

A Powerful Procedure for Pathway-based Meta-Analysis Using Summary Statistics Identifies 43 Pathways Associated with Type II Diabetes in European Populations

Han Zhang¹, William Wheeler², Paula L Hyland¹, Yifan Yang³, Jianxin Shi¹, Nilanjan Chatterjee^{4*}, Kai Yu^{1*}

¹ Division of Cancer Epidemiology and Genetics, National Cancer Institute, NIH, Bethesda, MD 20892, USA

² Information Management Services Inc., Calverton, MD 20904, USA

³ Department of Statistics, University of Kentucky, Lexington, KY 40508, USA

⁴ Department of Biostatistics, Bloomberg School of Public Health and Department of Oncology, School of Medicine, Johns Hopkins University, Baltimore, MD 21205, USA

* Corresponding Authors

Kai Yu (yuka@mail.nih.gov)

Nilanjan Chatterjee (nchatte2@jhu.edu)

1 Abstract

2 Meta-analysis of multiple genome-wide association studies (GWAS) has become an
3 effective approach for detecting single nucleotide polymorphism (SNP) associations
4 with complex traits. However, it is difficult to integrate the readily accessible SNP-
5 level summary statistics from a meta-analysis into more powerful multi-marker
6 testing procedures, which generally require individual-level genetic data. We
7 developed a general procedure called Summary based Adaptive Rank Truncated
8 Product (sARTP) for conducting gene and pathway meta-analysis that uses only
9 SNP-level summary statistics in combination with genotype correlation estimated
10 from a panel of individual-level genetic data. We demonstrated the validity and
11 power advantage of sARTP through empirical and simulated data. We conducted a
12 comprehensive pathway-based meta-analysis with sARTP on type 2 diabetes (T2D)
13 by integrating SNP-level summary statistics from two large studies consisting of
14 19,809 T2D cases and 111,181 controls with European ancestry. Among 4,713
15 candidate pathways from which genes in neighborhoods of 170 GWAS established
16 T2D loci were excluded, we detected 43 T2D globally significant pathways (with
17 Bonferroni corrected p-values < 0.05), which included the insulin signaling pathway
18 and T2D pathway defined by KEGG, as well as the pathways defined according to
19 specific gene expression patterns on pancreatic adenocarcinoma, hepatocellular
20 carcinoma, and bladder carcinoma. Using summary data from 8 eastern Asian T2D
21 GWAS with 6,952 cases and 11,865 controls, we showed 7 out of the 43 pathways
22 identified in European populations remained to be significant in eastern Asians at
23 the false discovery rate of 0.1. We created an R package and a web-based tool for

1 sARTP with the capability to analyze pathways with thousands of genes and tens of
2 thousands of SNPs.

3

4 **Author Summary**

5 As GWAS continue to grow in sample size, it is evident that these studies need to be
6 utilized more effectively for detecting individual susceptibility variants, and more
7 importantly, to provide insight into global genetic architecture of complex traits.
8 Towards this goal, identifying association with respect to a collection of variants in
9 biological pathways can be particularly insightful for understanding how networks
10 of genes might be affecting pathophysiology of diseases. Here we present a new
11 pathway analysis procedure that can be conducted using summary-level association
12 statistics, which have become the main vehicle for performing meta-analysis of
13 individual genetic variants across studies in large consortia. Through simulation
14 studies we showed the proposed method was more powerful than the existing state-
15 of-art method. We carried out a comprehensive pathway analysis of 4,713 candidate
16 pathways on their association with T2D using two large studies with European
17 ancestry and identified 43 T2D-associated pathways. Further examinations of those
18 43 pathways in 8 Asian studies showed that some pathways were trans-ethnically
19 associated with T2D. This analysis clearly highlights novel T2D-associated pathways
20 beyond what has been known from single-variant association analysis reported
21 from largest GWAS to date.

22

1 Introduction

2 Genome-wide association study (GWAS) has become a very effective way to identify
3 common genetic variants underlying various complex traits [1]. The most commonly
4 used approach to analyze GWAS data is the single-locus test, which evaluates one
5 single nucleotide polymorphism (SNP) at a time. Despite the enormous success of
6 the single-locus analysis in GWAS, proportions of genetic heritability explained by
7 already identified variants for most complex traits still remain small [2]. It is
8 increasingly recognized that the multi-locus test, such as gene-based analysis and
9 pathway (or gene-set) analysis, can be potentially more powerful than the single-
10 locus analysis, and shed new light on the genetic architecture of complex traits [3, 4].

11
12 The pathway analysis jointly tests the association between an outcome and SNPs
13 within a set of genes compiled in a pathway according to existing biological
14 knowledge [4]. Although the marginal effect of a single SNP might be too weak to be
15 detectable by the single-locus test, accumulated association evidence from all signal-
16 bearing SNPs within a pathway could be strong enough to be picked up by the
17 pathway analysis if this pathway is enriched with outcome-associated SNPs. Various
18 pathway analysis procedures have been proposed in the literature, with the
19 assumption that researchers could have full access to individual-level genotype data
20 [5-9]. In practice, pathway analysis usually utilizes data from a single resource with
21 limited sample size, as it can be challenging to obtain and manage individual-level
22 GWAS data from multiple resources. As a result, pathway analysis often fails to
23 identify new findings beyond what have already been discovered by the single-locus

1 tests. To maximize the chance of discovering novel outcome-associated variants by
 2 increasing sample size, a number of consortia have been formed to conduct single-
 3 locus meta-analysis on data across multiple GWAS [10-14]. The single-locus meta-
 4 analysis aggregates easily accessible SNP-level summary statistics from multiple
 5 studies. Similarly, the pathway-based meta-analysis [15-21] that integrates the
 6 same type of summary data across participating studies could provide us a greater
 7 opportunity for detecting novel pathway associations. Future association studies
 8 focusing on identified pathways would have a much-reduced multiple-comparison
 9 burden in searching for novel variants with main or complicated nonlinear joint
 10 effects on the outcome of interest.

11

12 In this paper, we developed a pathway-based meta-analysis procedure by extending
 13 the adaptive rank truncated product (ARTP) pathway analysis procedure [9], which
 14 was originally developed for analyzing individual-level genotype data. The new
 15 procedure, called Summary based ARTP (sARTP), accepts input from SNP-level
 16 summary statistics, with their correlations estimated from a panel of reference
 17 samples with individual-level genotype data, such as the ones from the 1000
 18 Genomes Project [22, 23]. This idea was initially used in conducting gene-based
 19 meta-analysis [24, 25] or conditional test [26]. As will be shown in the Results
 20 Section, sARTP usually has a power advantage over its competitors. In addition,
 21 sARTP is specifically designed for conducting pathway-based meta-analysis using
 22 SNP-level summary statistics from multiple studies. In real applications (e.g., the
 23 type 2 diabetes example described below), it is very common that different studies

1 could have genotypes measured or imputed on different sets of SNPs. As a result, the
2 sample size used in the pathway-based meta-analysis on each SNP can be quite
3 different. Ignoring the difference in sample sizes across SNPs in a pathway-based
4 meta-analysis would generate biased testing results.

5

6 Pathway analysis generally targets two types of null hypotheses [4], including the
7 competitive null hypothesis [15, 16, 18-20], i.e., the genes in a pathway of interest
8 are no more associated with the outcome than any other genes outside this pathway,
9 and the self-contained null hypothesis [17, 21], i.e., none of the genes in a pathway
10 of interest is associated with the outcome. The sARTP procedure focuses on the self-
11 contained null hypothesis, as our main goal is to identify outcome-associated genes
12 or loci. Also, as pointed out by [27], tests for the competitive null hypothesis often
13 assume that genotype measured at different genes are independent when evaluating
14 the association significance level. This assumption, which is generally invalid in
15 practice, is unnecessary for sARTP when testing the self-contained null hypothesis.
16 One may refer to [27] and [4] for more discussion and comparison of these two
17 types of hypotheses.

