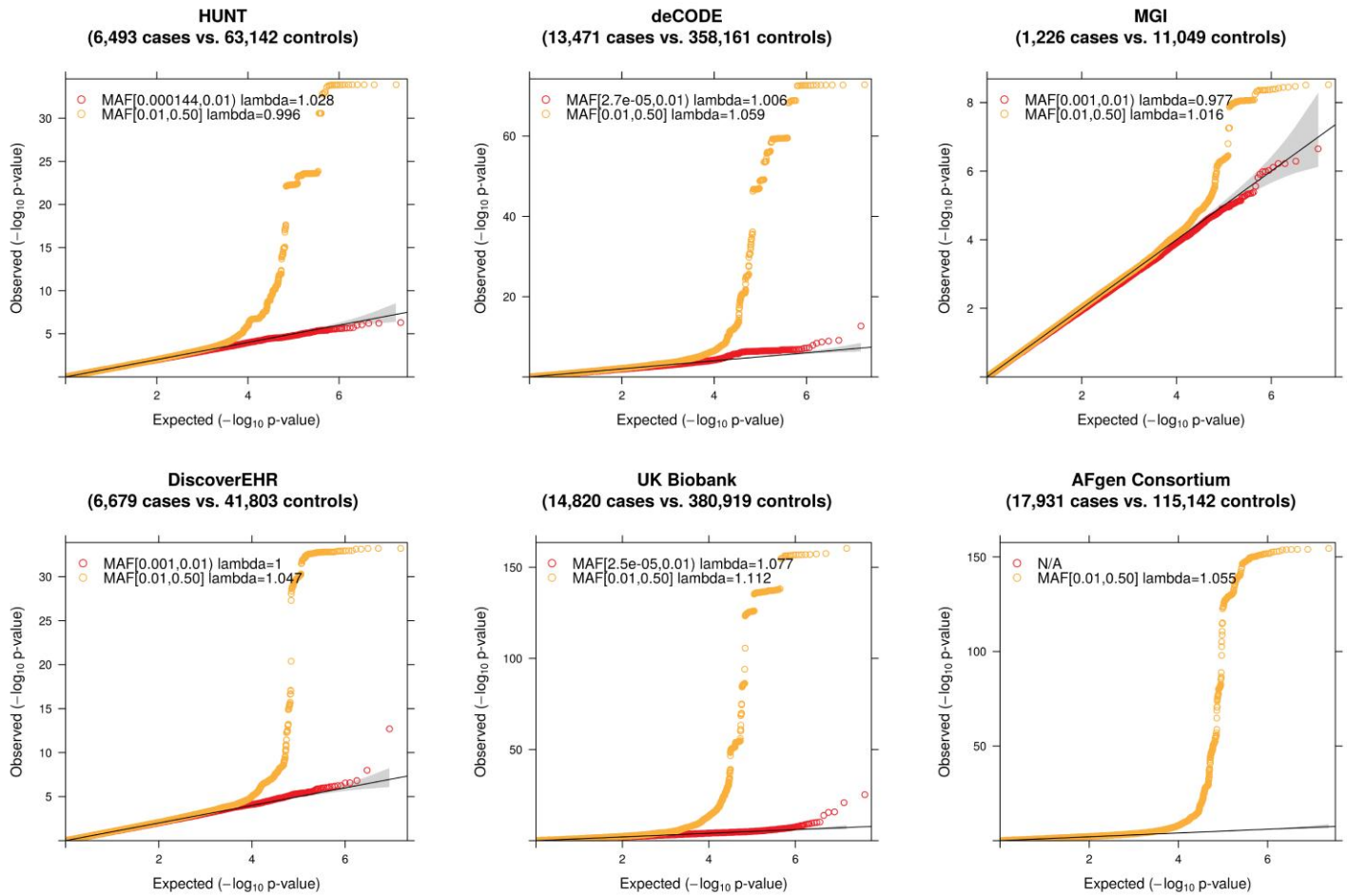


In the format provided by the authors and unedited.

