**Discovery and refinement of genetic loci associated with cardiometabolic risk using dense imputation maps**

**Participating studies**

**Stage 1 (discovery) cohorts**

**ALSPAC WGS.** The Avon Longitudinal Study of Parents and Children (ALSPAC) is a long-term health research project. More than 14,000 mothers enrolled during pregnancy in 1991 and 1992, and the health and development of their children has been followed in great detail ever since [1]. A random sample of 2,040 study participants was selected for WGS. The ALSPAC Genetics Advisory Committee approved the study and all participants gave signed consent to the study.

**TwinsUK WGS.** The Department of Twin Research and Genetic Epidemiology (DTR), is the UK's only twin registry of 11,000 identical and non-identical twins between the ages of 16 and 85 years [2]. The database used to study the genetic and environmental aetiology of age-related complex traits and diseases. The St Thomas's Hospital Ethics Committee approved the study and all participants gave signed consent to the study.

**ALSPAC GWA**: For ALSPAC, a total of 8,365 samples were genotyped in Illumina 550k. Besides the WGS samples, there were another 6,557 samples available [3].

**TwinsUK GWA**: For TwinsUK, there were another 2,575 samples that were unrelated to the sequence dataset (IBS>0.125) with genotypes on Illumina HumanHap300 or Illumina Human610 arrays [4]. Imputed TwinsUK data, although unrelated to those samples selected for WGS, did contain related individuals (mainly co-twins) which would require an association test that adjusts for the relatedness.

**1958 Birth Cohort.** Participants to the cohort have been followed-up regularly since birth with prospective information collected on a wide range of indicators related to health, health behaviour, lifestyle, growth and development. There have been 9 contacts with the participants since their birth (ages 7, 11, 16, 23, 33, 41, 45, 47, and 50 years). The biomedical survey at age 45 years included collection of blood samples and DNA from about 8000 participants. The survey was approved by the South East multicentre research ethics committee (MREC). There was an informed consent process conducted by the National Centre for Social Research [5].

**INGI-Val Borbera.** The INGI-Val Borbera population is a collection of 1,785 genotyped samples collected in the Val Borbera Valley, a geographically isolated valley located within the Appennine Mountains in Northwest Italy [6]. The valley is inhabited by about 3,000 descendants from the original population, living in 7 villages along the valley and in the mountains. Participants were healthy people 18-102 years of age that had at least one grandfather living in the valley. A standard battery of tests were performed by the laboratory of ASL 22 - Novi Ligure (AL), on sera from fasting blood collected in the morning. The project was approved by the Ethical committee of the San Raffaele Hospital and of the Piemonte Region. All participants signed an informed consent.

**INGI FVG.** The INGI Friuli Venezia Giulia (FVG) cohort comprised of about 1700 samples from six isolated villages covering a total area of 7858 km2 in a hilly part of Friuli-Venezia Giulia (FVG) county located in north-eastern Italy [7]. Genotyping and phenotypic data for 1590 samples are available. Participants were randomly selected people 3-92 years of age. People with age < 18 were excluded from analyses. Ethics approval was obtained from the Ethics Committee of the Burlo Garofolo children hospital in Trieste. Written informed consent was obtained from every participant to the study.

**INGI Carlantino.** Carlantino is a small village in the Province of Foggia in southern Italy. Genetic analyses of chromosome Y haplotypes as well as mitochondrial DNA show that Carlantino is a genetically homogeneous population and not only a geographically isolated village [8]. Participants were randomly selected in a range of 15 – 90 years of age. Genotyping and phenotypic data are available for 630 individuals. People with age < 18 were excluded from analyses. The local administration of Carlantino, the Health Service of Foggia Province, Italy, and

ethical committee of the IRCCS Burlo-Garofolo of Trieste approved the project. Written informed consent was obtained from every participant to the study.

**INCIPE.** For the INCIPE study, 6200 randomly chosen individuals, all Caucasians and at least 40 years of age as of 1 January 2006, received a letter inviting them to participate in the study. A total of 3870 subjects (62%) accepted and were enrolled. Two studies were included in the analysis: **1.** INCIPE1: Individuals genotyped on Affymetrix 500k**; 2.** INCIPE2: Individuals genotyped on HumanCoreExome-12v1**.** The ethics committees of the involved institutions approved the study protocol.

The **Ludwigshafen Risk and Cardiovascular Health (LURIC) study.** The LURIC study is a prospective study of more than 3,300 individuals of German ancestry in whom cardiovascular and metabolic phenotypes (CAD, MI, dyslipidaemia, hypertension, metabolic syndrome and diabetes mellitus) have been defined or ruled out using standardised methodologies in all study completed participants. A 10-year clinical follow-up for total and cause specific mortality has been completed [9]. From 1997 to 2002 about 3,800 patients were recruited at the Heart Center of Ludwigshafen (Rhein). Inclusion criteria were: German ancestry, clinical stability (except for acute coronary syndromes) and existence of a coronary angiogram. Exclusion criteria were: any acute illness other than acute coronary syndromes, any chronic disease where non-cardiac disease predominated and a history of malignancy within the last five years. The study was approved by the ethics review committee at the Landesärztekammer Rheinland-Pfalz in Mainz, Germany, and written informed consent was obtained from the participants.