18

19 The pathways defined in many public databases can consist of thousands of genes
20 and tens of thousands of SNPs. To make the procedure applicable to large pathways,
21 or pathways with high statistical significance, we implement sARTP with efficient
22 and parallelizable algorithms, and adopt the direct simulation approach (DSA) [28]
23 to evaluate the significance of the pathway association.

1

2 We demonstrated the validity and power advantage of sARTP through simulated
3 and empirical data. We applied sARTP to conduct a pathway-based meta-analysis on
4 the association between type 2 diabetes (T2D) and 4,713 candidate pathways
5 defined in the Molecular Signatures Database (MSigDB) v5.0. The analysis used SNP-
6 level summary statistics from two sources with European ancestry. One is generated
7 from the Diabetes Genetics Replication and Meta-analysis (DIAGRAM) consortium
8 [13], which consists of 12,171 T2D cases and 56,862 controls across 12 GWAS. The
9 other one is based on a T2D GWAS with 7,638 T2D cases and 54,319 controls that
10 were extracted from the Genetic Epidemiology Research on Aging (GERA) study [29,
11 30]. The novel T2D-associated pathways detected in the European population were
12 further examined in Asians using summary data generated by the Asian Genetic
13 Epidemiology Network (AGEN) consortium meta-analysis, which combined 8 GWAS
14 of T2D with a total of 6,952 and 11,865 controls from eastern Asian populations
15 [10].

16

17 **Material and Methods**

18 **The Pathway-based Meta-Analysis Procedure**

19 Here we describe the proposed method sARTP for assessing the association
20 between a dichotomous outcome and a pre-defined pathway consisting of J genes.
21 The same procedure can be applied to study a quantitative outcome with minor
22 modifications.

23

1 **Score Statistics and Their Variance-Covariance Matrix**

2 We assume we have data from L GWA studies, with each consisting of $n^{(l)}$ subjects,
 3 $l = 1, \dots, L$. Each gene in that pathway can contain one or multiple SNP(s), while any
 4 two genes may have some overlapped SNPs. For simplicity, we use superscript l to
 5 represent an individual study. For subject i in study l , $i = 1, \dots, n^{(l)}$, let $y_i^{(l)}$ be the
 6 dichotomous outcome (e.g., disease condition, case/control status) taking values
 7 from $\{0, 1\}$, and let $X_i^{(l)}$ be the vector of covariates to be adjusted for. The
 8 centralized genotypes of q SNPs within a pathway are presented as a vector
 9 $G_i^{(l)} = (g_{i1}^{(l)}, \dots, g_{iq}^{(l)})^T$ for subject i . We assume the following logistic regression model
 10 as the risk model

$$11 \quad \text{logit } P(y_i^{(l)} = 1 | X_i^{(l)}, G_i^{(l)}) = (X_i^{(l)})^T \alpha^{(l)} + (G_i^{(l)})^T \gamma, \quad i = 1, \dots, n,$$

12 Under the self-contained null hypothesis $H_0: \gamma = 0$, we denote the maximum
 13 likelihood estimate of $\alpha^{(l)}$ as $\hat{\alpha}^{(l)}$. Let $\hat{y}_i^{(l)} = 1 / (1 + \exp(-X_i^{(l)} \hat{\alpha}^{(l)}))$ and
 14 $u_i^{(l)} = \hat{y}_i^{(l)} (1 - \hat{y}_i^{(l)})$. The Rao's score statistic vector on γ , which is the sum of score
 15 vectors from L participating studies, follows the asymptotic multivariate normal
 16 distribution $N(0, V)$, where

$$17 \quad S = (S_t)_{q \times 1} = \sum_{l=1}^L \sum_{i=1}^{n^{(l)}} G_i^{(l)} (y_i^{(l)} - \hat{y}_i^{(l)}) \quad (1)$$

18 and

$$V = \sum_{l=1}^L \left(\sum_{i=1}^{n^{(l)}} u_i^{(l)} G_i^{(l)} (G_i^{(l)})^T - \sum_{i=1}^{n^{(l)}} u_i^{(l)} G_i^{(l)} (X_i^{(l)})^T \left(\sum_{i=1}^{n^{(l)}} u_i^{(l)} X_i^{(l)} (X_i^{(l)})^T \right)^{-1} \sum_{i=1}^{n^{(l)}} u_i^{(l)} X_i^{(l)} (G_i^{(l)})^T \right). \quad (2)$$

For study l , let $n_t^{(l)}$ be the number of subjects having their genotypes measured as

$H_t^{(l)}$ (or imputed) at SNP t , where $H_t^{(l)} = (g_{1t}^{(l)}, \dots, g_{n_t^{(l)}t}^{(l)})^T$. As pointed out by Hu, Berndt

(24) if the covariates and genotypes are uncorrelated or weakly correlated, the

covariance between scores at SNPs t and s can be approximated as

$$V_{ts} \approx \sum_{l=1}^L n_{ts}^{(l)} \bar{u}^{(l)} \overline{\text{Cov}(H_t^{(l)}, H_s^{(l)})} \approx \sum_{l=1}^L n_{ts}^{(l)} \rho_{ts} \sqrt{\bar{u}^{(l)} \overline{\text{Var}H_t^{(l)}}} \sqrt{\bar{u}^{(l)} \overline{\text{Var}H_s^{(l)}}}, \quad t, s = 1, \dots, q, \quad (3)$$

where $n_{ts}^{(l)}$ is the number of samples that have their genotypes available at both

SNPs in study l , $\bar{u}^{(l)} = (n^{(l)})^{-1} \sum_{i=1}^{n^{(l)}} u_i^{(l)}$, and $\overline{\text{Cov}(H_t^{(l)}, H_s^{(l)})} = (n^{(l)})^{-1} \sum_{i=1}^{n^{(l)}} g_{it}^{(l)} g_{is}^{(l)}$. Here, we

assume that the Pearson's correlation coefficient ρ_{ts} between two SNPs is the same

among all participating studies. This assumption is valid as long as subjects from all

studies are sampled from the same source population, or the population under

study is relatively homogeneous, such as a study of subjects with European ancestry

in the United States.

1 When only the summary statistics, i.e., the estimated marginal log odds ratios $\hat{\beta}_t^{(l)}$
 2 and their standard errors $\tau_t^{(l)}$ are available for each of the L studies, the score
 3 statistic at SNP t , defined by (1) can be approximated as

$$4 \quad S_t \approx \sum_{l=1}^L \left(\tau_t^{(l)} \right)^{-2} \hat{\beta}_t^{(l)}; t = 1, \dots, q. \quad (4)$$

5 Note that $n_t^{(l)} \widehat{\text{Var}} H_t^{(l)} \approx \left(\tau_t^{(l)} \right)^{-2}$, thus according to (3), we have

$$6 \quad V_{ts} \approx \sum_{l=1}^L \frac{n_{ts}^{(l)}}{\sqrt{n_t^{(l)} n_s^{(l)}}} \frac{\rho_{ts}}{\tau_t^{(l)} \tau_s^{(l)}}. \quad (5)$$

7 Assume that ρ_{ts} can be estimated from a public dataset (e.g., 1000 Genomes Project)

8 and the sample sizes $n_t^{(l)}$ and $n_{ts}^{(l)}$ are known, we can approximately recover the

9 variance-covariance matrix $V = (V_{ts})_{q \times q}$ of score statistics $S = (S_t)_{q \times 1}$. In cases when

10 we only have the SNP p-value p and its marginal log odds ratio $\hat{\beta}$, we can compute

11 its standard error as $\tau = |\hat{\beta}| / \sqrt{\chi_{1,p}^2}$, where $\chi_{1,p}^2$ is the quantile satisfying

12 $P(\chi_1^2 \geq \chi_{1,p}^2) = p$, with χ_1^2 representing a 1-df chi-squared random variable.

