

Supplementary Note

Inherited *GATA3* variants are associated with Ph-like childhood acute lymphoblastic leukemia and risk of relapse

¹Virginia Perez-Andreu, ²Kathryn G. Roberts, ³Richard C. Harvey, ¹Wenjian Yang, ⁴Cheng Cheng, ⁴Deqing Pei, ¹Heng Xu, ^{5,6}Julie Gastier-Foster, ¹Shuyu E, ^{1,7}Joshua Yew-Suang Lim, ³I-Ming Chen, ⁸Yiping Fan, ⁹Meenakshi Devidas, ¹⁰Michael J. Borowitz, ¹Colton Smith, ¹¹Geoffrey Neale, ¹²Esteban G. Burchard, ¹²Dara G. Torgerson, ¹³Federico Antillon Klussmann, ¹³Cesar Rolando Najera Villagran, ¹⁴Naomi J. Winick, ¹⁵Bruce M. Camitta, ¹⁶Elizabeth Raetz, ¹⁷Brent WoodM.D., ¹⁸Feng Yue, ¹⁶William L. Carroll, ¹⁹Eric Larsen, ²⁰W. Paul Bowman, ²¹Mignon L. Loh, ²²Michael Dean, ²³Deepa Bhojwani, ²³Ching-Hon Pui, ¹William E. Evans, ¹Mary V. Relling, ²⁴Stephen P. Hunger, ³Cheryl L. Willman, ²Charles G. Mullighan, ¹Jun J. Yang

¹Department of Pharmaceutical Sciences, St. Jude Children's Research Hospital, Memphis, TN, USA; ²Department of Pathology, St. Jude Children's Research Hospital, Memphis, TN, USA; ³Cancer Center, University of New Mexico, Albuquerque, NM, USA; ⁴Department of Biostatistics, St. Jude Children's Research Hospital, Memphis, TN, USA; ⁵Department of Pathology and Laboratory Medicine, Nationwide Children's Hospital, Columbus, OH, USA; ⁶Departments of Pediatrics, Ohio State University School of Medicine, Columbus, OH, USA; ⁷Department of Pediatrics, Yong Loo Lin School of Medicine, National University of Singapore, Singapore; ⁸Department of Computational Biology, St. Jude Children's Research Hospital, Memphis, TN, USA; ⁹Department of Epidemiology and Health Policy Research, University of Florida, Gainesville, FL, USA; ¹⁰Johns Hopkins Medical Institute, Baltimore, MD, USA; ¹¹Hartwell Center for Bioinformatics & Biotechnology, St. Jude Children's Research Hospital, Memphis, TN, USA; ¹²Department of Bioengineering & Therapeutic Science and Medicine, University of California at San Francisco, San Francisco, CA, USA; ¹³Unidad Nacional de Oncologia Pediatrica, Guatemala City, Guatemala; ¹⁴Department of Pediatric Hematology/Oncology, University of Texas Southwestern Medical Center, Dallas, TX, USA; ¹⁵Department of Pediatrics, Medical College of Wisconsin, Milwaukee, WI, USA; ¹⁶New York University Cancer Institute, New York, NY, USA; ¹⁷Department of Laboratory Medicine, University of Washington, Seattle, WA, USA; ¹⁸Ludwig Institute for Cancer Research, University of California at San Diego School of Medicine, La Jolla, CA, USA; ¹⁹Maine Children's Cancer Program, Scarborough, ME, USA; ²⁰Cook Children's Medical Center, Ft Worth, TX, USA; ²¹Department of Pediatrics, University of California at San Francisco, San Francisco, CA, USA; ²²Laboratory of Experimental Immunology, National Cancer Institute, Frederick, MD, USA; ²³Department of Oncology, St. Jude Children's Research Hospital, Memphis, TN, USA; ²⁴Children's Hospital Colorado, University of Colorado, Aurora, CO, USA.

Corresponding author

Jun J. Yang PhD
Dept. of Pharmaceutical Sciences, MS 313
St. Jude Children's Research Hospital
262 Danny Thomas Place, Memphis, TN 38105-3678
Email address: jun.yang@stjude.org
Phone: (901) 595-2517, FAX: (901) 595-8869

Acknowledgements for dbGaP datasets

MESA and the MESA SHARe project are conducted and supported by the National Heart, Lung, and Blood Institute (NHLBI) in collaboration with MESA investigators. Support for MESA is provided by contracts N01-HC-95159 through N01-HC-95169 and RR-024156. MESA, MESA Family, and the MESA SHARe project are conducted and supported by the National Heart, Lung, and Blood Institute (NHLBI) in collaboration with MESA investigators. Support for MESA is provided by contracts N01-HC-95159 through N01-HC-95169 and RR-024156.

Funding support for the Genome-Wide Association of Schizophrenia Study was provided by the National Institute of Mental Health (R01 MH67257, R01 MH59588, R01 MH59571, R01 MH59565, R01 MH59587, R01 MH60870, R01 MH59566, R01 MH59586, R01 MH61675, R01 MH60879, R01 MH81800, U01 MH46276, U01 MH46289 U01 MH46318, U01 MH79469, and U01 MH79470) and the genotyping of samples was provided through the Genetic Association Information Network (GAIN). The datasets used for the analyses described in this manuscript were obtained from the database of Genotype and Phenotype database (dbGaP) found at <http://www.ncbi.nlm.nih.gov/gap> through dbGaP (accession number phs000021.v1.p1). Samples and associated phenotype data for the Genome-Wide Association of Schizophrenia Study were provided by the Molecular Genetics of Schizophrenia Collaboration (Principal Investigator: P.V. Gejman, Evanston Northwestern Healthcare (ENH) and Northwestern University, Evanston, Illinois).

Funding support for the Whole Genome Association Study of Bipolar Disorder was provided by the US National Institute of Mental Health (NIMH) and the genotyping of samples was provided through the Genetic Association Information Network (GAIN). The datasets used for the analyses described in this manuscript were obtained from dbGaP (accession number phs000017.v1.p1). Samples and associated phenotype data for the Collaborative Genomic Study of Bipolar Disorder were provided by The NIMH Genetics Initiative for Bipolar Disorder. Data and biomaterials were collected in four projects that participated in the NIMH Bipolar Disorder Genetics Initiative. From 1991 to 1998, the principal investigators and co-investigators were: Indiana University, Indianapolis,

Indiana, U01 MH46282, J. Nurnberger, M. Miller and E. Bowman; Washington University, St. Louis, Missouri, U01 MH46280, T. Reich, A. Goate and J. Rice; Johns Hopkins University, Baltimore, Maryland, U01 MH46274, J.R. DePaulo Jr., S. Simpson and C. Stine; NIMH Intramural Research Program, Clinical Neurogenetics Branch, Bethesda, Maryland, E. Gershon, D. Kazuba and E. Maxwell. Data and biomaterials were collected as part of ten projects that participated in the NIMH Bipolar Disorder Genetics Initiative. From 1999 to 2003, the principal investigators and co-investigators were: Indiana University, Indianapolis, Indiana, R01 MH59545, J. Nurnberger, M.J. Miller, E.S. Bowman, N.L. Rau, P.R. Moe, N. Samavedy, R. El-Mallakh (University of Louisville, Louisville, Kentucky), H. Manji (Wayne State University, Detroit, Michigan), D.A. Glitz (Wayne State University), E.T. Meyer, C. Smiley, T. Foroud, L. Flury, D.M. Dick and H. Edenberg; Washington University, St. Louis, Missouri, R01 MH059534, J. Rice, T. Reich, A. Goate and L. Bierut; Johns Hopkins University, Baltimore, Maryland, R01 MH59533, M. McInnis, J.R. DePaulo Jr., D.F. MacKinnon, F.M. Mondimore, J.B. Potash, P.P. Zandi, D. Avramopoulos and J. Payne; University of Pennsylvania, Philadelphia, Pennsylvania, R01 MH59553, W. Berrettini; University of California, Irvine, California, R01 MH60068, W. Byerley and M. Vawter; University of Iowa, Iowa City, Iowa, R01 MH059548, W. Coryell and R. Crowe; University of Chicago, Chicago, Illinois, R01 MH59535, E. Gershon, J. Badner, F. McMahon, C. Liu, A. Sanders, M. Caserta, S. Dinwiddie, T. Nguyen and D. Harakal; University of California, San Diego, California, R01 MH59567, J. Kelsoe and R. McKinney; Rush University, Chicago, Illinois, R01 MH059556, W. Scheftner, H.M. Kravitz, D. Marta, A. Vaughn-Brown and L. Bederow; NIMH Intramural Research Program, Bethesda, Maryland, 1Z01MH002810-01, F.J. McMahon, L. Kassem, S. Detera-Wadleigh, L. Austin and D.L. Murphy.

Supplementary Table 1. Multivariate analysis of rs3824662 and rs3781093 for association with Ph-like ALL in COG AALL0232.

Gene	rs ID	Allele A	Allele B	Position (hg19)	SNP type	Ph-like ALL vs. non-Ph-like ALL				Ph-like ALL vs. non-ALL controls			
						P-value ¹	OR (95%, CI)	P-value ²	OR (95%, CI)	P-value ³	OR (95%, CI)	P-value ⁴	OR (95%, CI)
GATA3	rs3781093	T	C	8101927	Genotyped	2.62×10 ⁻⁷	3.09 (2.06-4.63)	0.306684	1.53 (0.66-3.51)	4.94×10 ⁻¹²	3.70 (2.61-5.25)	0.26976	1.58 (0.69-3.63)
GATA3	rs3824662	C	A	8104208	Genotyped	1.05×10 ⁻⁸	3.17 (2.12-4.74)	0.008	3.09 (1.32-7.25)	2.17×10 ⁻¹⁴	3.75 (2.65-5.30)	0.001	3.49 (1.61-7.55)

Abbreviations: OR, Odds ratio; CI, confidence interval.

Association of SNP genotype and Ph-like ALL was evaluated by logistic regression, after adjusting for genetic ancestry. Chromosomal locations are based on hg19.

P-value¹: Ph-like ALL vs. non-Ph-like ALL adjusting for genetic ancestry

P-value²: Ph-like ALL vs. non-Ph-like ALL adjusting for genetic ancestry and GATA3 SNPs genotype.

P-value³: Ph-like ALL vs. non-ALL controls adjusting for genetic ancestry

P-value⁴: Ph-like ALL vs. non-ALL controls adjusting for genetic ancestry and GATA3 SNPs genotype.

Supplementary Table 2. Clinical characteristics of patients included in this study by cohort.

	Children's Oncology Group Cohort		
Treatment Protocol	AALL0232	P9906	P9905
Number of patients	(n=511)	(n=215)	(n=889)
Race^a			
Asian	17 (3.0)	3 (1.3)	19 (2.1)
African-American	25 (5.0)	14 (6.5)	51 (5.7)
Hispanic	149 (29.1)	53 (24.6)	191 (21.5)
European-American	160 (31.3)	116 (54.0)	520 (58.5)
Sex			
Female	237 (46.4)	68 (31.6)	452 (50.8)
Male	273 (53.5)	147 (68.3)	437 (49.2)
Missing	1 (0.1)	NA	NA
Age at diagnosis, y			
<10	222 (43.5)	73 (34.0)	690 (77.6)
≥10	288 (56.4)	142 (66.0)	198 (22.3)
Missing	1 (0.1)	NA	1 (0.11)
Leucocyte count at diagnosis, μL			
<50 000	223 (43.7)	120 (56.0)	737 (82.9)
≥50 000	287 (56.2)	95 (44.1)	152 (17.1)
Missing	1 (0.1)	NA	NA
CNS status			
CNS3 or traumatic	11 (2.5)	28 (13.0)	0 (0.0)
CNS 1	397 (81.1)	163 (76.0)	795 (89.4)
CNS 2	81 (16.0)	24 (11.1)	93 (10.5)
Missing	22 (4.3)	NA	1 (0.11)

Data is presented as No. (%) unless otherwise indicated.

^aGenetic ancestry was determined by using STRUCTURE. Asians, African-Americans, Hispanics, and European-Americans were identified as >90%, >70%, >10% and higher than African ancestry, >95% of Asian, African, Native-American and European genetic ancestry, respectively.

NA: Not applicable.

Supplementary Table 3A. List of primers for Sanger sequencing of *GATA3* SNPs (rs3824662 and rs3781093).

<i>GATA 3</i> SNPs	PCR Primers		Sanger sequencing Primers
	Forward	Reverse	
rs3824662	5'-TATCACCCCTCCCCACCA	5'-GGAAAGCCCCAGATCAA	5'-TATCACCCCTCCCCACCA
rs3781093	5'-TTCCTGTGCTCTGTTTCCTT	5'-GGCTCAGGATAAACAATG	5'-TTCCTGTGCTCTGTTTCCTT

Supplementary Table 3B. List of primers for real-time PCR of *GATA3* mRNA

<i>GATA3</i> Forward Primer	5'-TCACAAAATGAACGGACAGAACC-3'
<i>GATA3</i> Reverse Primer	5'-CAGCCTTCGCTTGGGCTTAAT-3'

Supplementary Table 4. Association of germline *JAK2* SNPs with somatic *JAK2* mutation

rsID	<i>P</i>-value¹	OR (95%, CI)	<i>P</i>-value²	OR (95%, CI)
<i>rs2149556</i>	0.1978	1.35 (0.84-2.15)	0.1166	1.41 (0.91-2.19)
<i>rs7864782</i>	0.2329	1.30 (0.83-2.03)	0.1703	1.34 (0.87-2.05)
<i>rs10815144</i>	0.2584	0.76 (0.48-1.21)	0.2044	0.75 (0.48-1.16)
<i>rs10124001</i>	0.2854	1.97 (0.56-6.94)	0.3740	1.70 (0.51-5.59)
<i>rs10119004</i>	0.4588	0.84 (0.53-1.32)	0.3488	0.81 (0.52-1.25)
<i>rs10974944</i>	0.6754	0.89 (0.54-1.48)	0.6591	0.89 (0.54-1.47)
<i>rs11793659</i>	0.7882	1.06 (0.66-1.71)	0.7556	1.07 (0.67-1.72)
<i>rs1327493</i>	0.8624	1.10 (0.35-3.42)	0.9095	1.06 (0.37-3.02)
<i>rs17425637</i>	0.9616	1.01 (0.62-1.64)	0.9505	1.01 (0.63-1.63)
<i>rs12340895</i>	0.9661	0.98 (0.59-1.63)	0.9390	0.98 (0.59-1.60)
<i>rs6476934</i>	0.9921	1.00 (0.23-4.40)	0.7303	1.28 (0.30-5.30)

Association between *JAK2* somatic lesion and *JAK2* germline SNPs was tested in the combined cohort (COG AALL0232, COG P9906 and COG P9905) and in the discovery non-ALL control group (N=6,661), by logistic regression after adjusting for genetic ancestry.