# Biobank-driven genomic discovery yields new insight into atrial fibrillation biology

Jonas B. Nielsen<sup>1,2,23</sup>, Rosa B. Thorolfsdottir<sup>3,4,23</sup>, Lars G. Fritsche<sup>1,5,6,7,23</sup>, Wei Zhou<sup>1,7,23</sup>, Morten W. Skov<sup>8,23</sup>, Sarah E. Graham<sup>1,2,23</sup>, Todd J. Herron<sup>9,23</sup>, Shane McCarthy<sup>10,23</sup>, Ellen M. Schmidt<sup>11,23</sup>, Gardar Sveinbjornsson<sup>3,23</sup>, Ida Surakka<sup>1,2</sup>, Michael R. Mathis<sup>12</sup>, Masatoshi Yamazaki<sup>13</sup>, Ryan D. Crawford<sup>7</sup>, Maiken E. Gabrielsen<sup>5,6</sup>, Anne Heidi Skogholt<sup>5,6</sup>, Oddgeir L. Holmen<sup>5,6,14</sup>, Maoxuan Lin<sup>1,2</sup>, Brooke N. Wolford<sup>1,7</sup>, Rounak Dey<sup>11</sup>, Håvard Dalen<sup>15,16,17</sup>, Patrick Sulem<sup>1,3</sup>, Jonathan H. Chung<sup>10</sup>, Joshua D. Backman<sup>10</sup>, David O. Arnar<sup>3,4,18</sup>, Unnur Thorsteinsdottir<sup>3,4</sup>, Aris Baras<sup>10</sup>, Colm O'Dushlaine<sup>10</sup>, Anders G. Holst<sup>8</sup>, Xiaoquan Wen<sup>11</sup>, Whitney Hornsby<sup>1</sup>, Frederick E. Dewey<sup>10</sup>, Michael Boehnke<sup>1,11</sup>, Sachin Kheterpal<sup>12</sup>, Bhramar Mukherjee<sup>11</sup>, Seunggeun Lee<sup>1,11</sup>, Hyun M. Kang<sup>11</sup>, Hilma Holm<sup>3</sup>, Jacob Kitzman<sup>2</sup>, Jordan A. Shavit<sup>1,19</sup>, José Jalife<sup>1,9,20</sup>, Chad M. Brummett<sup>12</sup>, Tanya M. Teslovich<sup>10</sup>, David J. Carey<sup>21</sup>, Daniel F. Gudbjartsson<sup>1,3,22</sup>, Kari Stefansson<sup>3,4</sup>, Gonçalo R. Abecasis<sup>6,11\*</sup>, Kristian Hveem<sup>5,6,15\*</sup> and Cristen J. Willer<sup>1,2,7\*</sup>

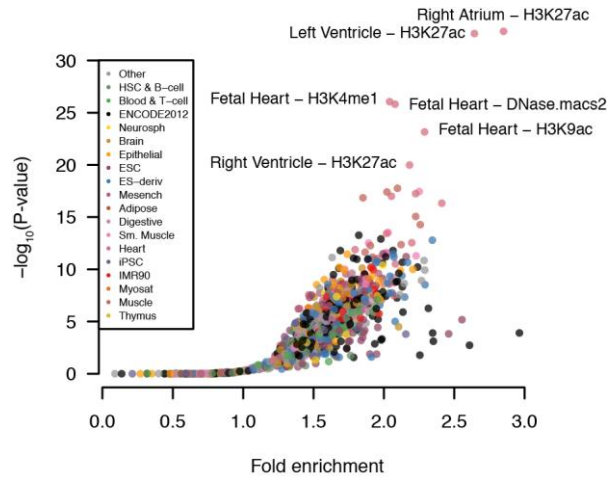
<sup>1</sup>Department of Internal Medicine, Division of Cardiovascular Medicine, University of Michigan, Ann Arbor, MI, USA. <sup>2</sup>Department of Human Genetics, University of Michigan, Ann Arbor, MI, USA. <sup>3</sup>deCODE genetics/Amgen, Inc., Reykjavik, Iceland. <sup>4</sup>Faculty of Medicine, University of Iceland, Reykjavik, Iceland. <sup>5</sup>HUNT Research Centre, Department of Public Health and General Practice, Norwegian University of Science and Technology, Levanger, Norway. <sup>6</sup>K.G. Jebsen Center for Genetic Epidemiology, Department of Public Health, Norwegian University of Science and Technology, Trondheim, Norway. <sup>7</sup>Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI, USA. <sup>8</sup>Laboratory of Molecular Cardiology, Department of Cardiology, The Heart Centre, Copenhagen University Hospital, Rigshospitalet, Copenhagen, Denmark. <sup>9</sup>Department of Internal Medicine, Center for Arrhythmia Research, University of Michigan, Ann Arbor, MI, USA. <sup>10</sup>Regeneron Genetics Center, Tarrytown, NY, USA. <sup>11</sup>Department of Biostatistics, Center for Statistical Genetics, University of Michigan, Ann Arbor, MI, USA. <sup>12</sup>Department of Anesthesiology, University of Michigan, Ann Arbor, MI, USA. <sup>13</sup>Medical Device Development and Regulation Research Center, The University of Tokyo, Tokyo, Japan. <sup>14</sup>Department of Cardiology, St. Olav's University Hospital, Trondheim, Norway. <sup>15</sup>Department of Medicine, Levanger Hospital, Nord-Trøndelag Hospital Trust, Levanger, Norway. <sup>16</sup>Department of Circulation and Medical Imaging, Norwegian University of Science and Technology, Trondheim, Norway. <sup>17</sup>Department of Cardiology, St. Olav's University Hospital, Trondheim University Hospital, Trondheim, Norway. <sup>18</sup>Department of Medicine, Landspítali - National University Hospital, Reykjavik, Iceland. <sup>19</sup>Department of Pediatrics and Communicable Diseases, University of Michigan, Ann Arbor, MI, USA. <sup>20</sup>Fundación Centro Nacional de Investigaciones Cardiovasculares (CNIC), Madrid, Spain. <sup>21</sup>Geisinger Health System, Danville, PA, USA. <sup>22</sup>School of Engineering and Natural Sciences, University of Iceland, Reykjavik, Iceland. <sup>23</sup>These authors contributed equally: Jonas B. Nielsen, Rosa B. Thorolfsdottir, Lars G. Fritsche, Wei Zhou, Morten W. Skov, Sarah E. Graham, Todd J. Herron, Shane McCarthy, Ellen M. Schmidt, Gardar Sveinbjornsson. \*e-mail: [goncalo@umich.edu](mailto:goncalo@umich.edu); [kristian.hveem@ntnu.no](mailto:kristian.hveem@ntnu.no); [cristen@umich.edu](mailto:cristen@umich.edu)