13

14 ***Combining Score Statistics for Pathway Analysis***

15 With recovered score statistics vector S and its variance-covariance matrix V , we

16 can conduct a pathway association test using the framework of the ARTP method.

17 The ARTP method first combines p-values of individual SNPs within a gene to form a

18 gene-based association statistic (i.e., the gene-level p-value), and then combines the

19 gene-level p-values into a final testing statistic for the pathway-outcome association.

1 In the original ARTP method, [9] proposed the use of a resampling-based method to
 2 evaluate the significance level of the pathway association test. Here we integrate the
 3 SNP-level score statistics into the ARTP framework and use DSA [28] to evaluate the
 4 significance level, which is much faster than the original ARTP algorithm [31]. Below
 5 is a brief summary of the improved ARTP algorithm.

6

7 First we obtain the p-values $p_{t_1}^{(0)}, \dots, p_{t_{q_j}}^{(0)}$ of q_j distinct SNPs in gene j as
 8 $p_t^{(0)} = P(\chi_1^2 \geq S_t^2 / V_u)$. Let $p_{j(1)}^{(0)}, \dots, p_{j(q_j)}^{(0)}$ be their order statistics such that
 9 $p_{j(1)}^{(0)} \leq \dots \leq p_{j(q_j)}^{(0)}$. For any predefined integer K and SNP-level cut points $c_1 < \dots < c_K$,
 10 we define the observed negative log product statistics for that gene at cut point c_k
 11 as

$$12 \quad w_{jk}^{(0)} = - \sum_{t=1}^{\min(q_j, c_k)} \log p_{j(t)}^{(0)}, \quad k=1, \dots, K.$$

13 We sample M copies of vectors of the score statistic from the null distribution
 14 $N(0, V)$ and convert each of them to be the tail probability of χ_1^2 as $p_{t_1}^{(m)}, \dots, p_{t_{q_j}}^{(m)}$,
 15 $m=1, \dots, M$, which are then used to calculate $w_{jk}^{(m)}$, $m=1, \dots, M$. The significance of
 16 $w_{jk}^{(0)}$ can be estimated as

$$17 \quad \xi_{jk}^{(0)} = \frac{\#\{w_{jk}^{(m)} \geq w_{jk}^{(0)}; m=1, \dots, M\}}{M+1}.$$

18

1 The ARTP statistic for testing association between gene j and the outcome is
 2 defined as $T_j^{(0)} = \min_{k=1, \dots, K} \xi_{jk}^{(0)}$. Note that for any $w_{jk}^{(m)}$, the set
 3 $\{w_{jk}^{(m')} : m' \in \{0, \dots, M\} \text{ and } m' \neq m\}$ forms its empirical null distribution. The
 4 significance of $w_{jk}^{(m)}$ therefore can be estimated as

$$5 \quad \xi_{jk}^{(m)} = \frac{\#\{w_{jk}^{(m')} \geq w_{jk}^{(m)}; m' \neq m \text{ and } m' = 1, \dots, M\}}{M+1}, \quad m = 1, \dots, M.$$

6 This idea, which was given by [32], can be used to avoid the computationally
 7 challenging nested two-layer resampling procedure for evaluating p-values. The p-
 8 value of $T_j^{(0)}$ can be readily calculated as

$$9 \quad z_j^{(0)} = \frac{\#\{T_j^{(m)} \leq T_j^{(0)} : m = 1, \dots, M\}}{M+1}, \quad j = 1, \dots, J.$$

10 where $T_j^{(m)} = \min_{k=1, \dots, K} \xi_{jk}^{(m)}$. $z_j^{(0)}$ is the estimated gene-level p-value for the association
 11 between the outcome and the j th gene. To obtain the pathway p-value, a similar
 12 procedure as above can be applied to combine already established gene-level p-
 13 values $z_j^{(0)}$, $j = 1, \dots, J$, through a set of K' gene-level cut points $d_1 < \dots < d_{K'}$. For
 14 simplicity, let $\zeta_k^{(0)}$ be the significance (p-value) of negative log product statistics
 15 defined on $z_j^{(0)}$, $j = 1, \dots, J$ at a specific cut point d_k , $k = 1, \dots, K'$. The ARTP statistic
 16 for the pathway association is defined as $T^{(0)} = \min_{k=1, \dots, K'} \zeta_k^{(0)}$. The top d_{k^*} genes, at
 17 which $\zeta_{k^*}^{(0)} = \min_{k=1, \dots, K'} \zeta_k^{(0)}$, can be regarded as the set of selected candidate genes that
 18 collectively convey the strongest pathway association signal.

1

2 In the following discussion, we will use the term sARTP to represent the proposed
3 pathway analysis procedure using the SNP-level summary statistics as input, and
4 reserve the term ARTP to represent the original ARTP procedure that requires the
5 individual-level genetic data. Both procedures adopt the DSA algorithm to accelerate
6 evaluating the significance level. When performing the pathway analysis in this
7 paper, we set SNP-level cut points as $(c_1, c_2) = (1, 2)$, i.e., gene-level association is
8 summarized by one or two most significant SNPs within each gene, and gene-level
9 cut points as $d_k = k \max(1, \lceil J/20 \rceil)$, $k = 1, \dots, 10$, where J is the number of genes in a
10 pathway, and $\lceil J/20 \rceil$ is the largest integer that is less or equal to $J/20$. We used
11 $M = 10^5$ DSA steps to assess the significance level of each pathway in the initial
12 screening. For pathways with estimated p-values $< 10^{-4}$, we further refined their p-
13 value estimates with $M = 10^7$ or 10^8 DSA steps.

14

15 ***Applying sARTP to meta-analysis result***

16 Many GWAS consortia usually publish their meta-analysis results by providing only
17 the combined results from the fixed effects model, rather than the summary
18 statistics from each participating study. We can apply sARTP to this meta-analysis
19 result directly, with some modifications. First, since the reported marginal log odds
20 ratios for each SNP by using the fixed effects inverse-variance weighting method is
21 given by

$$\hat{\beta}_t = \frac{\sum_{l=1}^L (\tau_t^{(l)})^{-2} \hat{\beta}_t^{(l)}}{\sum_{l=1}^L (\tau_t^{(l)})^{-2}},$$

with its standard error given by

$$\tau_t = \left(\sum_{l=1}^L (\tau_t^{(l)})^{-2} \right)^{-1/2}. \quad (6)$$

Based on (4), we can see $S_t \approx \tau_t^{-2} \hat{\beta}_t$. By assuming large sample sizes and certain conditions (see Appendix A), we can also approximate the covariance between S_t and S_s , which is given by (5), as

$$V_{ts} \approx \frac{n_{ts}}{\sqrt{n_t n_s}} \frac{\rho_{ts}}{\tau_t \tau_s}, \quad (7)$$

where $n_t = \sum_{l=1}^L n_t^{(l)}$, and $n_{ts} = \sum_{l=1}^L n_{ts}^{(l)}$. Thus, using just the meta-analysis result, without knowing summary statistics from each participating study, we can still obtain S_t exactly, and approximately recover V_{ts} . As a result, we can carry out the pathway-based meta-analysis based on the SNP-level meta-analysis result as if it were summary data from a single study. We call this approach the Meta-analysis based sARTP (MsARTP).