P-value¹ : comparing allele frequency between ALL with vs. without *JAK2* mutations

P-value² : comparing allele frequency between ALL with *JAK2* mutations vs. non-ALL controls

Abbreviations: OR, Odds Ratio; CI, confidence interval

Supplementary Table 5. Gene Set Enrichment Analysis of ALL cells overexpressing GATA3 vs. control

Nalm6				UOCB1			
Gene symbol	Rank Metric Score	Running ES	Core Enrichment	Gene symbol	Rank Metric Score	Running ES	Core Enrichment
<i>ANXA1*</i>	0.2595	0.3822	Yes	<i>ANXA1*</i>	3.9373	0.1921	Yes
<i>BCL6</i>	0.6874	0.1169	Yes	<i>CASP10</i>	0.2628	0.4673	Yes
<i>CASP10</i>	0.2619	0.3630	Yes	<i>GBP2*</i>	0.5453	0.3644	Yes
<i>CD99*</i>	0.4257	0.2728	Yes	<i>IL2RA*</i>	0.2100	0.5066	Yes
<i>CTHRC1*</i>	0.1918	0.4009	Yes	<i>LCP2</i>	0.4846	0.4372	Yes
<i>DOK4*</i>	0.2165	0.3922	Yes	<i>LYZ</i>	0.5054	0.4138	Yes
<i>ECM1*</i>	0.1871	0.4130	Yes	<i>MYBL1*</i>	0.2035	0.5128	Yes
<i>GBP2*</i>	0.6666	0.1664	Yes	<i>ANKRD28</i>	0.3527	0.4689	Yes
<i>IL2RA*</i>	0.1696	0.4312	Yes	<i>ANTXR2*</i>	1.6885	0.2743	Yes
<i>LCP2</i>	0.3432	0.3258	Yes	<i>CD300A*</i>	0.4056	0.4548	Yes
<i>LYZ</i>	0.3145	0.3479	Yes	<i>CDC42EP3*</i>	0.5181	0.3893	Yes
<i>MYBL1*</i>	0.4396	0.2415	Yes	<i>ENAM*</i>	0.2026	0.5223	Yes
<i>NRXN3*</i>	0.6131	0.2114	Yes	<i>MMRN1*</i>	0.2250	0.4844	Yes
<i>RAPGEF3*</i>	0.4032	0.3023	Yes	<i>PON2</i>	0.2216	0.4934	Yes
<i>SELL*</i>	0.8998	0.0662	Yes	<i>PSTPIP2*</i>	1.3373	0.3394	Yes
<i>TTN</i>	0.1702	0.4187	Yes	<i>S100Z*</i>	0.1939	0.5244	Yes
<i>ABCA9</i>	0.0429	0.3377	No	<i>SERPINA1</i>	0.2173	0.5017	Yes
<i>ABL1</i>	-0.0621	-0.0746	No	<i>STON2</i>	0.2425	0.4831	Yes
<i>AHR</i>	-0.0154	0.0598	No	<i>TBXAS1*</i>	0.2468	0.4733	Yes
<i>ANKRD28</i>	0.0419	0.3347	No	<i>BCL6</i>	-0.0679	0.0227	No
<i>ANTXR2</i>	0.0136	0.2152	No	<i>CD99</i>	-0.0369	0.0885	No
<i>ANXA4</i>	0.0272	0.2833	No	<i>CTHRC1</i>	-0.0742	0.0131	No
<i>ATP10A</i>	0.0513	0.3482	No	<i>DOK4</i>	-0.1730	-0.0245	No
<i>B3GNTL1</i>	-0.0376	-0.0216	No	<i>ECM1</i>	-0.0115	0.1645	No
<i>BAALC</i>	0.1278	0.4113	No	<i>NRXN3</i>	0.0116	0.2491	No
<i>BSPRY</i>	-0.0862	-0.1009	No	<i>RAPGEF3</i>	-0.0096	0.1706	No
<i>BST1</i>	-0.0242	0.0276	No	<i>SELL</i>	0.0753	0.4321	No
<i>C1QTNF4</i>	-0.0091	0.0879	No	<i>TTN</i>	0.0200	0.2757	No
<i>CA6</i>	0.0403	0.3357	No	<i>ABCA9</i>	-0.0732	0.0115	No
<i>CASP1</i>	-0.1761	-0.1232	No	<i>ABL1</i>	0.0543	0.3721	No
<i>CAV1</i>	0.0901	0.3956	No	<i>AHR</i>	-0.0090	0.1722	No
<i>CCL17</i>	0.0917	0.3925	No	<i>ANXA4</i>	0.0048	0.2253	No
<i>CCND2</i>	-0.0397	-0.0236	No	<i>ATP10A</i>	-0.0898	0.0054	No
<i>CD300A</i>	0.0414	0.3387	No	<i>B3GNTL1</i>	-0.1538	-0.0372	No
<i>CD302</i>	-0.0134	0.0677	No	<i>BAALC</i>	0.1046	0.4937	No
<i>CDC42EP3</i>	-0.0111	0.0788	No	<i>BSPRY</i>	-0.2403	-0.0231	No
<i>CEACAM6</i>	0.1038	0.4023	No	<i>BST1</i>	-0.0319	0.1004	No
<i>CFD</i>	0.0782	0.3889	No	<i>C1QTNF4</i>	0.0185	0.2704	No
<i>CFP</i>	0.0583	0.3540	No	<i>CA6</i>	0.1201	0.5037	No
<i>CHN1</i>	-0.3405	0.0038	No	<i>CASP1</i>	-0.1978	-0.0254	No
<i>CHN2</i>	0.0431	0.3355	No	<i>CAV1</i>	-0.0197	0.1393	No
<i>CHRNA1</i>	-0.2642	-0.0784	No	<i>CCL17</i>	-0.0075	0.1774	No
<i>CRADD</i>	0.0441	0.3308	No	<i>CCND2</i>	-0.0671	0.0213	No
<i>CSTA</i>	0.1191	0.4095	No	<i>CD302</i>	0.1744	0.5147	No
<i>CTDSPL</i>	-0.1187	-0.1181	No	<i>CEACAM6</i>	0.1266	0.5041	No
<i>CYYR1</i>	0.0535	0.3452	No	<i>CFD</i>	0.1061	0.4867	No
<i>DFNA5</i>	0.0440	0.3338	No	<i>CFP</i>	-0.0600	0.0324	No
<i>DPYD</i>	0.0729	0.3828	No	<i>CHN1</i>	-0.0538	0.0476	No
<i>DUSP6</i>	-0.0355	-0.0157	No	<i>CHN2</i>	-0.0228	0.1293	No
<i>EGFL7</i>	0.0858	0.3926	No	<i>CHRNA1</i>	0.0218	0.2801	No
<i>EMP1</i>	-0.0144	0.0640	No	<i>CRADD</i>	-0.1591	-0.0280	No
<i>ENAM</i>	0.0495	0.3431	No	<i>CSTA</i>	-0.0421	0.0771	No

Supplementary Table 5. Gene Set Enrichment Analysis of ALL cells overexpressing GATA3 vs. control

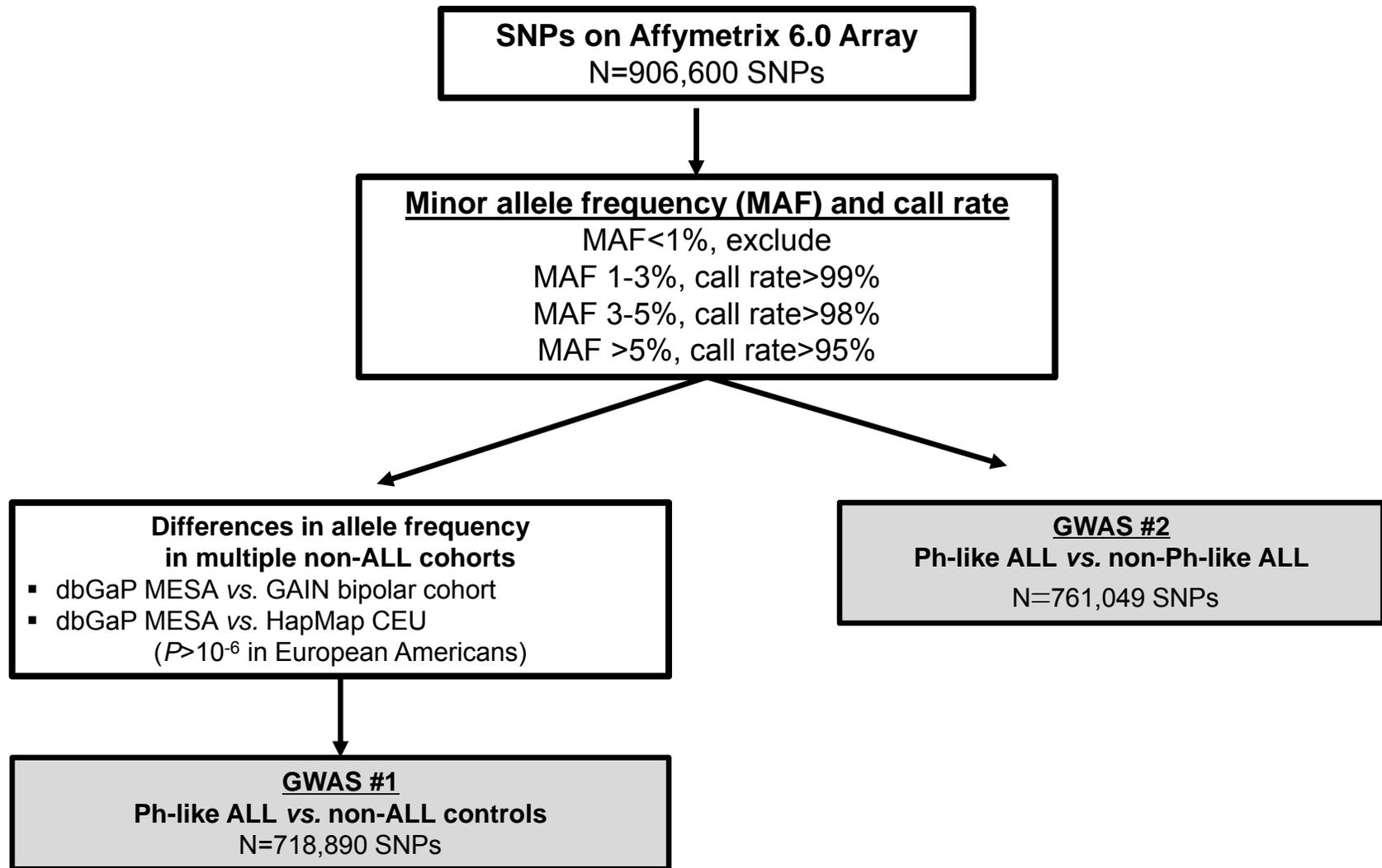
<i>FARP1</i>	0.0014	0.1468	No	<i>CTDSPL</i>	0.0202	0.2751	No
<i>FBXW7</i>	0.0151	0.2220	No	<i>CYYR1</i>	0.0955	0.4788	No
<i>FSCN1</i>	-0.0196	0.0448	No	<i>DFNA5</i>	0.1338	0.4976	No
<i>FUT7</i>	0.0688	0.3732	No	<i>DPYD</i>	0.1284	0.5012	No
<i>GADD45A</i>	-0.0849	-0.1047	No	<i>DUSP6</i>	-0.0181	0.1437	No
<i>GBP5</i>	-0.0004	0.1359	No	<i>EGFL7</i>	-0.0895	-0.0028	No
<i>GIMAP1</i>	0.0255	0.2758	No	<i>EMP1</i>	0.0000	0.2066	No
<i>GIMAP6</i>	0.0653	0.3710	No	<i>FARP1</i>	0.0020	0.2141	No
<i>GLIPR1</i>	-0.0380	-0.0201	No	<i>FBXW7</i>	-0.0333	0.0973	No
<i>GPR110</i>	0.0091	0.1925	No	<i>FSCN1</i>	-0.0065	0.1811	No
<i>GPR56</i>	-0.0186	0.0482	No	<i>FUT7</i>	-0.1438	-0.0410	No
<i>GYPC</i>	-0.2552	-0.0971	No	<i>GADD45A</i>	0.0112	0.2483	No
<i>HES1</i>	0.0324	0.3073	No	<i>GBP5</i>	0.1386	0.4992	No
<i>ID1</i>	-0.0214	0.0398	No	<i>GIMAP1</i>	0.0068	0.2330	No
<i>IFITM1</i>	0.0521	0.3474	No	<i>GIMAP6</i>	-0.0152	0.1520	No
<i>IFITM3</i>	-0.0207	0.0415	No	<i>GLIPR1</i>	0.0347	0.3219	No
<i>IGFBP7</i>	0.0565	0.3503	No	<i>GPR110</i>	0.1500	0.5087	No
<i>IPO11</i>	-0.0831	-0.1068	No	<i>GPR56</i>	-0.0421	0.0751	No
<i>KAZALD1</i>	0.0149	0.2219	No	<i>GYPC</i>	-0.3484	0.0131	No
<i>KBTBD8</i>	-0.0616	-0.0780	No	<i>HES1</i>	0.0354	0.3224	No
<i>KCNE3</i>	-0.0062	0.1034	No	<i>ID1</i>	0.0408	0.3350	No
<i>KLF9</i>	-0.0309	0.0016	No	<i>IFITM1</i>	-0.1198	-0.0350	No
<i>LIMS1</i>	0.0471	0.3390	No	<i>IFITM3</i>	0.0504	0.3609	No
<i>LST1</i>	0.0176	0.2340	No	<i>IGFBP7</i>	-0.2980	-0.0117	No
<i>MAPKAPK3</i>	-0.1170	-0.1339	No	<i>IPO11</i>	-0.0529	0.0474	No
<i>MCTP1</i>	0.0631	0.3675	No	<i>KAZALD1</i>	-0.2275	-0.0298	No
<i>MDFIC</i>	-0.0131	0.0689	No	<i>KBTBD8</i>	-0.1121	-0.0274	No
<i>MINA</i>	-0.2657	-0.0588	No	<i>KCNE3</i>	0.0494	0.3608	No
<i>MMP28</i>	0.0080	0.1856	No	<i>KLF9</i>	-0.1294	-0.0418	No
<i>MMRN1</i>	-0.0053	0.1083	No	<i>LIMS1</i>	0.1334	0.5031	No
<i>MS4A4A</i>	-0.0184	0.0481	No	<i>LST1</i>	-0.0744	0.0160	No
<i>MSRB3</i>	-0.0403	-0.0234	No	<i>MAPKAPK3</i>	-0.0757	0.0158	No
<i>MUC4</i>	0.0122	0.2088	No	<i>MCTP1</i>	0.0181	0.2703	No
<i>NFE2L2</i>	0.0156	0.2239	No	<i>MDFIC</i>	-0.0157	0.1510	No
<i>NPDC1</i>	0.0440	0.3370	No	<i>MINA</i>	0.0160	0.2639	No
<i>NT5E</i>	0.0310	0.3013	No	<i>MMP28</i>	-0.0127	0.1607	No
<i>NUDT4</i>	-0.0501	-0.0541	No	<i>MS4A4A</i>	-0.1490	-0.0391	No
<i>OLFML2A</i>	0.0851	0.3978	No	<i>MSRB3</i>	0.1058	0.4914	No
<i>OR7A5</i>	0.0071	0.1811	No	<i>MUC4</i>	0.0418	0.3359	No
<i>PELI1</i>	-0.0601	-0.0782	No	<i>NFE2L2</i>	-0.0133	0.1589	No
<i>PHACTR1</i>	-0.1256	-0.1077	No	<i>NPDC1</i>	0.0222	0.2800	No
<i>PON2</i>	0.0334	0.3103	No	<i>NT5E</i>	-0.1546	-0.0305	No
<i>PRX</i>	0.0097	0.1947	No	<i>NUDT4</i>	-0.3040	0.0019	No
<i>PSTPIP2</i>	0.0383	0.3286	No	<i>OLFML2A</i>	-0.0050	0.1866	No
<i>PTPN14</i>	0.1094	0.4045	No	<i>OR7A5</i>	0.0306	0.3101	No
<i>RASSF8</i>	0.0760	0.3877	No	<i>PELI1</i>	0.1256	0.5087	No
<i>ROBO3</i>	-0.0014	0.1303	No	<i>PHACTR1</i>	0.0404	0.3353	No
<i>ROBO4</i>	0.0079	0.1859	No	<i>PRX</i>	0.0332	0.3178	No
<i>RRAS</i>	-0.0565	-0.0714	No	<i>PTPN14</i>	0.0041	0.2225	No
<i>S100A8</i>	-0.1265	-0.0994	No	<i>RASSF8</i>	-0.0377	0.0878	No
<i>S100Z</i>	-0.0623	-0.0705	No	<i>ROBO3</i>	0.0867	0.4590	No
<i>SCHIP1</i>	0.0829	0.3969	No	<i>ROBO4</i>	0.0173	0.2682	No
<i>SERPINA1</i>	-0.0707	-0.0891	No	<i>RRAS</i>	-0.0357	0.0913	No
<i>SH3BP5</i>	-0.0021	0.1259	No	<i>S100A8</i>	-0.1936	-0.0319	No
<i>SLC2A5</i>	-0.2737	-0.0393	No	<i>SCHIP1</i>	-0.0822	0.0082	No