**CBR: Cambridge BioResource:** CBR is a collection of pseudo-anonymised DNA samples from 8,000 healthy blood donors that has been established in 2008 and 2010 by the NIHR funded Cambridge Biomedical Research Centre in collaboration with NHS Blood and Transplant for use in genotype-phenotype association studies [10]. Four thousand donors each were enrolled during 2007 and 2009. Full blood counts (FBCs) were obtained from EDTA anticoagulated samples of blood drawn from the pouches of the donation collection sets. FBCs performed on an ABX Pentra 60 automated haematology analyser (ABX Diagnostics, Montpellier, France) or on a Sysmex XE-2100. For the purpose of calibration measurements, 500 blood samples were performed on both the Beckman-Coulter and Sysmex instruments. Measurements were performed between 16-24 hours after phlebotomy.

**HELIC-MANOLIS (HA).** The HELIC (Hellenic Isolated Cohorts; www.helic.org) MANOLIS (Minoan Isolates) collection focuses on Anogia and surrounding Mylopotamos villages. Recruitment of this population-based sample was primarily carried out at the village medical centres. All individuals were older than 17 years and had to have at least one parent from the Mylopotamos area. The study includes biological sample collection for DNA extraction and lab-based blood measurements, and interview-based questionnaire filling. The phenotypes collected include anthropometric and biometric measurements, clinical evaluation data, biochemical and haematological profiles, self-reported medical history, demographic, socioeconomic and lifestyle information. The Harokopio University Bioethics Committee approved the study and informed consent was obtained from every participant.

**HELIC-Pomak (HP).** The HELIC (Hellenic Isolated Cohorts; www.helic.org) Pomak collection focuses on the Pomak villages, a set of isolated mountainous villages in the North of Greece. Recruitment of this population-based sample was primarily carried out at the village medical centres. The study includes biological sample collection for DNA extraction and lab-based blood measurements, and interview-based questionnaire filling. The phenotypes collected include anthropometric and biometric measurements, clinical evaluation data, biochemical and haematological profiles, self-reported medical history, demographic, socioeconomic and lifestyle information. The Harokopio University Bioethics Committee approved the study and informed consent was obtained from every participant.

**TEENAGE.** Participants were drawn from the TEENAGE (TEENs of Attica: Genes and Environment) study. A random sample of 857 adolescent students attending public secondary schools located in the wider Athens area of Attica in Greece were recruited in the study from 2008 to 2010. Our sample comprised 707 (55.9% females) adolescents of Greek origin aged 13.42 ± 0.88 years. Details of recruitment and data collection have been described elsewhere [11]. Prior to recruitment all study participants gave their verbal assent along with their parents'/guardians' written consent forms. Harokopio University Bioethics Committee and the Greek Ministry of Education, Lifelong Learning and Religious Affairs approved the study.

**Women's Health Initiative (WHI**): WHI is one of the largest (n=161,808) studies of women's health ever undertaken in the U.S [12]. There are two major components of WHI: (1) a Clinical Trial (CT) that enrolled and randomized 68,132 women ages 50 – 79 into at least one of three placebo-control clinical trials (hormone therapy, dietary modification, and calcium/vitamin D); and (2) an Observational Study (OS) that enrolled 93,676 women of the same age range into a parallel prospective cohort study. A diverse population including 26,045 (17%) women from minority groups was recruited from 1993-1998 at 40 clinical centers across the U.S. The design has been published [13,14]. For the CT and OS participants enrolled in WHI and who had consented to genetic research, DNA was extracted by the Specimen Processing Laboratory at the Fred Hutchinson Cancer Research Center (FHCRC) from specimens that were collected at the time of enrolment in to the study (between 1993 and 1998).

**Stage 2 cohorts**

**London Life Sciences Prospective Population Study (LOLIPOP).** LOLIPOP is an ongoing community cohort of 25,372 individuals aged 35-75 years, recruited in West London, UK to study the environmental and genetic factors that contribute to cardiovascular disease among UK Indian Asians. The study includes both European and Indian Asian subjects. For the current study, only white individuals were included in the primary meta-analysis. Three studies were included in the analysis: (1). LOLIPOP - EWA: European whites from the general population, genotyped on Affymetrix 500K arrays. (2). LOLIPOP - EWP: European whites from the general population, genotyped on Perlegen custom array. (3). LOLIPOP - EW610: European whites from the general population, genotyped on Illumina Human610 array.

