

# Quality control and conduct of genome-wide association meta-analyses

Thomas W Winkler<sup>1</sup>, Felix R Day<sup>2</sup>, Damien C Croteau-Chonka<sup>3,4</sup>, Andrew R Wood<sup>5</sup>, Adam E Locke<sup>6</sup>, Reedik Mägi<sup>7</sup>, Teresa Ferreira<sup>8</sup>, Tove Fall<sup>9,10</sup>, Mariaelisa Graff<sup>11</sup>, Anne E Justice<sup>11</sup>, Jian'an Luan<sup>2</sup>, Stefan Gustafsson<sup>9</sup>, Joshua C Randall<sup>12</sup>, Sailaja Vedantam<sup>13–15</sup>, Tsegaselassie Workalemahu<sup>16</sup>, Tuomas O Kilpeläinen<sup>17</sup>, André Scherag<sup>18,19</sup>, Tonu Esko<sup>7,13–15</sup>, Zoltán Kutalik<sup>20–22</sup>, Iris M Heid<sup>1,27</sup>, Ruth J F Loos<sup>23–25,27</sup> & the Genetic Investigation of Anthropometric Traits (GIANT) Consortium<sup>26</sup>

<sup>1</sup>Department of Genetic Epidemiology, Institute of Epidemiology and Preventive Medicine, University of Regensburg, Regensburg, Germany. <sup>2</sup>Medical Research Council (MRC) Epidemiology Unit, Institute of Metabolic Science, Addenbrooke's Hospital, Cambridge, UK. <sup>3</sup>Department of Genetics, University of North Carolina, Chapel Hill, North Carolina, USA. <sup>4</sup>Channing Division of Network Medicine, Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts, USA. <sup>5</sup>Genetics of Complex Traits, University of Exeter Medical School, University of Exeter, Exeter, UK. <sup>6</sup>Department of Biostatistics and Center for Statistical Genetics, University of Michigan, Ann Arbor, Michigan, USA. <sup>7</sup>Estonian Genome Center, University of Tartu, Tartu, Estonia. <sup>8</sup>Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, UK. <sup>9</sup>Department of Medical Sciences, Molecular Epidemiology and Science for Life Laboratory, Uppsala University, Uppsala, Sweden. <sup>10</sup>Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden. <sup>11</sup>Department of Epidemiology, University of North Carolina, Chapel Hill, North Carolina, USA. <sup>12</sup>Wellcome Trust Sanger Institute, Cambridge, UK. <sup>13</sup>Divisions of Endocrinology and Genetics and Center for Basic and Translational Obesity Research, Boston Children's Hospital, Boston, Massachusetts, USA. <sup>14</sup>Program in Medical and Population Genetics, Broad Institute, Cambridge, Massachusetts, USA. <sup>15</sup>Department of Genetics, Harvard Medical School, Boston, Massachusetts, USA. <sup>16</sup>Department of Nutrition, Harvard School of Public Health, Boston, Massachusetts, USA. <sup>17</sup>The Novo Nordisk Foundation Center for Basic Metabolic Research, Section of Metabolic Genetics, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark. <sup>18</sup>Institute for Medical Informatics, Biometry and Epidemiology (IMIBE), University Hospital of Essen, University of Duisburg-Essen, Essen, Germany. <sup>19</sup>Clinical Epidemiology, Integrated Research and Treatment Center, Center for Sepsis Control and Care (CSCC), Jena University Hospital, Jena, Germany. <sup>20</sup>Department of Medical Genetics, University of Lausanne, Lausanne, Switzerland. <sup>21</sup>Institute of Social and Preventive Medicine (IUMSP), Centre Hospitalier Universitaire Vaudois (CHUV), Lausanne, Switzerland. <sup>22</sup>Swiss Institute of Bioinformatics, Lausanne, Switzerland. <sup>23</sup>The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, New York, USA. <sup>24</sup>The Mindich Child Health and Development Institute, Icahn School of Medicine at Mount Sinai, New York, New York, USA. <sup>25</sup>The Genetics of Obesity and Related Metabolic Traits Program, Icahn School of Medicine at Mount Sinai, New York, New York, USA. <sup>26</sup>A full list of members is available in the **Supplementary Note**. <sup>27</sup>These authors jointly supervised this work. Correspondence should be addressed to R.J.F.L. ([ruth.loos@mssm.edu](mailto:ruth.loos@mssm.edu)) or I.M.H. ([iris.heid@klinik.uni-regensburg.de](mailto:iris.heid@klinik.uni-regensburg.de)).

Published online 24 April 2014; doi:10.1038/nprot.2014.071

**Rigorous organization and quality control (QC) are necessary to facilitate successful genome-wide association meta-analyses (GWAMAs) of statistics aggregated across multiple genome-wide association studies. This protocol provides guidelines for (i) organizational aspects of GWAMAs, and for (ii) QC at the study file level, the meta-level across studies and the meta-analysis output level. Real-world examples highlight issues experienced and solutions developed by the GIANT Consortium that has conducted meta-analyses including data from 125 studies comprising more than 330,000 individuals. We provide a general protocol for conducting GWAMAs and carrying out QC to minimize errors and to guarantee maximum use of the data. We also include details for the use of a powerful and flexible software package called EasyQC. Precise timings will be greatly influenced by consortium size. For consortia of comparable size to the GIANT Consortium, this protocol takes a minimum of about 10 months to complete.**

## INTRODUCTION

### Background

The genome-wide association study (GWAS) approach has been extremely successful in pinpointing the association of common genetic variants with diseases or disease-related quantitative phenotypes<sup>1,2</sup>. However, given the small sizes of the expected effect under a polygenic model, individual GWASs are generally too small to provide the necessary power to detect single-nucleotide polymorphism (SNP) associations while accounting for the multiple number of independent tests. Therefore, the genetics community has widely adopted the approach of combining summary statistics from multiple GWASs into a single meta-analysis to increase the statistical power of the analysis by augmenting the effective sample size<sup>3,4</sup>. These GWAMAs collate data from GWASs conducted around the world and thus require an enormous organizational effort to ensure effective communication, standardization of analytical procedures and coordination at both the study-specific level and the meta-analysis level, followed by rigorous QC during the meta-analysis process. Although a QC protocol for individual GWASs has been described before<sup>5</sup>,

a comprehensive protocol describing state-of-the-art procedures to conduct and perform QC of large-scale GWAMAs is currently lacking.

The typical GWAMA approach is to design a standardized analysis plan centrally and share it with the individual study partners who will perform the GWAS according to the designated analysis plan. More specifically, the study analysts conduct study-specific GWA QC for each SNP, and they impute the genome-wide SNP array data. Next, they compute association statistics for each SNP, including effect size estimates with standard errors (or odds ratios with corresponding confidence intervals for case-control samples), allele frequencies, sample size, and *P* values, and they provide these summary statistics to the meta-analyses centers. Typically, data on the individual participants, alongside phenotype and genome-wide SNP genotype information, are not shared in order to guarantee anonymity of study participants and to conform to strict data-sharing policies. The unavailability of individual participant data at the meta-analyses centers creates unique analytical challenges for QC, requiring specific statistical

and graphical tools to track errors in the study-specific analysis from the available aggregated data.

Study-specific data issues that need to be detected at the meta-analysis stage include file-naming errors (e.g., female-specific files labeled as male-specific), erroneous SNP genotype data (e.g., flipped alleles, duplicate SNPs and bad imputation quality), and association issues stemming from incorrect analysis models (e.g., improper model adjustments, population stratification and unaccounted relatedness of individuals). Although some errors impede the meta-analysis (e.g., file formatting errors), others (e.g., incorrect trait transformations and flipped alleles) limit the full contribution of a study to the meta-analysis and thus lower the power of the meta-analysis or inflate the number of false positives (type I errors, e.g., unaccounted population stratification). Issues that inflate the number of type I errors should be avoided with higher priority than issues that increase the number of false negatives (type II errors), which negatively affect the statistical power of the meta-analysis. False positives could set researchers onto the wrong track, leading them to spend time and money on misguided follow-up studies, whereas missed genetic signals can be expected to emerge in a subsequent, larger GWAMA.

A typical GWAMA involves two stages: (i) a discovery stage, in which meta-analyzed GWAS data are used to select promising variants; and (ii) a follow-up stage, in which analyses are performed on data derived either from *de novo* genotyping or from existing genome-wide data (*in silico*). This protocol focuses on the discovery stage. Although *in silico* follow-up data can generally be treated similarly to discovery GWAS data for QC purposes, *de novo* genotyped data need to be checked with a particular focus on SNP strand issues, call rate, Hardy-Weinberg equilibrium (HWE)<sup>5</sup> or other technical steps related to the particular genotyping technology applied.

In recent years, GWAMAs have become more and more complex. First, GWAMAs can extend from simple analysis models to more complex models including stratified<sup>6</sup> and interaction<sup>7,8</sup> analyses. Second, beyond imputed genome-wide SNP arrays, new custom-designed arrays such as Metachip<sup>9</sup>, ImmunoChip<sup>10</sup> and ExomeChip<sup>11</sup> are increasingly integrated into meta-analyses. Because of differing SNP densities, strand annotations, builds of the genome and the presence of low-frequency variants, data from such arrays require additional processing and QC steps (also outlined in this protocol by using the example of the Metachip). Finally, GWAMAs involve an ever-increasing number of studies. Up to a 100 studies were involved in recent GWAMAs<sup>12–17</sup>, often involving 1,000–2,000 study-specific files. Increasing the scale and complexity of GWAMAs increases the likelihood of errors by study analysts and meta-analysts, underscoring the need for more extensive and automated GWAMA QC procedures.

We present a pipeline model that provides GWAMA analysts with organizational instruments, standard analysis practices and statistical and graphical tools to carry out QC and to conduct GWAMAs. The protocol is accompanied by an R package, EasyQC, a flexible, user-friendly software that implements this GWAMA QC pipeline and can accommodate additional and alternative steps.

### Development of the protocol

Our protocol was developed by analysts from the GIANT Consortium, which is one of the largest global collaborations to

study complex traits and diseases, currently including up to 125 studies into the meta-analysis. Established in 2006, GIANT has accumulated a lot of experience with GWAMAs. Four rounds of analyses have already been conducted, with each round incorporating new studies and chip technologies<sup>13,15,18–20</sup>. Our work illustrates the increasing complexity of GWAMAs: we deal with multiple phenotypes (e.g., height, body mass index (BMI), waist-hip ratio (WHR), waist and hip circumference (WC and HIP), the latter three also with adjustment for BMI ( $WHR_{adjBMI}$ ,  $WC_{adjBMI}$  and  $HIP_{adjBMI}$ ) and body fat percentage), multiple SNP platforms (genome-wide SNP and Metachip arrays), multiple analysis models (without and with adjustment for BMI, interaction with smoking status and with physical activity, sex- and age-stratified analyses and various dichotomizations of the BMI distribution<sup>6,21</sup>), including imputed and unimputed SNP data, and an ever-increasing number of studies per meta-analysis (16 initially and up to 125 in the current analyses). Our ongoing analyses include more than 1,500 GWAS input files, necessitating an efficient QC pipeline. The size and experience of the GIANT Consortium provides an ideal basis for the development of a GWAMA protocol. The protocol and tools can readily be applied by other consortia using aggregated statistics for meta-analysis, studying other quantitative traits and using other statistical models or other genotyping platforms. We have incorporated all QC steps that proved to be helpful during our GIANT work and that have been known to be efficient in the work of other consortia. We have also developed special tools to conduct meta-level QC and to handle the particularly large number of files.

### Limitations

First, this protocol has been developed for human genomic data. Although some aspects can be applied to non-human data, a detailed protocol for other species is beyond the scope of this protocol.

Second, even a perfect protocol for the meta-analysis of aggregated statistics cannot fully compensate for not having access to individual participant data, which would guarantee standardized QC and analyses across studies. Advantages and disadvantages of meta-analyses using individual participant data are summarized in the ‘Comparison with other approaches’ section below. However, ethically motivated restrictions to sharing genome-wide genotype and phenotype data currently limit the realization of individual participant GWAMAs, which is the reason why the aggregated-statistics GWAMA—as described here—is currently the most widely applied approach.