Supplementary Table 5. Gene Set Enrichment Analysis of ALL cells overexpressing *GATA3* vs. control

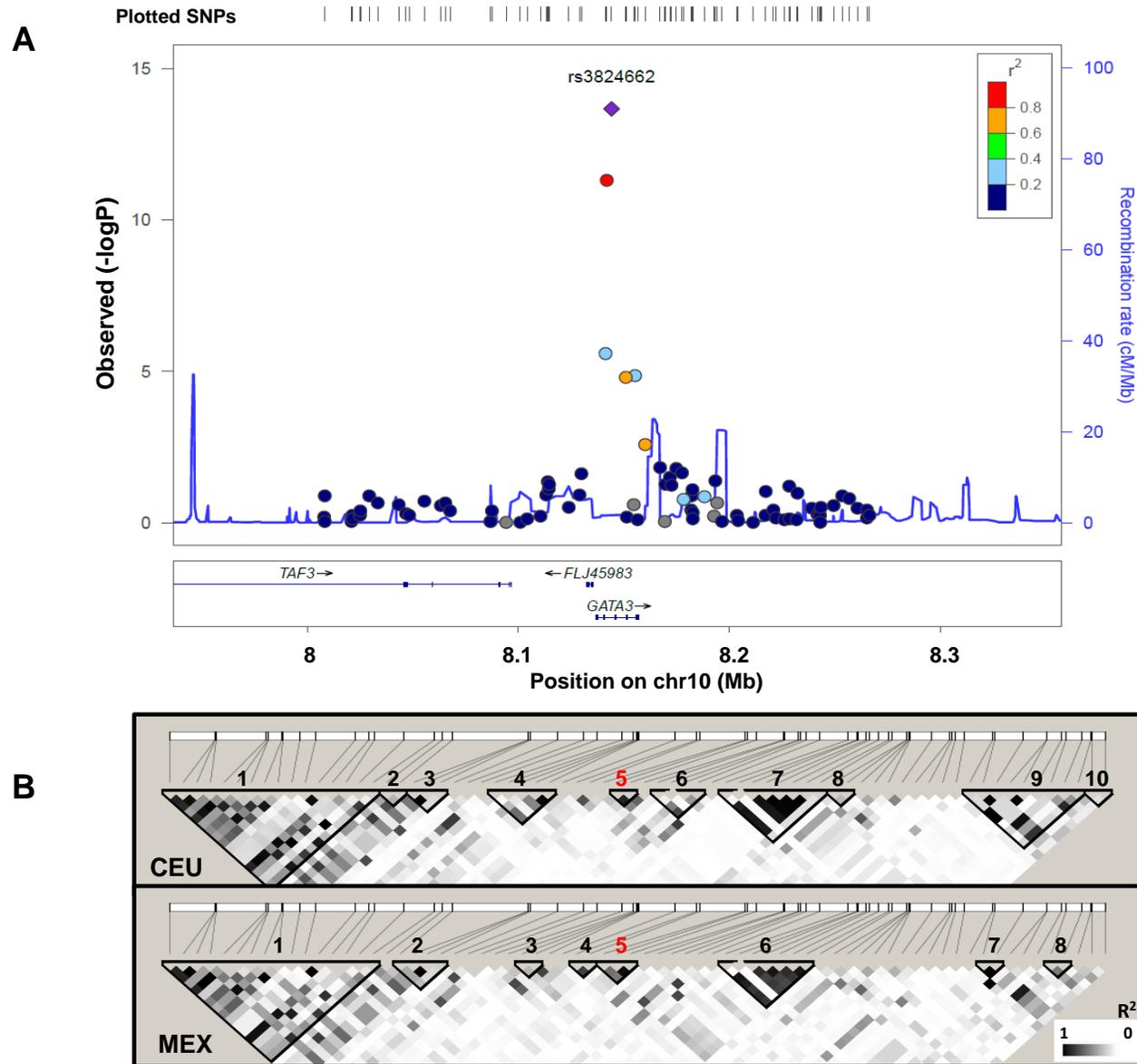
<i>SLC39A10</i>	-0.2173	-0.1092	No	<i>SH3BP5</i>	0.0360	0.3230	No
<i>SLC44A1</i>	0.0134	0.2154	No	<i>SLC2A5</i>	0.0072	0.2344	No
<i>SOCS2</i>	-0.0105	0.0811	No	<i>SLC39A10</i>	-0.0762	0.0191	No
<i>SPARC</i>	-0.1224	-0.1135	No	<i>SLC44A1</i>	-0.0042	0.1897	No
<i>SPON1</i>	0.0445	0.3295	No	<i>SOCS2</i>	-0.1358	-0.0379	No
<i>STAB1</i>	0.0545	0.3455	No	<i>SPARC</i>	-0.2845	-0.0232	No
<i>STON2</i>	-0.1032	-0.1246	No	<i>SPON1</i>	-0.1067	-0.0241	No
<i>SUSD3</i>	0.1036	0.4098	No	<i>STAB1</i>	0.0628	0.3966	No
<i>SV2A</i>	0.0249	0.2734	No	<i>SUSD3</i>	-0.1684	-0.0291	No
<i>TBXAS1</i>	0.0678	0.3752	No	<i>SV2A</i>	-0.0668	0.0189	No
<i>THBS1</i>	0.0369	0.3266	No	<i>THBS1</i>	0.0370	0.3249	No
<i>TMEM154</i>	-0.0787	-0.1026	No	<i>TMEM154</i>	0.1572	0.5112	No
<i>TTYH2</i>	-0.1903	-0.1156	No	<i>TTYH2</i>	-0.1322	-0.0389	No
<i>UACA</i>	-0.1186	-0.1267	No	<i>UACA</i>	0.1635	0.5115	No
<i>UPP1</i>	0.0371	0.3251	No	<i>UPP1</i>	-0.0896	0.0015	No

GSEA analysis comparing ALL cells ectopically overexpressing *GATA3* vs. cells transduced with empty vector, with PAM-based Ph-like signature (upregulated) as the *a priori* gene set.

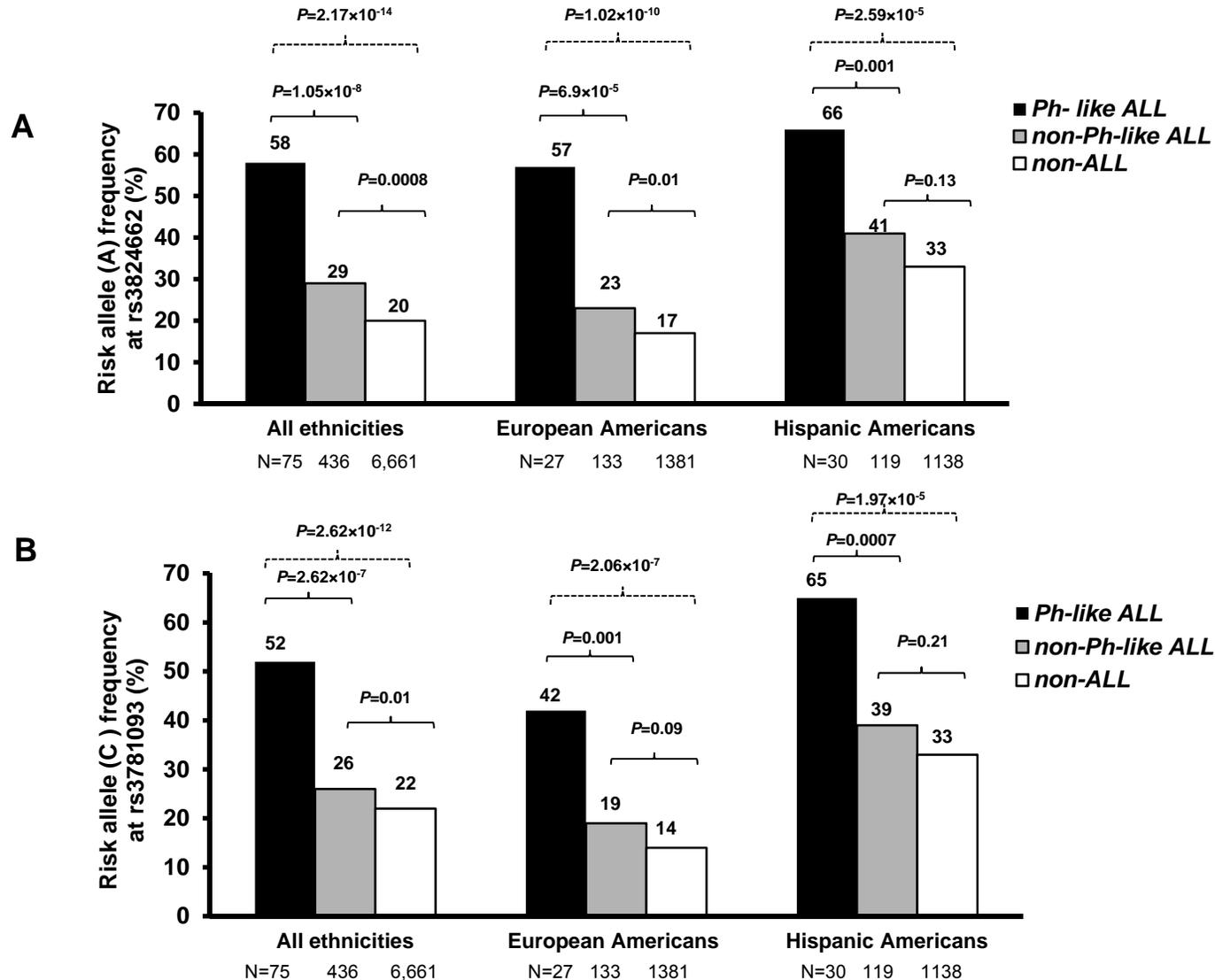
*Denotes genes that were upregulated in patients carrying the A allele at rs3824662 compared with those not carrying the A allele in both COG AALL0232 and COG P9906.



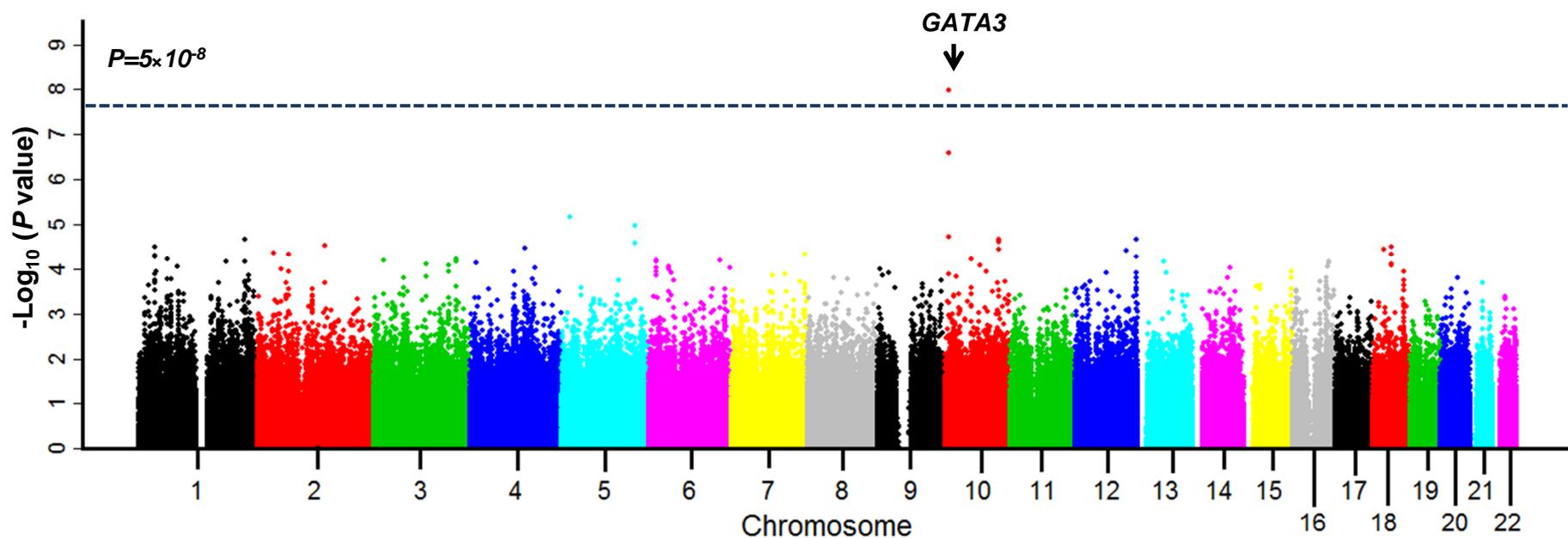
Supplementary Figure 1. Flow chart of SNP quality control/filtering in the discovery GWAS. SNPs were filtered on the basis of allele frequency and call rate, as detailed in “*Genotyping and Quality Control*”.