**UK BioBank (UKBB).** The UK Biobank resource (www.ukbiobank.ac.uk) is a large-scale prospective study consisting of ~500,000 samples aged between 40-69 from all over the UK. The aim of the resource is to advance the diagnosis and treatment, as well as the prevention of serious illnesses, such as heart diseases. The purpose of considering this age group is that in the following decades, the participants are at an increased risk of developing a wide range of diseases, including diabetes, heart disease and stroke and thus allows investigations to be carried out for such conditions. The genotyping and imputation procedures are described here: https://biobank.ctsu.ox.ac.uk/crystal/docs/genotyping_sample_workflow.pdf . Our sample comprised 98,880 (54.3% females) individuals with mean age of 57.2.

**Rotterdam Study cohort I (RS-I).** The Rotterdam Study is an ongoing prospective population-based cohort study, focused on chronic disabling conditions of the elderly. The study comprises an outbred ethnically homogenous population of Dutch Caucasian origin. The rationale of the study has been described in detail elsewhere [15]. In summary, 7,983 men and women aged 55 years or older, living in Ommoord, a suburb of Rotterdam, the Netherlands, were invited to participate in the first phase. Fasting blood samples were taken during the participant's third visit to the research center.

**Rotterdam Study cohort II (RS-II).** The Rotterdam Study cohort II prospective population-based cohort study comprises 3,011 residents aged 55 years and older from the same district of Rotterdam. The rationale and study design of this cohort is similar to that of the RS-I [15]. The baseline measurements, including the fasting HDL measurements, took place during the first visit. The Rotterdam Study has been approved by the Medical Ethics Committee of the Erasmus MC and by the Ministry of Health, Welfare and Sport of the Netherlands, implementing the "Wet Bevolkingsonderzoek: ERGO (Population Studies Act: Rotterdam Study)". All participants provided written informed consent to participate in the study and to obtain information from their treating physicians [15].

**Lifelines.** LifeLines is a multidisciplinary prospective population-based cohort study examining, in a unique three-generation design, the health and health-related behaviours of 167729 persons living in the northeast region of The Netherlands. It employs a broad range of investigative procedures in assessing the biomedical, socio-demographic, behavioural, physical and psychological factors that contribute to the health and disease of the general population, with a special focus on multimorbidity and complex genetics. The LifeLines study has been approved by the review board of the University Medical Center, Groningen, and adheres to the principles expressed in the Declaration of Helsinki. All study participants provided written informed consent.

**The Precocious Coronary Artery Disease Study (PROCARDIS)** consists of coronary artery disease (CAD) cases and controls from four European countries (UK, Italy, Sweden and Germany). CAD (defined as myocardial infarction,

acute coronary syndrome, unstable or stable angina, or need for coronary artery bypass surgery or percutaneous coronary intervention) was diagnosed before 66 years of age and 80% of cases had a sibling fulfilling the same criteria for CAD. Subjects with self-reported non-European ancestry were excluded. Among the "genetically-enriched" CAD cases, 70% had suffered myocardial infarction (MI).

### Phenotype harmonization

Phenotype harmonization across studies consisted of identification of outliers, best fit trait transformation (natural log, inverse normal, square root, inverse, or non-transformed) within study, and adjustment for potential confounders, including age, $age^2$, gender, and body mass index (BMI), and batch effect (instrument and date of measurement), depending on the study and trait. Covariates were fit into a linear regression model and only those significantly associated with the traits were included as adjustment variables in the final model. When inverse-normal transformation was used, the samples were divided into males and females for transformation and covariates adjustment, separately. The use of standardization as the last step of the phenotype harmonization facilitated meta-analysis and cross-traits examination of effect sizes.

### Lipids

Lipids measurement methods were as following: for ALSPAC, plasma levels of TC, HDL and TG were measured with enzymatic colorimetric assays (Roche) on a Hitachi Modular P Analyser. LDL was derived from the following formula: TC- (HDL+TG/2.19); for TwinsUK, Enzymatic colorimetric assays were used to measure serum levels of TC, HDL and TG were measured using three analysing devices (Cobas Fara; Roche Diagnostics, Lewes, UK; Kodak Ektachem dry chemistry analysers (Johnson and Johnson Vitros Ektachem machine, Beckman LX20 analysers, Roche P800 modular system)); for 1958BC, serum TG,
TC and HDL were measured in serum by Olympus model AU640 autoanalyser in a central lab in Newcastle. Enzymatic colorimetric determination GPO-PAP method was used to determine TG, CHOD-PAP method for TC and for HDL; for INGI-VB, lipids were measured using HITACHI 917 ROCHE and Unicel Dx-C 800 BECKMAN devices; for INGI-FVG and INGI-Carl, lipids were measured using BIOTECNICA BT-3000 TARGA chemistry analyser; for INCIPE, enzymatic determination of TC and TG was performed on Dimension RxL apparatus (Siemens Diagnostics). HDL cholesterol was determined by the homogeneous method; LDL cholesterol by the Friedewald formula (Friedewald et al. 1972); for LURIC, TC and TG were obtained by ß-quantification from serum and measured enzymatically using WAKO reagents on a WAKO 30R analyser (Neuss, Germany). LDL and HDL were measured after separating lipoproteins with a combined ultracentrifugation- precipitation method; for HELIC Manolis and HELIC Pomak and Teenage, TC, HDL, TG were assessed using enzymatic colorimetric assays and while LDL levels were calculated according to Friedewald equation (Friedewald et al. 1972). For WHI, HDL, LDL, and TG measurements were performed at the University of Minnesota by standard biochemical methods on the Roche Modular P Chemistry analyzer (Roche Diagnostics): HDL was measured in serum by the HDL-C plus third generation direct method; TG was measured in serum by Triglyceride GB reagent, and total cholesterol (TC) was measured in serum by a cholesterol oxidase method. LDL was calculated in serum specimens having a TG value < 400 mg/dl according to the formula of Friedewald et al. Based on the LDL-lowering effects of statins, we estimated the pretreatment LDL value for individuals on lipid-lowering medication by dividing treated LDL values by 0.75.