**Supplementary Figure 1**

**Quantile–quantile plots for genome-wide single-variant association analyses for the six contributing study cohorts.**

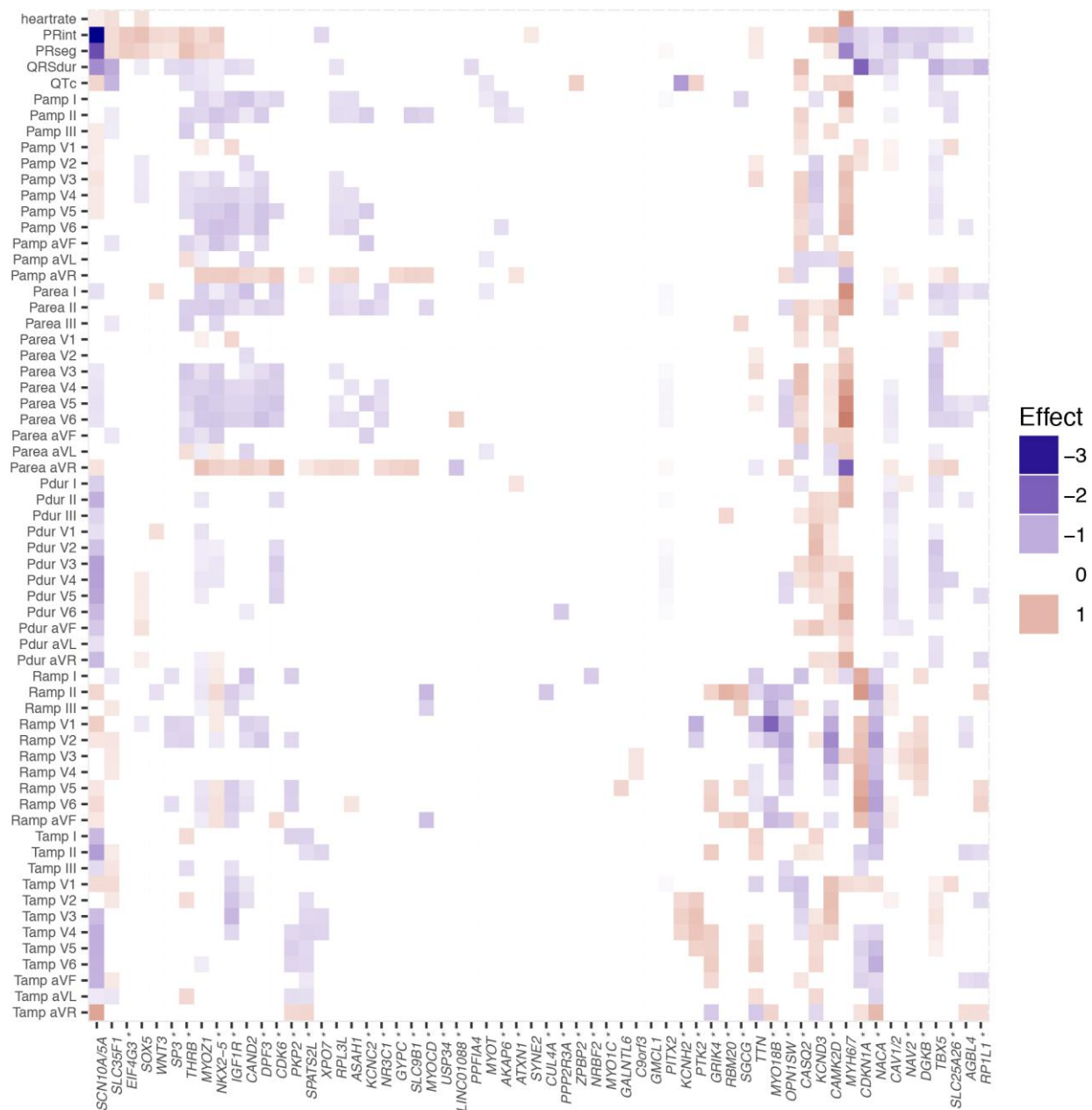
Markers are stratified by minor allele frequency below versus above 0.01. A genomic control factor of 1.38 was applied to the deCODE association results. Dots indicate observed  $P$  values ( $-\log_{10}(P \text{ value})$ ) compared with those expected by chance under the null hypothesis (no association). The black line indicates the identity (no association) with corresponding 95% confidence intervals.