14

However, to apply the MsARTP, we need additional sample size information n_t and n_{st} in order to properly estimate the variance-covariance matrix defined by (7). If the same set of SNPs are studied by all participating studies, we have $n_t = n_s = n_{st}$,

1 and the approximation (7) becomes $V_{ts} \approx \frac{\rho_{ts}}{\tau_t \tau_s}$, i.e., we can obtain the estimated
2 variance-covariance matrix without knowing n_t and n_{st} . But in most applications,
3 not all GWAS choose the same SNP genotyping array, even after the imputation
4 using the same reference genomes. As a result, the SNP coverage, i.e., the set of SNPs
5 evaluated in each participating study can be quite different. In those situations, we
6 need to know the SNP coverage information in each participating study in order to
7 obtain n_s and n_{st} . We will show in the Results Section that using MsARTP with an
8 inappropriate uniform coverage assumption (i.e., $n_{ts} = n_t = n_s$), which is commonly
9 made by many multi-locus approaches, can lead to inflated type I error.

10
11 Given SNP-level summary statistics from each participating study, we can either
12 apply sARTP directly, or first conduct a SNP-level meta-analysis, and then apply
13 MsARTP to the meta-analysis result. These two approaches use the same score
14 statistics, and different but consistent estimates for the variance-covariance matrix.
15 Numeric experiments in the Results Section suggest that these two approaches
16 generate very similar pathway p-values.

17

18 **Study Materials**

19 ***Pathway and Gene Definition***

20 We downloaded definitions for 4,716 human and murine (mammalian) pathways
21 (gene sets) from the MSigDB v5.0 (C2: curated gene sets). Genomic definitions for

1 genes were downloaded from Homo sapiens genes NCBI36 and reference genome
2 GRCh37.p13 using the Ensemble BioMart tool.

3

4 ***DIAGRAM Study***

5 The DIAGRAM (DIABetes Genetics Replication And Meta-analysis) consortium
6 conducted a large-scale GWAS meta-analysis to characterize the genetic architecture
7 of T2D [13]. We downloaded the summary statistics generated by the DIAGRAMv3
8 (Stage 1) GWAS meta-analysis from www.diagram-consortium.org [13]. The meta-
9 analysis studied 12 GWAS with European ancestry consisting of 12,171 cases and
10 56,862 controls. Up to 2.5 million autosomal SNPs with minor allele frequencies
11 (MAFs) larger than 1% were imputed using CEU samples from Phase II of the
12 International HapMap Project. Study-specific covariates were adjusted in testing
13 T2D-SNP association under an additive logistic regression model [13]. SNP-level
14 summary statistics from each GWAS were first adjusted for residual population
15 structure using the genomic control (GC) method [33], and then combined in the
16 fixed effects meta-analysis.

17

18 We sorted 2.5 million autosomal SNPs by their corresponding meta-analysis sample
19 sizes in Figure S1, which shows that there are two major groups of SNPs with equal
20 sample sizes. One group of 469,985 SNPs (19.0%) had 12,171 cases and 56,862
21 controls, which included all the available samples in the meta-analysis; another
22 group of 1,431,361 SNPs (57.9%) had 9,580 cases and 53,810 controls. Since the
23 calculation of covariance V_{ts} in (7) relies on n_{ts} , the number of samples having

1 genotypes available at both SNP s and SNP t , in order to obtain an accurate
 2 estimate of n_{ts} , we focused on these two groups of SNPs, which in combination had a
 3 total of 1,901,346 SNPs. For any two SNPs in this reduced set, it is certain
 4 $n_{ts} = \min(n_t, n_s)$. The Pearson's correlation coefficients ρ_{ts} were estimated using an
 5 external reference panel consisting of genotypes on 503 European subjects (CEU,
 6 TSI, FIN, GBR, and IBS) from the 1000 Genomes Project (Phase 3, v5, 2013/05/02).

8 ***GERA Study***

9 We assembled a GWAS on T2D from the Genetic Epidemiology Research on Adult
 10 Health and Aging (GERA, dbGaP Study Accession: phs000674.v1.p1). The GERA
 11 project includes a cohort of over 100,000 adults who are members of the Kaiser
 12 Permanente Medical Care Plan, Northern California Region, and participating in the
 13 Kaiser Permanente Research Program on Genes, Environment, and Health (RPGEH).
 14 From the GERA data, we compiled a GWAS with 7,638 T2D cases and 54,319
 15 controls (subjects without T2D) who self-reported to be non-Hispanic White
 16 Europeans in the RPGEH survey. We performed the genotype imputation with
 17 IMPUTE2 [34] using CEU reference samples from Phase II of the International
 18 HapMap Project. After removing SNPs with low imputation quality ($r^2 < 0.3$), we
 19 ended up with 2.4 million SNPs for further analysis. In the single-locus analysis, we
 20 adjusted for the categorized body mass index (BMI) provided in the downloaded
 21 dataset (adding a category for missing BMI), gender, year of birth (in five-year
 22 categories), a binary indicator on whether or not a participant was diagnosed with
 23 cancer (includes malignant tumors, neoplasms, lymphoma and sarcoma), and the

1 top five eigenvectors for the adjustment of population stratification. In the following
2 discussion, we refer this assembled T2D GWAS as the GERA study.

3

4 When analyzing the SNP-level summary data from the GERA study, the Pearson's
5 correlation coefficients ρ_{ts} were estimated using an external reference panel
6 consisting of genotypes on 503 European subjects from the 1000 Genomes Project.

7

8 ***AGEN-T2D Study***

9 The Asian Genetic Epidemiology Network (AGEN) consortium carried out a meta-
10 analysis by combining eight GWAS of T2D with a total of 6,952 cases and 11,865
11 controls from eastern Asian populations [10]. The meta-analysis was conducted
12 with the fixed effect model. We obtained SNP-level summary statistics on 2.6 million
13 imputed and genotyped autosomal SNPs from AGEN, and used this summary data to
14 evaluate whether pathway associations identified in European populations remain
15 to be present in Asians. We adopted an external reference panel consisting of 312
16 eastern Asian subjects (103 from CHB, 105 from CHS, and 104 from JPT) from the
17 1000 Genomes Project for the variance-covariance matrix estimation in the pathway
18 analysis.

19

20 **Results**

21 **Simulation Studies**

1 Firstly, we conducted a simulation study to evaluate the empirical size of sARTP and
 2 MsARTP. Secondly, we compared empirical powers of different strategies for
 3 carrying out pathway-based meta-analysis that integrated summary statistics from
 4 multiple studies. We also evaluated whether results from sARTP were consistent
 5 with the ones from MsARTP. Thirdly, we compared our method to the recently
 6 developed method aSPUsPath [8] that can be used for pathway-based meta-analysis.
 7 We used the R package, aSPU (version 1.39), with the default settings given in [8, 17]
 8 to conduct the aSPUsPath test.

9

Empirical Size of sARTP and MsARTP

11 To evaluate the empirical size of sARTP and MsARTP, we conducted a simulation
 12 study by using individual-level GWAS data of the pathway
 13 PUJANA_BREAST_CANCER_WITH_BRCA1_MUTATED_UP (including 728 SNPs in 50
 14 genes) from the GERA study. We picked 12,000 samples randomly for this
 15 experiment. By keeping their genotypes unchanged, we randomly assigned 6,000
 16 subjects as cases and the remaining as controls to generate 500,000 datasets. We
 17 split each dataset into three case-control studies, each with 2,000 cases and 2,000
 18 controls. To mimic the scenario when not all studies have their genotypes measured
 19 on the same set of SNPs (such as the one occurred in the DIAGRAM and AGEN data),
 20 we assumed that each case-control study had genotypes measured on only half of
 21 SNPs in the pathway. For each generated dataset that consisted of three case-control
 22 studies, we applied sARTP to the SNP-level summary data obtained from each case-
 23 control study, and MsARTP to the meta-analysis result based on the three case-

control studies, with the variance-covariance matrix estimated by an external reference panel (with 503 European reference samples from the 1000 Genomes Project), or an internal reference panel (with 500 samples randomly selected from the GERA data).

