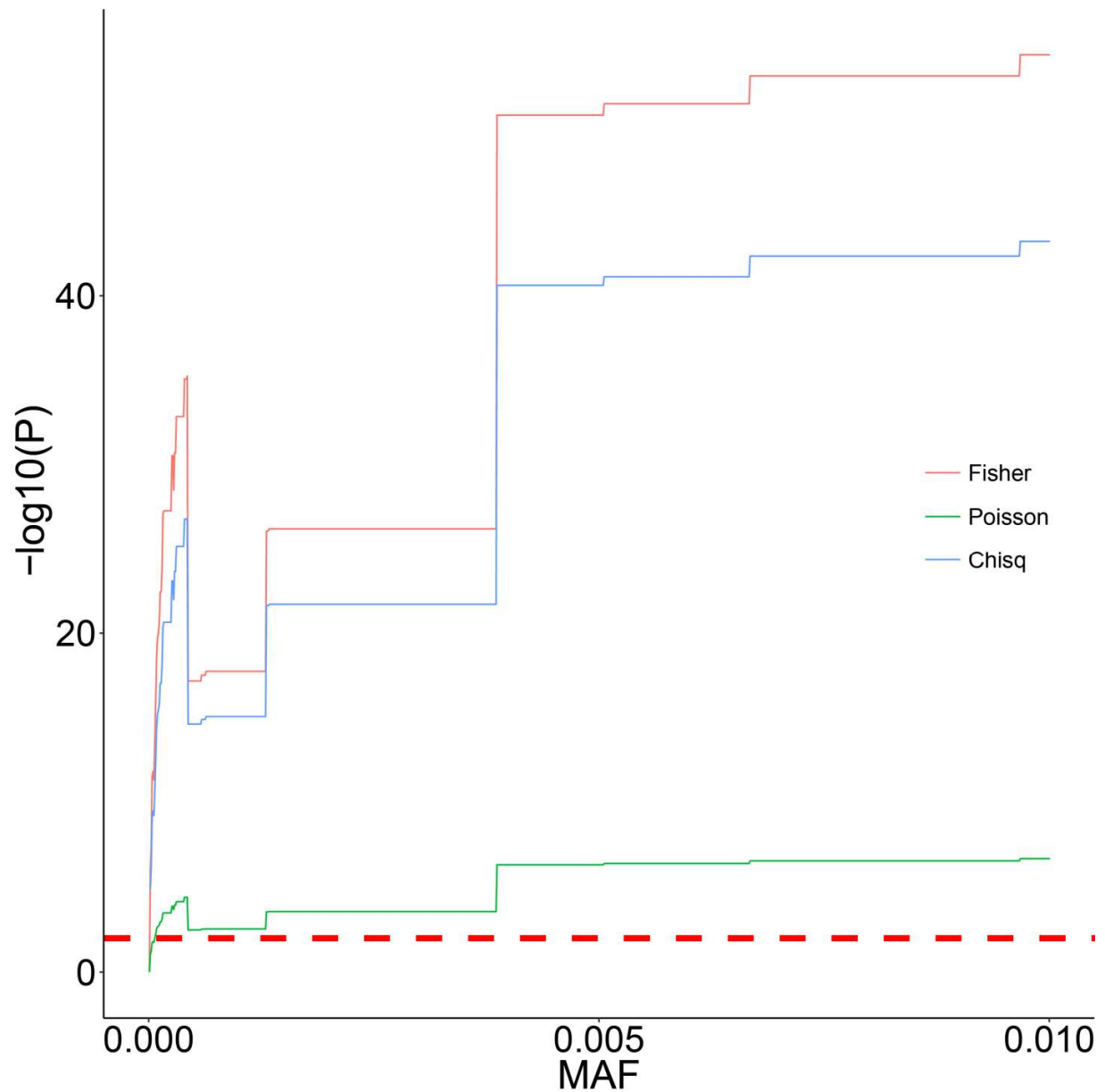
Supplement Table S4. FH-FHRs fusion proteins identified in 400 aHUS patients