**Supplementary Figure 2**

**Enrichment of atrial fibrillation-associated risk variants in regulatory elements across 127 Roadmap Epigenomics tissue groups.**

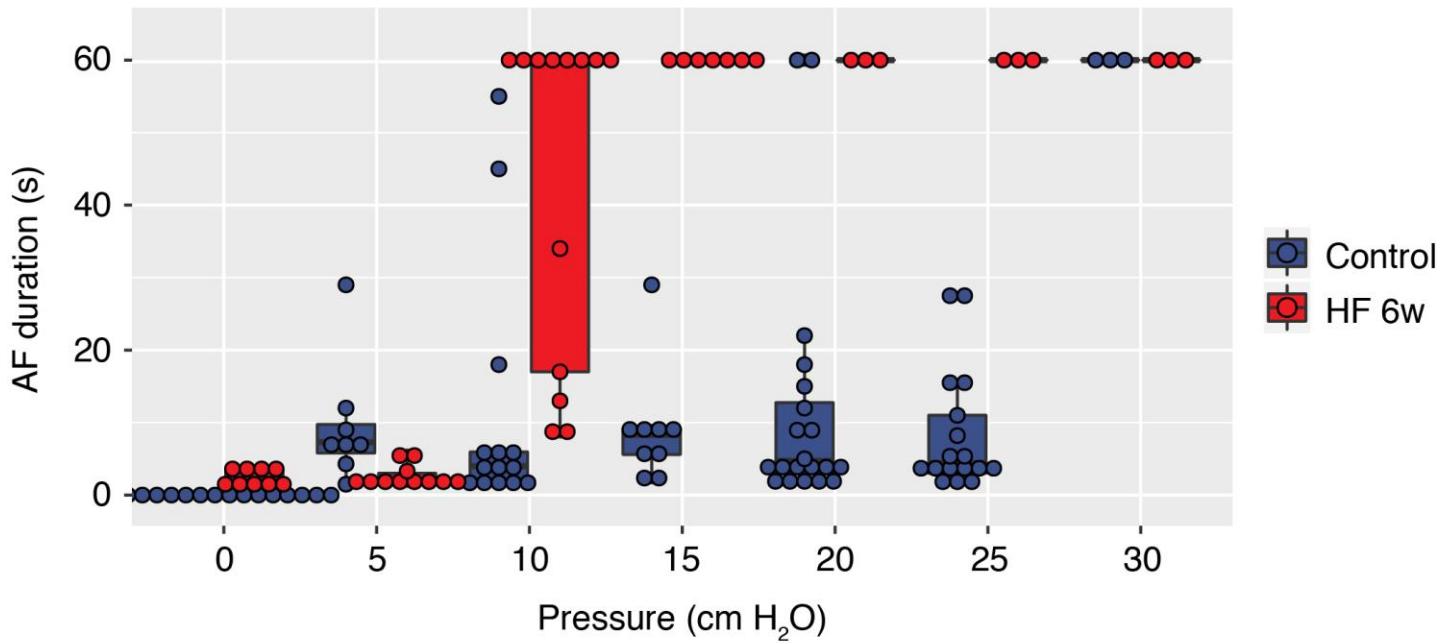
A total of 785 combinations of regulatory features and tissues were examined. *P* values and fold enrichment were estimated using GREGOR. The most statistically significant findings comprised an overlap with H3K27 in right atrium and left ventricle along with H3K4me1 and DNase sites in fetal heart.