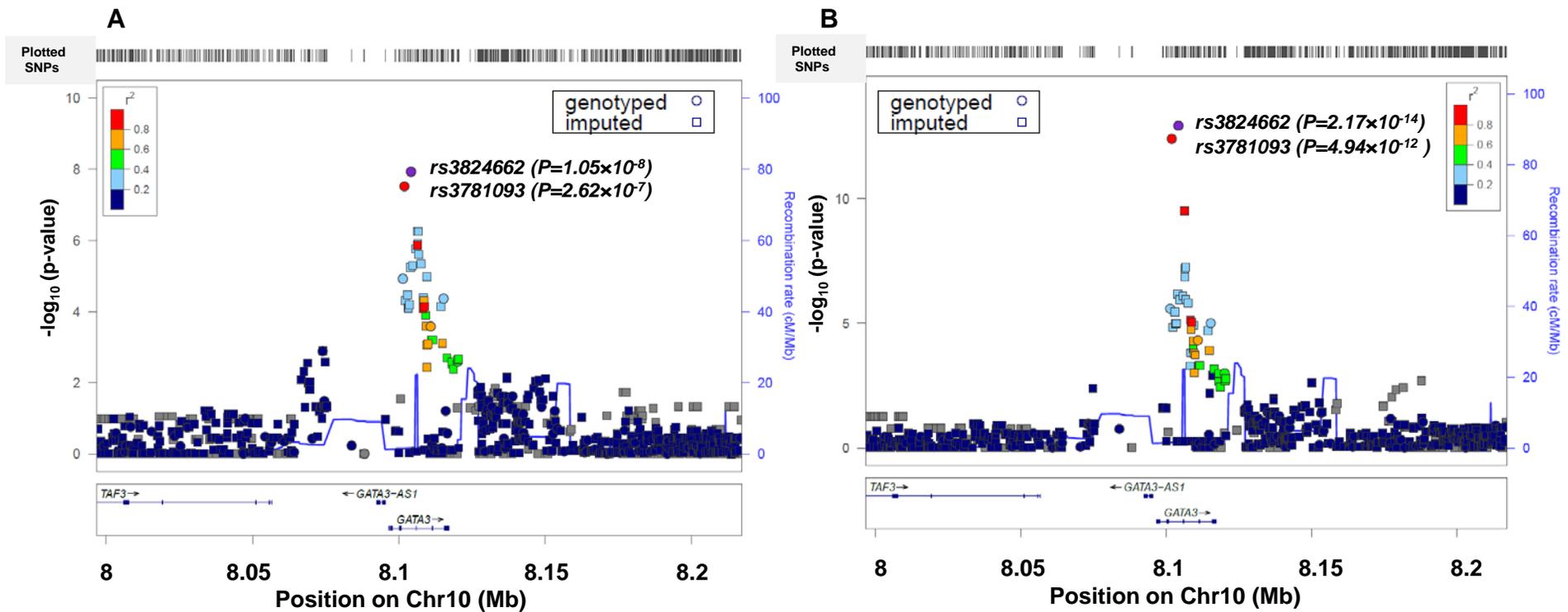
Supplementary Figure 2. Association results and linkage disequilibrium (LD) at the 10p14 locus. Panel A illustrates association signals in the discovery GWAS (Ph-like ALL vs. non-ALL controls). The negative logarithm of the P value (left axis) and recombination rate (right axis) are plotted for a 250 Kb window at the 10p14 locus, using LocusZoom (*Bioinformatics* 26: 2336). Color indicates LD (r^2) with rs3824662 in the HapMap CEU samples and chromosome position is based on hg18. In **panel B**, LD at this locus is depicted based on r^2 in HapMap CEU and MEX cell lines, and the plots were constructed using the HaploView software. *GATA3* SNPs (rs3824662 and rs3781093) were in high LD with each other (CEU: $r^2=0.94$ and $D'=1$; MEX: $r^2=0.90$ and $D'=0.95$) within LD block # 5 in both CEU and MEX populations.



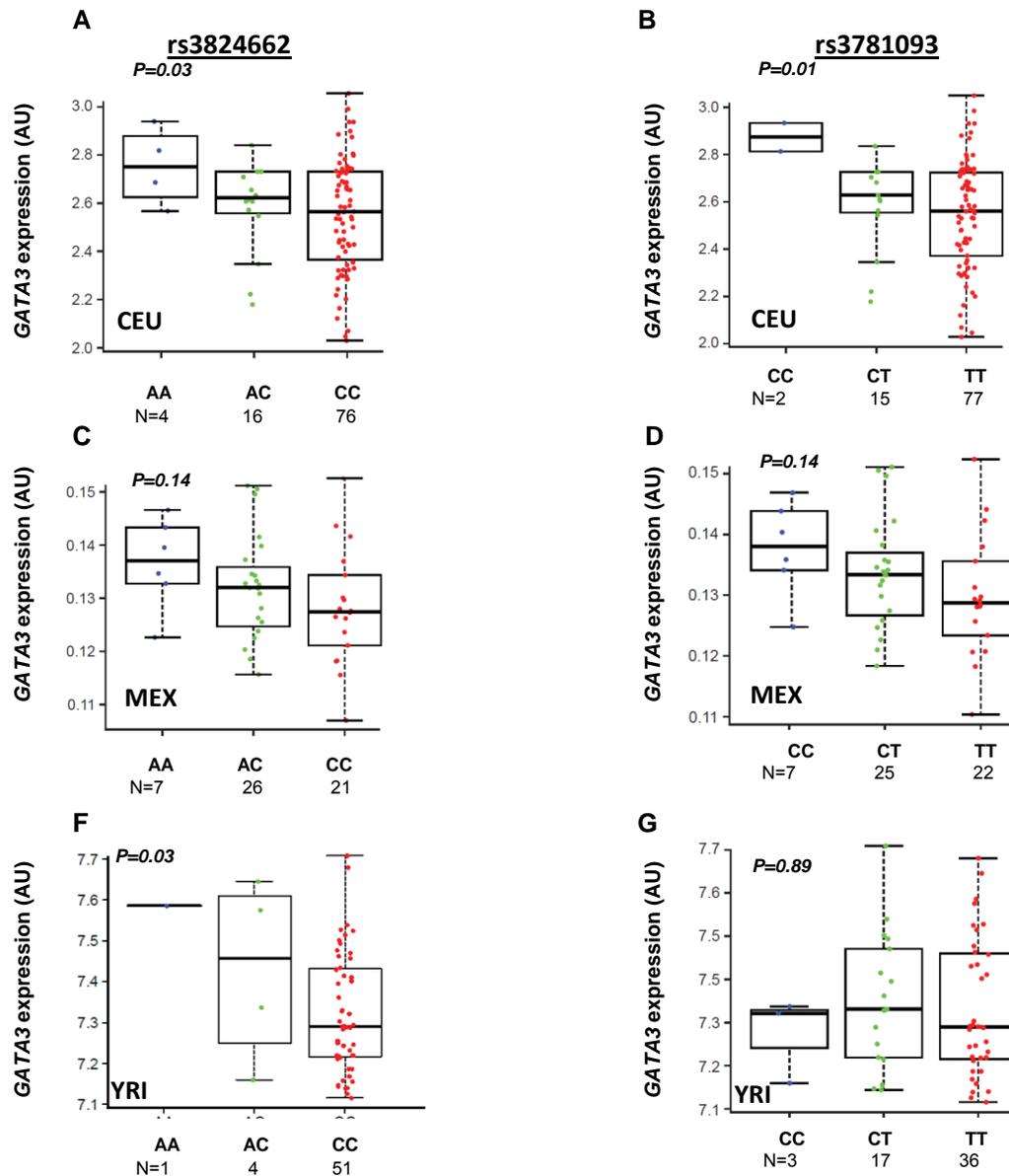
Supplementary Figure 3. Relationship between *GATA3* SNPs (rs3824662 and rs3781093) and Ph-like ALL by ethnicity in the COG AALL0232 cohort. The A allele at rs3824662 (**Panel A**) was over-represented in Ph-like ALL relative to non-Ph-like ALL and non-ALL controls. This association was true within the European Americans (>95% European genetic ancestry) or Hispanic Americans (>10% Native American genetic ancestry and Native American ancestry > African genetic ancestry). Similar association was confirmed for the risk allele (C) at the *GATA3* SNP rs3781093 (**Panel B**). Genetic ancestry was determined by using STRUCTURE (version 2.2.3) with HapMap CEU, YRI, CHB/JPT, and indigenous Native Americans as reference populations.



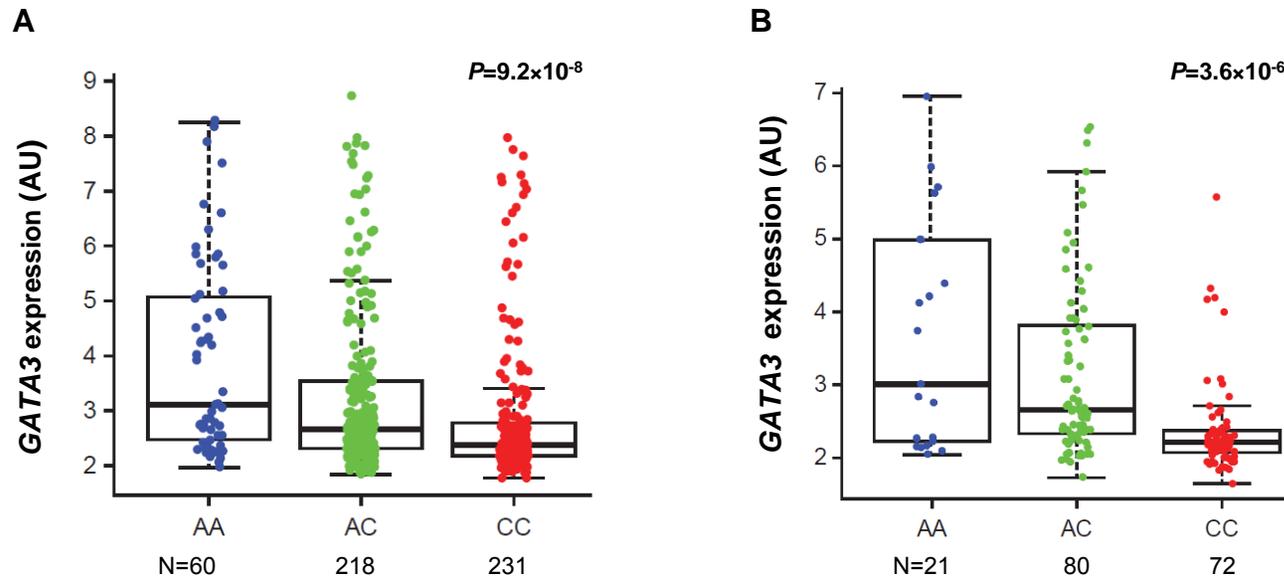
Supplementary Figure 4. GWAS of Ph-like ALL by comparing allele frequency between Ph-like ALL vs. non-Ph-like ALL. The association between genotype and Ph-like was evaluated using logistic regression model for 761,049 SNPs in 75 ALL cases with Ph-like gene expression profile and 436 ALL cases without this expression signature. P -values ($-\log_{10} P$, y axis) were plotted against respective chromosomal position of each SNP (x axis). Points above the blue horizontal line indicate SNPs achieving the genome-wide significant threshold ($P < 5 \times 10^{-8}$). Gene symbol was indicated for the *GATA3* locus at 10p14.



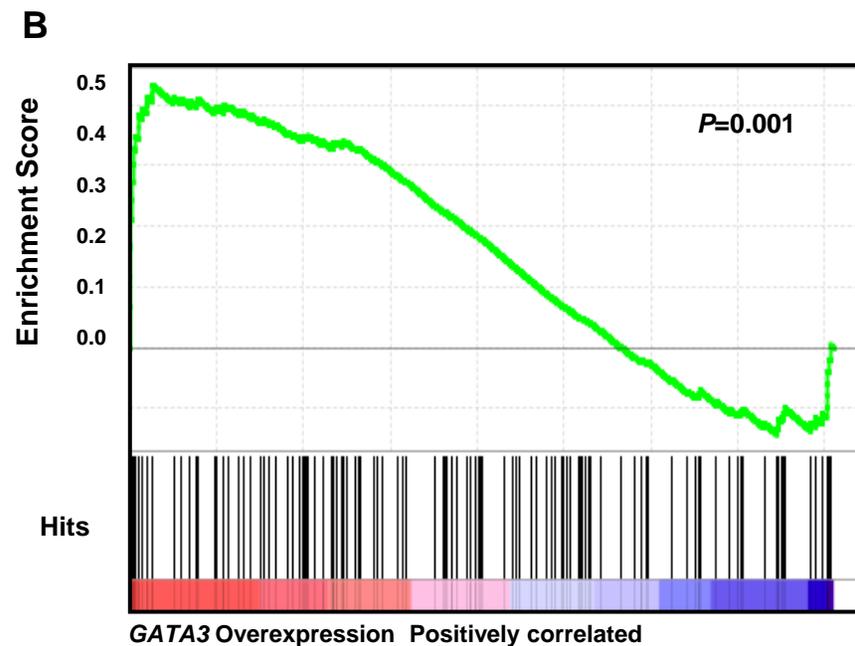
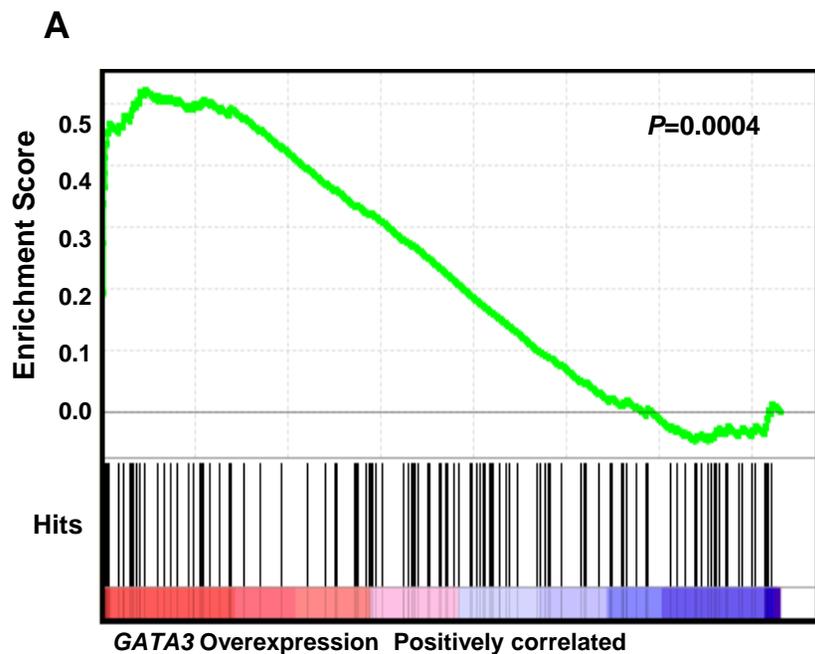
Supplementary Figure 5. Association with Ph-like ALL for imputed SNPs on chr10p14. Genotype was imputed for a 10 Mb region on 10p14 (chr10:60,523-10,060,447, hg19), using MaCH-Admix 2.0.185 with 1,000 Genome data set as references. Association was tested by comparing genotype frequency between Ph-like ALL (N=75) vs. non-Ph-like ALL (N=436, Panel **A**) and between Ph-like ALL (N=75) vs. non-ALL controls (N=6,661, Panel **B**) at 37,493 imputed or directly genotyped SNPs. Shown here are the association results for a 220 Kb window centered around rs3824662 with 727 SNPs spanning chr10:7,996,666-8,217,164 (hg19), and the plots are constructed using LocusZoom (*Bioinformatics* 26: 2336). Color indicates LD (r^2) with rs3824662 in the HapMap CEU population.