Patient ID	Fusion Protein (inferred based on MLPA)
1	FH SCR 1-18 + FHR1 SCR 4-5
2	FH SCR 1-18 + FHR1 SCR 4-5
3	FH SCR 1-18 + FHR1 SCR 4-5
4	FH SCR 1-19 + FHR1 SCR 5
5	FH SCR 1-19 + FHR1 SCR 5
6	FHR1 SCR 1-2 + FH SCR 18-20
7	FHR1 SCR 1-2 + FH SCR 18-20
8	FHR1 SCR 1-3 + FH SCR 19-20
9	FHR3 SCR 1-4 + FHR4 SCR 4-5
10	FHR3 SCR 1-4 + FHR4 SCR 4-5
11	FHR3 SCR 1-4 + FHR4 SCR 4-5
12	FHR3 SCR 1-4 + FHR4 SCR 4-5
13	FHR3 SCR 1-4 + FHR4 SCR 4-5
14	Complex (FH and FHR1 involved)
15	Complex (FH, FHR3, FHR1, FHR4 and FHR2 involved)

Supplement Table S5. Rare variants with $0.1\% < \text{MAF (NFE)} < 1\%$ identified in *CFH*, *CD46*, *C3*, *CFI* and *CFB* in 400 aHUS patients

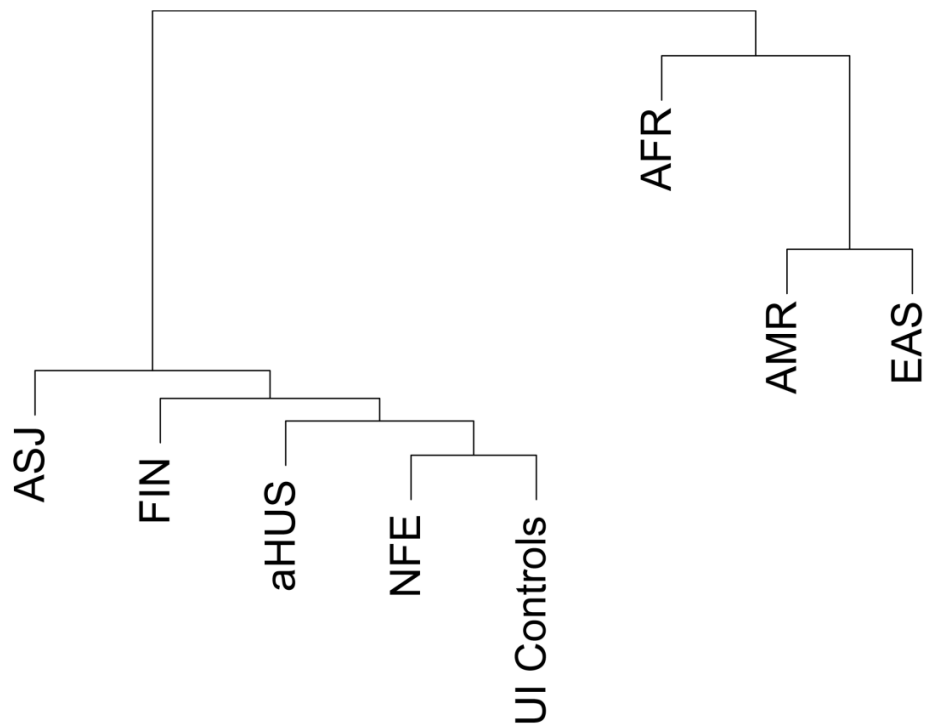
Gene	HGVS	dbSNP	aHUS MAF	Control MAF	NFE MAF	Max MAF	Max Pop	Pathogenicity
<i>CFH</i>	c.2850G>T; p.Gln950His	rs149474608	0.63%	0.75%	0.59%	1.78%	ASJ	B
<i>CFH</i>	c.2867C>T; p.Thr956Met	rs145975787	0.13%	0.33%	0.17%	0.17%	NFE	LB
<i>CFI</i>	c.1657C>T; p.Pro553Ser	rs113460688	0.38%	0.33%	0.27%	0.27%	NFE	LB
<i>CFI</i>	c.1322A>G; .Lys441Arg	rs41278047	0.75%	0.50%	0.24%	4.77%	ASJ	B
<i>CFI</i>	c.782G>A; p.Gly261Asp	rs112534524	0.13%	0.00%	0.19%	0.46%	ASJ	LB
<i>CFI</i>	c.782G>A; p.Gly261Asp	rs112534524	0.13%	0.00%	0.19%	0.46%	ASJ	LB
<i>CFB</i>	c.1697A>C; p.Glu566Ala	rs45484591	3.00%	0.00%	1.00%	2.10%	ASJ	VUS
<i>C3</i>	c.4855A>C; p.Ser1619Arg	rs2230210	0.63%	0.25%	0.22%	0.22%	NFE	VUS
<i>C3</i>	c.2203C>T; p.Arg735Trp	rs117793540	0.13%	0.17%	0.25%	1.24%	ASJ	B
<i>C3</i>	c.463A>C; p.Lys155Gln	rs147859257	0.75%	0.58%	0.54%	0.54%	NFE	LB