**Supplementary Figure 3**

**Heat map showing the effects of atrial fibrillation (AF) variants on electrocardiogram (ECG) traits in sinus rhythm ECGs, excluding AF cases.**

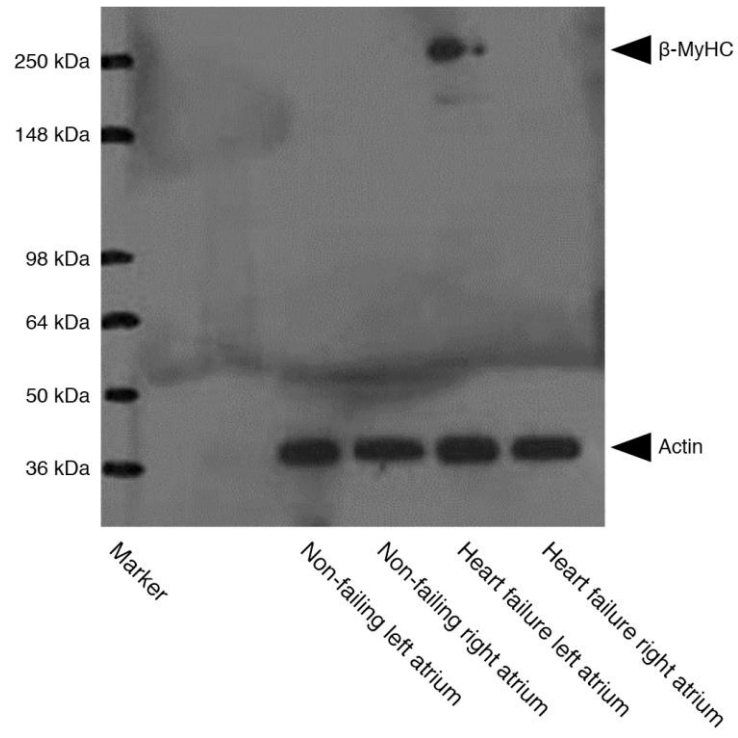
Sinus rhythm ECG measurements were available for 62,974 Icelandic individuals without diagnosis of AF. Each column shows the estimated effect of the AF risk allele on various ECG traits. The effect of each variant, annotated with the locus gene name, is scaled with the log AF odds ratio. Novel variants are marked with an asterisk. Red represents a positive effect of the AF risk allele on the ECG variable, and blue represents negative effect. The effect is shown only for significant associations after adjusting for multiple testing with a false discovery rate procedure for each variant. Non-significant associations are white in the heat map. Sixty of 111 variants with at least one association are shown. *P* values and effect estimates were obtained using BOLT-LMM. For readability, selected highly correlated lead-specific time duration ECG variables (P interval,  $r^2 > 0.51$ ; PR segment,  $r^2 > 0.46$ ; QRS duration,  $r^2 > 0.47$ ; and T duration,  $r^2 > 0.16$ ) have been omitted from the plot. A complete set of association results is provided in Supplementary Table 12. Print, PR interval; PRseg, PR segment; QRSdur, QRS interval duration; Pamp, P-wave amplitude; Parea, P-wave area; Pdur, P-wave duration; Ramp, R-wave amplitude; Tamp, T-wave amplitude.



**Supplementary Figure 4**

**Relationship between left atrium pressure and duration of atrial fibrillation (AF) following burst pacing of rabbit hearts.**

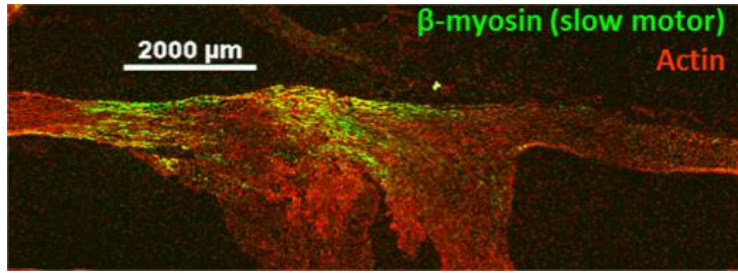
This is an extended version of Fig. 4b showing all individual data points. Heart failure (HF) hearts ( $n = 4$ ) developed long-lasting AF (>60 s) when intra-atrial pressure was increased to 10 cm H<sub>2</sub>O. Control hearts ( $n = 4$ ) did not develop long-lasting AF until intra-atrial pressure was increased to 30 cm H<sub>2</sub>O. Each individual measurement (represented by a dot) is superimposed on box plots showing the median (horizontal black lines), interquartile range (upper and lower box borders), and interquartile range  $\times 1.5$  (vertical black lines) of AF duration.