For phenotype harmonization, extra care was given to the TwinsUK cohorts given there was random efforts of different dates of visits and different instrumental measurements. For ALSPAC and other discovery and replication cohorts, the same phenotype protocol was used. Inverse normal transformation was applied to all cohorts. For each cohort, the residuals with confounding variables regressed out were standardized so that the phenotype had a mean of 0 and a standard deviation of 1.

### Hematologic traits

For ALSPAC, HGB were measured with Hemocue Hb201+ analyser. For all eight hematologic traits in TwinsUK and all other discovery and replication cohorts, the traits were measured with Beckman Coulters, except for WHI, where HGB, HCT, WBC, and PLT were determined at local laboratories using automated hematology cell counters and standardized quality assurance procedures (Margolis et al. 2005). Different phenotype transformation protocol was applied to the eight hematologic phenotypes: inverse normal transformation for HGB, PCV, PLT,

square root for MCH, natural log for WBC, and no transformation for MCHC, MCV, RBC. For each trait of each cohort, the residuals with confounding variables regressed out were standardized so that the phenotype has a mean of 0 and a standard deviation of 1.

### C-reactive protein

CRP was measured by high-sensitivity immunologic assay in all participating cohorts. CRP measurement methods are as following: for ALSPAC, CRP was measured by Latex enhanced assay; for TwinsUK, CRP was measured by automated particle-enhanced immunoturbidimetric assay (Roche UK, Welwyn Garden City, UK); for 1958BC, CRP antigen levels were measured by high sensitivity nephelometric assay using latex particles coated with monoclonal antibodies to human CRP in the BN Prospec protein analyzer (Dade Behring, Marburg, Germany). For HELIC-MANOLIS and HELIC-Pomak, CRP was measured using an immunoturbidimetric assay on a COBAS 8000 analyser (Roche). For WHI, CRP was measured using a latex-particle enhanced immunoturbidimetric assay kit (Roche Diagnostics, Indianapolis, IN). For FHS and the rest of discovery cohorts, CRP was measured in fasting serum samples using various versions of high-sensitivity assay, mostly the Dade Behring BN100. For phenotype harmonizaiton, abnormal values of CRP, defined as <0.1mg/L or >10mg/L, were excluded. For TwinsUK, the phenotype harmonization was conducted for WGS and GWA samples separately. Inverse normal transformation was applied to the full dataset without gender specific transformation. Regression test found no significant effects of dates of visits or analysers. BMI was not included as a covariate. For each trait of each cohort, the residuals with confounding variables regressed out were standardized so that the phenotype has a mean of 0 and a standard deviation of 1.

### Genome-wide association analyses

We present Q-Q plots from the association analyses on each trait in each cohort in Supplementary Figure 1. They are created with no LD filtration as per common practices and overall show no genomic inflation as measured by the λ values. After further inspection of the study-trait pairs showing the earliest deviation from the null, we found that to be a results of single large regions of association, the removal of which diminishes all signal and results in observed p-values following the theoretical null distribution. In particular, a single region on chromosome 11 (spanning ~5Mb) appears to cause inflation in the Q-Q plot for association to MCV in the HELIC (Pomak) sample. This population has been previously been reported to have unusually long haplotypes in the same region for the corresponding allele frequencies (~1.8Mb), which suggest positive selection could have acted on them[16]. Similarly, a single long region is responsible for the early deviations in the QQ plots in HELIC Pomac for MCH and MCHC (same region), and TwinsUK WGS and INCIPE2 for Uric Acid.

### GARFIELD enrichment analysis

### Functional annotation

The genome-wide discovery meta-analysis p-value (denoted later by p) distribution of each trait was used in order to assess whether we observe enrichment in various functional and regulatory features using genic annotations, chromatin states, DNaseI hypersensitive sites, transcription factor (TF) binding sites, FAIRE-seq elements and histone modifications. To do this, a set of independent SNVs for each phenotype was selected and annotated such that each variant was said to have a certain feature if it itself or if any of its LD proxies fell into an appropriate region. A "n-fold" enrichment score was computed to quantify the observed enrichment at various GWAS p-value cut-offs and finally permutations were performed matching on MAF, distance to nearest TSS and number of LD proxies in order to assess the significance of the observed enrichment. The matching is introduced to control for genomic features of the data, which could otherwise lead to biased results.