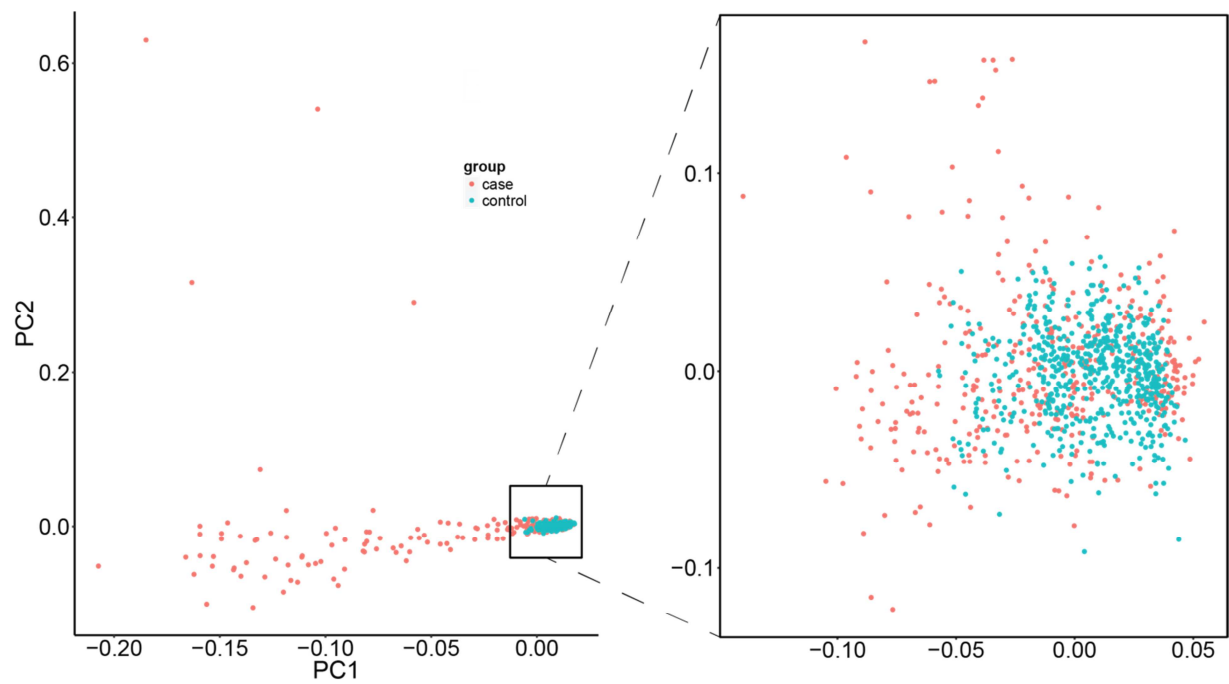
SUPPLEMENTAL FIGURES



Supplement Figure S1. Rare variants in the *CFH* gene are significantly more abundant in the non-Finnish European (NFE) subpopulation as compared to the Finnish (FIN) subpopulation from gnomAD due to population stratification.