**Supplementary Figure 5**

**Western blotting for MYH7 expression (β-MyHC protein) indicates MYH7 expression exclusively in the remodeled heart failure left atrium.**

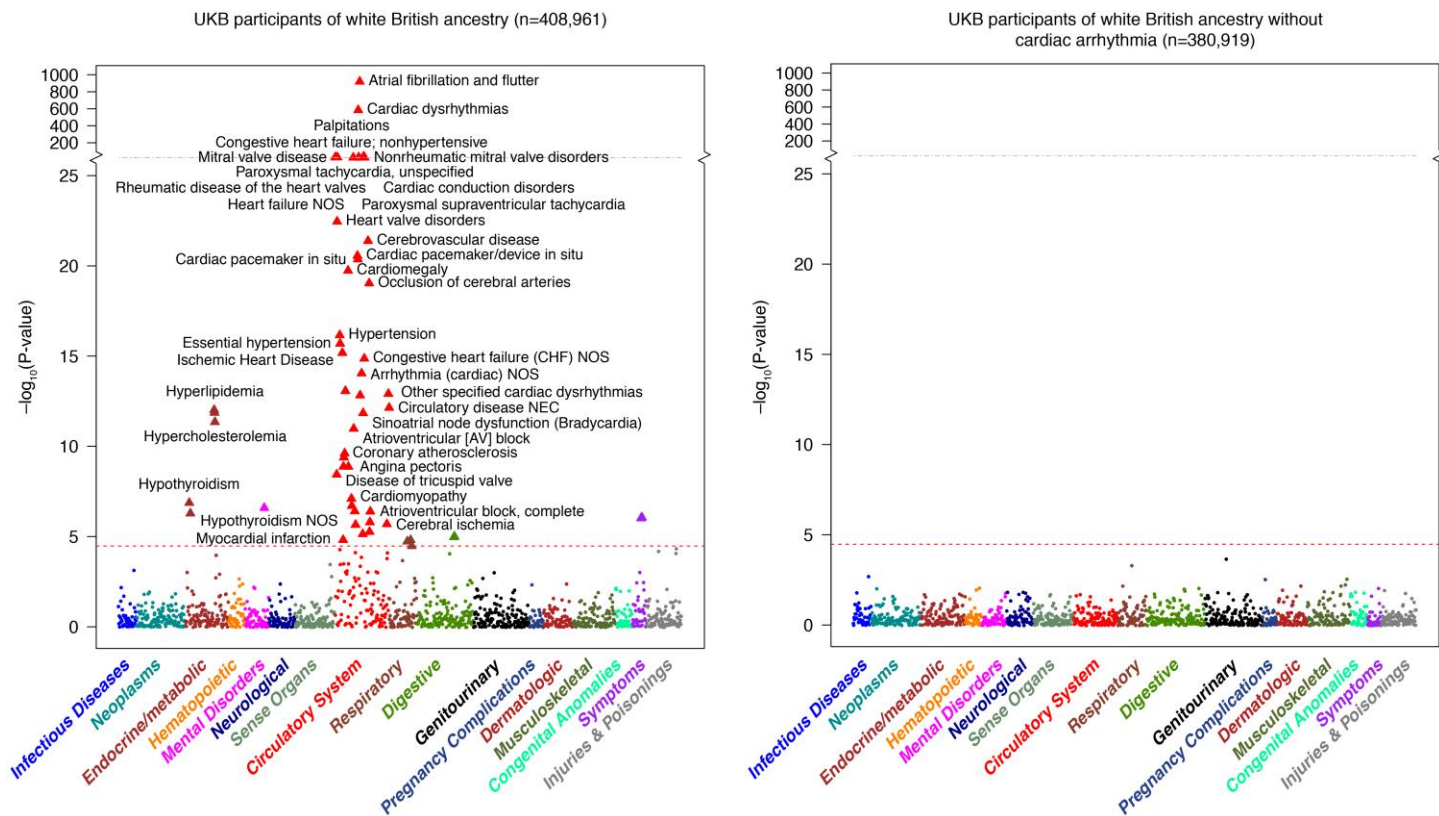
Uncropped version of Fig. 4c.



**Supplementary Figure 6**

**Immunostaining and confocal microscopy reveal heterogeneous MYH7 expression in the heart failure left atrium.**

Green represents MYH7 expression ( $\beta$ -myosin), and red represents actin filaments.

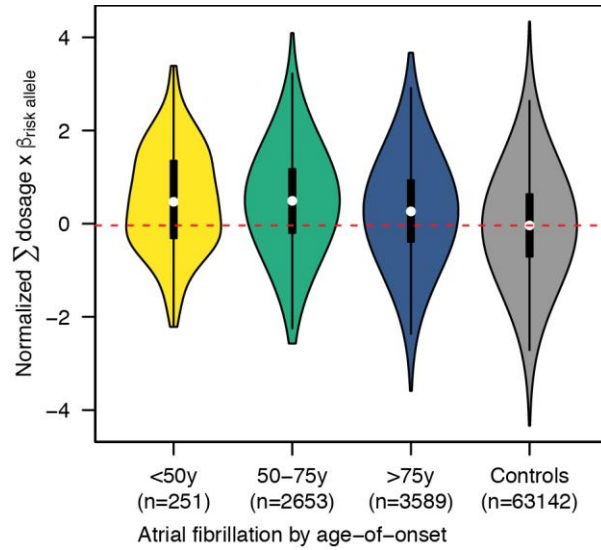


Supplementary Figure 7

**Association between atrial fibrillation polygenic risk score ( $n = 142$  markers) and 1,494 ICD-based traits in UK Biobank participants of white British ancestry.**

Association tests were performed using a logistic regression adjusted for sex and birth year. The horizontal dotted red line represents a  $P$ -value threshold of significance based on Bonferroni correction ( $P < 0.05/1,494 = 3.3 \times 10^{-5}$ ). Some labels have been omitted on the left plot (see Supplementary Table 15 for details on association results).