### Applications of the protocol

Generally, this protocol assumes that the study analysts have quality-controlled their study data regarding phenotype and genotype, as well as accounted for ethnicity, race and familial relatedness. For these steps, there are standardized protocols available<sup>5</sup>. It also assumes that they have imputed their genome-wide SNP array data—ideally with a prespecified common reference panel—to ensure a common SNP panel across all studies or that they have data from an unimputed custom genotype array available.

This protocol specifically focuses on the discovery stage of a GWAMA, but it can be readily applied to the follow-up stage as well. Imputed *in silico* follow-up data can be treated in a similar way, as the imputed genome-wide SNP array data, nonimputed

*in silico* or *de novo* genotyped data described here can be treated like the Metachip data with regard to the cleaning of call rate, HWE and strand issues.

Although this protocol has been developed for quantitative phenotypes and HapMap-imputed or typed common autosomal genetic variants, it can be extended to 1000 Genomes-imputed variants, dichotomous phenotypes, rare variants, gene-environment interaction (G×E) analyses and to sex-chromosomal variants. A summary of directly applicable protocol steps or steps requiring adaptation is given in **Table 1**. As 1000 Genomes-imputed data extend to a larger SNP panel and include structural variants (SV) and insertions or deletions (indels), the allele coding and harmonization of marker names require special considerations: (i) additional allele codes (other than 'A', 'C', 'G' or 'T') are needed for indels and SVs (e.g., 'I' and 'D' for insertions and deletions); and (ii) to account for the fact that some SVs and indels map to the same genomic position as SNPs, the identifier format 'chr<chromosome>:<position>' would introduce duplicates. Therefore, the identifier format needs to be amended (e.g., to 'chr<chromosome>:<position>:[snp|indel]', which adds the type to the format).

For dichotomous traits, the effective sample size needs to be computed by  $N_{\text{eff}} = 2/(1/N_{\text{Cases}} + 1/N_{\text{Controls}})$ , an expression that balances the number of cases with the number of controls. Custom-array data require checks of genotype quality per case status. The analysis is usually performed by using logistic instead of linear regression, providing beta estimates and standard errors that enable the implementation of the same meta-analysis methods. The minor allele count (MAC) cutoff requires more consideration: it depends on the logistic regression-based test used and on the ratio between the number of cases and controls<sup>22</sup>.

For rare and low-frequency variants, more refined considerations regarding the minimal sample size or the minimally acceptable MAC cutoff per file are required. The comparability of the study frequencies with reference data such as HapMap or 1000 Genomes is of limited use, as Exomechip or custom-made chips focusing on rare variants and low frequency tend to include novel or population-specific variants. Often, the single-variant analyses are complemented by gene-based burden tests requiring special consideration. For single-variant analyses, most of the protocol steps described herein are directly applicable.

Results for analyses models that include an interaction term can also be quality-controlled by this protocol. The main SNP effect estimates can be treated like the SNP effects without interaction. The interaction effect estimates need to be cleaned and meta-analyzed in addition. This objective can be achieved in the same manner as the main effect estimates or by implementing alternate methods<sup>23</sup>. As the analysis of the interaction between SNP and the environment is more and more included into GWAMA efforts, this approach will be of increased importance.

Analyses with sex-chromosomal variants require some special considerations, especially in men. We assume that study partners have quality-controlled their data regarding rare gonosomal aberrations (X0, XXX and XYY). The potential errors in coding variants in men include differences in the coding of X-chromosomal variants (either 0|1 or 0|2 for men) or erroneous coding of pseudo-autosomal variants (should be 0|1|2). Separating the QC by X-, Y- and pseudo-autosomal variants in men can be grasped by deflated or inflated beta estimates (and

thus standard errors) in the SE-N (i.e., inverse of the median standard error versus the square root of the sample size) plot. Generally, sex-chromosomal variants should be cleaned and analyzed in men and women separately.

### Comparison with other approaches

Over the past 6 years, >100 large phenotype-driven consortia of genetic association studies have emerged<sup>1</sup>. Most of these consortia follow a similar framework for QC and data 'sanity checks,' as outlined here<sup>24</sup>.

Some consortia, such as the Uric Acid (UA) Consortium, follow slightly modified procedures, whereby study-specific QC metrics, generated by GWASToolbox<sup>25</sup>, were collected next to summary-level association statistics<sup>26</sup>. This approach enables the easy detection of basic data problems even before the results are shared, but at the same time it poses an extra burden on the analysts, and its implementation does not help the necessity of meta-level checks. The Chronic Kidney Disease Genetics (CKDGen) Consortium omits filtering data on the basis of poor imputation quality<sup>27</sup>, whereas most consortia, including GIANT, delete badly imputed variants from the meta-analysis.

Whereas most GWAMAs meta-analyze study-specific statistics, where study analysts have provided GWAS results to the meta-analysis center, the Psychiatric Genomic Consortium (PGC) conducts a meta-analysis of individual participant data, as both the individual-level genotype and phenotype data of all participating studies are deposited centrally<sup>28</sup>. This approach has the following advantages: (i) central QC: genotype and phenotype data can be modeled and quality-controlled centrally, eliminating the need for subsequent troubleshooting; (ii) standardized study-specific analyses: fewer analysts are involved and the utilization of the same imputation and association analysis software is guaranteed; and (iii) flexibility: more complex and comprehensive statistical analyses can be conducted without burdening a large number of study analysts. However, our GWAMA approach has also advantages compared with the meta-analysis of individual participant data: (i) gathering experts: the more analysts involved, the more the network can profit from the accumulated expertise; (ii) local know-how: local study analysts know their study better than a central team of meta-analysts; and (iii) compliance: ethically motivated restrictions may limit the sharing of genome-wide genotype and phenotype data owing to the risk of participant identification may inhibit the study contribution<sup>29–31</sup>. In summary, the framework presented in this protocol reflects the currently most widely applied GWAMA conduct and QC approach.

### Experimental design

**Organizational aspects of the conduct of a typical GWAMA (Steps 1–6).** The typical GWAMA starts with setting up logistics aimed at achieving a smooth communication between participating partners, analysts and principal investigators, which limits the burden for study analysts so as to ensure a timely delivery of results to the meta-analysis team.

Once the study partners have been identified, general rules for the collaboration can be issued in a 'memorandum of understanding' to set out the guidelines of confidentiality, data access, publication of results and authorship. Subsequently, collaborators and analysts are invited to join task groups and regular teleconference calls.

**TABLE 1** | Expandability of the protocol to 1000 Genomes imputed data, dichotomous traits, rare variants, SNP × E interactions, and X-chromosomal variants.

Procedure steps	Step no.	1000 Genomes	Dichotomous trait	Rare variant analyses	SNP × E interaction	Analyses of the sex chromosomes
Setting up logistics of meta-analysis	1–3	DA	DA	DA	DA	DA
Collecting aggregated statistics per study	4–6	DA	DA	DA	DA	DA
File-level QC	7–18	AA, allow for indels and SVs; adjust SNP name harmonization	AA, calculate $\beta = \ln(\text{OR})$ ; filter on effective $N_e$ , adjust MAC cutoff	AA, adjust MAC cutoff	AA, add checks on $\beta_{\text{gxe}}$	AA, extra checks for pseudo-autosomal variants in men
Identification of analytical issues by the SE-N plot	19,20	DA	DA	AA, adjust $c$	AA, add checks on $\beta_{\text{gxe}}$	AA, separately for men and women, extra considerations of coding errors in male X and Y
Identification of analytical issues by the P-Z plot	21,22	DA	DA	DA	AA, add checks on $\beta_{\text{gxe}}$	DA
Identification of problems with allele frequencies or strand	23,24	AA, use 1000 Genomes allele frequencies as reference	AA, limit checks to control group	AA, update allele frequency reference	DA	AA, separately for men and women
Identification of population stratification	25,26	DA	DA	DA	AA, add $\lambda_{\text{GC}}$ from $P_{\text{gxe}}$	AA, use autosomal variant to compute $\lambda_{\text{GC}}$
Meta-analysis	27,28	DA	DA	DA for single variant analyses	AA, add meta-analysis of beta G × E	DA
Meta-analysis QC—compare results from two analysts	29,30	DA	DA	DA	DA	DA
Meta-analysis QC—identify analytical issues by calculating the study level $\lambda_{\text{GC}}$	31,32	DA	DA	DA	AA, add $\lambda_{\text{GC}}$ from $P_{\text{gxe}}$	AA, use autosomal variant to compute $\lambda_{\text{GC}}$
Finalizing meta-analysis	33	DA	DA	DA	DA	DA

Abbreviations:  $\beta$ , SNP effect on outcome;  $\beta_{\text{gxe}}$ , SNP × E interaction effect; AA, applicable with adaptation; DA, directly applicable; E, environment; OR, odds ratio;  $P$ , association  $P$  value;  $P_{\text{gxe}}$ , SNP × E interaction  $P$  value; indels, insertions or deletions; SVs, structural variants.

An analysis plan is designed centrally by the meta-analysts to describe the standardized analyses to be performed ‘locally’ and to detail phenotype transformation (e.g., to deal with non-normal phenotype distributions and to enable comparability across studies), genotype handling, imputation requirements and

association analysis methods (statistical model, adjustment and stratification). Where possible and reasonable, software scripts are provided to every participating study group to minimize the potential of errors and to alleviate the analysis burden for the study analyst. The analysis plan also defines the required



## Box 1 | Study-specific GWAS results; columns as requested by GIANT

Stated are the columns requested by GIANT from the study partners for each GWAS to ensure uniform study-specific files:

'MarkerName': character string; the SNP identifier of the marker analyzed

'Strand': a single character '-' or '+'; strand on which the alleles are reported

'Chr': character; chromosome

'Pos': integer; base position of the SNP

'N': positive integer; the effective number of subjects analyzed

'Effect\_allele': a single upper case character 'A', 'C', 'G' or 'T'; the allele associated with phenotypic traits (corresponding to change in beta estimates)

"Other\_allele": a single upper-case character 'A' 'C' 'G' or 'T'; indicating the other (non-effect) allele

"EAF": numeric; effect allele frequency (range 0–1)

"BETA": numeric; estimate of the effect size

"SE": numeric; estimated standard error on the estimate of the effect size

"P": numeric; significance of the variant association, uncorrected for genomic control.

### Only for genotyped data

"P\_HWE": numeric; Exact HWE *P* value for the sample analyzed

"Callrate": numeric; Call rate for this SNP across all subjects. Perfectly genotyped (100%) data will have a Callrate = 1.000

### Only for imputed data

'Information\_type': integer; code indicating the type of data in the 'Information' column (i.e., the type of the imputation and analysis software used):

0: if the SNP was not tested by using imputation or genotyping uncertainty, in which case the following column 'Information' should be missing (e.g., for directly genotyped SNPs)

1: for a MACH-imputed SNP, whereas the following column 'Information' either contains 'r2\_Hat' from MACH2DAT/MACH2QTL OR 'INFO' from PLINK (if have used PLINK for the association with MACH-imputed SNP data)

2: if the following column 'Information' contains 'proper\_info' from SNPTTEST

3: for a PLINK-imputed SNP, i.e., the following column 'Information' contains 'INFO' from PLINK (if the SNP was imputed using PLINK as well)

4: if the following 'Information' column contains 'rSqHat' from QUICKTEST

'Information': numeric; A value (range 0–1; PLINK values can exceed 1) corresponding to the information content output from the association testing (according to the data type specified in the 'InformationType' column above)

aggregated association statistics (e.g., SNP identifier, effect allele, allele frequency, beta estimate, standard error, sample size, call rate or imputation quality and *P* value) and details the format in which they need to be submitted (**Box 1**). In the design of the analysis plan, the decision regarding whether or not to provide detailed and lengthy guidelines, possibly including even software codes, needs to be weighed against providing a short and comprehensive—but potentially more error-prone—description. The less-standard the requested analyses, the more details that need to be provided. A general analysis plan format cannot be provided, but the GIANT analysis plan can serve as an example that has worked and that has been improved through several rounds of meta-analyses (**Supplementary Manual**). The analysis plan is discussed with the study collaborators and then sent out to each study analyst, including a deadline and server access details for data upload.

When data from all studies have been uploaded to a password-secured file server, a data freeze ensures the integrity of the data for all meta-analysts, regardless of download time (**Supplementary Fig. 1**).