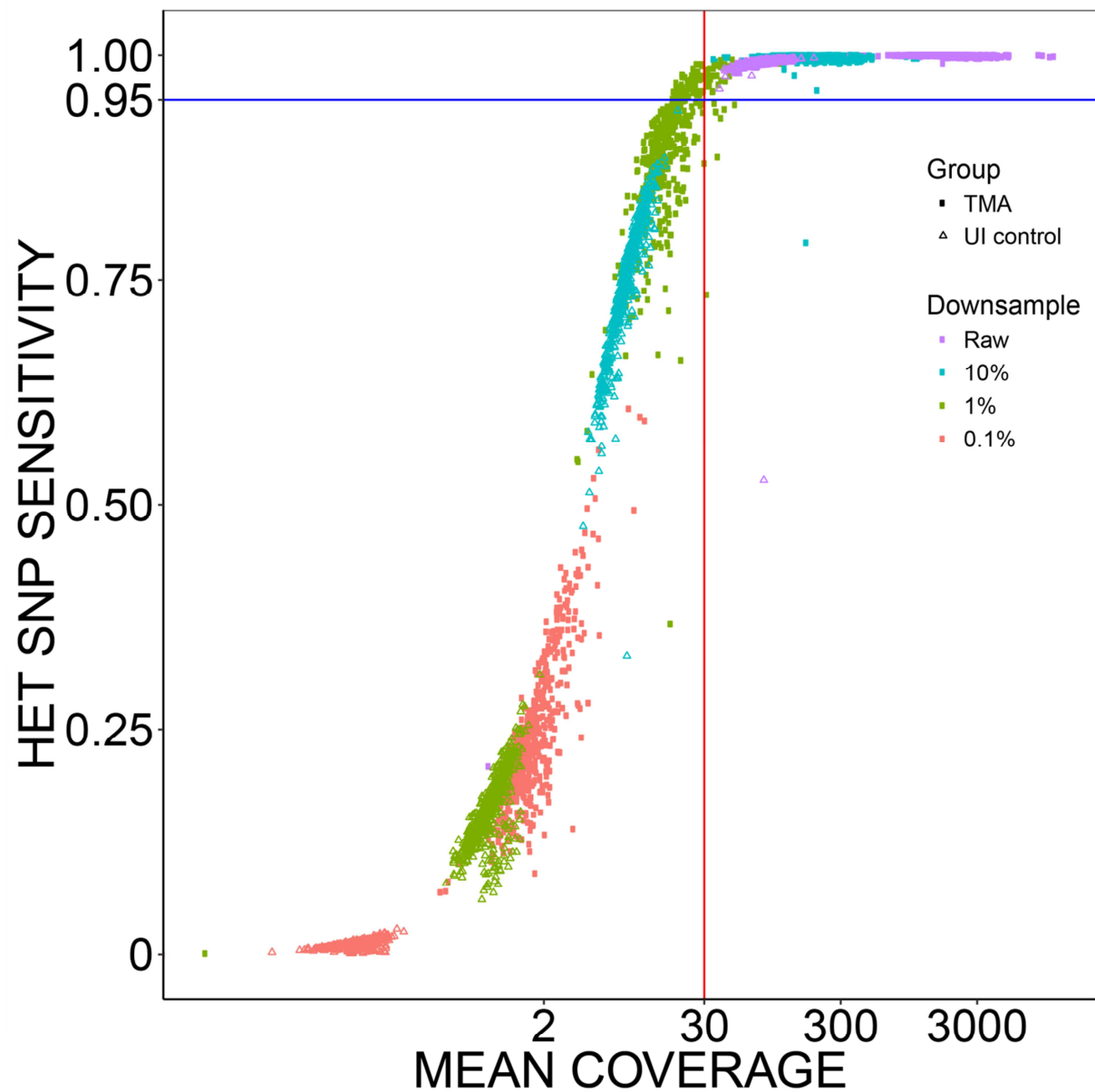
Three algorithms are used to test for enrichment and demonstrate that the modified Poisson exact test is least sensitive to population stratification. P values are shown as curves: red curve, Fisher's exact test; green curve, Poisson exact test; blue curve, Chi-square test; red dashed line, P value=0.05.