### Supplementary Figure 8

#### Polygenic risk score distributions for atrial fibrillation-associated variants stratified by age of onset of disease.

Results are based on the HUNT Study only. White dots represent the median, black boxes represent interquartile ranges, black whiskers are the interquartile range times 1.5, and the colored areas show the probability density of the data. The horizontal red dotted line represents the median score for controls.

**Supplementary Note** for the paper “*Biobank-driven genomic discovery yields new insight into atrial fibrillation biology*” by Nielsen J.B., Thorolfsdottir R.B. et. al.

## ***Extended description of selected study cohorts***

### **The HUNT Study**

#### *Sample ascertainment and phenotype definition*

The Nord-Trøndelag Health Study (HUNT) is a population-based health survey conducted in the county of Nord-Trøndelag, Norway. Individuals were included at three different time points during approximately 20 years (HUNT1 [1984-1986], HUNT2 [1995-1997] and HUNT3 [2006-2008]).<sup>1</sup> At each time point, the entire adult population ( $\geq 20$  years) was invited to participate by completing questionnaires, attending clinical examinations and interviews. Participation rates have generally been high: 89.4% (n = 77,212), 69.5% (n = 65 237) and 54.1% (n = 50 807) in HUNT1, HUNT2 and HUNT3, respectively.<sup>1</sup> Taken together, the health studies included information from over 120,000 different individuals from Nord-Trøndelag. Biological samples including DNA have been collected for approximately 70,000 participants. Atrial fibrillation was defined based on ICD-10 codes collected from local hospitals and out-patient clinics between 1999-2016. Cases were defined as individual with one or more ICD-9 or ICD-10 codes specific for atrial fibrillation ("I48" or "427.3") whereas controls were all individuals without a code specific for atrial fibrillation.

#### *Genotyping, quality control, and imputation*

In total, DNA from 71,860 HUNT samples was genotyped using one of three different Illumina HumanCoreExome arrays (HumanCoreExome12 v1.0, HumanCoreExome12 v1.1 and UM HUNT Biobank v1.0). We excluded samples that failed to reach a 99% call rate, had contamination  $> 2.5\%$  as estimated with BAF Regress,<sup>2</sup> large chromosomal copy number variants, lower call rate of a technical duplicate pair and twins, gonosomal constellations other than XX and XY, or whose inferred sex contradicted the reported gender. Samples that passed quality control were analyzed in a second round of genotype calling following the Genome Studio quality control protocol described elsewhere.<sup>3</sup> Genomic position, strand orientation and the reference allele of genotyped variants were determined by aligning their probe sequences against the human genome (Genome Reference Consortium Human genome build 37 and revised Cambridge Reference Sequence of the human mitochondrial DNA; <http://genome.ucsc.edu>) using BLAT.<sup>4</sup> PLINK v1.90<sup>5</sup> was then used to exclude variants if their probe sequences could not be perfectly mapped, cluster separation was  $< 0.3$ , Gentrain score  $< 0.15$ , showed deviations from Hardy Weinberg equilibrium in unrelated samples of European ancestry with p-value  $<$

0.0001), had a call rate < 99%, or another assay with higher call rate genotyped the same variant. Ancestry of all samples was inferred by projecting all genotyped samples into the space of the principal components of the Human Genome Diversity Project (HGDP) reference panel (938 unrelated individuals; downloaded from <http://csg.sph.umich.edu/chaolong/LASER/>),<sup>6,7</sup> using PLINK. Recent European ancestry was defined as samples that fell into an ellipsoid spanning exclusively European populations of the HGDP panel. The different arrays were harmonized by reducing to a set of overlapping variants and excluding variants that showed frequency differences > 15% between data sets, or that were monomorphic in one and had MAF > 1% in another data set. The resulting genotype data were phased using Eagle2 v2.3.<sup>8</sup> Imputation was performed on samples of recent European ancestry using Minimac3 (v2.0.1, <http://genome.sph.umich.edu/wiki/Minimac3>)<sup>9</sup> and a merged reference panel that was constructed by combining the Haplotype Reference Consortium panel (release version 1.1)<sup>10</sup> and a local reference panel based on 2,202 whole-genome sequenced HUNT study participants. A maximal set of relatively unrelated individuals (kinship coefficient < 0.0884) was chosen using KING<sup>11</sup> and FastIndep<sup>12</sup>.