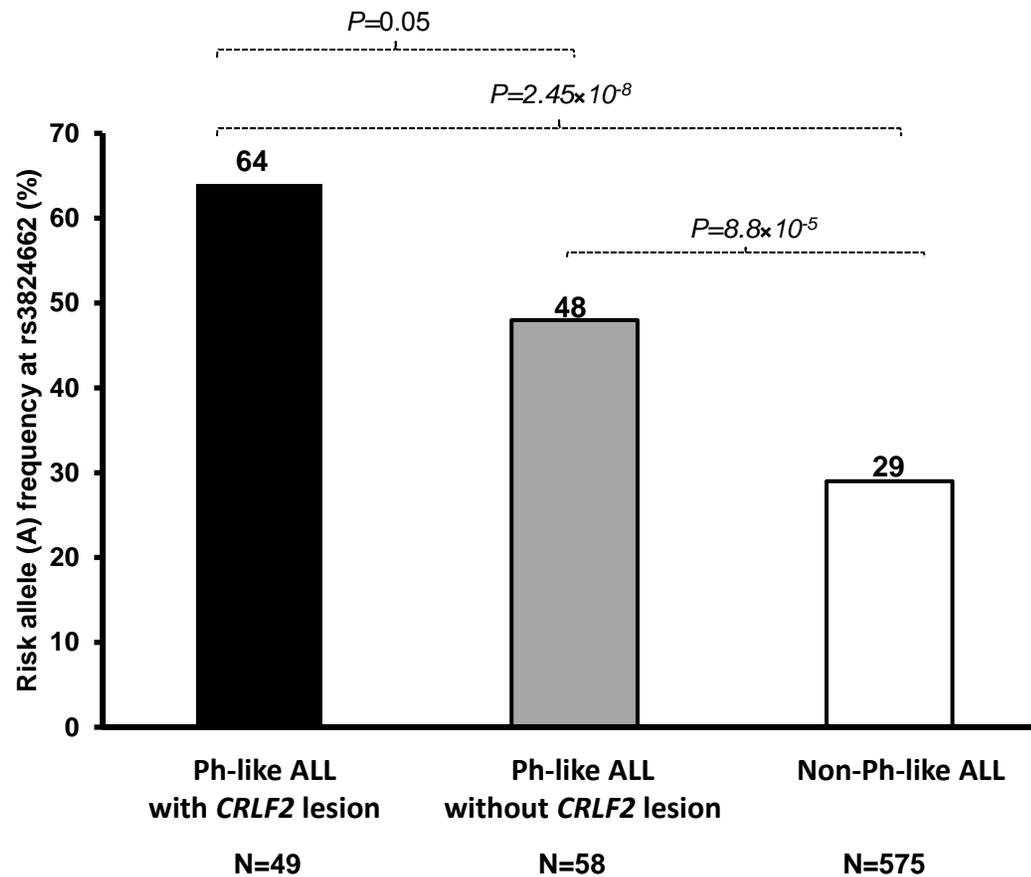
Supplementary Figure 6. eQTL analyses of rs3824662 and rs3781093 in diverse HapMap populations. *GATA3* SNP rs3824662 risk allele (A) and the rs3781093 (C) were associated with higher *GATA3* mRNA in 96 unrelated lymphoblastoid cell lines from the HapMap CEU population, using the publicly available gene expression data set GSE5859 (A and B, respectively). Similar trend was observed in 54 unrelated HapMap MEX cell lines for which *GATA3* expression was evaluated by real time-PCR (Panels C and D for rs3824662 and rs3781093, respectively). Genotype-expression association in 56 unrelated YRI samples is represented in panel E (rs3824662) and panel F (rs3781093), using gene expression data set GSE7851. Genotype-expression association was evaluated using a linear regression model adjusting for ancestry as appropriate. AU, arbitrary units. Boxes include data between the twenty-fifth and the seventy-fifth percentiles.



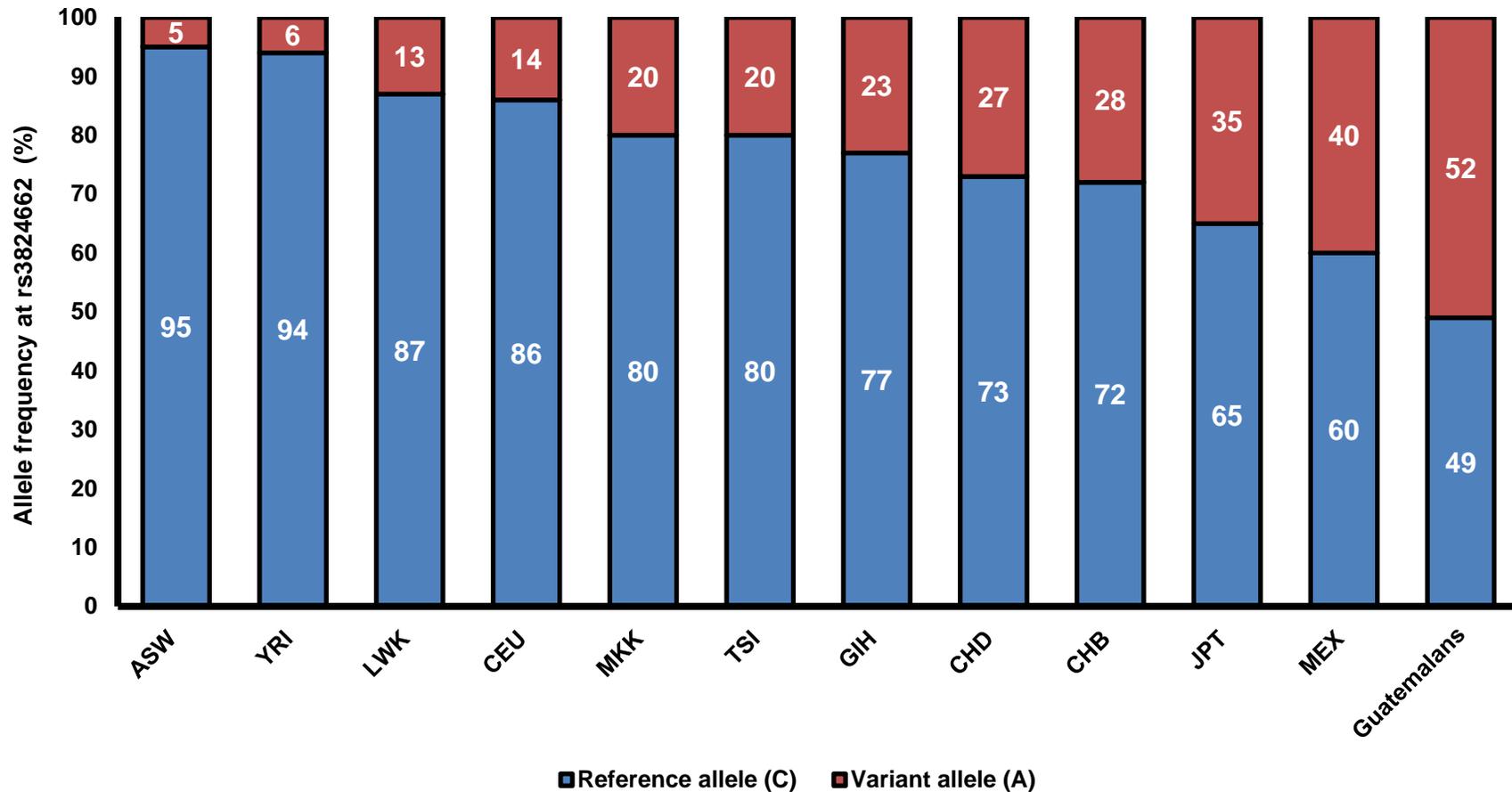
Supplementary Figure 7. rs3824662 genotype was associated with GATA3 expression in ALL blasts. GATA3 SNP rs3824662 risk allele (the A allele) was associated with higher GATA3 mRNA expression in diagnostic ALL blasts from 511 children in the COG AALL0232 cohort (**A**) and 173 children in the COG P9906 cohort (**B**). Genotype-expression association was evaluated using a linear regression model, adjusting genetic ancestry as appropriate. AU, arbitrary unit. Boxes include data between the twenty-fifth and the seventy-fifth percentiles.



Supplementary Figure 8. Enrichment of Ph-like signature genes in ALL cell lines ectopically overexpressing *GATA3* compared with those transduced with control vectors, using the Gene Set Enrichment Analysis (GSEA). GSEA tested the upregulation of Ph-like ALL genes (i.e., PAM-based Ph-like signature) in ALL cell line UOCB1 (A) and Nalm6 (B) after ectopic overexpression of *GATA3* gene. *P* value was based permutations.

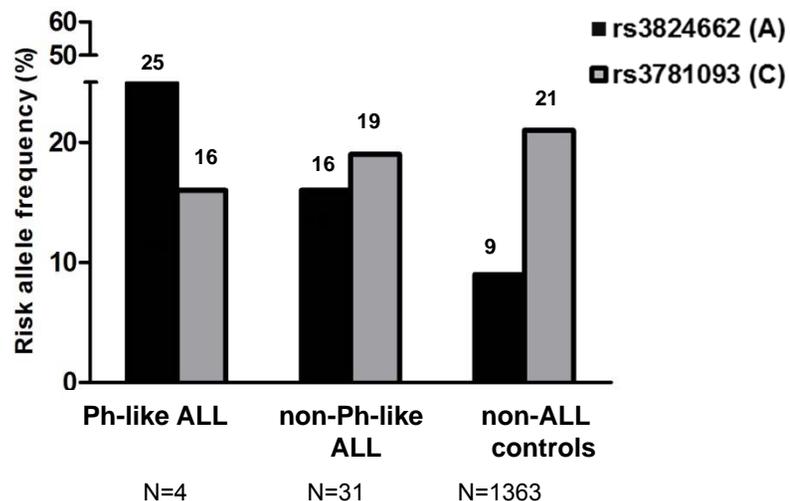


Supplementary Figure 9. Risk allele frequency at *GATA3* SNP rs3824662 in Ph-like patients (according to *CRLF2* status) and non-Ph-like patients. Combined cohort includes COG AALL0232 and COG P9906 (N=682), and *P* values were estimated by logistic regression test after adjusting for genetic ancestry.

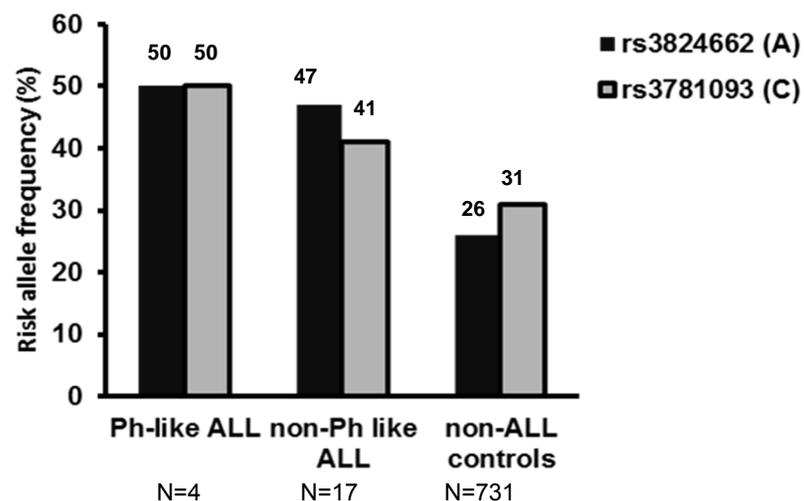


Supplementary Figure 10. GATA3 SNPs rs3824662 allele frequency in worldwide populations. The frequency of the Ph-like ALL-related allele (A) is shown in an ascending order for HapMap populations and Native Americans in Guatemala. **Population descriptors:** **ASW:** African ancestry in Southwest USA, **CEU:** Utah residents with Northern and Western European ancestry from the CEPH collection, **CHB:** Han Chinese in Beijing, China, **CHD:** Chinese in Metropolitan Denver, Colorado, **GIH:** Gujarati Indians in Houston, Texas, **JPT:** Japanese in Tokyo, Japan, **LWK:** Luhya in Webuye, Kenya, **MEX:** Mexican ancestry in Los Angeles, California, **MKK:** Maasai in Kinyawa, Kenya, **TSI:** Tuscan in Italy, **YRI:** Yoruban in Ibadan, Nigeria. **Guatemalan:** 65 unrelated Guatemalan individuals with high Native American ancestry.