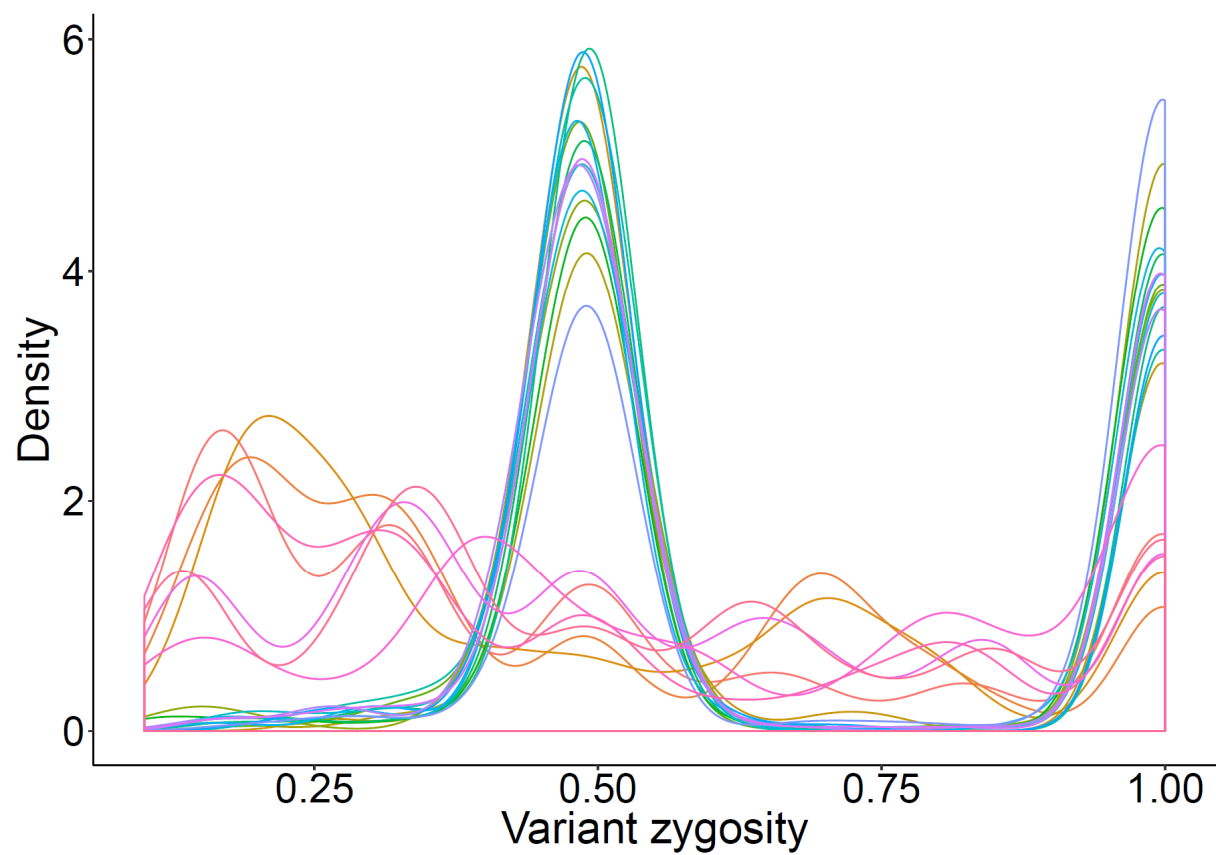
Supplement Figure S3. Population stratification within aHUS cases and UI controls was used to remove outliers.