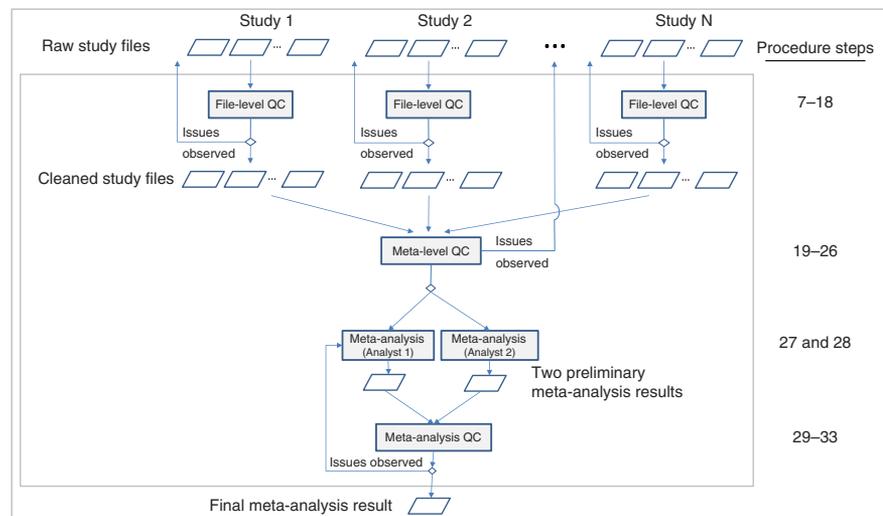
The complete turnaround time for consortia comparable in size to GIANT (>100 studies in meta-analysis) is, at minimum, ~10 months: 2 months to set up the logistics and to develop the analysis plan, 2 months to collect the data after the analysis plan has been sent out and 6 months to perform QC and meta-analysis.

**QC workflow.** The workflow involves three QC steps: file-level QC (Steps 7–18), meta-level QC (Steps 19–26) and meta-analysis QC (Steps 29–32). The file-level QC tackles formatting issues that can be checked independently on each study file. In the meta-level QC, the study-specific statistics are compared across studies or with reference panels to detect errors in the analyses that cannot be identified by examining the study files individually. The meta-analysis QC works on the level of already aggregated meta-analysis results and helps remove or flag suspicious SNP results. The workflow and the three QC steps are presented in **Figure 1**.

**File-level QC (Steps 7–18).** This stage involves 'cleaning' (deleting poor quality data) and 'checking' (providing summaries to judge data quality) data. Thresholds for what data to remove are typically defined *a priori* (e.g., by this protocol). Although data checking should ascertain that there are no issues left, it often reveals further issues, which require recleaning and rechecking. A few QC iterations may be needed before all files are fully cleaned and ready for meta-analyses. Which SNPs or study files are to be removed depends on how much the improvement in data quality weighs against loss of data. On the one hand, the stricter the QC, the more SNPs or study files are removed and thus the lower the coverage or sample size (and thus power). On the other hand, the more relaxed the QC requirements, the larger the coverage and sample size at the expense of data quality, which also decreases power.



**Figure 1** | Workflow of the QC and the meta-analysis. A typical GWAMA includes four major stages. The first stage is file-level QC (Steps 7–18), which includes the QC of each study file to ensure validity. This stage involves file cleaning (e.g., adjustments of column headings, file format changes, SNP exclusions based on certain criteria or adding columns) and file checks (e.g., checking overall characteristics of the file or the number of SNP exclusions), usually in an iterative manner. Typically, this task is divided by study among analysts of the meta-analysis team. Files that pass the file-level QC are labeled as ‘CLEANED’. Any issues observed with particular files should be clarified with the respective study analyst directly. Second, the meta-level QC (Steps 19–26) addresses the comparison of file-specific statistics across files in order to depict study-specific issues that are yet undetected. In case issues of specific studies cannot be resolved centrally, the relevant study analyst should be contacted for clarification. Third, meta-analysis (Steps 27 and 28) is the stage at which the meta-analysis is actually conducted, a task typically performed by two analysts independently. Finally, meta-analysis QC (Steps 29–33) involves checking the meta-analysis results and includes the comparison of the two meta-analyses performed by the different analysts and the QC of the meta-analysis result.



In case issues of specific studies cannot be resolved centrally, the relevant study analyst should be contacted for clarification. Third, meta-analysis (Steps 27 and 28) is the stage at which the meta-analysis is actually conducted, a task typically performed by two analysts independently. Finally, meta-analysis QC (Steps 29–33) involves checking the meta-analysis results and includes the comparison of the two meta-analyses performed by the different analysts and the QC of the meta-analysis result.

Clearly, monomorphic SNPs or SNPs with missing (e.g., missing *P* value, beta estimate or alleles) or nonsensical information (e.g., alleles other than A, C, G or T, *P* values or allele frequencies >1 or <0 or standard errors = 0, infinite beta estimates or standard errors) are of no help to the meta-analysis and need to be removed. Systematically missing values or errors can point toward analysis problems; thus, such data calls into question the correctness of the data and should be discussed with the study analyst. A large number of monomorphic SNPs can also point to study-specific array problems.

If a study includes a low number of individual participants, its summary statistics can be unstable (e.g., zero or infinite standard errors, zero *P* values or extremely large beta estimates), which might drive the meta-analysis toward detecting false positives. This risk pertains particularly to low-frequency variants. The detection of false positives due to the low statistical power of the meta-analysis can be avoided by requiring a minimum sample size per study and a minimum number of minor alleles contributing to a SNP for each participating study. For example, in meta-analyses performed by the GIANT Consortium, SNPs were removed from the study file if the number of individuals informative for the SNP was <30 or the MAC was (computed as  $2 \times \text{MAF} \times N$ , with MAF being the minor allele frequency)  $\leq 6$ .

Imputed genotype data are often filtered on the basis of the imputation quality. For example, in the GIANT Consortium,

poorly imputed SNPs were removed according to a threshold that depended on the imputation method and on the imputation quality metric (Table 2). Arguably, however, SNPs with poor imputation quality can be retained in the meta-analysis<sup>27</sup>: on the one hand, a badly imputed SNP can be considered a random, nondifferential error in the genotype (i.e., not systematically prioritizing one genotype and independent of the phenotype), and thus it will not tend to create a false signal and, on the other hand, a study with the SNP badly imputed will neither contribute to a true signal nor mask it. Filtering poorly imputed SNPs has the advantage that no nonsensical results will unduly decrease the statistical significance of truly informative data.

Sex-chromosomal and autosomal SNPs require different genotype models, and therefore they are often studied separately from each other. To focus on autosomal SNPs and consistent genotype models across studies in its analyses, the GIANT Consortium has removed any sex-chromosomal SNPs.

SNP identifiers often differ between arrays and/or imputation reference panels and, therefore, they often differ between studies. Their harmonization across studies is pivotal to the meta-analysis. For example, a SNP that is assigned to two different SNP identifiers (e.g., rs123 in half of the studies and rs17614680 in the other half) will appear as two different SNPs in the meta-analysis output, with the total sample size split across the two SNPs; a true signal might, therefore, be missed because of loss of statistical

**TABLE 2** | Imputation quality metrics for different combinations of imputation and analysis software packages as observed in GIANT.

		Association software				
			PLINK	SNPtest (–expected)	QUICKtest	Other
Imputation	MACH	0.3 (r2_hat)	0.3 (INFO)	0.4 (proper_info) <sup>a</sup>	0.3 (rSqHat)	0.3 (rSqHat)
Software	IMPUTE	–	–	0.4 (proper_info) <sup>a</sup>	0.3 (rSqHat)	–
	PLINK	–	0.8 (INFO)	–	–	–

<sup>a</sup>Newer versions of SNPtest output a column called ‘info’ instead of ‘proper\_info’.



power. For HapMap-imputed studies, a unique SNP identifier can be generated by combining the SNP's genetic positions to generate the format 'chr<chromosome>:<position>'. However, for some arrays (e.g., MetaboChip), not all SNPs map to a standard reference panel. In such cases, the DNA probe sequences need to be mapped to the reference genome build of interest to arrive at a common chromosome and position, which can then be used to generate the SNP identifier. This procedure will also remove SNPs that do not map uniquely to the genome. Maps with unique SNP identifiers and genomic positions (for several different genome builds) for several commercial arrays are freely available for download (<http://www.well.ox.ac.uk/~wrayner/strand/>).

**Meta-level QC (Steps 19–26).** This stage consists of the cross-study comparison of statistics to identify study-specific problems. This QC stage compensates for not having the individual participant data of each study available to the meta-analyst. We recommend that the following plots be included in the GWAMA QC protocol.

*The SE-N plot (Steps 19 and 20).* Several types of analytical problems can be identified by depicting, for each study file, the inverse of the median standard error of the beta estimates across all SNPs against the square root of the sample size. The inverse proportionality between the median standard error and the square root of the sample size derives from the fact that the sampling variance of a linear regression–derived beta estimate of a specific SNP  $j$  depends on the variance of the phenotype,  $\text{Var}(Y)$ , the variance of the SNP genotype,  $\text{Var}(X_j)$ , and the sample size  $N_j$ :

$$SE_j^2 = \text{Var}(\beta_j) = \frac{\text{Var}(Y)}{N_j \cdot \text{Var}(X_j)}$$

If the regression model is adjusted, then  $\text{Var}(Y)$  reflects the variance of the residuals. Thus, the average of the standard errors across all SNPs will reflect the sample size. Assuming that the sample size for a given SNP is close enough to the maximum sample size for all SNPs,  $N_j = N$ , the median of the standard errors across all  $m$  SNPs ( $j = 1 \dots m$ ) can be written as median  $(SE_j) = (\sqrt{\text{Var}(Y)})/\sqrt{N} \times \text{median}(1/\sqrt{\text{Var}(X_j)})$  and therefore

$$c\sqrt{\text{Var}(Y)} \cdot \frac{1}{\text{median}(SE_j)} = \sqrt{N} \quad (1)$$

with

$$c = \text{median}\left(\frac{1}{\sqrt{\text{Var}(X_j)}}\right)$$

The constant  $c$  can be computed per study file incorporating the genotype frequencies (for genotyped variants) or the genotype dosages and imputation quality (for imputed variants), and it will depend on the individuals' ethnicity, genotyping platform, imputation reference panel and imputation quality. By ignoring the uncertainty from the imputation,  $c$  can be approximated by

$$c \sim \text{median}\left(\frac{1}{\sqrt{2\text{MAF}_j(1-\text{MAF}_j)}}\right) \quad (2)$$

However, the computation of  $c$  per study is not ideal for comparing the studies with each other. Differences in the MAF distribution

**TABLE 3 | SE-N Plot calibration factors for various genotyping platforms, imputation reference panels, and ethnicities.**

Genotyping platform	Imputation reference panel	Ethnicity	Calibration factor (c)
GWAS chip	HapMap	CEU	1.75
GWAS chip	HapMap	YRI	1.83
GWAS chip	1000 Genomes	ALL	8.86
MetaboChip	–	EUR	1.93
MetaboChip	–	AFR	2.18

The calibration factors were estimated from the publically available HapMap and 1000 Genomes reference data. Only autosomal and non-monomorphic SNPs were used in the estimation. The MetaboChip  $c$  factors were estimated from 179,000 overlapping SNPs from 1000 Genomes reference data frequencies.

between any individual study and the reference would not be detected. For several standard platforms, imputation panels and ethnicities, these approximate  $c$  values to be used in the SE-N plot are given in **Table 3**. For other platforms, panels or ethnicities,  $c$  is to be computed from a reference study or the imputation reference panel.