Based on results from the 500,000 generated datasets, this simulation study showed that both sARTP and MsARTP, using the internal or external reference samples, can well control their empirical sizes (Table 1). Given the same reference panel, the p-values estimated from sARTP and MsARTP are highly consistent (Pearson's correlation coefficient > 0.99). Furthermore, the p-values of sARTP (or MsARTP) estimated with an external or internal reference panel are also very consistent (Pearson's correlation coefficient > 0.99). More numeric experiments demonstrating the validity of sARTP under the null are described in Appendix B.

Table 1: Empirical sizes of the sARTP, MsARTP, and MsARTP-u procedures

		Size				
Reference		0.05	0.01	0.005	0.001	0.0005
sARTP	External ^a	0.050	0.0093	0.0040	0.00078	0.00044
	Internal ^b	0.046	0.0087	0.0042	0.00074	0.00040
MsARTP	External	0.048	0.0093	0.0040	0.00076	0.00041
	Internal	0.048	0.0084	0.0042	0.00074	0.00041
MsARTP-u	External	0.082	0.018	0.0081	0.0013	0.00064
	Internal	0.094	0.022	0.011	0.0016	0.00081

Empirical sizes are estimated based on 500,000 datasets simulated from the GERA data.

^aUsing 503 European samples from the 1000 Genomes Project as an external reference;

^bUsing 500 samples from the GERA data as an internal reference.

1 To demonstrate the importance of knowing n_t and n_{ts} when applying MsARTP to
 2 the meta-analysis result, we analyzed each simulated dataset using MsARTP
 3 assuming the uniform coverage ($n_{ts} = n_t = n_s$). We called this approach MsARTP-u. It
 4 is clear from Table 1 that MsARTP-u assuming the uniform coverage suffers from
 5 inflated type I errors with either the internal or external reference panel.

6

7 ***Empirical Power of sARTP and MsARTP for Pathway-based Meta-analysis***

8 We conducted a set of simulation studies to compare the power of different
 9 strategies to carry out pathway analysis when SNP-level summary statistics were
 10 available from multiple studies. We considered a hypothetical pathway consisting of
 11 50 genes randomly selected from chromosome 17, each with 20 randomly chosen
 12 SNPs. The joint genotype distribution at the 20 SNPs within each gene was defined
 13 by the observed genotypes in the GERA study. We further assumed that all genes in
 14 that pathway are independent. This assumption is unnecessary for sARTP and
 15 MsARTP, but it was introduced for simplifying the simulation. For the risk model, we
 16 assumed the first \mathcal{M} ($\mathcal{M} = 5, 10, 15$) genes were associated with the outcome.
 17 Within each outcome-associated gene, we picked the SNP with its MAF closest to the
 18 median MAF level within the gene to be functional. We considered the following risk
 19 model

$$20 \quad \text{logit } P(y=1 | g_1^*, \dots, g_{\mathcal{M}}^*) = \alpha + \sum_{l=1}^{\mathcal{M}} \gamma_l^* g_l^*, \quad (8)$$

21 where g_l^* is the genotype (encoded as 0, 1, or 2 according to counts of minor alleles)
 22 at the functional SNP within gene l . Under this model, γ_l^* is also the marginal log

odds ratio for the l th functional SNP [9]. Given the sample sizes of cases and controls, and the MAF of the l th functional SNP, γ_l^* was chosen such that the theoretical power of the trend test to detect the l th functional SNP is equal to \mathcal{P} ($\mathcal{P} = 0.3, 0.4$), with 0.05 as the targeted type I error rate. For every pair of $(\mathcal{M}, \mathcal{P})$, we generated 1,000 datasets, each consisting of three case-control studies, with the same sample size and SNP coverage configurations used for evaluating the empirical size. Given the genotype distribution in the general population, individual-level genotype data for a case-control study can be generated according to the assumed risk model (8).

10

We assumed that only SNP-level summary statistics from each of the three studies were available. For each simulated dataset, we applied sARTP and MsARTP, using either an internal or external reference panel to estimate the variance-covariance matrix. The sARTP and MsARTP approaches integrate association evidence across SNP-level summary statistics, which are obtained by pooling information from all participating studies on individual SNPs. As a comparison, we also considered a naïve approach, in which we first applied sARTP to analyze the summary statistics from each study separately, and then combined the three pathway p-values with Fisher's method. This naïve approach could be useful when the researchers do not have access to the SNP-level summary data but the pathway p-values from individual studies. The empirical powers are compared at the type I error level of 0.05, and are summarized in Table 2. It is obvious that the pathway-based meta-analysis using sARTP, with either the internal or external reference panel, have

1 almost the same level of power as the MsARTP method. It is also evident that both
 2 sARTP and MsARTP are more powerful than the naïve approach, which suggests
 3 that it is always be beneficial to have the SNP-level summary statistics from each
 4 participating study, or SNP-level meta-analysis result when conducting a pathway
 5 analysis.

6

7 Given the SNP coverage information, the MsARTP method is a valid pathway
 8 association test that has well controlled type I error and similar power to the sARTP
 9 method. In the following analysis, either sARTP or MsARTP is chosen depending on
 10 the type of available data. For the sake of simplicity, we always label the chosen
 11 procedure as sARTP.

12

Table 2. Power comparisons under the type I error rate of 0.05 when analyzing data from three studies

\mathcal{P}^a	\mathcal{M}^b	Internal reference ^c			External reference ^d		
		sARTP	MsARTP	Fisher	sARTP	MsARTP	Fisher
0.3	5	0.165	0.170	0.110	0.170	0.167	0.105
	10	0.405	0.402	0.229	0.399	0.401	0.221
	15	0.573	0.578	0.334	0.564	0.561	0.323
0.4	5	0.292	0.293	0.162	0.295	0.297	0.154
	10	0.642	0.637	0.363	0.640	0.635	0.362
	15	0.858	0.858	0.574	0.855	0.856	0.561

For every pair of \mathcal{P} and \mathcal{M} , the empirical powers are computed from 1,000 simulated datasets at the level of 0.05. Each dataset contains three studies. The pathway consists of 50 independent genes, each with 20 SNPs. Fisher's method is used to combine the three pathway p-values obtained by applying sARTP to the SNP-level summary data from each of three studies separately.

^a The theoretical power of the single-locus trend test on the functional SNP under the type I error rate of 0.05, given the sample sizes of cases and controls, and the MAF of the functional SNP;

^b The number of genes including the functional SNPs;

^c Using 500 samples from the GERA data as an internal reference;

^d Using 503 European samples from the 1000 Genomes Project as an external reference.

1 *Power Comparison between sARTP and aSPUsPath*

2 Since the aSPUsPath method in the current aSPU package cannot handle summary
3 data from multiple studies, or meta-analysis results from studies with varied SNP
4 coverage, we focused on the scenario with just one study, and adopted the similar
5 simulation strategy as the one used by [8] to compare the power between sARTP
6 and aSPUsPath. We simulated haplotypes on a set of SNPs within a gene in the
7 general population using the algorithm of Wang and Elston (35). Then the joint
8 genotypes on a subject can be formed by randomly pairing two haplotypes. In brief,
9 we first chose the MAF for each SNP by randomly sampling a value from the uniform
10 distribution $U(0.1, 0.4)$. Then for the set of SNPs in a gene we sampled a latent
11 vector $Z = (z_1, \dots, z_q)^T$ from a multivariate normal distribution with a covariance
12 matrix $\text{Cov}(z_i, z_j) = \rho^{|i-j|}, 1 \leq i, j \leq q$, where ρ was sampled from the uniform
13 distribution $U(0, 0.8)$ for a given gene. We randomly picked 50% of the SNPs and
14 converted their simulated z_i into minor and major alleles (coded as 0, 1), with the
15 cuts chosen for each z_i such that the resultant minor allele has its frequency defined
16 by the specified MAF. For the remaining SNPs, we used the same algorithm to
17 dichotomize $-z_i$ into minor and major alleles. This created a more realistic
18 haplotype structure such that a haplotype can consist of a mixture of minor and
19 major alleles. Genotypes on SNPs from different genes were generated
20 independently.