Left panel, distribution of patients and controls prior to sample removal; right panel, distribution of cases and controls after removing outliers.



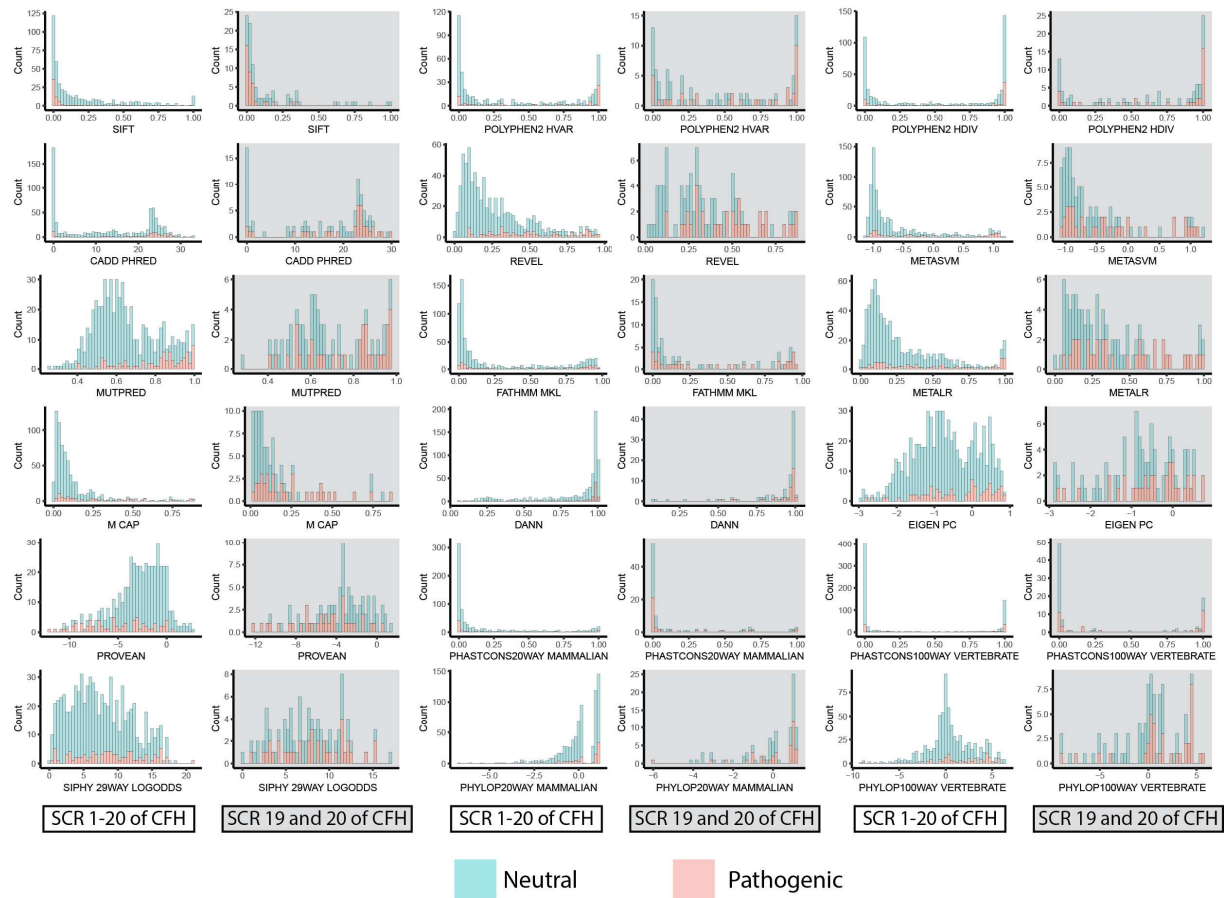
Supplement Figure S4. Mean coverage is correlated with theoretical heterozygous SNP sensitivity (HET SNP sensitivity) in patients and controls.