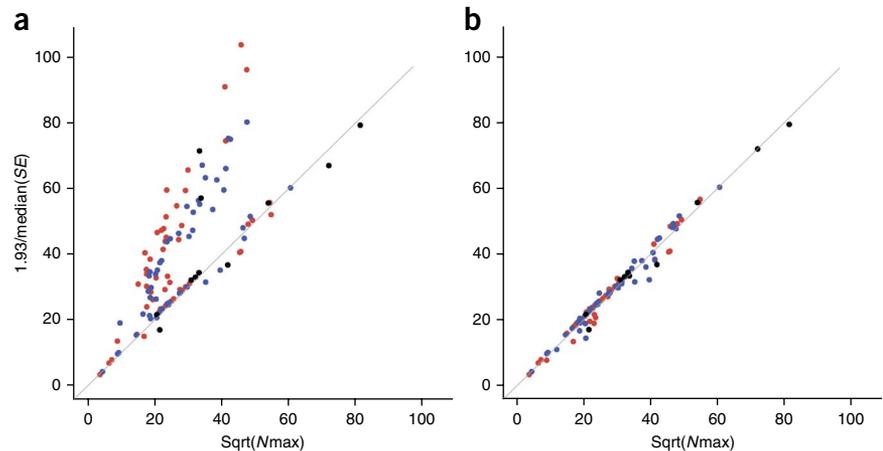
The study-specific data points of the SE-N plot will tend to describe a straight line. However, studies will deviate from the overall trend, if:

- the study's phenotypic variance differs from other studies, which might be explained by a different study design or special study population;
- the study's MAFs differ from other studies, which might be explained by a diverging genotyping platform, reference panel for the imputation or a different ethnicity;
- the study's SNP imputation qualities differ from those of other studies, which might reflect errors in the imputation or a different reference panel;
- the study's effective sample size differs from the stated sample size, which might be due to unaccounted relatedness between study participants or miscoded sample size;
- the study analyst has used a different statistical test; or
- the study analyst has mis-specified the phenotype transformation or the regression model, which results in a different phenotype variance or residual variance (**Fig. 2**; ANTICIPATED RESULTS; **Supplementary Fig. 2**).

*The P-Z plot (Steps 21 and 22).* Analytical problems related to the study-specific computation of beta estimates, standard errors or  $P$  values can also be revealed by a study-specific scatter plot that, for each SNP, compares the reported  $P$  values with the  $P$  values computed from the  $Z$ -statistics based on reported beta estimate and standard error ( $Z$ -statistics =  $\beta_j/\text{SE}(\beta_j)$ ) (**Fig. 3**; ANTICIPATED RESULTS; **Supplementary Fig. 3**).

*The effect allele frequency (EAF) plot (Steps 23 and 24).* Plotting the reported EAFs against a reference set, such as from the HapMap<sup>32</sup> or 1000 Genomes<sup>33</sup> projects, or from one specific study, can help visualize patterns that pinpoint strand issues, allele miscoding or the inclusion of individuals whose self-reported ancestry did not match their genetic ancestry (**Fig. 4**; ANTICIPATED RESULTS). A strand mismatch or allele miscoding may severely reduce statistical power. If, for example, a study (or several studies) reports alleles on the '–' instead of the '+' strand, which cannot be corrected for 'palindromic' A/T or C/G

**Figure 2** | SE-N plots to reveal issues with trait transformations. **(a,b)** SE-N plots to detect issues with trait transformations contrasting the study-specific standard errors with sample sizes for GIANT studies typed on Metabochip and tested for association with  $HIP_{adjBMI}$  ( $N = 81,000$ ). **(a)** Before QC: a number of studies (in fact, the majority of studies) revealed errors by clustering above the identity line. **(b)** After QC: the same plot after having gone back to the relevant study analysts and having resolved all trait-transformation issues. Different colors for the points in the plot indicate men-specific (blue), women-specific (red) or sex-combined (black) association results.



SNPs, a true signal will be diminished, abolished or even reversed. Although comparison of allele frequencies across studies will not detect strand issues or allele miscoding for SNPs with MAF close to 0.5, this comparison will be informative for most SNPs.

**The lambda-N plot (Steps 25 and 26).** Population stratification can either inflate or deflate association  $P$  values and can be grasped by the genomic control (GC) inflation factor ( $\lambda_{GC}$ )<sup>34</sup>. As  $\lambda_{GC}$  increases with sample size in the case of polygenic phenotypes<sup>35</sup>, plotting  $\lambda_{GC}$  versus sample size per study file identifies inflated  $\lambda_{GC}$  and thus potential problems with population stratification (Fig. 5; ANTICIPATED RESULTS). In the GIANT Consortium, analysts of studies with  $\lambda_{GC} > 1.1$  are contacted and asked to revisit their analyses (e.g., adjusting for principal components) and results.

**Meta-analysis and QC of meta-analysis output (Steps 27–32).**

The meta-analysis combines the study-specific association results to obtain an overall estimate of the association and its  $P$  value. The inverse-variance weighted meta-analysis using the fixed-effects model is most commonly used for GWAMAs (e.g., implemented in METAL<sup>36</sup>). The  $Q$ -statistic and  $I^2$  measures test and estimate between-study heterogeneity<sup>22,37</sup>. For SNPs with pronounced heterogeneity ( $I^2 > 75\%$ ), the effect estimation benefits from a random-effects meta-analysis<sup>38</sup>. An alternative approach for deriving overall  $P$  values is the sample size–weighted  $z$ -score meta-analysis<sup>39</sup>. This approach is used when beta estimates or standard errors are not available, or when the meta-analyzed traits are on a different scale (e.g., blood-level data measured in different laboratories or differences in trait transformation) at the cost of losing power.

Meta-analyses are conducted by two meta-analysts independently, each uploading the results and log files onto the server (Supplementary Fig. 1). Results are compared by using (i) the log files that specify the study files included and the meta-analysis parameters set in the software program; (ii) descriptive statistics (min, median and max) of sample size and number of SNPs included in meta-analysis results; and (iii) correlation and scatter plots of  $P$  values. Differences between the two analyses are resolved until agreement is reached.

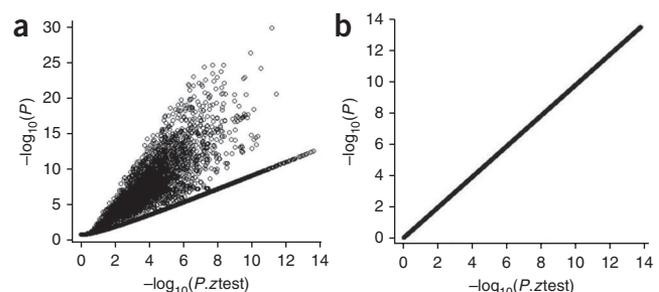
To evaluate whether the statistics of the meta-analyzed effect are inflated owing to population stratification accumulated across studies or owing to unaccounted relatedness, the  $\lambda_{GC}$  is computed for the meta-analysis result (complementing the file-specific  $\lambda_{GC}$  values, see above). A high value ( $\lambda_{GC} > 1.1$ ) might be due to (i) an excess of association signals in large GWAMAs for

highly polygenic traits<sup>14</sup>, (ii) residual population stratification per study file accumulated across studies, (iii) relatedness between individuals across strata (when the study-specific analyses have been performed separately by strata) or (iv) related subjects across studies, which are more likely to occur in very large GWAMAs. In the third case, a meta-analysis across strata per study can be conducted and a study-specific  $\lambda_{GC} > 1.1$  might provide insight into inflation requiring contact of the study analyst. Generally, we recommend applying the  $\lambda_{GC}$  correction at the file level and at the meta-analysis level (double GC correction), but very large GWAMAs ( $> 200,000$  individuals) on highly polygenic traits (e.g., height) may opt to omit the second GC correction (single GC correction)<sup>35</sup>.

Finally, when all issues are resolved, one of the analysts shares the final results with the analysis task group (Supplementary Fig. 1). The final results file will be used for all subsequent steps, including SNP selection for top hit identification and/or follow-up.

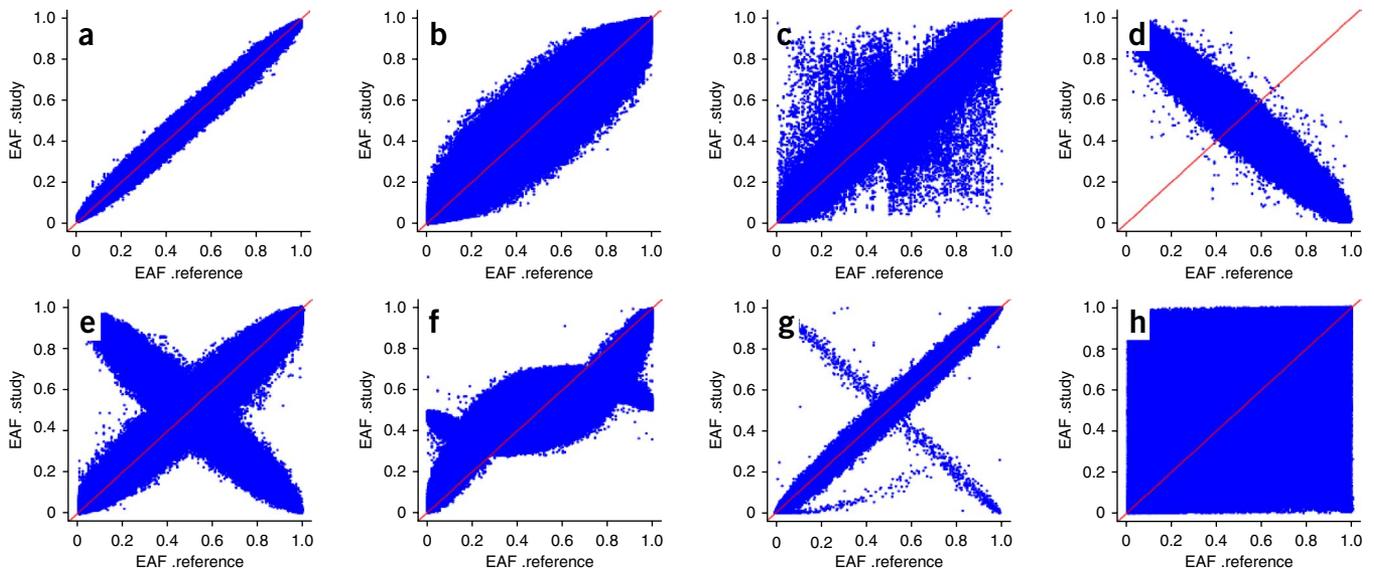
**Special considerations for custom array data instead of genome-wide SNP array data.**

The GIANT Consortium has worked on data genotyped with Metabochip, which is a custom genotyping array that contains ~195,000 replication and fine-mapping SNPs chosen from GWAMAs of metabolic, cardiovascular and anthropometric traits<sup>9</sup>. Although many of the QC steps for HapMap-imputed SNP data can be directly applied to Metabochip and other customized genotype arrays, some steps need to be



**Figure 3** | P-Z plot to reveal analytical issues with beta, standard error and  $P$  values. Plots to reveal issues with beta estimates, standard errors and  $P$  values for **(a)** an uncleaned study file showing severe deviations from the identity line; and **(b)** the cleaned data set showing perfect concordance. The plots compare  $P$  values reported in the association result file with  $P$  values calculated from  $Z$ -statistics ( $P.ztest$ ) derived from the reported beta and standard error from an example GIANT file. The uncleaned study file contained a large number of highly significant but erroneous (reported)  $P$  values.

## PROTOCOL

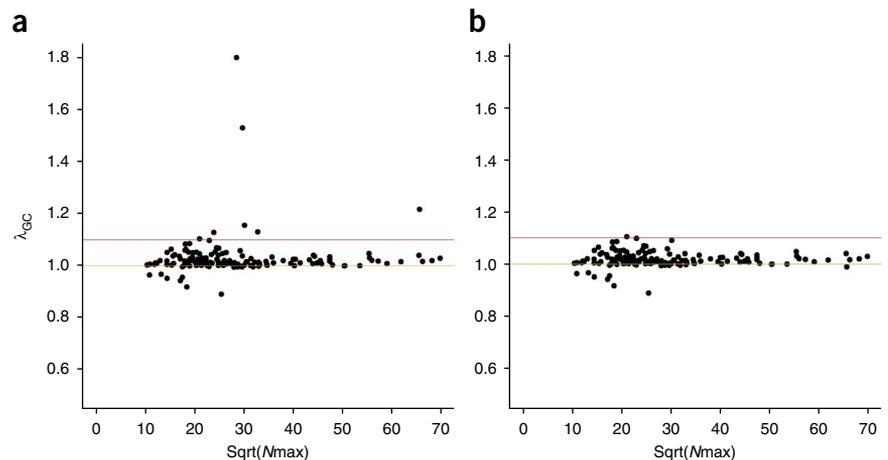


**Figure 4** | Different patterns of allele frequencies in the EAF plot. These different patterns have been observed during the QC checks performed by the GIANT analysts. In the graphs, the observed (study-specific) allele frequencies reported on the y axis are plotted against the expected (HapMap or 1000 Genomes) allele frequencies reported on the x axis. (a–c) These plots represent data from studies in which allele frequencies and strand annotation are correct but participants exhibit different ancestries compared with the reference, which includes mostly samples of European ancestry. (a) A study in which data are relatively consistent with the reference. (b) A study in which participants had slightly different ancestry to the reference, resulting in a thicker band across the diagonal. (c) A study involving participants of non-European ancestry resulting in substantial deviation from the reference. (d–h) These plots pertain to studies with errors in coding the effect allele, the effect allele frequency and/or strand annotation. (d) A study in which the wrong allele was consistently labeled as the effect allele. (e) A study in which a fraction of the effect alleles was mis-specified, e.g., from stating the MAF instead of the effect allele frequency, or from incorrectly assigning the strand owing to data management or wrong strand reference (sometimes specific to ‘palindromic’ SNPs A/T or C/G). (f–h) Studies with other data management or analytical errors in calculating the allele frequencies.

adjusted, which are summarized in the following section and given in the protocol as an alternative route to using HapMap-imputed SNP data: (i) to control genotype quality instead of imputation quality, a filter on call rate and deviation from HWE is required; (ii) to identify strand and allele frequency errors, other references may be required because some genotyped SNPs may not be available in the HapMap reference data; and (iii) to perform GC correction for chips that are designed to cover multiple traits, the calculation of the  $\lambda_{GC}$  needs to be limited to a subset of SNPs that are chosen from a trait that is uncorrelated with the trait of interest. This  $\lambda_{GC}$  is then to be applied to all SNPs on the array. For example, GIANT limits the  $\lambda_{GC}$  computation for Metabochip data to the 4,427 QT-interval SNPs, as the QT-interval is uncorrelated with the GIANT traits, as recommended by the Metabochip designers<sup>9</sup>.