**Genome functional annotation maps**

a) Genic annotations were obtained from GENCODEv13 and variants were split into categories defined as follows: Intron, Exon (coding), 3'UTR, 5'UTR, Downstream genetic variant (within 5KB of the end of a gene), Upstream genetic variant (within 5KB of the start of a gene).

b) DNaseI data was obtained from ENCODE (ftp://ftp.ebi.ac.uk/pub/databases/ensembl/encode/integration_data_jan2011) the NIH Roadmap Epigenomics Project (http://www.genboree.org/EdaccData/Current-Release/experiment-sample/Chromatin_Accessibility/) on all available cell types. DHS data was processed following DHS data processing protocol described in an ENCODE study[17].

c) Chromatin states, Histone modifications, TF binding sites and FAIRE-seq elements were obtained from ENCODE on the Tier1 and Tier2 cell types.

**Data Processing**

To remove possible biases due to linkage disequilibrium (LD) or dependence between variants we compute the $r^2$ between all SNVs within 1Mb windows and consider $r^2$ of less than 0.1 between two variants to mean (approximate) independence. Next, from the full set of genetic variants for each phenotype, we create an independent set of SNVs, where in order to keep all possible GWAS signals we sequentially find and retain the next most significant (lowest p) variant independent of all other variants in our independence set. Then we annotate each independent SNV and consider it as overlapping a functional element if (1) the SNV itself resides in such a genomic region or (2) at least one of its proxies in LD ($r^2 \geq 0.8$) and within 500Kb with it does. We include the latter as the association of a SNV in GWAS potentially tags the effect of other variants, which could underlie the observed association signal.

**Quantifying enrichment**

To find the enrichment of GWAS signals within a given annotation, we calculate fold enrichment as the fraction of variants that fall in that annotation and have p less than threshold T, divided by the fraction of total number of variants in that annotation.

Specifically,

$$Fold\ Enrichment(t) = \frac{N_a^t}{N^t} \bigg/ \frac{N_a}{N}$$

where N denotes the total number of variants, $N_a$ - the total number of variants that fall in the annotation of interest, $N^t$ - the total number of variants with p less than threshold t, $N_a^t$ - the number of variants with p less than t that fall in the annotation of interest.

**Statistical testing**

We consider that perhaps variants with specific annotations may be more likely to be at specific positions of the genome. To account for this we use permutation testing, where we shuffle the p associated to each variant in our independence set in such a way as to match SNVs according to MAF, distance to nearest TSS and number of LD proxies ($r^2 >= 0.8$) they have. Specifically, we bin variants according to 5 quantiles of MAF, number of LD proxies and distance to nearest TSS, resulting in 125 bins overall. We then permute variants within each bin separately.

**Multiple testing**

We test at 95% significance level and correct for multiple testing by applying a Bonferroni correction for the effective number of distinct annotation used, which we determine by adapting an approach proposed in Galwey, 2009 [18].

## Fine-mapping of regions with multiple causal variants

Results from conditional analysis revealed 85% of associated loci as having a single independent association signal, showing that for the vast majority of cases our assumption of a single causal variant per region is appropriate. For the remaining cases this assumption is expected to result in fine-mapping of the stronger of the underlying association signals. The additional rounds of conditional analysis and further fine-mapping on the conditional analysis summary statistics are performed in order to further fine-map these secondary and potentially weaker signals. To provide evidence for this, we further extracted all regions from our conditional analysis results (Supplementary Table 4) for which conditional analysis identified multiple independent signals. For these 13 regions (11 loci) we additionally performed analysis with FINEMAP[19] with 2 or 3 causal variants (numbers obtained from conditional analysis). Supplementary Table 7 shows a comparison of these results to the ones obtained from our fine-mapping pipeline, and in particular highlight that in all cases the top variant identified from our initial analyses is identified with this more complex approach. Furthermore, for 3 of the regions the more complex fine-mapping strategy would not result in successful fine-mapping (>20 credible set configurations) despite clearly showing our top variant as the most likely causal variant. This is likely to happen because FINEMAP is unable to clearly fine-map the second causal variant in those regions, which we have also observed in our two stage approach.

## Supplementary Figure and Table Legends

**Supplementary Figure 1. Q-Q plots and $\lambda$ values for association analysis results.** Each of the 20 panels shows the Q-Q plots and $\lambda$ values for all cohorts at a given trait, where cohorts are presented with different colours.

**Supplementary Figure 2. Enrichment of GWAS variants in regulatory features.** All figures show fold enrichment (at $10^{-5}$ GWAS significance threshold) for a particular regulatory element and celltype for a given trait, if it is significant, except (**a**) and (**c**), which show the most significant statistic per tissue. Traits and features/celltypes with no enrichment have been removed from figures. (**a**) ENCODE and Roadmap Epigenomics DNaseI hypersensitive sites (peaks). (**b**) GENCODE genic elements. (**c**) ENCODE and Roadmap Epigenomics DNaseI hypersensitive sites (hotspots). (**d**) ENCODE transcription factor binding sites. (**e**) ENCODE segmentation states. (**f**) ENCODE footprints. (**g**) Histone modifications in six Tier1 and Tier2 ENCODE cell lines. (**h**) ENCODE FAIRE-seq elements.