A

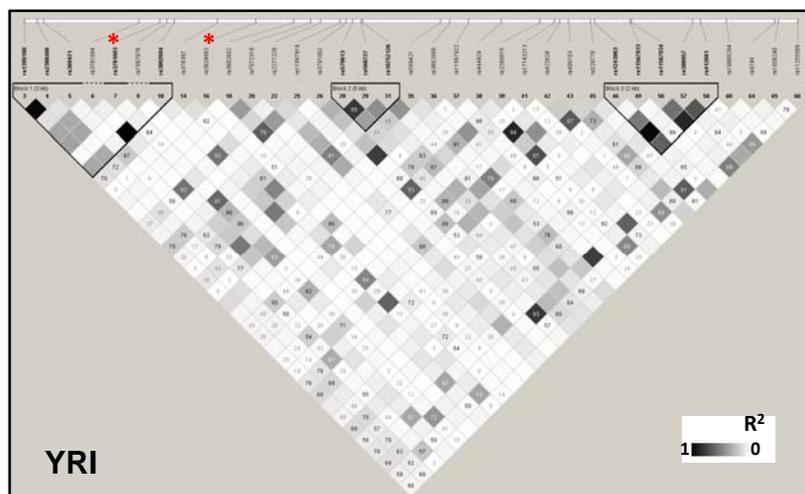


B



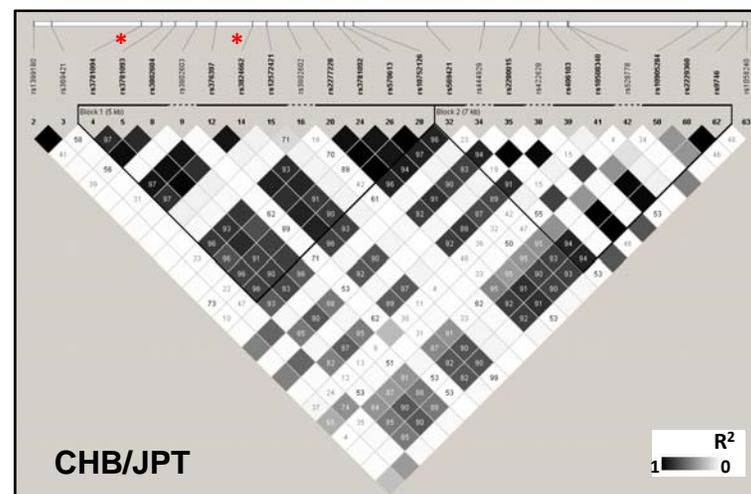
C

Chr10; 8,137,066.....Chr10:8,157,058



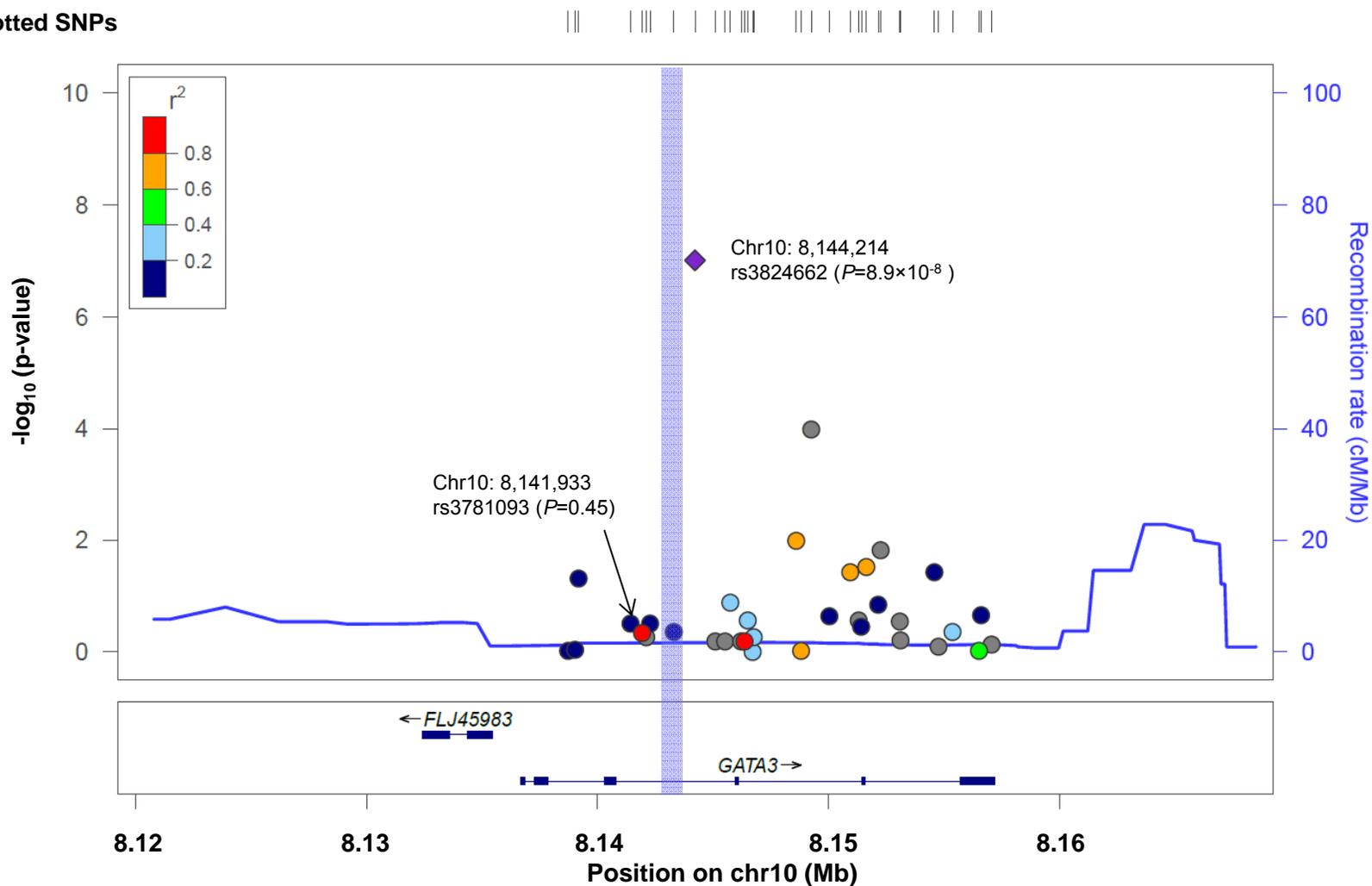
D

Chr10; 8,138,647.....Chr10:8,157,058

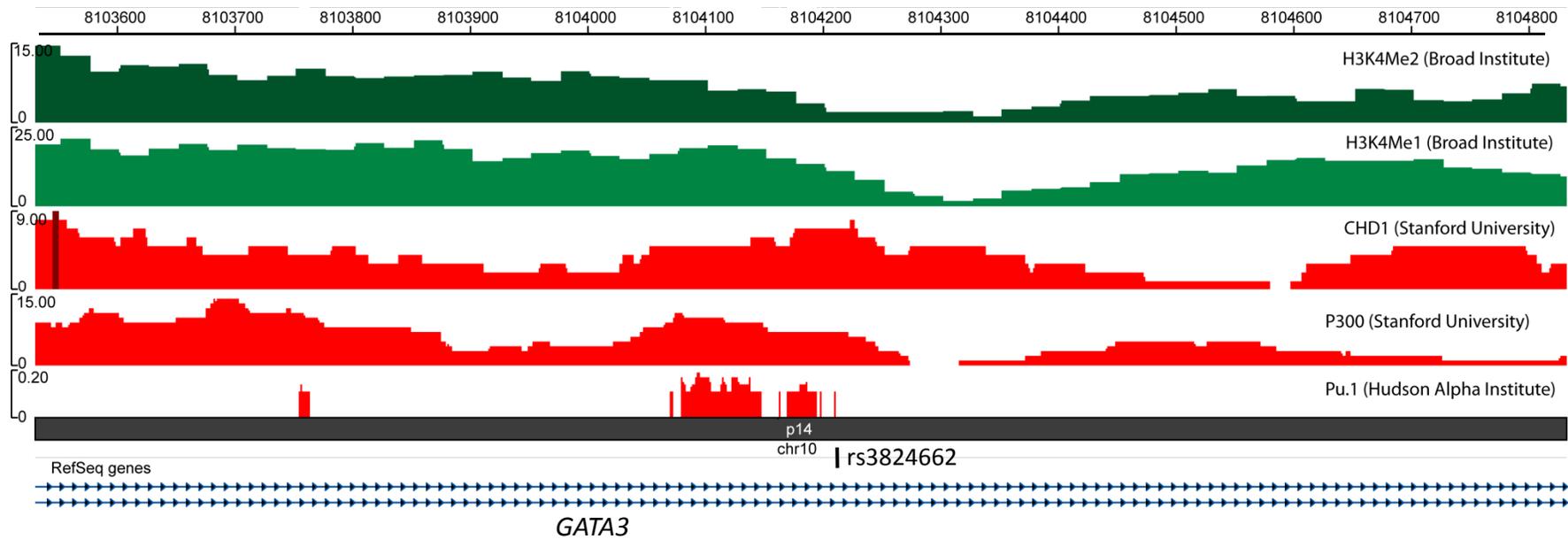


Supplementary Figure 11. Association of *GATA3* SNPs with Ph-like ALL in African and Asian populations. Panels A and B illustrate the allele frequency at rs3824662 and rs3781093 for Ph-like ALL, non-Ph-like ALL, and non-ALL controls with >70% African ancestry (A) and those with >90% Asian ancestry (B). In panels C and D, LD is depicted based on r^2 in unrelated HapMap YRI (C) and CHB/JPT (D) samples, and the plots were constructed using the HaploView software. rs3824662 and rs3781093 are in high LD with each other in CHB/JPT population ($r^2=0.97$ and $D'=1.0$) but not in YRI population ($r^2=0.006$ and $D'=0.16$). Chromosome position is based on hg18.

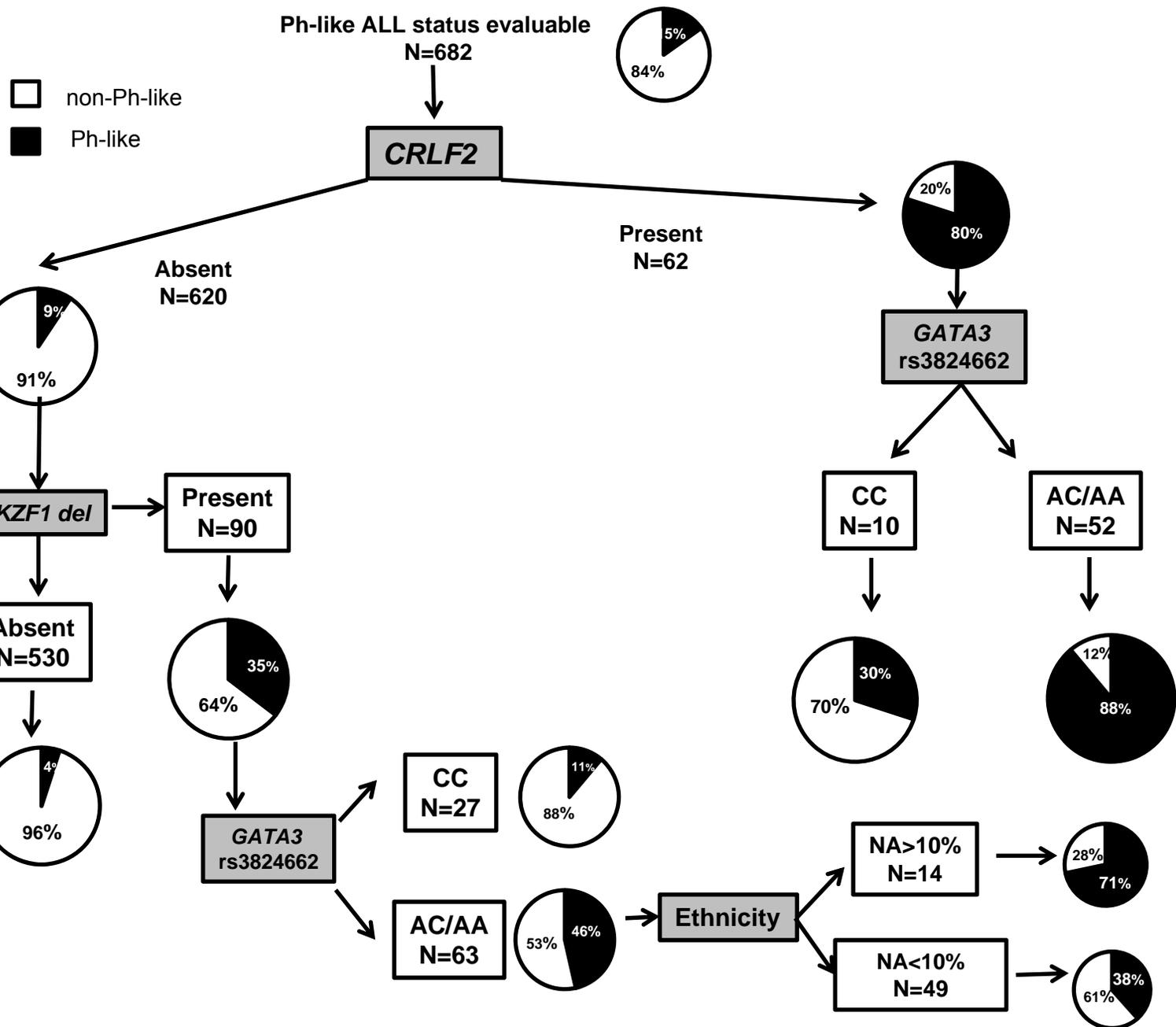
Plotted SNPs



Supplementary Figure 12. Association of *GATA3* variants with local DNase hypersensitivity. DNase hypersensitivity was obtained for 70 HapMap YRI samples from a previously published data set (Nature 482:390), and genotype at 35 SNPs was retrieved from the 1,000 Genomes data set. rs3824662 (chr10; 8,144,214, hg18) showed the strongest association with the local DNase sensitivity window (chr10 8,144,000-8,144,100; shaded region), as determined by a linear regression test.

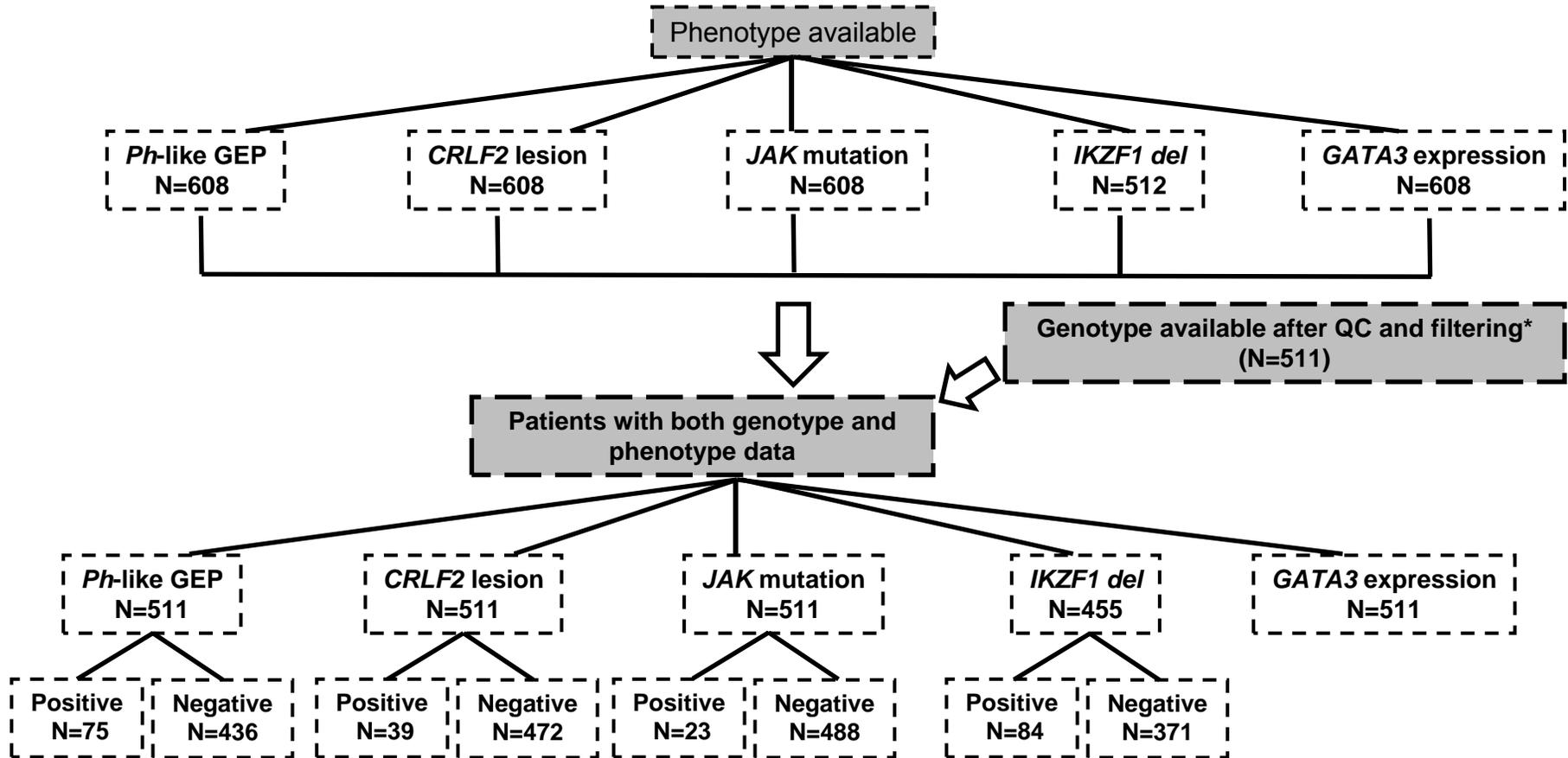


Supplementary Figure 13. Enhancer signal within the genomic region encompassing rs3824662 in the ENCODE data set. ENCODE data was queried for possible regulatory activities around rs3824662, focusing on histone marks and transcription factor binding. In GM12878 (ENCODE Tier 1), H3K4Me1 and H3K4Me2 signals indicate enhancer activity at rs3824662 and the reduction of ChIP-seq reads correspond to possible occupancy by transcription factor (P300 and PU.1). Consistent enrichment of CHD1-binding signal (chromatin remodeling activity) was also noted, plausibly in close proximity to the methylated H3K4. Graph is constructed using WashU Epigenome Browser (Nat Methods 10:375).