**Software.** By using the standard, open-source and freely available software R<sup>40</sup> and the graphical R package ‘Cairo’, we created

**Figure 5** | Lambda-N plot to reveal issues with population stratification. (a,b) Plot to detect issues with population stratification contrasting the study-specific  $\lambda_{GC}$  with sample sizes for GIANT studies typed on Metabochip and tested for association with  $HIP_{adjBMI}$  ( $N = 81,000$ ). (a) Before QC: a number of studies displayed high  $\lambda_{GC}$  values. (b) After QC: the same plot in a after having gone back to the study analyst and having resolved all issues. The orange line indicates the optimal  $\lambda_{GC} = 1$ . Dots above the red line, which visualizes the threshold  $\lambda_{GC} = 1.1$ , represent problematic studies.



a pipeline for completing this protocol into a downloadable GWAMA-QC R package called EasyQC. We provide code application directly in the procedure steps. The general basic usage is described in **Box 2**. Minimum system requirements are described in the MATERIALS section.

We provide a number of template scripts that enable to conduct multiple procedure steps at once: (i) the EasyQC scripts ‘1\_filelevel\_qc.gwa.ecf’ and ‘1\_filelevel\_qc.metabochip.ecf’ to perform file-level QC (Steps 7–18); (ii) the EasyQC script ‘2\_metalevel\_qc.ecf’ to perform meta-level QC (Steps 19–26); (iii) the METAL script ‘3\_metaanalysis.metal.txt’ to perform the meta-analysis (Steps 27 and 28); (iv) the EasyQC script ‘4\_metaanalysis\_qc.compare.ecf’ to compare two meta-analysis results for meta-analysis QC (Steps 29 and 30); (v) the

## Box 2 | Easy QC programming

Generally EasyQC is started by calling the EasyQC function at the R prompt and by using an ecf file as the parameter:

```
> library(EasyQC)
> EasyQC("/path2ecffile/examplescript.ecf")
```

Every data input/output (I/O) and the conducted pipeline are defined in the ecf file.

EasyQC's ecf files are modularized and each step can be conducted separately.

An ecf file consists of two parts: a header or config-section at the beginning that defines data I/O by using the DEFINE and EASYIN functions; this is followed by a scripting interface, which defines the QC steps being executed.

Structure of an ecf file:

### Header to define I/O:

```
[DEFINE, EASYIN]
```

### Scripting interface with EasyQC function steps:

```
[CLEAN, GETNUM, ADDCOL ...]
```

Several example scripts and templates that combine multiple steps described in this protocol are available from <http://www.genepi-regensburg.de/easyqc/>.

R script '4\_metaanalysis\_qc.compare\_logfiles.r' to compare two meta-analysis log files with regard to included and excluded files for meta-analysis QC (Step 30); and (vi) the EasyQC script '4\_metaanalysis\_qc.studymeta.ecf' to perform study-specific meta-analyses for meta-analysis QC (Steps 31 and 32).

Parts of the EasyQC template scripts and single EasyQC functions can also be included in other existing QC pipelines. This task can be accomplished by removing functions from the scripting interface of the template scripts (**Box 2**).

Future studies, such as GWAMAs using 1000 Genomes-imputed data, will have an increased number of variants and will include additional genetic structures such as indels or SVs. The EasyQC software is specifically designed to handle large data sets,

and it can thus be used for larger SNP panels. With regard to memory requirement, EasyQC requires a minimum of 30 GB random access memory (RAM) for 1000 Genomes-imputed data (~40M SNPs) for the file-level QC, which is the protocol part requiring the largest memory. Alternatively, the file-level QC steps can be parallelized by splitting the data into smaller parts, e.g., into chromosomes or into overlapping segments of 5 Mb, as recommended for 1000 Genomes imputation. To handle indels and SVs, adjustments to the scripts, such as allowing for 'I' (insertion) and 'D' (deletion) alleles, are needed and can be made directly to the provided EasyQC scripts. To this end, the EasyQC package is under active development, and future updates will include scripts tailored to 1000 Genomes data.

## MATERIALS

### EQUIPMENT

#### Data

- Allele frequency reference panels. For HapMap-imputed GWAS data: HapMap 'CEU' frequencies as given in 'AlleleFreq\_HapMap\_CEU.v2.txt.gz'. For typed MetaboChip data: 1000 Genomes 'EUR' frequencies as given in 'AlleleFreq\_1,000G\_EUR\_MetaboChip.v1.txt.gz'. Both files are available from the relevant website of the Department of Genetic Epidemiology, University of Regensburg <http://www.genepi-regensburg.de/easyqc/>
- SNP identifier reference panel for marker harmonization: The file 'SNPID\_to\_ChrPosID.b36\_v2.txt.gz' (available from <http://www.genepi-regensburg.de/easyqc/>) maps ~9.1 million known different SNP-IDs (column 'SNPID', which contains different versions of rs-IDs from b35, b36 or b37, as well as array-specific marker names such as 'SNP\_1\_12345') to ~4.8 million unique ChrPosIDs (column 'ChrPosID'). It can be used to harmonize SNP identifier names between HapMap-imputed or MetaboChip data (Step 15). It does not include sex-chromosomal SNPs. Please see the **Supplementary Methods** for a description of the file creation

- QT-interval SNPs for GC correction of typed MetaboChip data and only for traits that are not correlated with the QT interval: 'QTSNPs\_AEL\_TW.txt' (available from <http://www.genepi-regensburg.de/easyqc/>)
- Multiple summary-level association result files

#### Software

- R statistical software (<http://cran.r-project.org/>)
- EasyQC R package (<http://www.genepi-regensburg.de/easyqc/>)
- METAL Meta-analysis software (<http://www.sph.umich.edu/csg/abecasis/metal/>)
- Template R, EasyQC and METAL scripts that can be used to conduct multiple procedure steps are available from <http://www.genepi-regensburg.de/easyqc/>

#### Hardware

- Computer workstation or server with Unix or Linux operating system
- Minimum memory requirements: To perform file-level QC (which is the most memory-intensive step owing to evaluating unfiltered data) with HapMap-imputed data (~2.8 M SNPs), at least 4 GB of RAM should be available

## PROCEDURE

### Setting up logistics of meta-analysis ● TIMING ~2 months

1| Identify GWAS partners and lay out rules of cooperation (memorandum of understanding). Form task groups and set up phone meetings.

2| Develop a GWAS analysis plan (**Supplementary Manual**), including instructions on phenotype transformation, analysis models, covariate adjustment, stratification, use of reference panels for imputation and formatting of data submissions.

## PROTOCOL

3| Set up an sftp site that will be used to collect and securely store the data and to organize and label directories and subdirectories in a logical, self-explanatory manner (**Supplementary Fig. 1**).

### Collecting aggregated statistics per study ● **TIMING** ~2 months

4| Send out the analysis plan and allow 2 months for the collaborators to provide the data.

5| In the meantime, prepare file cleaning instructions and a meta-analysis plan.

6| When all files are available (or at least files from >80% of studies), freeze the data (i.e., protect it from further changes (**Supplementary Fig. 1**)) and start conducting the file-level QC.

### File-level QC ● **TIMING** ~2 months

▲ **CRITICAL** The following file-level QC tasks (Steps 7–18) can be grouped by study and assigned to a set of analysts. Check whether the format and variable names included in the study file match the requested format and columns. The following example uses the terminology of the GIANT format described in **Box 1** and assumes that data are provided by study collaborators in tabular (TAB)-delimited text files with missing values being indicated by ' . '.

7| To check variable names and format in EasyQC, define the requested columns and format by using the DEFINE and the EASYIN functions in the ecf header (for more information on how to use EasyQC, see **Box 2**) by using option A for imputed data or option B for genotyped MetaboChip data.

#### (A) Defining columns and format for imputed data

(i) Type the following commands:

```
DEFINE --acolIn
MarkerName;Strand;Chr;Pos;N;Effect_allele;Other_allele;EAF;Information_type;
Information;BETA;SE;P

--acolInClasses
character;character;character;integer;integer;character;character;numeric;
numeric;numeric;numeric;numeric;numeric

--strMissing .

--strSeparator TAB

EASYIN --fileIn /path2input/study.gwa.file1.txt
EASYIN --fileIn /path2input/study.gwa.file2.txt
...
```

#### (B) Defining columns and format for genotyped MetaboChip data

(i) Type the following commands:

```
DEFINE --acolIn
MarkerName;Strand;Chr;Pos;N;Effect_allele;Other_allele;EAF;
P_HWE;Callrate;BETA;SE;P

--acolInClasses
character;character;character;integer;integer;character;character;numeric;
numeric;numeric;numeric;numeric;numeric

--strMissing .

--strSeparator TAB

EASYIN --fileIn /path2input/study.metaboChip.file1.txt
EASYIN --fileIn /path2input/study.metaboChip.file2.txt
...
```

8| (Optional) If column names were incorrectly labeled (e.g., if the analyst used 'Pvalue' instead of 'P'), change the column names centrally, as this is more time-efficient. If any of the requested columns cannot be clearly allocated or are even missing, consult the study analyst for clarification or, if needed, ask for re-upload. EasyQC will only start to iterate over the defined input files if their headings and format match the requested columns and the requested format. Minor changes to the requested format (e.g., renaming column names or using a different delimiter) can be handled by EasyQC directly through small adjustments in the ecf header.