#### *Association analysis*

We performed testing for association with AF using a generalized mixed model including covariates birth year, sex, genotype batch, and principal components (PC) 1-4 as implemented in SAIGE.<sup>13</sup> PCs were computed using PLINK. Additional filters applied to the analysis included minor allele count  $\geq 10$  and imputation  $r^2 \geq 0.3$ .

#### **Michigan Genomics Initiative (MGI)**

DNA from blood samples of surgical patients at the University of Michigan Health System was genotyped on a customized Illumina HumanCore Exome array. Genotypes of the Haplotype Reference Consortium were imputed into the phased MGI genotypes, resulting in dense mapping at over 20 million variants. Atrial fibrillation cases were derived from electronic health records for patients with at least 1 encounter with ICD-9 billing code 427.31 ('Atrial fibrillation'). We performed a genome-wide association analysis for Atrial fibrillation in 1,226 cases and 11,049 controls of European ancestry (excluding 1<sup>st</sup> and 2<sup>nd</sup> degree relatives). Exclusion criteria for controls included individuals with ICD-9 codes 426-427.99. We used the Firth bias-corrected logistic likelihood ratio test with adjustment for age, sex, and principal components 1-4.

#### **DiscovEHR Collaboration Cohort**

The DiscovEHR human genetics cohort analyzed here included 58,124 consented enrollees of European ancestry from the ongoing MyCode Community Health Initiative of the Geisinger Health System ("DiscovEHR study"). Participants were recruited from outpatient primary care and specialty clinics, the cardiac catheterization laboratory, and from patient

populations referred for bariatric and abdominal vascular surgery between 2007 and 2016. Clinical laboratory measurements, International Classification of Diseases, Ninth Revision (ICD-9) disease diagnosis codes, medications, and procedural codes were extracted from the electronic health records (EHR) recording a median of 15 years of clinical care.

Cases were defined as DiscovEHR participants with at least one electronic health record problem list entry or at least two diagnosis code entries for two separate clinical encounters on separate calendar days for ICD-10 I48: Atrial fibrillation and flutter. Corresponding controls were defined as individuals with no EHR diagnosis code entries (problem list or encounter codes) for ICD-10 I48.\*.

Aliquots of DNA were sent to Illumina for genotyping on the Human OmniExpress Exome Beadchip. All individuals of European ancestry, as determined using principal components (PC) analysis, were imputed to the HRC Reference Panel using the Michigan Imputation Server. Markers with imputation  $r^2 > 0.3$  were carried forward for analysis with no filtering for MAF. BOLT-LMM was used to analyze BGEN dosage files, and variants were tested for association with atrial fibrillation under an additive genetic model, adjusting for sex, age, age<sup>2</sup>, and the first four PCs of ancestry; additionally, a genetic relatedness matrix (calculated using variants with MAF > 1%, per-genotype missing data rate < 1%, and HWE P-value <  $1 \times 10^{-15}$ ) was included as a random-effects variable in the model.

### ***Supplementary references***

1. Krokstad, S. *et al.* Cohort profile: The HUNT study, Norway. *Int. J. Epidemiol.* **42**, 968–977 (2013).
2. Jun, G. *et al.* Detecting and estimating contamination of human DNA samples in sequencing and array-based genotype data. *Am. J. Hum. Genet.* **91**, 839–848 (2012).
3. Guo, Y. *et al.* Illumina human exome genotyping array clustering and quality control. *Nat. Protoc.* **9**, 2643–2662 (2014).
4. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
5. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience* **4**, 7 (2015).
6. Wang, C. *et al.* Ancestry estimation and control of population stratification for sequence-based association studies. *Nat. Genet.* **46**, 409–415 (2014).

7. Li, J. Z. *et al.* Worldwide human relationships inferred from genome-wide patterns of variation. *Science* **319**, 1100–1104 (2008).
8. Loh, P.-R. *et al.* Reference-based phasing using the Haplotype Reference Consortium panel. *Nat. Genet.* **48**, 1443–1448 (2016).
9. Das, S. *et al.* Next-generation genotype imputation service and methods. *Nat. Genet.* **48**, 1284–1287 (2016).
10. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat. Genet.* **48**, 1279–1283 (2016).
11. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies. *Bioinforma. Oxf. Engl.* **26**, 2867–2873 (2010).
12. Ma, C., Blackwell, T., Boehnke, M., Scott, L. J. & GoT2D investigators. Recommended joint and meta-analysis strategies for case-control association testing of single low-count variants. *Genet. Epidemiol.* **37**, 539–550 (2013).
13. Zhou, W. *et al.* Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *bioRxiv* 212357 (2017). doi:10.1101/212357