**Supplementary Figure 3. Fine mapping and annotation.** Panels show the regional association locuszoom plot, the PP statistics from the fine-mapping methods, the CATO and deltaSVM scores, VEP genic annotations and overlap of regulatory annotations found significant (coloured in blue) from GARFIELD enrichment analysis. Circles sizes and colours for all scores have been scaled with respect to score type (i.e. PP, CATO or deltaSVM) and numbers have been plotted below each circle.

**Supplementary Figure 4. Histograms of residuals for non-inverse normalized traits.** A panel is shown for each trait-cohort pair, where the trait histogram is shown in light blue and the (assumed) standard normal distribution density is shown in purple.

**Supplementary Table 1. Study descriptives.**

**Supplementary Table 2. Phenotype preparation protocols.** Table summarizing trait transformation, exclusions and covariates for all the traits studied in this study.

**Supplementary Table 3. Association statistics for new loci in discovery and replication.**

**Supplementary Table 4. GENCODE, ENCODE and Roadmap epigenomics annotations used for enrichment analysis with software GARFIELD.** Category refers to Genic, Histone Modifications, Chromatin States, FAIRE, Peaks, Hotspots, Footprints and TFBS (transcription factor binding site); Tissue, Cell type and Type give further information on the annotations, where Type shows the elements within category (e.g. Exon, 3'UTR for Genic elements).

**Supplementary Table 5. Enrichment of cardio-metabolic traits in 1005 GENCODE, ENCODE and Roadmap Epigenomics annotations at the $10^{-5}$ and $10^{-8}$ GWAS significance thresholds**. Each row corresponds to a significantly enriched annotation for a given phenotype (Trait) at a given GWAS P-value threshold (PThresh), where FE denotes the fold enrichment, EmpPval the empirical p-value for significance of the observed enrichment, N - the total number of pruned variants, NThresh - number of variants at that threshold, NAnnot - number of variants overlapping the given annotation, NAnnotThresh - number of variants at threshold which overlap the annotation. Category, Type, Tissue and Celltype provide further information on the annotation under consideration (where available). The table is sorted in descending order according to the FE column.

**Supplementary Table 6. Fine-mapping results.**

**Supplementary Table 7. FINEMAP analysis with a relaxed assumption of multiple causal variants per locus.** Comparison between current fine-mapping approach (assuming single causal variant per region and performing fine-mapping on conditional analyses when necessary) and using FINEMAP with at most N causal variants, where N is the number of conditionally independent signals per region obtained from conditional analysis. The table contains all regions we have fine-mapped which have N>1 from conditional analyses. In red we have highlighted the top variant from our round 0 analysis (unconditional) and in blue the top variant from our round 1 analysis (fine-mapping after conditioning on the best round 0 variant).

**Supplementary References**

1. Golding, J., Pembrey, M., Jones, R. & Team, A.S. ALSPAC--the Avon Longitudinal Study of Parents and Children. I. Study methodology. *Paediatr Perinat Epidemiol* **15**, 74-87 (2001).
2. Moayyeri, A., Hammond, C.J., Hart, D.J. & Spector, T.D. The UK Adult Twin Registry (TwinsUK Resource). *Twin Res Hum Genet*, 1-6 (2012).
3. Bonnelykke, K. *et al.* Meta-analysis of genome-wide association studies identifies ten loci influencing allergic sensitization. *Nat Genet* **45**, 902-6 (2013).
4. Soranzo, N. *et al.* A genome-wide meta-analysis identifies 22 loci associated with eight hematological parameters in the HaemGen consortium. *Nature Genetics* **41**, 1182-1190 (2009).
5. Power, C. & Elliott, J. Cohort profile: 1958 British birth cohort (National Child Development Study). *Int J Epidemiol* **35**, 34-41 (2006).
6. Traglia, M. *et al.* Heritability and Demographic Analyses in the Large Isolated Population of Val Borbera Suggest Advantages in Mapping Complex Traits Genes. *PLoS ONE* **4**, e7554 (2009).
7. Esko, T. *et al.* Genetic characterization of northeastern Italian population isolates in the context of broader European genetic diversity. *Eur J Hum Genet* **21**, 659-65 (2013).
8. Sala, C. *et al.* Variation of hemoglobin levels in normal Italian populations from genetic isolates. *Haematologica* **93**, 1372-5 (2008).
9. Winkelmann, B.R. *et al.* Rationale and design of the LURIC study--a resource for functional genomics, pharmacogenomics and long-term prognosis of cardiovascular disease. *Pharmacogenomics* **2**, S1-73 (2001).
10. Dendrou, C.A. *et al.* Cell-specific protein phenotypes for the autoimmune locus IL2RA using a genotype-selectable human bioresource. *Nat Genet* **41**, 1011-5 (2009).
11. Ntalla, I. *et al.* Body composition and eating behaviours in relation to dieting involvement in a sample of urban Greek adolescents from the TEENAGE (TEENs of Attica: Genes & Environment) study. *Public Health Nutr* **17**, 561-8 (2014).
12. Design of the Women's Health Initiative clinical trial and observational study. The Women's Health Initiative Study Group. *Control Clin Trials* **19**, 61-109 (1998).
13. Anderson, G.L. *et al.* Implementation of the Women's Health Initiative study design. *Ann Epidemiol* **13**, S5-17 (2003).
14. Hays, J. *et al.* The Women's Health Initiative recruitment methods and results. *Ann Epidemiol* **13**, S18-77 (2003).
15. Hofman, A. *et al.* The Rotterdam Study: 2014 objectives and design update. *Eur J Epidemiol* **28**, 889-926 (2013).
16. Panoutsopoulou, K. *et al.* Genetic characterization of Greek population isolates reveals strong genetic drift at missense and trait-associated variants. *Nat Commun* **5**, 5345 (2014).
17. Thurman, R.E. *et al.* The accessible chromatin landscape of the human genome. *Nature* **489**, 75-82 (2012).
18. Galwey, N.W. A new measure of the effective number of tests, a practical tool for comparing families of non-independent significance tests. *Genet Epidemiol* **33**, 559-68 (2009).
19. Benner, C. *et al.* FINEMAP: efficient variable selection using summary data from genome-wide association studies. *Bioinformatics* **32**, 1493-501 (2016).