### ? TROUBLESHOOTING

**9|** Filter monomorphic SNPs. Exclude and count SNPs with allele frequency = 0 or = 1. In EasyQC, this can be done by using the 'CLEAN' function:

```
CLEAN --rcdClean (EAF==0)|(EAF==1) --strCleanName numDrop_Monomorph
```

**10|** Filter SNPs with missing values. Exclude and count all SNPs with missing alleles, *P* value, beta estimate, standard error, allele frequency or sample size. In EasyQC, this can be done by using the 'CLEAN' function:

```
CLEAN --rcdClean is.na(Effect_allele) --strCleanName numDrop_Missing_EA
CLEAN --rcdClean is.na(Other_allele) --strCleanName numDrop_Missing_OA
CLEAN --rcdClean is.na(P) --strCleanName numDrop_Missing_P
CLEAN --rcdClean is.na(BETA) --strCleanName numDrop_Missing_BETA
CLEAN --rcdClean is.na(SE) --strCleanName numDrop_Missing_SE
CLEAN --rcdClean is.na(EAF) --strCleanName numDrop_Missing_EAF
CLEAN --rcdClean is.na(N) --strCleanName numDrop_Missing_N
```

**11|** Filter SNPs with nonsense values. Exclude and count all SNPs with alleles other than 'A', 'C', 'G' or 'T'; *P* values <0 or >1; negative or infinite standard errors ( $\leq 0$  or equal to infinity); and infinite beta estimates or allele frequencies <0 or >1. In EasyQC, this can be done by using the 'CLEAN' function:

```
CLEAN --rcdClean !(Effect_allele%in%c('A', 'C', 'G', 'T')) --strCleanName
numDrop_invalid_EA
CLEAN --rcdClean !(Other_allele%in%c('A', 'C', 'G', 'T')) --strCleanName
numDrop_invalid_OA
CLEAN --rcdClean P<0|P>1 --strCleanName numDrop_invalid_P
CLEAN --rcdClean SE<=0|SE==Inf --strCleanName numDrop_invalid_SE
CLEAN --rcdClean abs(BETA)==Inf --strCleanName numDrop_invalid_BETA
CLEAN --rcdClean (EAF<0)|(EAF>1) --strCleanName numDrop_invalid_EAF
```

**12|** Filter SNPs on the basis of allele frequency and sample size. Exclude and count SNPs with a sample size <30. Add a column called MAC, defined as two times the sample size times MAF, and exclude and count all SNPs with MAC values  $\leq 6$ . In EasyQC, these steps can be performed by using the following EasyQC code:

```
CLEAN --rcdClean N<30 --strCleanName numDrop_Nlt30
ADDCOL --rcdAddCol 2*pmin(EAF,1-EAF)*N --colOut MAC
CLEAN --rcdClean MAC<=6 --strCleanName numDrop_MAClet6
```

**13|** Filter SNPs on the basis of genotype quality. Use option A for imputed data or option B for genotyped MetaboChip data:

**(A) Filtering SNPs in imputed data**

(i) Filter SNPs due to nonsense or missingness. Exclude and count SNPs with missing Information\_type, genotyped SNPs (indicated by Information\_type = 0) with an imputation quality <1 (Information <1) and imputed SNPs (Information\_type !=0) with missing imputation quality. In EasyQC, this can be done by using the 'CLEAN' function:

```
CLEAN --rcdClean is.na(Information_type)
--strCleanName numDrop_MissingInformationType
CLEAN --rcdClean Information_type==0&Information<1
--strCleanName numDrop_Genotyped_LowInformation
CLEAN --rcdClean (Information_type!= 0)&(is.na(Information))
--strCleanName numDrop_Imputed_MissingInformation
```

## PROTOCOL

- (ii) Filter SNPs on imputation quality. Exclude and count SNPs with low imputation quality by using a threshold that depends on the imputation and association software used (**Table 2**). In EasyQC, this can be done with the 'CLEAN' function:

```
CLEAN
--rcdClean
(Information_type!=0&Information<0.3) | (Information_type==2&Information<0.4) |
(Information_type==3&Information<0.8)
--strCleanName numDrop_LowInformation
```

### (B) Filtering SNPs in genotyped MetaboChip data

- (i) Filter SNPs due to nonsense or missingness. Exclude and count SNPs with missing per-SNP call rates; missing HWE *P* values (*P*<sub>HWE</sub>); and call rate or *P*<sub>HWE</sub> values <0 or >1. In EasyQC, this can be done with the 'CLEAN' function:

```
CLEAN --rcdClean is.na(Callrate) --strCleanName numDrop_MissingCallrate
CLEAN --rcdClean is.na(P_HWE) --strCleanName numDrop_MissingPhwe
CLEAN --rcdClean Callrate<0|Callrate>1 --strCleanName numDrop_InvalidCallrate
CLEAN --rcdClean P_HWE<0|P_HWE>1 --strCleanName numDrop_InvalidPhwe
```

- (ii) Filter SNPs on the basis of low call rate and SNPs violating the HWE. Exclude and count SNPs with call rates <0.95 and SNPs with *P*<sub>HWE</sub> values <10<sup>-6</sup>. In EasyQC, this can be done with the 'CLEAN' function:

```
CLEAN --rcdClean Callrate<0.95 --strCleanName numDrop_LowCallrate
CLEAN --rcdClean P_HWE <1e-6 --strCleanName numDrop_LowHwe
```

- 14|** Filter and count SNPs on sex chromosomes. Keep the sex-chromosomal SNPs in a separate file for optional subsequent analyses. In EasyQC, this can be done with the 'CLEAN' function:

```
CLEAN --rcdClean !Chr%in%c(1:22,NA) --strCleanName numDropSNP_ChrXY
--blnWriteCleaned 1
```

If the chromosomal information is missing in the input file, all SNPs on the sex chromosomes will be excluded by the next step.

**15|** Harmonize SNP identifiers. To maximize the overlap in the number of SNPs between the study files and to ensure a proper meta-analysis, create a unique SNP-ID called ChrPosID, which uses the unique format 'chr<chr>:<position>' (e.g., 'chr10:104207431', which only uses genetic positions on build 36). We propose two alternative approaches for this SNP-ID harmonization. Use option A for studies that lack information on genetic positions (columns 'Chr' and 'Pos'). Option A was implemented in GIANT meta-analyses, as the genetic positions were not available in many of the studies (in particular in those that contributed to earlier rounds of analyses). For future studies, we recommend using option B for which Chr and Pos are requested from each collaborator to allow compiling the ChrPosID from the provided information. Option B is the preferable, more generic approach that easily handles novel genotyping arrays (e.g., Exomechip), imputation reference panels (e.g., 1000 Genomes) or genome builds that are not depicted by the provided reference panel 'SNPID\_to\_ChrPosID.b36\_v2.txt.gz' (**Supplementary Methods**).

#### (A) Creating ChrPosID if genetic positions are not available in the study file

- (i) Create a SNP identifier reference panel. Create a reference file that can be used to remap different versions of SNP names to unique ChrPosIDs (see **Supplementary Methods** for detailed descriptions on how to create such a reference file). To analyze HapMap-imputed or MetaboChip data on genome build 36, use the provided SNP identifier reference panel 'SNPID\_to\_ChrPosID.b36\_v2.txt.gz' (**Supplementary Methods**).
- (ii) Add the unique ChrPosID to the study file by merging the study file column MarkerName with the reference file column SNPID. In EasyQC, this can be done by using the 'RENAMEMARKER' function:

```
RENAMEMARKER --colInMarker MarkerName
              --fileRename /path2reffiles/SNPID_to_ChrPosID.b36_v2.txt.gz
              --colRenameOldMarker SNPID
              --colRenameNewMarker ChrPosID
```

- (iii) Check the format of existing ChrPosIDs. To avoid formatting errors with existing ChrPosIDs in study files, remove all spaces from the SNP names (i.e., transform 'chr10: 104207431' to 'chr10:104207431') and add the character string 'chr' at the beginning of the SNP name in case it was forgotten (i.e., transform '10:104207431' to 'chr10:104207431'). In EasyQC, correcting the format of mislabeled ChrPosID SNPs can be performed by using the following commands:

```
EDITCOL --rcdEditCol gsub(" ", "", ChrPosID) --colEdit ChrPosID

EDITCOL --rcdEditCol ifelse(regexpr(":", ChrPosID) == 2 | regexpr(":", ChrPosID) == 3,
paste("chr", ChrPosID, sep = ""), ChrPosID)

--colEdit ChrPosID
```

### (B) Creating ChrPosID if genetic positions are available in the study file

- (i) Generate ChrPosID directly from the provided Chr and Pos columns by horizontally concatenating the string 'chr', column Chr, character ':' and column Pos. This approach requires genetic positions to be given in the study file. In EasyQC, this can be done with the 'ADDCOL' function:

```
ADDCOL --rcdAddCol paste("chr", Chr, ":", Pos, sep = "") --colOut ChrPosID
```

- 16|** Filter duplicate SNPs. To use the best candidate, exclude the duplicate with the smaller sample size. In EasyQC, this can be done with the 'CLEANDUPLICATES' function:

```
CLEANDUPLICATES --colInMarker ChrPosID --strMode samplesize --colN N
```

- 17|** Save cleaned files. Add the prefix "CLEANED." to the filename, save the cleaned file and use "." as missing character. In EasyQC, this can be done with the 'WRITE' function:

```
WRITE --strPrefix CLEANED. --strMissing . --strMode gz
```

- 18|** To perform a file-level QC check, prepare a summary for each study file. Count and check the number of SNPs in the cleaned file and the number of exclusions for each procedure step. An example list of report variables is given in **Supplementary Table 1**. The number of SNPs in the cleaned file should be >2.2 million for GWAS data (if imputed to a HapMap II reference panel) and >100,000 for Metachip data. Major departures from these expected values, generally large numbers of exclusions or any exclusions due to missing or nonsense values (Steps 10, 11, 13A(i), 13B(i)), may indicate systematic issues with the file; consult the study analyst to clarify. When using EasyQC, open the generated summary report in Excel. The report is automatically written to the output path and carries the file extension '.rep'. It contains one row per input file and the QC variables, to be checked, in columns (**Supplementary Table 1**).

### Meta-level QC ● TIMING ~2 months

- 19|** Identify analytical issues by the SE-N plot. To check for issues with trait transformation, the coded sample size or file-naming, calculate the median standard error and maximum sample size of every input and produce a plot of  $c/\text{median}(\text{SE})$  versus  $\sqrt{\text{max}(N)}$  (one point for each file; **Fig. 2**). The proportionality constant  $c$  depends on the genotyping platform or the imputation reference panel (**Table 3**). Find values for  $c$  for standard platforms and panels in **Table 3** (i.e., use 1.93 for typed Metachip data or 1.75 for HapMap II-imputed GWAS data). For platforms or panels other than those given in **Table 3**, the value of  $c$  needs to be computed *de novo* by equation (2) for one study with the respective platform or for the imputation reference panel; this  $c$  can then be applied to the other studies. In EasyQC, calculate the statistics and create the

## PROTOCOL

plot by using the 'CALCULATE' and 'RPLOT' functions (note that here and below, items prefaced by a '#' are comments with instructions and are ignored by the shell):

```
CALCULATE --rcdCalc max(N,na.rm=T) --strCalcName Nmax
CALCULATE --rcdCalc median(SE,na.rm=T) --strCalcName SEmedian
RPLOT --rcdRPlotX sqrt(Nmax)
      --rcdRPlotY [c]/SEmedian
      --arcAdd2Plot abline(a=0,b=1,col="orange")
      --strAxes zeroequal
      --strPlotName SEN-PLOT
# Please replace [c] at --rcdRPlotY with the respective value from Table 3.
```

**20|** Check whether the points follow the identity line. In case any points clearly deviate from the diagonal, consult the study analyst to clarify trait transformation, sample-size coding and file-naming (**Fig. 2; ANTICIPATED RESULTS**). Studies with unaccounted relatives show deviation from the identity line, as the effective sample size is different from the actual sample size, but whether unaccounted relatedness is the reason for an observed deviation should be confirmed after consultation with the analyst.

**21|** Identify analytical issues by using the P-Z scatter plot. To check for problems with beta estimates, standard errors and *P* values, create plots comparing *P* values (on the  $-\log_{10}$  scale) calculated from a *Z*-statistic ( $Z = \beta/SE(\beta)$ ) with the *P* values directly provided by study partners (**Fig. 3 and Supplementary Fig. 3**). In EasyQC, this can be done by using the 'PZPLOT' function:

```
PZPLOT --colBeta BETA --colSe SE --colPval P
```

**22|** Check whether the points follow the identity line. In case any points clearly deviate from the diagonal, consult study analyst (**Fig. 3; ANTICIPATED RESULTS**).

**23|** Identify problems with allele frequencies or strand. To check for strand and allele frequency issues, plot the allele frequency of each SNP and for each file against a reference allele frequency (one plot for each file) (**Fig. 4 and Supplementary Fig. 4**). For HapMap-imputed GWAS data, plot allele frequencies against publically available HapMap allele frequencies, which are reported in the reference file 'AlleleFreq\_HapMap\_CEU.v2.txt.gz'. For genotyped MetaboChip data, plot allele frequencies against publically available 1000 Genomes allele frequencies, which are reported in the reference file 'AlleleFreq\_1,000G\_EUR\_MetaboChip.v1.txt.gz'. In EasyQC, the AFCHECK function can be used to create these plots (please replace [reffile] in the following code with the respective reference file name):

```
AFCHECK --colInMarker ChrPosID
        --colInStrand Strand
        --colInA1 Effect_allele
        --colInA2 Other_allele
        --colInFreq EAF
        --fileRef /path2reffiles/[reffile]
        --colRefMarker ChrPosID
        --colRefA1 A1
        --colRefA2 A2
        --colRefFreq Freq1
        --blnMetalUseStrand 1
```

# Replace the path to the reference and the reference-file name at --fileRef

**24** The frequencies should be distributed along the identity line. Check whether there are patterns (**Fig. 4**; ANTICIPATED RESULTS) that indicate problems with strand or allele frequencies. If you observe such patterns, contact the study analyst to clarify the issue. To define the problem more precisely, it can be helpful to provide the collaborator with a list of outlying SNPs (i.e., SNPs with allele frequencies that deviate by >20% from the reference population) and mismatching SNPs (i.e., SNPs with alleles that do not match the reference, such as AC in the study population versus AT in the reference population). The AFCHECK function automatically saves the lists of outlying or mismatching SNPs to the output path (files indicated by suffix 'AFCHECK.outlier.txt' and 'AFCHECK.mismatch.txt'). In case of problems, it can also be helpful to check the summary report variables indicated by 'AFCHECK.[variablename]' (**Supplementary Table 2**).