21

1 Given the number of genes (20, 50, or 80) in a pathway, the proportion of genes (5%,
2 10%, 20%, and 30%) associated with the outcome, and a chosen common value for
3 all log odds ratios (γ^*) in the risk model (8), we repeated the following steps to
4 generate 1,000 case-control studies, with each consisting of 1,000 cases and 1,000
5 controls. First, the number of SNPs within each gene was randomly chosen from 10
6 to 100. Second, for each randomly selected outcome-associated gene, we randomly
7 picked a functional SNP. Third, we use the aforementioned algorithm of Wang and
8 Elston (35) to generate the individual-level genotype data for a case-control study
9 according to the specified risk model. We also considered the situation where all γ^*
10 in the risk model (8) had the same magnitude but different directions. More
11 precisely, when generating a case-control study at the third step, we defined the risk
12 model (8) by randomly choosing the direction of each log odds ratio to be positive
13 or negative with equal probability. Furthermore, we considered a more complex
14 scenario where each outcome-associated gene had one or two functional SNPs, each
15 with equal probability.

16

17 All simulation results are given in Table S1 and Table S2. It is clear that sARTP are
18 generally more powerful than aSPUsPath, especially when the signal-to-noise ratio
19 (the proportion of genes including a functional SNP) is relatively low. The two types
20 of tests tend to have comparable performance when the signal-to-noise ratio
21 increases to 30%, although it is uncommon for a candidate pathway to have such a
22 high signal-to-noise ratio in real applications. For example, among the 4,713
23 candidate pathways analyzed in the next section, only 4.2% and 0.9% of the

1 pathways have over 20% and 30% of their genes that are likely to contain
2 association signals (i.e., with gene-level p-values < 0.05).

3

4 From Table S1 and Table S2, we also notice that the advantage of sARTP over
5 aSPUsPath is more evident if not all minor alleles of the functional SNPs are
6 deleterious (or protective) variants (i.e., γ^* in the risk model (8) are not all positive).
7 This is expected, as the sARTP approach does not take the effect direction of the
8 minor allele at each SNP into consideration, while aSPUsPath integrates a set of
9 candidate statistics, including the one similar to the burden test that assumes all
10 minor alleles are either deleterious or protective. When this assumption is not valid,
11 the inclusion of the burden test statistic in aSPUsPath is unlikely to enhance the
12 power, but certainly would increase the multiple-testing penalty.

13

14 **Evaluation of sARTP using Data from T2D Studies**

15 To demonstrate the consistency between results obtained by sARTP using SNP-level
16 summary statistics and the ones by ARTP using individual-level genotype data, we
17 compared pathway analysis results from three different procedures on the 4,713
18 candidate pathways using the GERA GWAS data. Details on how those 4,713
19 pathways were pre-processed are given in the Results of T2D Pathway Analysis
20 Section. We applied sARTP to the SNP-level summary statistics generated from the
21 GERA study, using either an internal or an external reference panel. We also
22 obtained the pathway p-values by directly applying the ARTP method to the
23 individual-level GERA GWAS data. Figure 1 shows the comparison among p-values

- 1 from these three analyses, and demonstrates that all three approaches can generate
- 2 very consistent results.

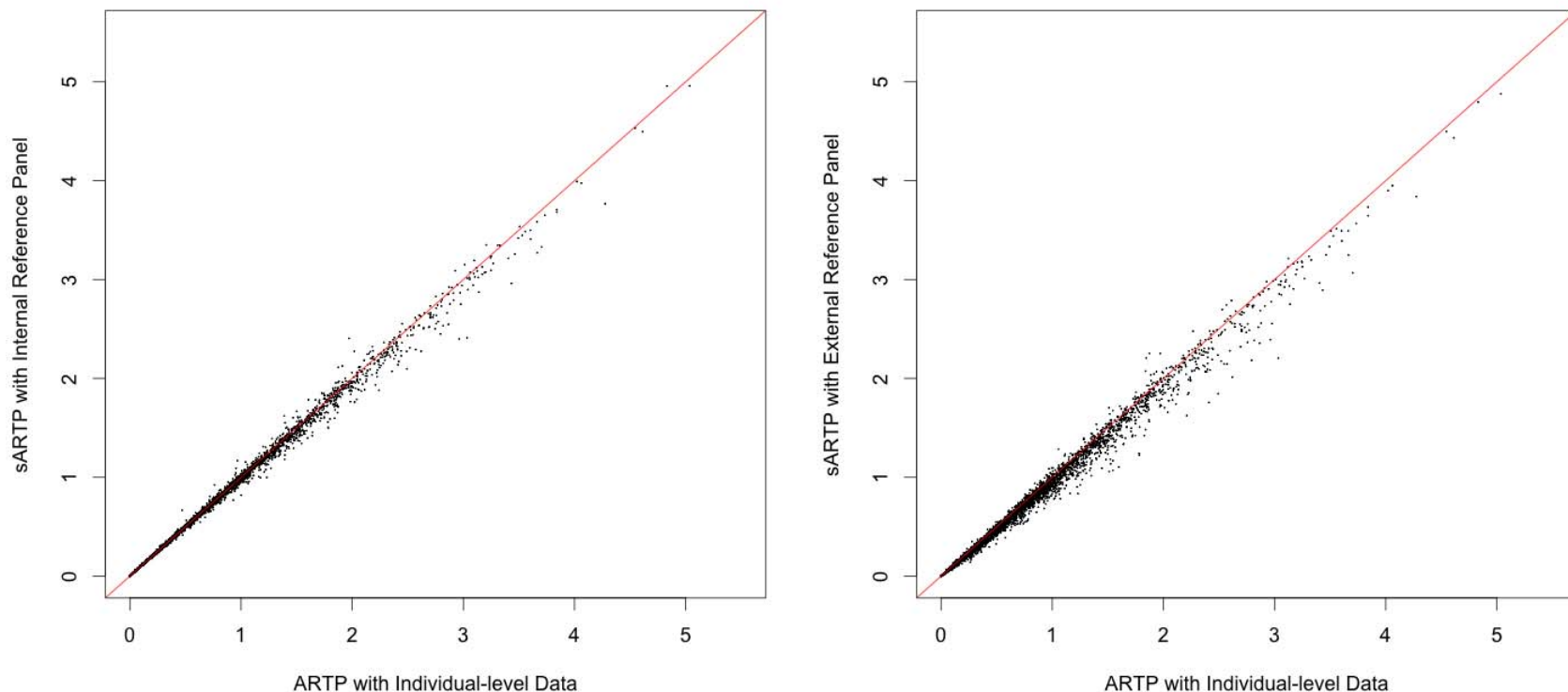


Figure 1. Comparisons of p-values from three types of pathway analyses on the GERA data

Based on the GERA data, 4,713 pathways are analyzed in three different ways. Pathway p-values obtained by ARTP using the GERA individual-level genetic data (x-axis) are compared with the ones obtained by sARTP using summary statistics in combination with the internal reference panel that consists of 500 randomly selected GERA samples (left), and the ones using the summary statistics in combination with the external reference panel that consists of 503 European subjects from the 1000 Genomes Project (right).

1 **Results of T2D Pathway Analysis**

2 *Findings from the European Populations*

3 Since our goal was to identify new susceptibility loci for T2D through the pathway
4 analysis, we excluded 170 high evidence T2D associated SNPs that were either listed
5 in [13] or found from the GWAS Catalog satisfying the following three conditions
6 simultaneously: (1) were investigated by GWAS of samples with European ancestry;
7 (2) had reported p-values $<10^{-7}$ on the initial study; and (3) were replicated on
8 independent studies. We excluded 195 SNPs that has their single-locus testing p-
9 values less than 10^{-7} in either DIAGRAM or GERA data to ensure that the pathway
10 analysis result was not driven by a single SNP. In addition, we further excluded
11 genes within a $\pm 500\text{kb}$ region from each of the removed SNPs to eliminate potential
12 association signals that could be caused by linkage disequilibrium (LD) with the
13 index SNPs.