Random down sampling demonstrated a quick drop of HET SNP sensitivity when mean coverage is below 30X. Most raw sequencing data (purple) from patients (dots) and controls (triangles) is good. Two low quality samples were excluded. Blue line, 95% of HET SNP sensitivity, Red line, 30X of mean coverage.



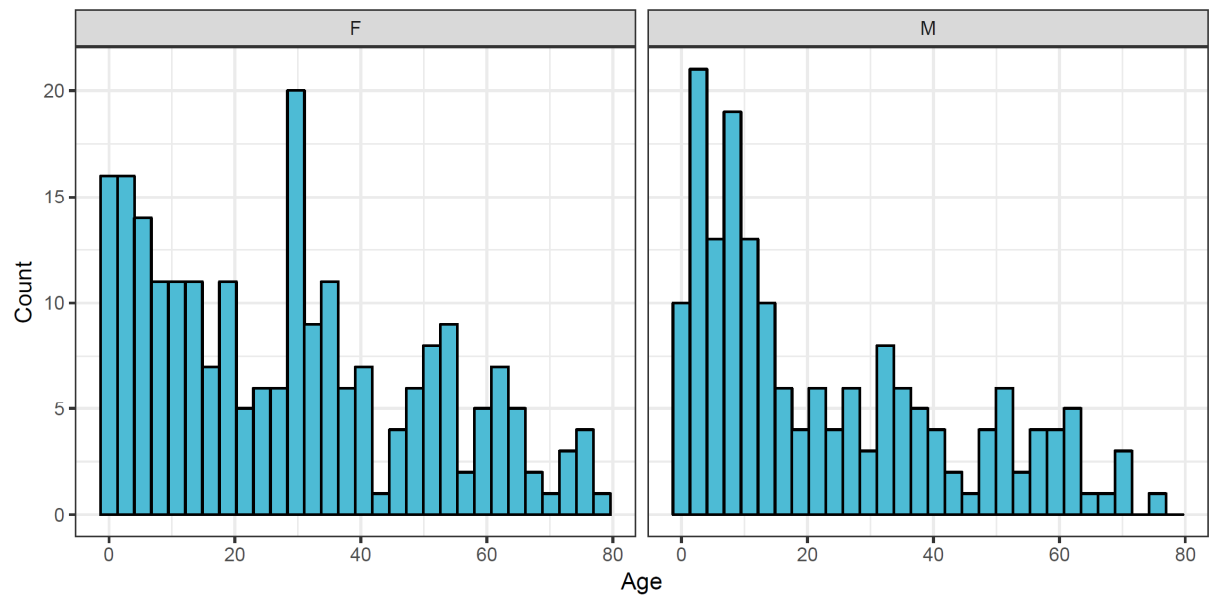
Supplement Figure S5. Six samples with a shifted ratio of ref/alt reads were excluded from the study cohort.

High quality samples are expected to have peaks at 0.5 and 1.0.