**25** Identify population stratification. Calculate  $\lambda_{GC}$  for each study file, without applying the GC correction at this stage, by using all SNPs for imputed GWAS data. For custom chip data, use only a subset of SNPs that are not associated with the outcome of interest. In GIANT, 4,425 QT-interval SNPs (defined in 'QTSNPS\_AEL\_TW.txt') were used to derive the  $\lambda_{GC}$  for typed Metabochip data. To get an overview of the  $\lambda_{GC}$  values across all studies and to identify studies with high  $\lambda_{GC}$ , produce a plot of  $\lambda_{GC}$  values versus the maximum sample sizes (**Fig. 5**). In EasyQC, the calculation of the  $\lambda_{GC}$  and the plotting can be done with the 'GC' and the 'RPLOT' functions:

```
GC      --colPval P
        --blnSuppressCorrection 1
        # --fileGcSnps /path2reffiles/QTSNPS_AEL_TW.txt
        # --colInMarker ChrPosID
        # --colGcSnpsMarker ChrPosID
        # Uncomment the last three parameters for Metabochip data
RPLOT  --rcdRPlotX Nmax
        --rcdRPlotY Lambda.P.GC
        --arcAdd2Plot abline(h=1,col='orange');abline(h=1.1,col='red')
        --strAxes lim(0,NULL,0,NULL)
        --strPlotName GC-PLOT
```

**26** Examine the plot and check whether  $\lambda_{GC}$  is above 1.1 in any of the individual studies. If this is the case, go back to the relevant study analyst to clarify potential issues with population stratification, unaccounted relatedness or duplicated samples included in the analyses (**Fig. 5**; ANTICIPATED RESULTS). The summary report table created by EasyQC might be helpful to identify studies that exhibit high  $\lambda_{GC}$  (variable 'GC.P.Lambda', **Supplementary Table 2**).

### Meta-analysis ● TIMING ~0.5 months

**27** Prepare scripts for an inverse variance-weighted meta-analysis by using a fixed-effects model with METAL, as follows: for QC, we recommend that two analysts perform the meta-analysis independently. The two analysts should ensure that the order in which the studies are read into METAL is the same, because the first study defines the allele coding directions and the following studies are compared with this study. We advise running METAL with the following column definitions and options in the METAL script:

```
# Input columns:
MARKER ChrPosID
ALLELE Effect_allele Other_allele
EFFECT BETA
STDERRLABEL SE
FREQLABEL EAF
PVALUE P
STRAND Strand
CUSTOMVARIABLE N
LABEL N AS N
```



## PROTOCOL

```
# Metal Options:
SCHEME STDERR
WEIGHT N
USESTRAND ON
AVERAGEFREQ ON
MINMAXFREQ ON
VERBOSE OFF

GENOMICCONTROL ON
# GENOMICCONTROL LIST /path2reffiles/QTSNPs_AEL_TW.txt
# Use the latter for Metabochip data!

PROCESS /path2cleanedfiles/CLEANED.study1.file1.txt.gz
PROCESS /path2cleanedfiles/CLEANED.study1.file2.txt.gz
# .
PROCESS /path2cleanedfiles/CLEANED.study1.fileM.txt.gz
PROCESS /path2cleanedfiles/CLEANED.study2.file1.txt.gz
# .
PROCESS /path2cleanedfiles/CLEANED.studyN.fileM.txt.gz

OUTFILE metalout .TBL

ANALYZE HETEROGENEITY
```

To correct for file-specific population stratification, 'GENOMICCONTROL' should be set to 'ON', as this will apply GC correction to each study file. For Metabochip studies, the 'GENOMICCONTROL LIST' parameter can be used to limit the calculation of the  $\lambda_{GC}$  to the subset of QT-interval SNPs. An alternative to using METAL for the GC correction by study file during the meta-analysis is provided by the EasyQC function 'GC' (see the EasyQC manual provided on the EasyQC website for further details). Implementation of this function can be added to the file-level QC to correct study-specific standard errors and  $P$  values in the same way METAL does. To add metrics that measure between-study heterogeneity, use the command 'ANALYZE HETEROGENEITY' at the end of the METAL script file. We provide template METAL scripts, which include the described options and commands ('3\_metaanalysis.metal').

**28|** Perform the inverse variance-weighted meta-analysis and create a METAL log file by using the following command from the command line:

```
metal 3_metaanalysis.metal > metalout_log.txt
```

### Meta-analysis QC ● TIMING ~1.5 months

**29|** Compare results from two meta-analysts. For each of the two meta-analysis results, calculate descriptive statistics of  $P$  values and sample sizes (length, number of missing values, minimum, maximum, median, mean and standard deviation) and the meta-level  $\lambda_{GC}$  (again, restrict calculation of the  $\lambda_{GC}$  to QT-interval SNPs for Metabochip results) and check the values for discrepancies. To compare the meta-analyzed  $P$  values directly, merge the two data sets, create a scatter plot of  $P$  values (on the  $-\log_{10}$  scale) and calculate their Spearman correlation coefficient. In EasyQC, the calculation of the statistics, as well as the merging of the data sets and the creation of the plot, can be done by using the following 'EasyQC' code:

```
DEFINE      --acolIn MarkerName;P.value;N
            --acolInClasses character;numeric;numeric
```

```

EASYIN      --fileIn /path2metalresults/metalout.analyst1.TBL --fileInTag A1
EASYIN      --fileIn /path2metalresults/metalout.analyst1.TBL --fileInTag A2
START EASYQC
EVALSTAT    --colStat P.value
EVALSTAT    --colStat N
GC          --colPval P.value
            --blnSuppressCorrection 1
            #--fileGcSnps /path2reffiles/QTSNPs_AEL_TW.txt
            #--colInMarker MarkerName
            #--colGcSnpsMarker ChrPosID
            # Uncomment last three parameters for Metabochip data
MERGEEASYIN --colInMarker MarkerName
CALCULATE   --rcdCalc cor(P.value.A1,P.value.A2,method="spearman",use="pairwise.
            complete.obs")
            --strCalcName corr_Pvals
PLOT        --rcdSPlotX -log10(P.value.A1)
            --rcdSPlotY -log10(P.value.A2)
            --arcAdd2Plot abline(a=0,b=1,col='orange')
STOP EASYQC

```

The summary report table created by EasyQC contains the descriptive values, the  $\lambda_{GC}$  and the correlation coefficient (**Supplementary Table 3**).

**30|** Examine the calculated values and the scatter plot to check for discrepancies between the two meta-analysis results. All summary statistics should be identical, and the *P* values should lie on the identity line. Most discrepancies observed between meta-analysts are usually explained by different file inclusions in the meta-analysis. To get a quick overview on the files included in the meta-analysis of each analyst, run the R script '4\_metaanalysis\_qc.compare\_logfiles.r'. This action takes the two meta-analysis log files as inputs and creates a table that can be used to compare file inclusions.

**31|** (Optional) Identify analytical issues by calculating the study-level  $\lambda_{GC}$ . If the verified and agreed-on meta-analysis result displays a large meta-level  $\lambda_{GC}$  (>1.1, check the  $\lambda_{GC}$  calculated by Step 29), conduct one meta-analysis for each study (e.g., pooling strata-specific files per study) and calculate the study-level  $\lambda_{GC}$ . An inflated study-level  $\lambda_{GC}$  might pinpoint unaccounted relatedness or overlap of samples across the strata of the study; it can also pinpoint errors as simple as misnaming the strata files (e.g., one file is labeled as 'men', the other as 'women', but the 'men' file was uploaded twice). A substantial fraction of the inflated meta-level  $\lambda_{GC}$  might be explained by such study-specific issues. In EasyQC, the study-specific meta-analysis and the calculation of the study-level  $\lambda_{GC}$  can be performed by using the following EasyQC code:

```

DEFINE      --acolIn ChrPosID;Effect_allele;Other_allele;BETA;SE
            --acolInClasses character;character;character;numeric;numeric
EASYIN      --fileIn /path2cleanedfiles/CLEANED.study1.file1.txt --fileInTag 1
EASYIN      --fileIn /path2cleanedfiles/CLEANED.study1.file2.txt --fileInTag 2
START EASYQC
MERGEEASYIN --colInMarker ChrPosID

```



## PROTOCOL

```
METAANALYSIS --acolBETAs BETA.1;BETA.2
               --acolSEs SE.1;SE.2
               --acolA1s Effect_allele.1;Effect_allele.2
               --acolA2s Other_allele.1;Other_allele.2
               --colOutBeta betaPooled
               --colOutSe sePooled
               --colOutP pPooled

GC             --colPval pPooled
               --blnSuppressCorrection 1
               #--fileGcSnps /path2reffiles/QTSNPs_AEL_TW.txt
               #--colInMarker ChrPosID
               #--colGcSnpsMarker ChrPosID
               # Uncomment last three parameters for metabochip data

STOP EASYQC
```

The summary report table created by EasyQC contains the study-level  $\lambda_{GC}$ .

**32|** Check the study-level  $\lambda_{GC}$  and consult the relevant study analyst in case of a study-level  $\lambda_{GC} > 1.1$ . If the study analyst then flags analytical errors, re-analysis of the study data is needed and Steps 7–32 have to be repeated for the affected files.

**33|** Finalize the meta-analysis. After passing all meta-analysis quality checks, upload the final meta-analysis results file to the ftp site and freeze the upload directory (**Supplementary Fig. 1**). Use the agreed-on result files to extract significant SNPs, to create plots (e.g., Manhattan or QQ plots), and to perform any further evaluation. If a replication of the findings using independent follow-up data is planned, all steps of the PROCEDURE can be repeated for the follow-up meta-analysis.

### ? TROUBLESHOOTING

#### Step 8

It is likely that data from some studies may have been uploaded in a format that differs from the requested one. If the format of an input file does not match the requested format, EasyQC stops with an error message before it starts to iterate over all input files. Issues such as completely missing columns may require contacting the study analyst. Some obvious problems, such as different column names (e.g., 'Pvalue' instead of 'P'), different column separators (e.g., ',' instead of TAB) or missing characters (e.g., 'NaN' instead of '.') can instead be fixed by EasyQC directly (by overwriting the DEFINE parameters at the respective EASYIN statement):

```
EASYIN --fileIn /home/fileWithDifferentFormat.txt
--acolIn
MarkerName;Strand;Chr;Pos;N;Effect_allele;Other_allele;EAF;Information_type;
Information;BETA;SE;Pvalue
--acolInClasses
character;character;character;integer;integer;character;character;numeric;numeric;
numeric;numeric;numeric;numeric
--acolNewName
MarkerName;Strand;Chr;Pos;N;Effect_allele;Other_allele;EAF;Information_type;
Information;BETA;SE;P
--strMissing NaN
--strSeparator COMMA
```

## ● TIMING

The timing of the whole QC and GWAMA pipeline depends on the number of studies involved and also on the experience of the analysts. The estimates reported below are based on the assumption that an existing pipeline of QC and meta-analysis is available (as given by this protocol). The original GIANT conduct and QC has taken longer because of the exploratory nature of the effort. The estimates provided are realistic, as they are given by experienced meta-analysts. For a consortium of comparable size to GIANT's, we estimate the timing to be as follows:

Steps 1–3, setting up logistics of meta-analysis: ~2 months

Steps 4–6, collecting aggregated statistics per study: ~2 months

Steps 7–18, file-level QC: ~2 months

Steps 19–26, meta-level QC: ~2 months

Steps 27 and 28, meta-analysis: ~0.5 months

Steps 29–33, meta-analysis QC: ~1.5 months

## ANTICIPATED RESULTS

### Meta-level QC: identification of analytical issues by the SE-N plot (Steps 19 and 20)

In the case of an inverse normal transformed phenotype, forcing the phenotype into the standard normal distribution,  $N(0,1)$ , the data points on the SE-N plot should tend to describe a straight line on the diagonal (i.e., the identity line). **Figure 2a** illustrates a major deviation of a cluster of GIANT studies from the identity for  $HIP_{adjBMI}$  in the initial round of meta-level QC.