14
15 We conducted three types of pathway-based meta-analyses using sARTP, including
16 the one using the DIAGRAM SNP-level summary statistics, the one using the GERA
17 SNP-level summary statistics, and the pathway meta-analysis combining SNP-level
18 summary statistics from both DIAGRAM and GERA studies. When applying the
19 pathway-based meta-analysis to a single gene, we refer to this as the gene-level
20 meta-analysis. We used the external reference panel of 503 Europeans from the
21 1000 Genomes Project to estimate the variance-covariance matrix.

22

Before performing a pathway analysis, we applied LD filtering to remove redundant SNPs. For any two SNPs with their pairwise squared Pearson's correlation coefficient > 0.9 estimated from the external reference panel from the 1000 Genomes Project, we removed the one with a smaller value defined as, $2f(1-f)n_0n_1n^{-1}$, where n_0 and n_1 are numbers of controls and cases, $n = n_0 + n_1$, and f is the MAF based on the reference panel. This value is proportional to the non-centrality parameter of the trend test statistic at a given SNP. We also excluded SNPs with $MAF < 1\%$. After all SNP filtering steps, we had a total of 4,713 pathways for the analysis. The summary of the number of genes and SNPs used in each pathway analysis is given in Figure S2.

11

The DIAGRAM study had a genomic control inflation factor $\lambda_{GC} = 1.10$ based on the published meta-analysis result. The assembled GERA T2D GWAS had $\lambda_{GC} = 1.08$. When conducting the pathway analysis on each of two studies, we adjusted the inflation by using the corresponding $\sqrt{\lambda_{GC}}$ to rescale the standard error of estimated log odds ratio at each SNP. The single-locus meta-analysis combining results from DIAGRAM and GERA datasets had an inflation factor $\lambda_{GC} = 1.067$ after each study had adjusted for its own inflation factor. We further adjusted this inflation in the pathway and gene-level meta-analysis when combining SNP-level summary statistics from both studies using formulas (4) and (5).

21

1 The Q-Q plots of gene-level and pathway p-values are given in Figure 2. Gene-level
 2 p-value Q-Q plots based on the three analyses show no sign of inflation with their
 3 λ_{GC} close to 1.0, but suggest that there are enriched gene-level association signals at
 4 the tail end. The pathway p-value Q-Q plots, on the other hand, shift away from the
 5 diagonal identify line and have much higher λ_{GC} , which suggests that T2D
 6 associated genes are preferably included in pathways under study. In fact, it can be
 7 seen from Figure S3 that a gene with a smaller gene-level meta-analysis p-value
 8 tends to be included in more pathways, even though the 4,713 pathways collected
 9 from MSigDB v5.0 are not specifically defined for the study of T2D.

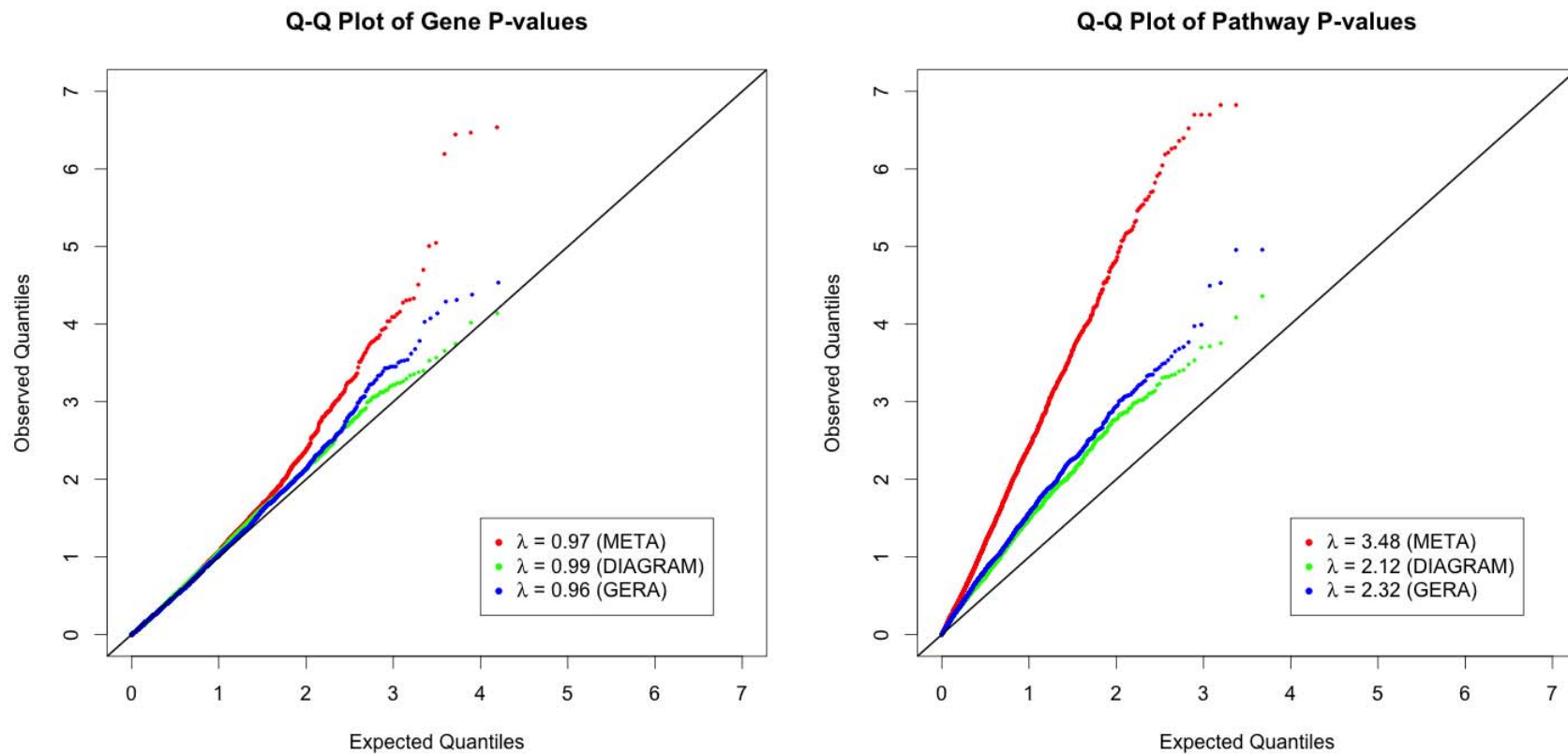


Figure 2. Q-Q plots of gene-level and pathway p-values based on the sARTP procedure on the DIAGRAM study, the GERA study, and the two studies combined

Left: Q-Q plots of gene-level p-values on 15,946 genes based on the sARTP gene-based analysis of the DIAGRAM study (DIAGRAM), the GERA study (GERA), and the two studies combined (META)
 Right: Q-Q plots of pathway p-values on 4,713 pathways based on the sARTP pathway analysis of the DIAGRAM study (DIAGRAM), the GERA study (GERA), and the two studies combined (META)

1

2 Figure 2 illustrates that the gene and pathway level signal from the GERA study
3 tends to be slightly stronger than that from the DIAGRAM study. The main reason is
4 that the DIAGRAM summary result had gone through two rounds of inflation
5 adjustments, with the first round done at each participating study, and the second
6 round on the meta-analysis result. Also, its second round adjustment ($\lambda_{GC} = 1.10$) is
7 larger than the one applied to the GERA study ($\lambda_{GC} = 1.08$). Adjusting for λ_{GC} in the
8 pathway analysis could be too conservative, since some proportion of the inflation
9 can be caused by the real polygenic effect. A less conservative adjustment could be
10 possible, but it might not be adequate. More discussions on this issue are given in
11 the Discussion Section.