**Supplementary Figure 1**: **Q-Q plots and λ values for association analysis results.** Each of the 20 panels shows the Q-Q plots and λ values for all cohorts at a given trait, where cohorts are presented with different colours.

**Supplementary Figure 2**: **Enrichment of GWAS variants in regulatory features.** All figures show fold enrichment (at $10^{-5}$ GWAS significance threshold) for a particular regulatory element and celltype for a given trait, if it is significant, except (a) and (c), which show the most significant statistic per tissue. Traits and features/celltypes with no enrichment have been removed from figures. (a) ENCODE and Roadmap Epigenomics DNaseI hypersensitive sites (peaks). (b) GENCODE genic elements. (c) ENCODE and Roadmap Epigenomics DNaseI hypersensitive sites (hotspots). (d) ENCODE transcription factor binding sites. (e) ENCODE segmentation states. (f) ENCODE footprints. (g) Histone modifications in six Tier1 and Tier2 ENCODE cell lines. (h) ENCODE FAIRE-seq elements.

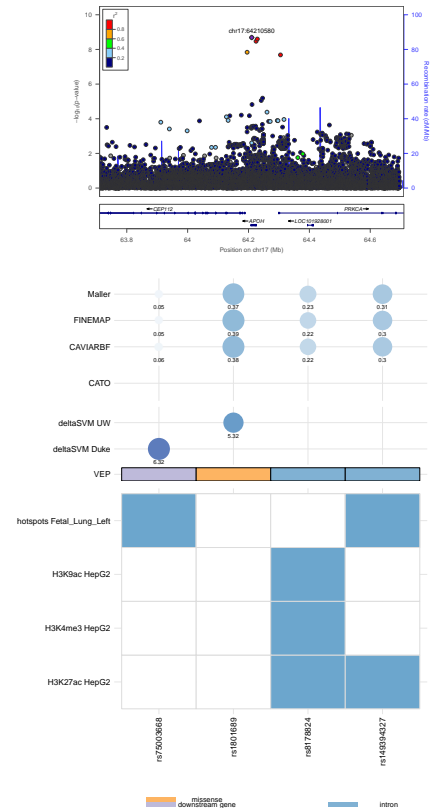Locus 1, TC; Locus 1, LDL; Locus 2, CRP; Locus 3, CRP; Locus 4, PLT; Locus 5, PLT
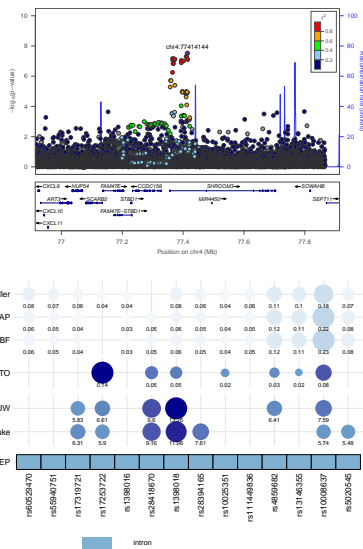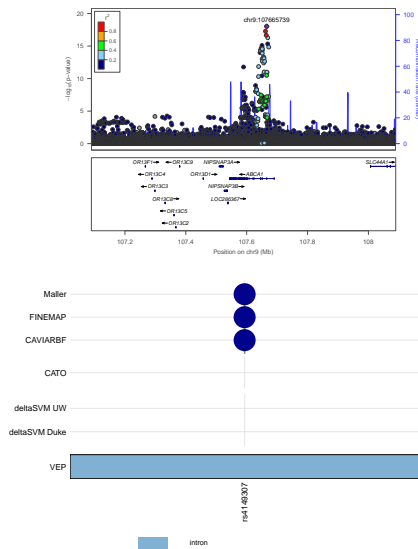
Locus 44, PLT