To investigate the reason for this deviation, we surveyed the way each study analyst performed the phenotype: whether the analyst adjusted the phenotype for age, age squared ( $age^2$ ), study-specific covariates and BMI by sex according to the analysis plan and then subjected it to the inverse normal transformation, again separately by sex. This survey revealed that the studies in the cluster above the identity line 'first' (instead of 'last') applied the inverse normal transformation and then adjusted the phenotype for the covariates; a few studies had done the adjustment and/or transformation in men and women combined (instead of by sex) and separated the data by sex afterward.

Subsequent explorations revealed that the SE-N plot identified this problem for phenotypes adjusted for BMI (such as  $HIP_{adjBMI}$ ), but not the BMI-unadjusted phenotypes, as the adjustment for BMI after the inverse normal transformation had disrupted the  $N(0,1)$  distribution of the phenotype (**Supplementary Fig. 2a**). Further explorations revealed that such type of trait transformation issue would result in a loss of power (QQ plot; **Supplementary Fig. 2b**) and in estimates biased toward the null (**Supplementary Fig. 2c**).

Other transformation errors that we were able to identify by using the SE-N plot (not shown) include (i) lack of inverse normal transformation; (ii) stratification by sex conducted after the adjustment and inverse normal transformation; and (iii) miscoded sample size (e.g., stating the full sample size rather than the sample size used for the analysis).

### Meta-level QC: identification of analytical issues by the P-Z scatter plot (Steps 21 and 22)

Occasionally, for a large proportion of SNPs, we observed a discrepancy between the  $P$  value reported by an analysis software and the  $P$  value calculated manually from the  $Z$ -statistic on the basis of the reported beta estimates and standard errors ( $Z = \beta/SE(\beta)$ ). In the GIANT Consortium, we observed such discrepancies caused by the '—score' option in the SNPtest software. The P-Z plots can detect such issues (**Fig. 3a**), and we resolved these issues by asking the study analyst to reanalyze the data using the requested and (in our case) correct '—expected' option (**Fig. 3b**). Panels of such plots for each file in the meta-analysis can provide a quick overview across files and studies (**Supplementary Fig. 3**).

### Meta-level QC: identification of problems with allele frequencies or strands (Steps 23 and 24)

Heterogeneity in allelic patterns may be observed when study allele frequencies are plotted against a reference set, either derived from HapMap, 1000 Genomes or the meta-analysis mean allele frequency. **Figure 4** shows patterns observed in data submitted to the GIANT Consortium. Deviations from the reference frequencies are expected for studies of different ancestry and for studies that have incorrectly coded effect alleles, allele frequencies or strand. Creating a panel displaying such plots for each study file at once provides a quick overview and can identify studies with any of the above issues (**Supplementary Fig. 4**).

### Meta-level QC: identification of population stratification (Steps 25 and 26)

To detect studies with population substructure, the file-specific  $\lambda_{GC}$  can be plotted against the square root of the sample size (**Fig. 5**). Study analysts providing  $\lambda_{GC}$  values  $>1.1$  should be contacted to provide reanalysis (e.g., including principal components in their analysis model). The EasyQC report may help identify which studies have high  $\lambda_{GC}$  values (**Supplementary Table 2**).

Note: Any Supplementary Information and Source Data files are available in the online version of the paper.

**ACKNOWLEDGMENTS** This work was supported by grants from the German Federal Ministry of Education and Research (BMBF) (01ER1206 for I.M.H.); the Leenaards Foundation and the Swiss National Science Foundation (31003A-143914 for Z.K.); the US National Institutes of Health (DK078150, T32 HL007427 for D.C.C.-C.; R01DK075787 for T.E.); the UK Medical Research Council (MRC; U106179471, U106179472 for F.R.D.); the European Research Council (SZ-245 50371-GLUCOSEGENES-FP7-IDEAS-ERC for A.R.W.); the Targeted Financing from the Estonian Ministry of Science and Education (SF0180142s08 for T.E.); the Development Fund of the University of Tartu (SP1GVARENG for T.E.); the European Regional Development Fund to the Centre of Excellence in Genomics (EXCEGEN, 3.2.0304.11-0312 for T.E.); and FP7 (313010 for T.E.). We are also thankful for the GIANT Consortium and the many participating research groups that have allowed us to develop this protocol.

**AUTHOR CONTRIBUTIONS** T.W.W., F.R.D., D.C.C.-C., A.R.W., A.E.L., R.M., T. Ferreira, T.O.K., A.S., T.E., Z.K., I.M.H. and R.J.F.L. comprised the writing group. T.W.W., F.R.D., D.C.C.-C., A.R.W., A.E.L., R.M., T. Ferreira, T.O.K., A.S., T.E. and Z.K. were involved in the pipeline and procedure development. T.W.W., F.R.D., D.C.C.-C., A.R.W., A.E.L., R.M., T. Ferreira, T. Fall, M.G., A.E.J., J.L., S.G., J.C.R., S.V., T.W., T.O.K., A.S., T.E. and Z.K. were the analysts contributing to the QC of the recent GIANT papers.

**COMPETING FINANCIAL INTERESTS** The authors declare no competing financial interests.

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>.

- Hindorf, L.A. *et al.* Potential etiologic and functional implications of genome-wide association loci for human diseases and traits. *Proc. Natl. Acad. Sci. USA* **106**, 9362–9367 (2009).
- McCarthy, M.I. & Hirschhorn, J.N. Genome-wide association studies: past, present and future. *Human Mol. Genet.* **17**, R100–R101 (2008).
- Hirschhorn, J.N. & Gajdos, Z.K. Genome-wide association studies: results from the first few years and potential implications for clinical medicine. *Annu. Rev. Med.* **62**, 11–24 (2011).
- Visscher, P.M., Brown, M.A., McCarthy, M.I. & Yang, J. Five years of GWAS discovery. *Am. J. Hum. Genet.* **90**, 7–24 (2012).
- Anderson, C.A. *et al.* Data quality control in genetic case-control association studies. *Nat. Protoc.* **5**, 1564–1573 (2010).
- Randall, J.C. *et al.* Sex-stratified genome-wide association studies including 270,000 individuals show sexual dimorphism in genetic loci for anthropometric traits. *PLoS Genet.* **9**, e1003500 (2013).
- Surakka, I. *et al.* A genome-wide screen for interactions reveals a new locus on 4p15 modifying the effect of waist-to-hip ratio on total cholesterol. *PLoS Genet.* **7**, e1002333 (2011).
- Manning, A.K. *et al.* A genome-wide approach accounting for body mass index identifies genetic variants influencing fasting glycemic traits and insulin resistance. *Nat. Genet.* **44**, 659–669 (2012).
- Voight, B.F. *et al.* The metabochip, a custom genotyping array for genetic studies of metabolic, cardiovascular, and anthropometric traits. *PLoS Genet.* **8**, e1002793 (2012).
- Cortes, A. & Brown, M.A. Promise and pitfalls of the ImmunoChip. *Arthritis Res. Ther.* **13**, 101 (2011).
- Huyghe, J.R. *et al.* Exome array analysis identifies new loci and low-frequency variants influencing insulin processing and secretion. *Nat. Genet.* **45**, 197–201 (2013).
- Teslovich, T.M. *et al.* Biological, clinical and population relevance of 95 loci for blood lipids. *Nature* **466**, 707–713 (2010).
- Heid, I.M. *et al.* Meta-analysis identifies 13 new loci associated with waist-hip ratio and reveals sexual dimorphism in the genetic basis of fat distribution. *Nat. Genet.* **42**, 949–960 (2010).
- Lango Allen, H. *et al.* Hundreds of variants clustered in genomic loci and biological pathways affect human height. *Nature* **467**, 832–838 (2010).
- Speliotes, E.K. *et al.* Association analyses of 249,796 individuals reveal 18 new loci associated with body mass index. *Nat. Genet.* **42**, 937–948 (2010).
- Scott, R.A. *et al.* Large-scale association analyses identify new loci influencing glycemic traits and provide insight into the underlying biological pathways. *Nat. Genet.* **44**, 991–1005 (2012).
- Schunkert, H. *et al.* Large-scale association analysis identifies 13 new susceptibility loci for coronary artery disease. *Nat. Genet.* **43**, 333–338 (2011).
- Loos, R.J. *et al.* Common variants near *MC4R* are associated with fat mass, weight and risk of obesity. *Nat. Genet.* **40**, 768–775 (2008).
- Willer, C.J. *et al.* Six new loci associated with body mass index highlight a neuronal influence on body weight regulation. *Nat. Genet.* **41**, 25–34 (2009).
- Lindgren, C.M. *et al.* Genome-wide association scan meta-analysis identifies three loci influencing adiposity and fat distribution. *PLoS Genet.* **5**, e1000508 (2009).
- Berndt, S.I. *et al.* Genome-wide meta-analysis identifies 11 new loci for anthropometric traits and provides insights into genetic architecture. *Nat. Genet.* **45**, 501–512 (2013).
- Cochran, W.G. The combination of estimates from different experiments. *Biometrics* **10**, 101–129 (1954).
- Manning, A.K. *et al.* Meta-analysis of gene-environment interaction: joint estimation of SNP and SNP × environment regression coefficients. *Genet. Epidemiol.* **35**, 11–18 (2011).
- de Bakker, P.I. *et al.* Practical aspects of imputation-driven meta-analysis of genome-wide association studies. *Hum. Mol. Genet.* **17**, R122–R128 (2008).
- Fuchsberger, C., Taliun, D., Pramstaller, P.P., Pattaro, C. & CKDGen Consortium. GWAToolbox: an R package for fast quality control and handling of genome-wide association studies meta-analysis data. *Bioinformatics* **28**, 444–445 (2012).
- Kottgen, A. *et al.* Genome-wide association analyses identify 18 new loci associated with serum urate concentrations. *Nat. Genet.* **45**, 145–154 (2013).
- Kottgen, A. *et al.* New loci associated with kidney function and chronic kidney disease. *Nat. Genet.* **42**, 376–384 (2010).
- Schizophrenia Psychiatric Genome-Wide Association Study Consortium. Genome-wide association study identifies five new schizophrenia loci. *Nat. Genet.* **43**, 969–976 (2011).
- Knoppers, B.M., Dove, E.S., Litton, J.E. & Nietfeld, J.J. Questioning the limits of genomic privacy. *Am. J. Hum. Genet.* **91**, 577–578: author reply 579 (2012).
- Gymrek, M., McGuire, A.L., Golan, D., Halperin, E. & Erlich, Y. Identifying personal genomes by surname inference. *Science* **339**, 321–324 (2013).
- Visscher, P.M. & Hill, W.G. The limits of individual identification from sample allele frequencies: theory and statistical analysis. *PLoS Genet.* **5**, e1000628 (2009).
- International HapMap Consortium. *et al.* Integrating common and rare genetic variation in diverse human populations. *Nature* **467**, 52–58 (2010).
- Genomes Project Consortium. *et al.* A map of human genome variation from population-scale sequencing. *Nature* **467**, 1061–1073 (2010).
- Devlin, B. & Roeder, K. Genomic control for association studies. *Biometrics* **55**, 997–1004 (1999).
- Yang, J. *et al.* Genomic inflation factors under polygenic inheritance. *Eur. J. Hum. Genet.* **19**, 807–812 (2011).
- Willer, C.J., Li, Y. & Abecasis, G.R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* **26**, 2190–2191 (2010).
- Higgins, J.P., Thompson, S.G., Deeks, J.J. & Altman, D.G. Measuring inconsistency in meta-analyses. *BMJ* **327**, 557–560 (2003).
- DerSimonian, R. & Laird, N. Meta-analysis in clinical trials. *Control. Clin. Trials* **7**, 177–188 (1986).
- Whitlock, M.C. Combining probability from independent tests: the weighted Z-method is superior to Fisher’s approach. *J. Evol. Biol.* **18**, 1368–1373 (2005).
- R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing (2013).

© 2014 Nature America, Inc. All rights reserved.