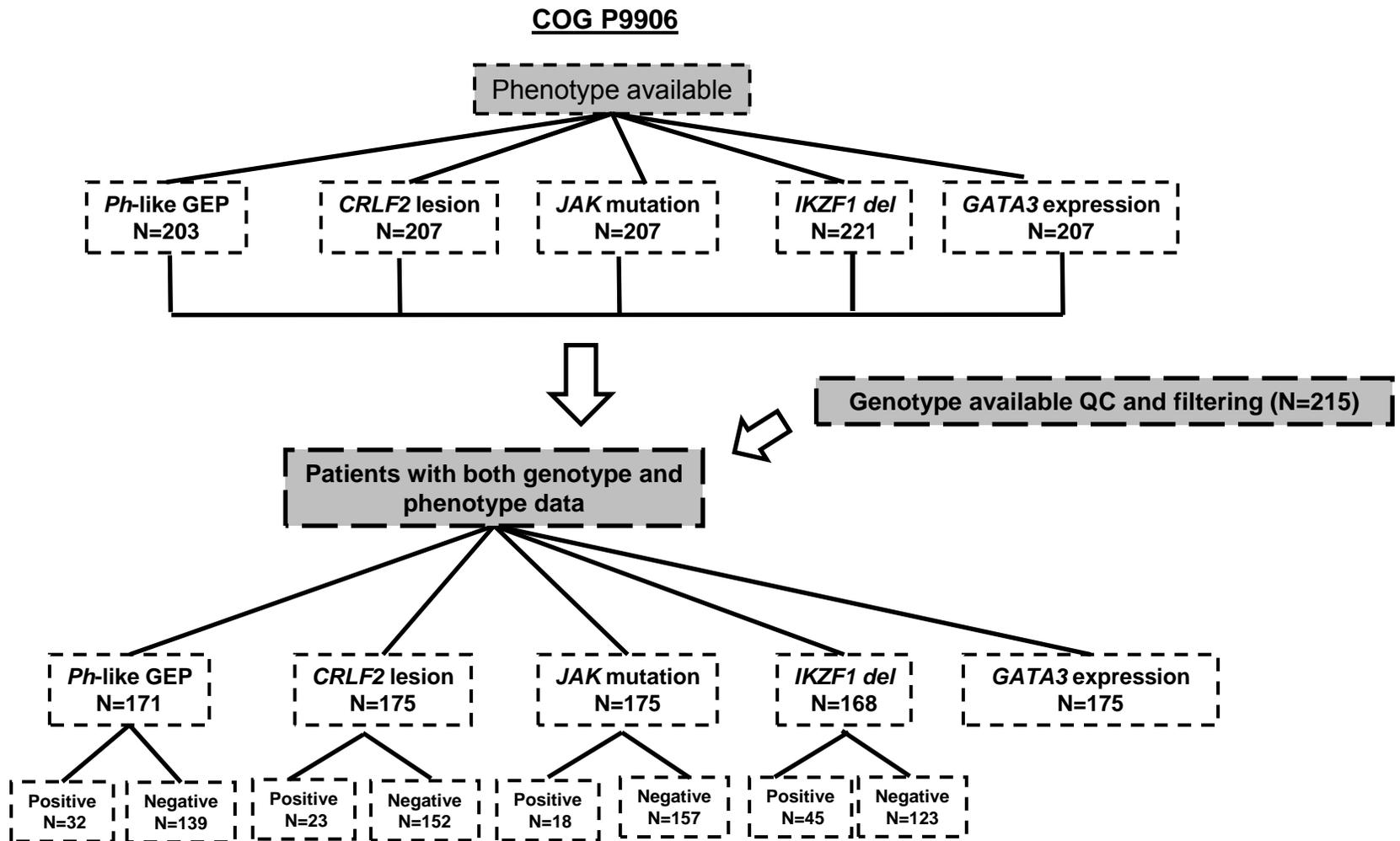


Supplementary Figure 14. Classification and regression tree (CART) analysis of predictors for Ph-like ALL. 682 patients with Ph-like ALL status evaluable from the COG AALL0232 and COG P9906 protocols were included. CART analysis was performed using the rpart function in R software, with *CRLF2*, *JAK*, *IKZF1* lesion, *GATA3* SNP, and genetic ancestry included in the model building process. NA, Native American genetic ancestry.

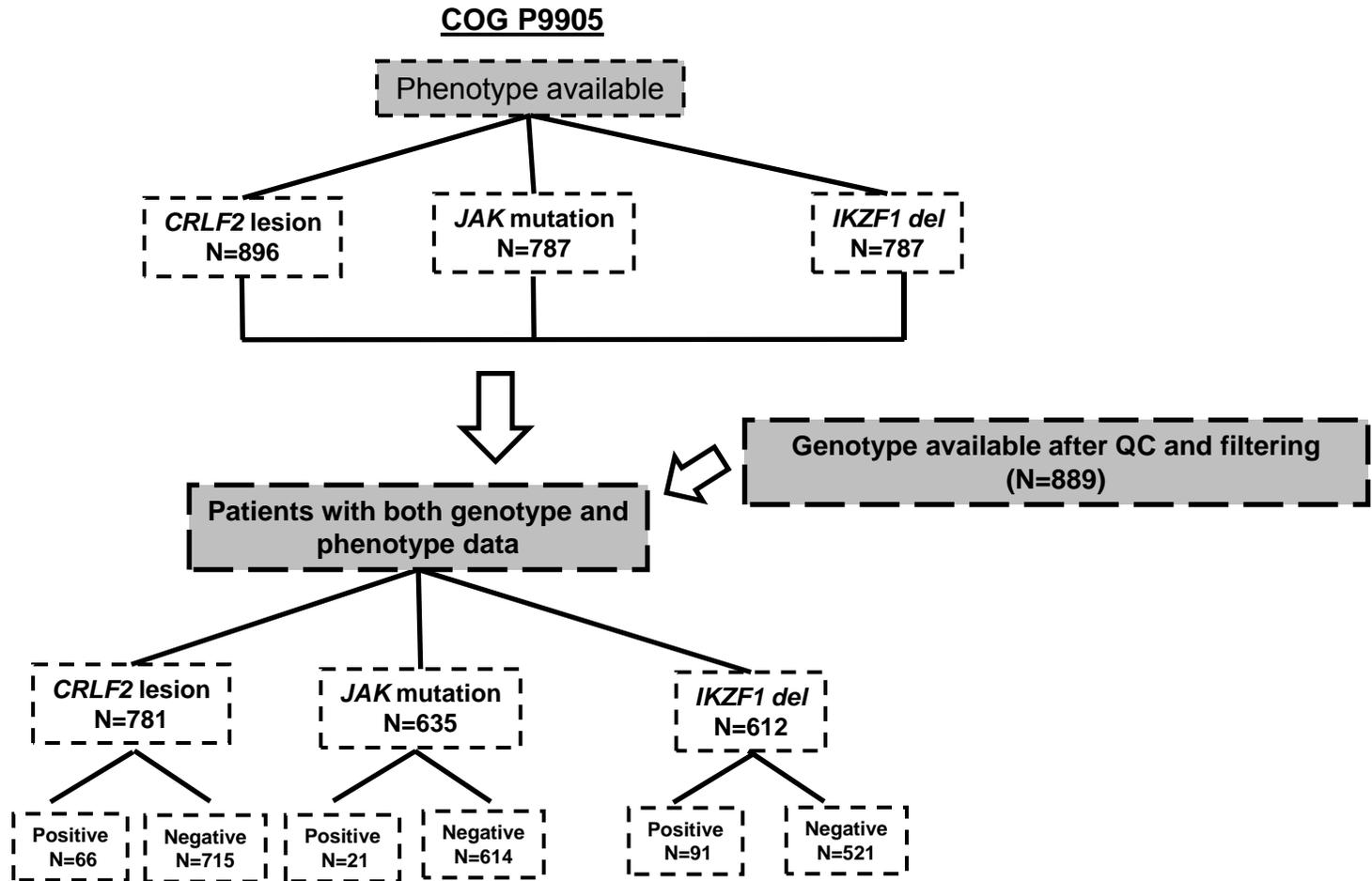
COG AALL0232



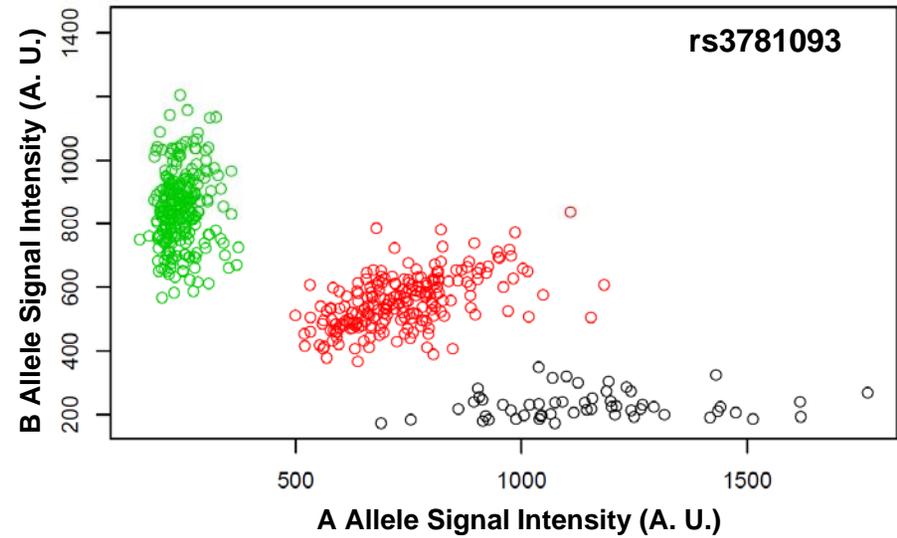
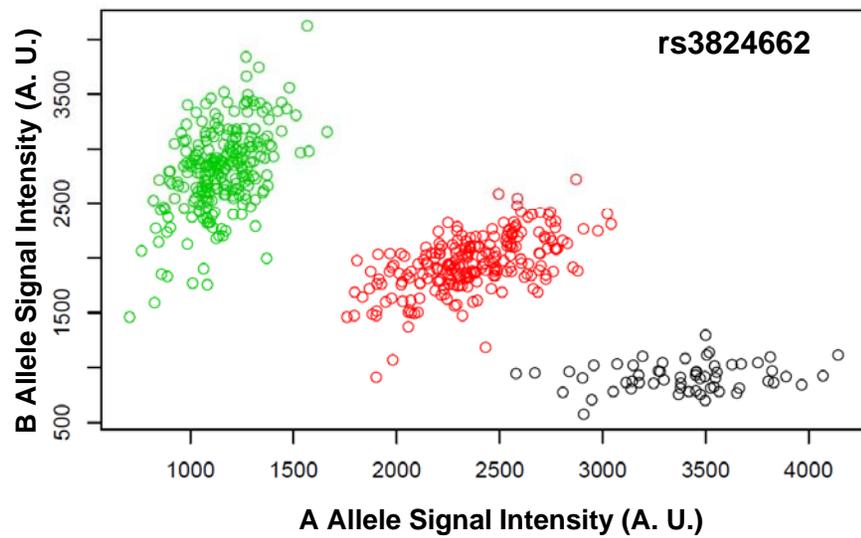
Supplementary Figure 15. Detailed description of patients tested for Ph-like, somatic *IKZF1* deletion, *CRLF2* lesion, *JAK* mutation in COG AALL0232 discovery GWAS group. *Of 550 sample genotyped, 12 samples were removed due to poor call rate or mismatch, and 27 Ph+ ALL were included in a separate analysis.



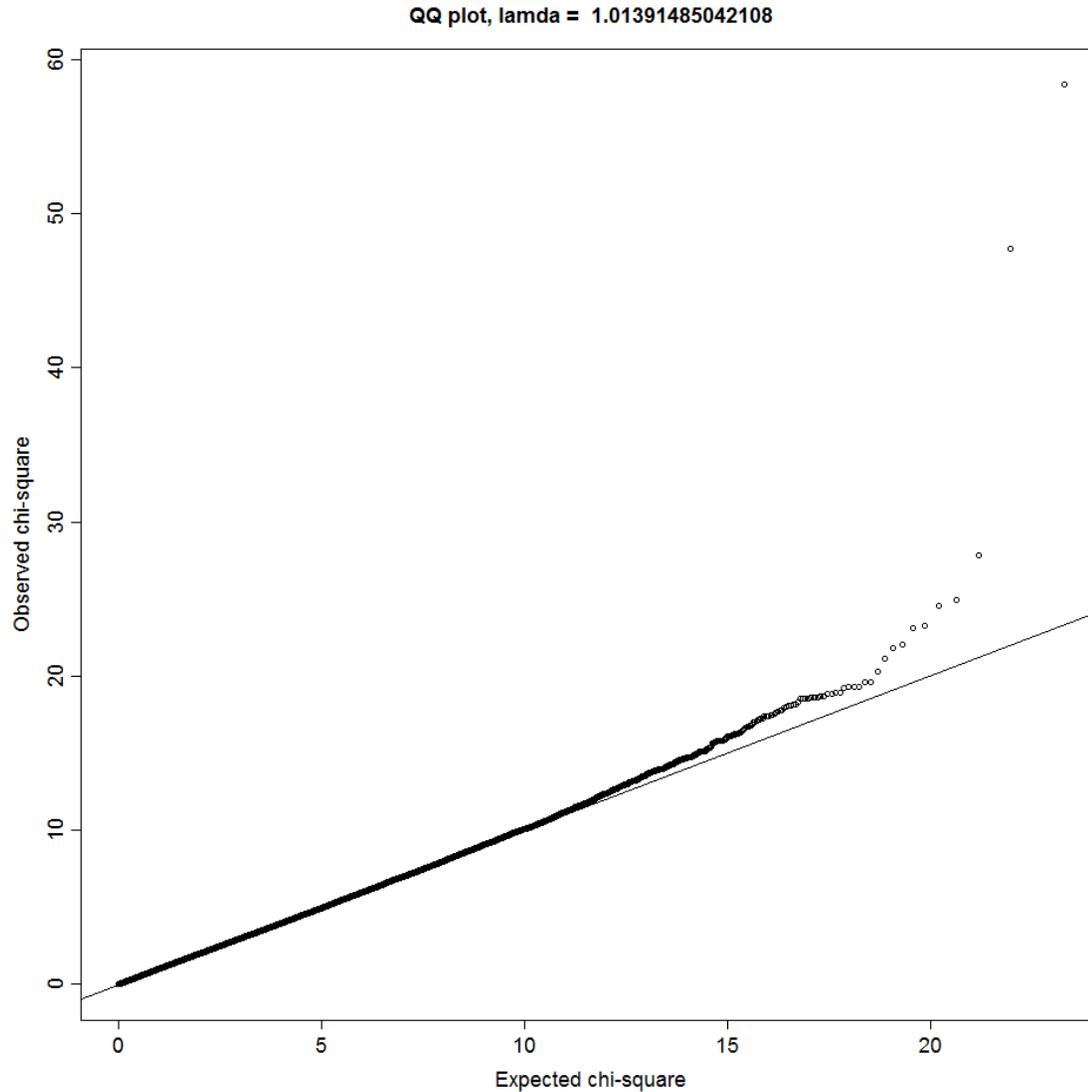
Supplementary Figure 16. Detailed description of patients tested for Ph-like, somatic *IKZF1* deletion, *CRLF2* lesion, *JAK* mutation in COG P9906 replication group.