Locus 45, PCV

Locus 46a, TG

Locus 46b, TC

Locus 47, PLT
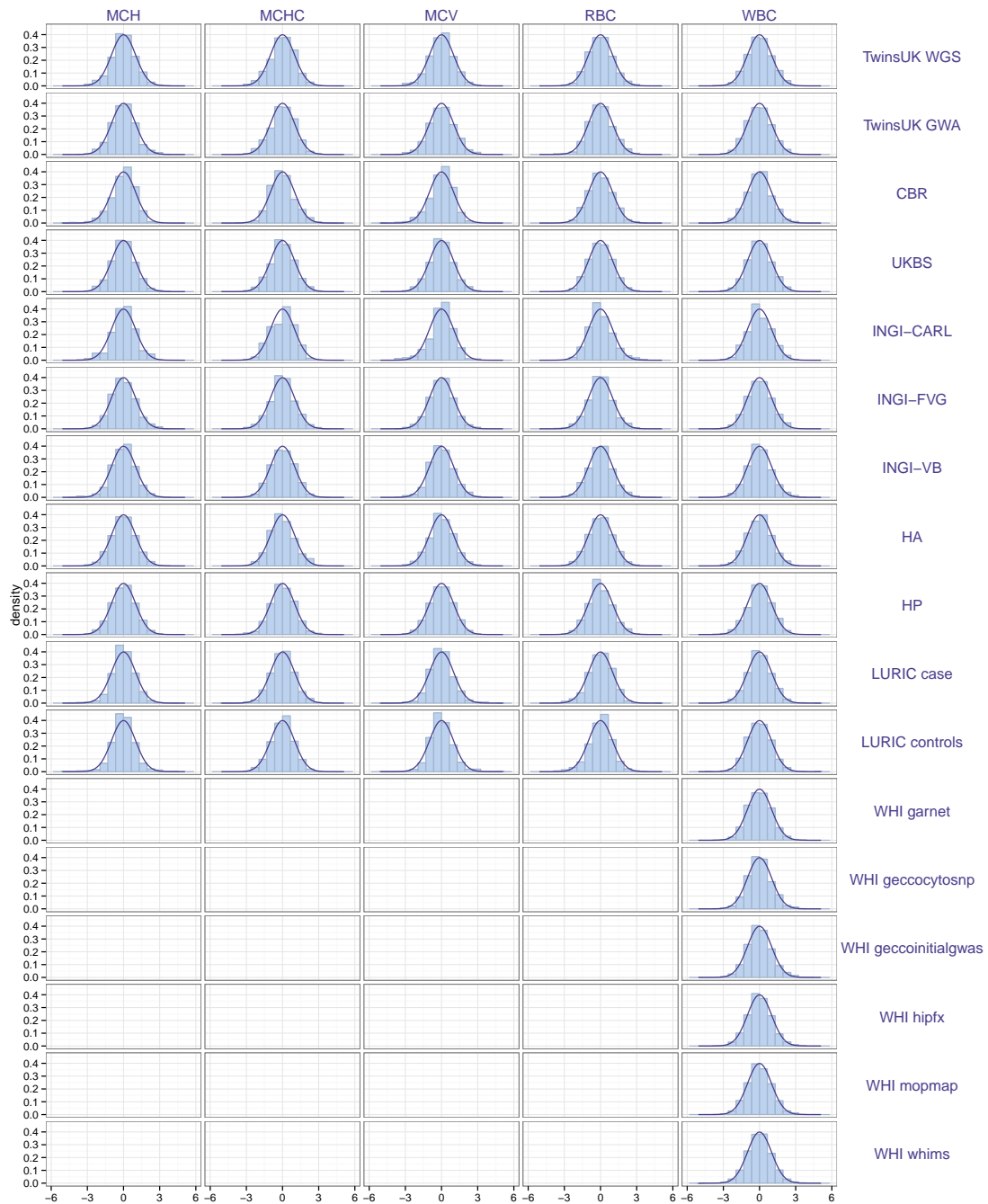
Locus 48, Glucose

**Supplementary Figure 3**: **Fine mapping and annotation.** Panels show the regional association locuszoom plot, the PP statistics from the fine-mapping methods, the CATO and deltaSVM scores, VEP genic annotations and overlap of regulatory annotations found significant (coloured in blue) from GARFIELD enrichment analysis. Circles sizes and colours for all scores have been scaled with respect to score type (i.e. PP, CATO or deltaSVM) and numbers have been plotted below each circle.

**Supplementary Figure 4**: **Histograms of residuals for non-inverse normalized traits.** A panel is shown for each trait-cohort pair, where the trait histogram is shown in light blue and the (assumed) standard normal distribution density is shown in purple.