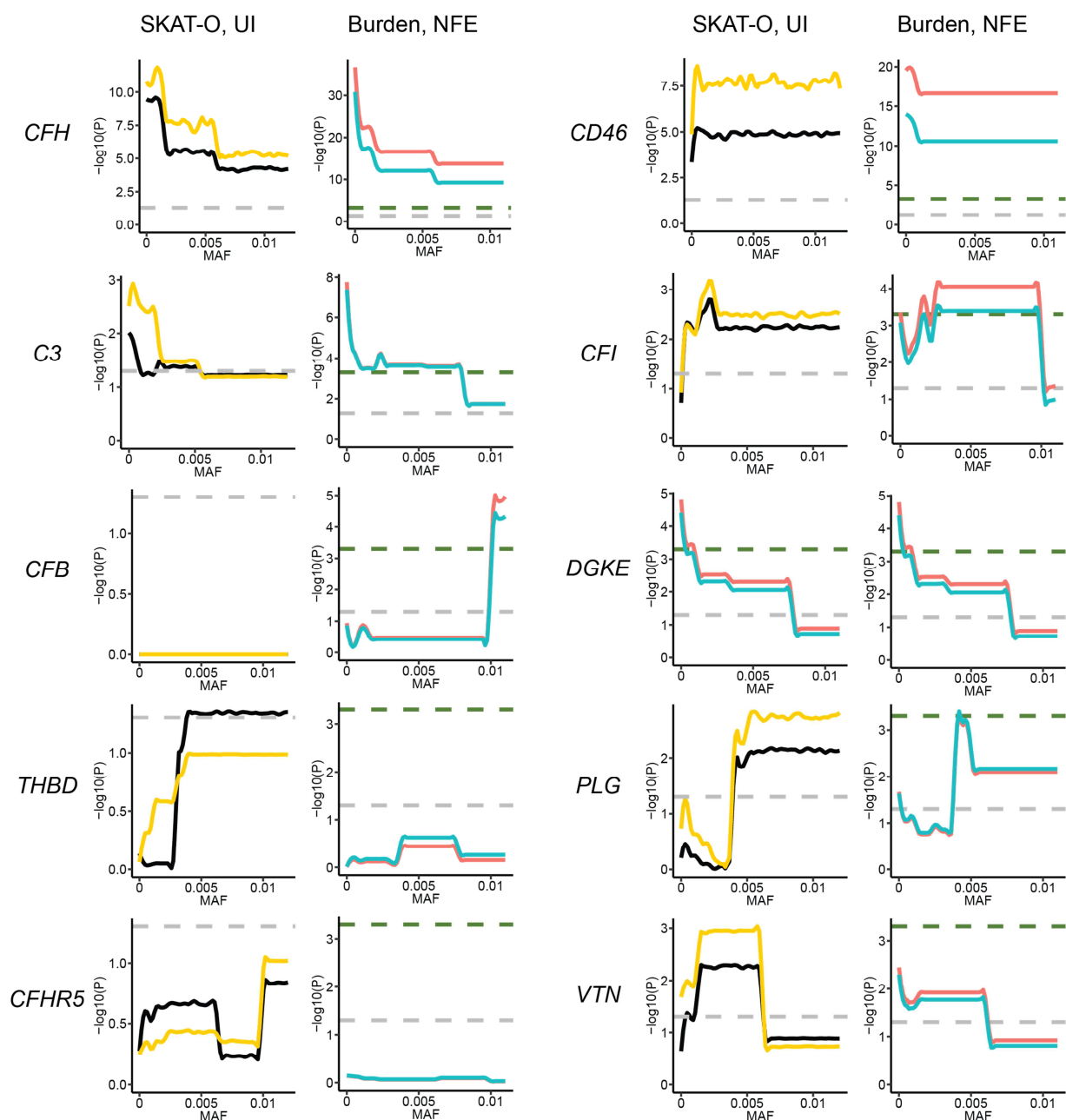
Supplement Figure S6. Prediction score distribution of neutral and pathogenic variants in *CFH*.

Neutral variants from gnomAD and pathogenic variants from the literature were both filtered by $MAF < 0.1\%$. 18 different tools were used to perform in-silico prediction on these variants to compare neutral and pathogenic variants across the entire gene (white background) or restricting the analysis to SCR19-20 (grey background). Note the heavily mixed distribution of neutral and pathogenic variants for all tools, indicating their ineffectiveness in predicting variant effect in *CFH* for aHUS.



Supplement Figure S7. Median age of female patients is significantly lower than that of male patients

Female: left panel, 27.3 years; Male: right panel, 14.8 years; Mann-Whitney U test $P = 0.011$



Supplement Figure S8: Enrichment of ultra-rare variants 'contaminates' the result of the association analysis when MAF thresholds are set higher.

The minor allele frequency threshold (cut off) was increased in a stepwise fashion to select variants for the analyses. Sets of p values are shown as curves: black curve, SKAT-O test adjusting for population stratification in UI controls; yellow curve, SKAT-O test without adjusting in UI controls; red curve, Fisher's exact test in NEF controls; blue curve, Poisson exact test in NEF controls; grey dashed line, $P < 0.05$; green dashed line, $P < 0.0005$.