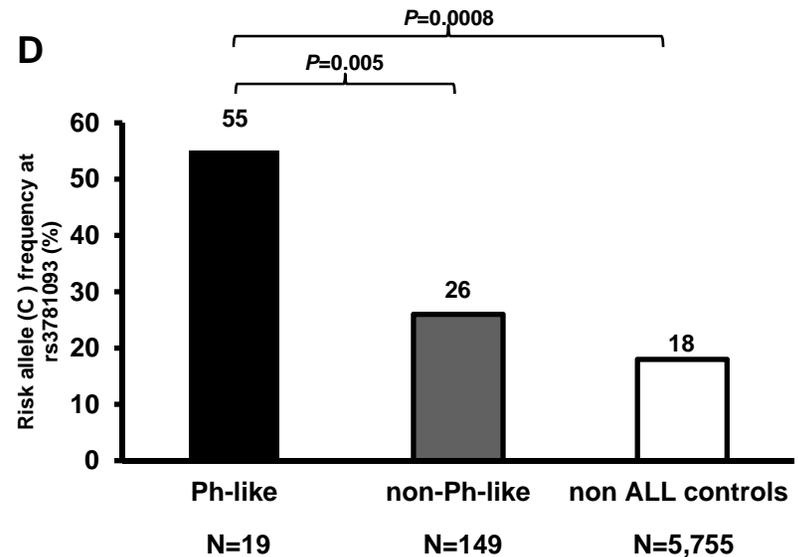
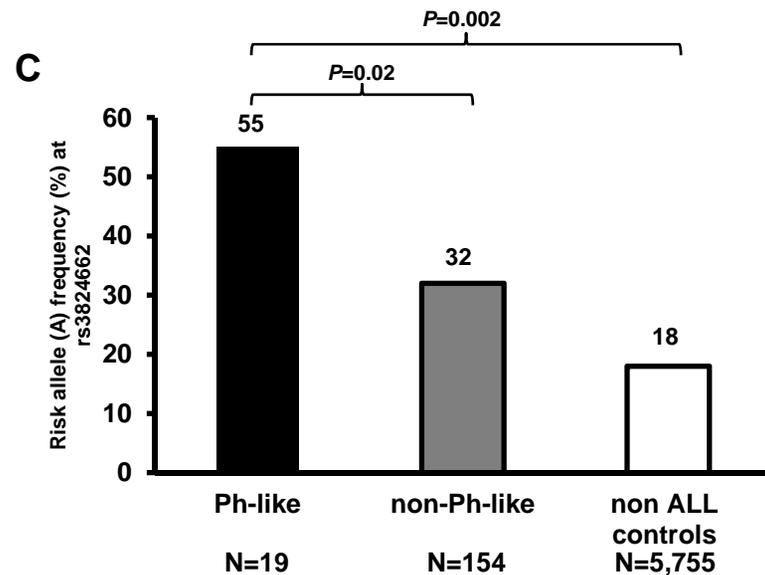
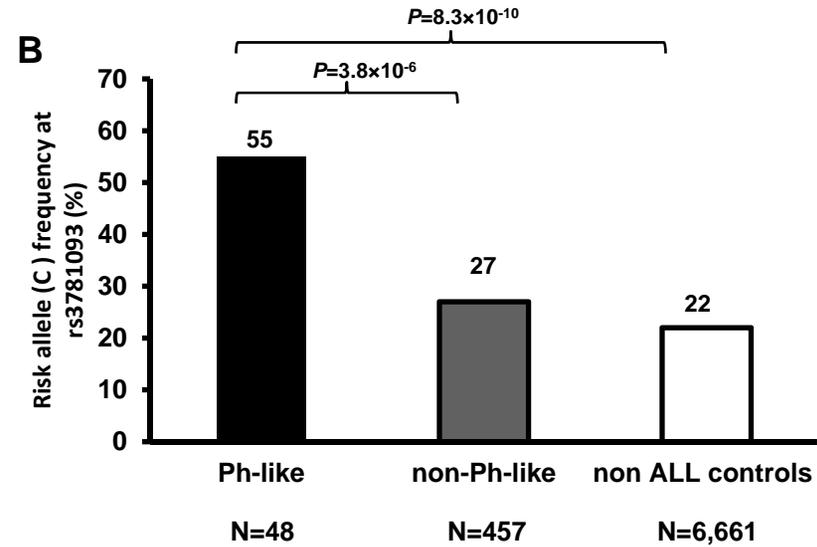
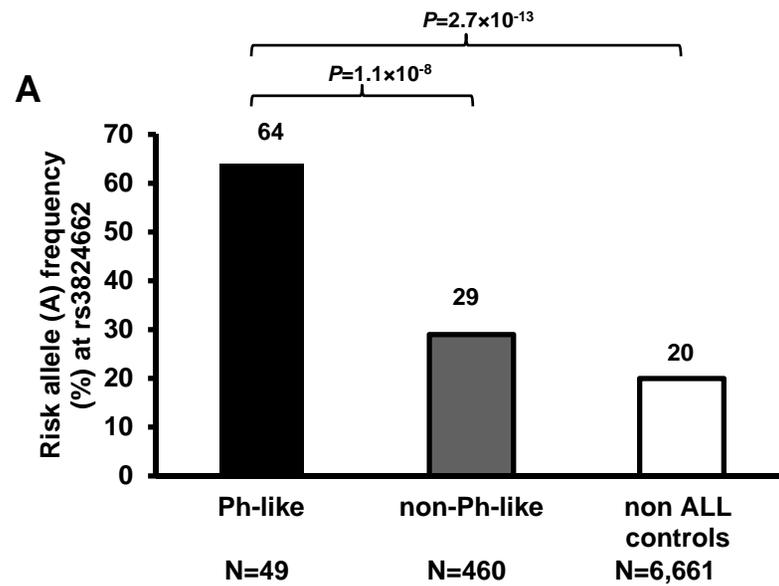
Supplementary Figure 17. Detailed description of patients tested for Ph-like, somatic *IKZF1* deletion, *CRLF2* lesion, *JAK* mutation in COG P9905 group.



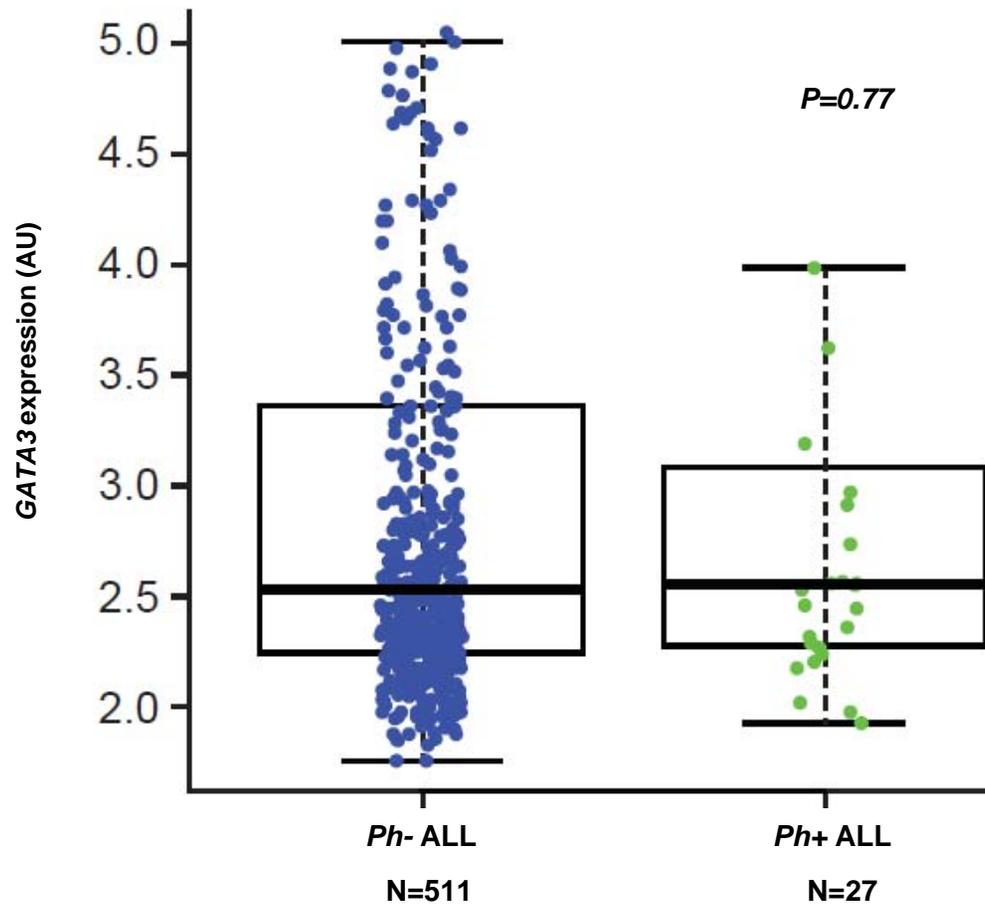
Supplementary Figure 18. Cluster plots of rs3824662 and rs3781093. Signal intensity for A and B alleles at both SNPs was based on theta value of 538 samples in the COG AALL0232 cohort, using Affymetrix Genotyping Console. At both SNPs, samples with AA, AB, and BB genotype clearly clustered into distinct groups, indicating high-quality genotype calls. AU, arbitrary units.



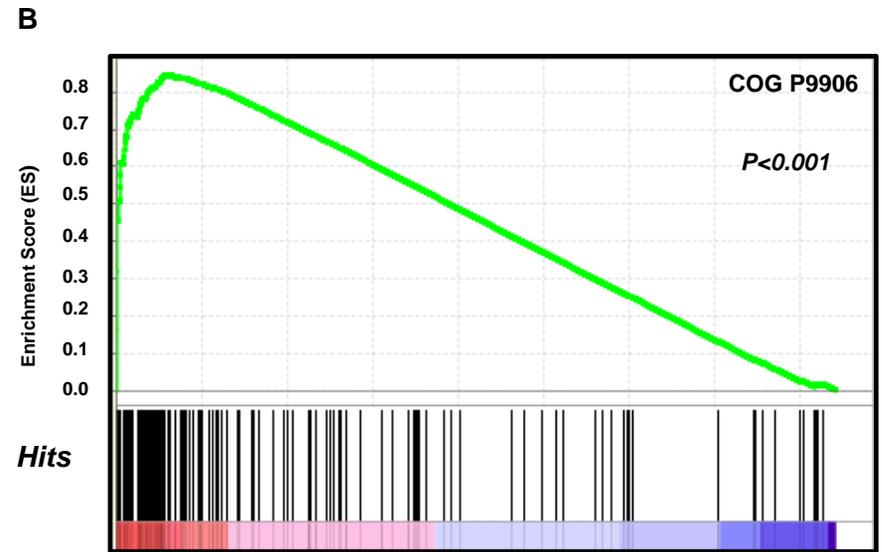
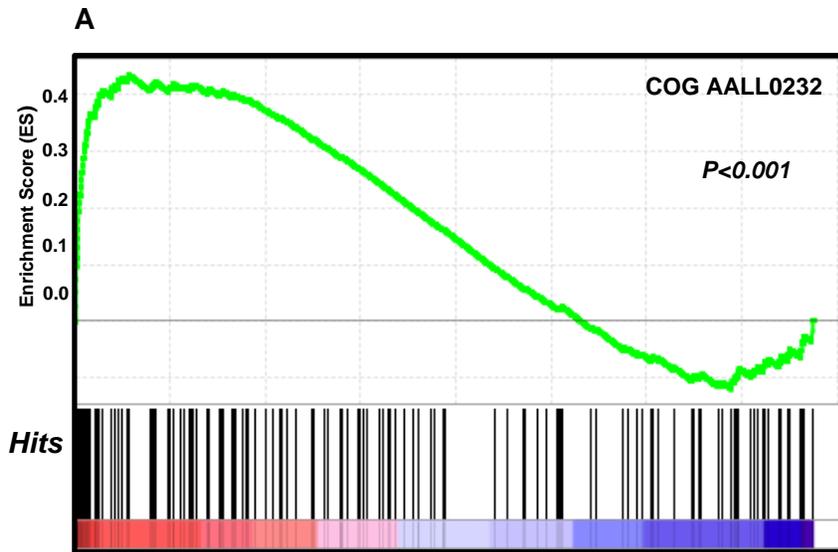
Supplementary Figure 19. Quantile-quantile (Q-Q) plot of logistic regression test for GWAS. The negative logarithm of the observed (y axis) and the expected (x axis) P value is plotted for each SNP (dot), and the black line indicates the null hypothesis of no true association. Deviation from the expected P value distribution is evident only in the tail area ($\lambda=1.01$), suggesting that population stratification was adequately controlled by adjusting for genetic ancestry.



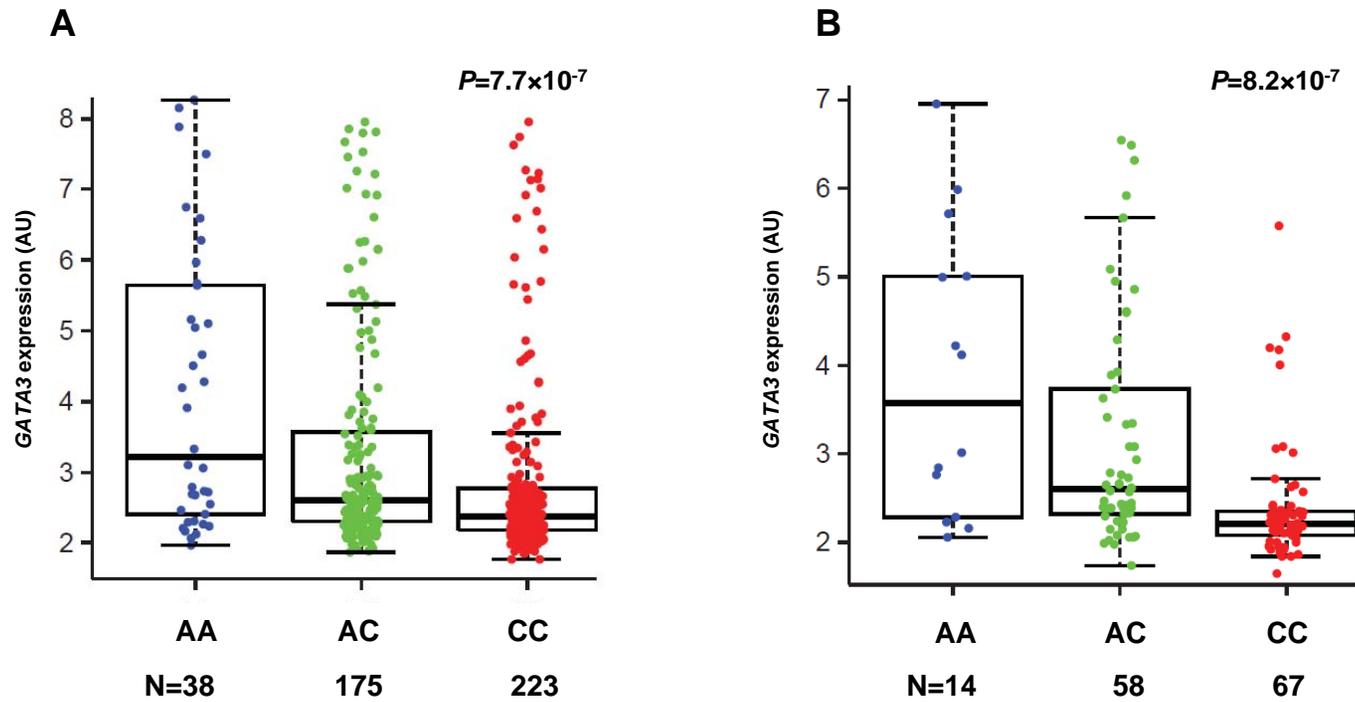
Supplementary Figure 20. Association between *GATA3* SNPs and Ph-like ALL identified on the basis of ROSE clustering. *GATA3* SNPs rs3824662 (A) and rs3781093 (B) were associated with Ph-like ALL in the discovery GWAS group (COG AALL0232 and dbGAP-MESA). The association was also validated in the replication group (COG P9906 and independent non-ALL controls, Panels C and D). *P*-values were estimated by the logistic regression test after adjusting for ancestry.



Supplementary Figure 21. GATA3 expression in Ph-positive ALL vs. Ph-negative ALL. GATA3 expression was quantified by Affymetrix U133A array in diagnostic bone marrow in COG AALL0232 and association with Ph+ status was tested using a logistic regression model adjusting for ancestry as appropriate. AU, arbitrary units. Boxes include data between the twenty-fifth and the seventy-fifth percentiles.



Supplementary Figure 22. Enrichment of Ph-like signature in genes differentially expressed by rs3824662 genotype in ALL blasts. We first identified genes for which expression was associated with *GATA3* rs3824662 genotype (AA+AC vs. CC) in COG AALL0232 and COG P9906 cohorts. Over-representation of the Ph-like gene signature in those affected by *GATA3* SNP genotype was then evaluated using GSEA and *P* value was estimated based on permutations.



Supplementary Figure 23. rs3824662 was associated with GATA3 expression in non-Ph-like ALL. The A allele at rs3824662 was associated with higher GATA3 expression in non-Ph-like ALL cases in the COG AALL0232 (N=436; $P=7.7 \times 10^{-7}$; **Panel A**) and COG P9906 cohorts (N=139; $P=8.2 \times 10^{-7}$; **Panel B**), indicating direct influence of SNP on GATA3 transcription. Genotype-expression association was evaluated using a linear regression model adjusting by ancestry as appropriate. AU, arbitrary units. Boxes include data between the twenty-fifth and the seventy-fifth percentiles.