12

13 Based on the pathway meta-analysis on a total of 4,713 pathways, we identified 43
14 significant pathways with p-values less than 1.06×10^{-5} , the family-wise significant
15 threshold based on the Bonferroni correction. Their pathway meta-analysis results
16 as well as results from individual studies are summarized in Table 3. More detailed
17 results on each of 43 significant pathways are given in the Figures S6-S48, and
18 Supplemental Data. There are a total of 15,946 unique genes in all 4,713 pathways.
19 The top 50 genes with smallest gene-level p-values based on the gene meta-analysis
20 are listed in Table S3. Because of the LD filtering, a gene belonging to two pathways
21 might end up with slightly different sets of SNPs. To remove this ambiguity, we
22 obtained the gene-level p-values by conducting a gene-level meta-analysis on each
23 gene separately.

1
2
3

Table 3. Summary of 43 significant pathways detected by the pathway meta-analysis based on the DIAGRAM and GERA studies

Pathway	META ^a	DIAGRAM ^b	GERA ^c
SCHLOSSER_SERUM_RESPONSE_UP ^d	2.50E-08	2.92E-04	1.77E-03
PENG_RAPAMYCIN_RESPONSE_DN ^{ef}	1.50E-07	1.68E-03	2.08E-04
YAGI_AML_WITH_T_8_21_TRANSLOCATION ^{ef}	1.50E-07	4.33E-03	4.46E-04
PATIL_LIVER_CANCER ^d	2.00E-07	4.35E-05	3.89E-03
PUJANA_CHEK2_PCC_NETWORK ^{ef}	2.00E-07	1.18E-02	3.39E-03
STEIN_ESRRA_TARGETS ^{ef}	2.00E-07	9.37E-04	8.39E-04
STEIN_ESRRA_TARGETS_UP ^{ef}	3.00E-07	6.38E-03	1.02E-04
WANG_CISPLATIN_RESPONSE_AND_XPC_UP ^{ef}	4.00E-07	6.75E-03	1.18E-01
CADWELL_ATG16L1_TARGETS_DN ^e	4.35E-07	1.45E-03	9.59E-03
SONG_TARGETS_OF_IE86_CMV_PROTEIN ^d	5.30E-07	7.88E-04	3.20E-05
CASORELLI_ACUTE_PROMYELOCYTIC_LEUKEMIA_DN ^{ef}	5.50E-07	2.48E-02	5.70E-03
RIZ_ERYTHROID_DIFFERENTIATION ^d	6.15E-07	2.33E-02	2.31E-02
BORCZUK_MALIGNANT_MESOTHELIOMA_UP ^d	6.50E-07	7.61E-03	2.52E-02
HILLION_HMGA1_TARGETS ^e	9.00E-07	3.39E-01	1.10E-05
KEGG_MATURITY_ONSET_DIABETES_OF_THE_YOUNG ^d	1.14E-06	1.68E-02	3.58E-04
HOLLEMAN_ASPARAGINASE_RESISTANCE_BALL_DN ^e	1.22E-06	3.42E-02	1.67E-03
PUJANA_BRCA1_PCC_NETWORK ^{ef}	1.50E-06	1.69E-02	5.11E-03
HOSHIDA_LIVER_CANCER_SUBCLASS_S3 ^d	1.95E-06	6.14E-03	5.83E-03
GRAESSMANN_APOPTOSIS_BY_DOXORUBICIN_DN ^{ef}	2.00E-06	9.57E-04	6.41E-04
REACTOME_REGULATION_OF_BETA_CELL_DEVELOPMENT ^d	2.26E-06	4.81E-02	1.15E-03
PUJANA_BREAST_CANCER_WITH_BRCA1_MUTATED_UP ^d	2.48E-06	7.61E-03	2.39E-02
BLALOCK_ALZHEIMERS_DISEASE_UP ^d	2.50E-06	3.52E-02	3.26E-02
GOBERT_OLIGODENDROCYTE_DIFFERENTIATION_UP ^d	2.85E-06	3.68E-02	5.37E-04
MCBRYAN_PUBERTAL_BREAST_4_5WK_DN ^d	2.95E-06	2.20E-01	1.10E-05
REACTOME_REGULATION_OF_GENE_EXPRESSION_IN_BETA_CELLS ^d	3.11E-06	2.81E-02	2.33E-03
SANSOM_APC_TARGETS_DN ^{ef}	3.25E-06	3.08E-01	2.84E-03
NABA_MATRISOME ^d	3.45E-06	4.46E-02	1.30E-02
PUJANA_BRCA2_PCC_NETWORK ^d	4.65E-06	1.25E-02	6.32E-02

Pathway	META ^a	DIAGRAM ^b	GERA ^c
KEGG_TYPE_II_DIABETES_MELLITUS ^d	4.85E-06	6.38E-02	9.93E-04
LINDGREN_BLADDER_CANCER_CLUSTER_1_DN ^d	5.50E-06	5.20E-02	4.36E-03
ROPERO_HDAC2_TARGETS ^e	6.04E-06	2.11E-02	1.41E-03
KEGG_INSULIN_SIGNALING_PATHWAY ^d	6.20E-06	2.65E-02	4.26E-03
CHEN_PDGF_TARGETS ^d	6.36E-06	2.48E-03	1.04E-02
REACTOME_INTEGRATION_OF_ENERGY_METABOLISM ^e	6.50E-06	6.11E-02	1.06E-04
PETROVA_ENDOTHELIUM_LYMPHATIC_VS_BLOOD_UP ^d	6.60E-06	1.79E-01	8.72E-03
REACTOME_PPARA_ACTIVATES_GENE_EXPRESSION ^d	6.70E-06	4.97E-02	1.69E-02
AGUIRRE_PANCREATIC_CANCER_COPY_NUMBER_UP ^e	6.90E-06	1.73E-01	1.97E-04
DACOSTA_UV_RESPONSE_VIA_ERCC3_UP ^d	7.45E-06	2.20E-02	4.25E-02
TOYOTA_TARGETS_OF_MIR34B_AND_MIR34C ^d	8.10E-06	1.20E-01	2.67E-03
HOLLEMAN_ASPARAGINASE_RESISTANCE_ALL_DN ^e	8.43E-06	5.90E-02	3.22E-03
REACTOME_CLASS_I_MHC_MEDIATED_ANTIGEN_PROCESSING_PRESENTATION ^e	8.45E-06	5.67E-02	2.00E-02
DODD_NASOPHARYNGEAL_CARCINOMA_DN ^{ef}	1.00E-05	2.86E-03	1.71E-04
REACTOME_MEMBRANE_TRAFFICKING ^e	1.04E-05	6.43E-03	4.52E-04

The 43 pathways are identified among 4,713 candidate pathways for having their pathway meta-analysis p-values less than the $<1.06 \times 10^{-5}$, the Bonferroni correction threshold.

^a P-values based on summary statistics combined from the DIAGRAM and GERA studies;

^b P-values based on summary statistics from the DIAGRAM study;

^c P-values based on summary statistics from the GERA study;

^d Pathways that do not contain genes in the 17q21 region;

^e Pathways that contain at least one gene in the 17q21 region;

^f Pathways that remain globally significant after excluding genes in the 17q21 region.

1 From Table 3, we can notice that some identified pathways have relatively weak
2 association signals from each of the two studies, but have very significant p-values
3 based on the pathway meta-analysis on the two studies combined. For example, the
4 pathway RIZ_ERYTHROID_DIFFERENTIATION has p-values of 0.0233 and 0.0231
5 based on DIAGRAM and GERA studies, respectively. Combining these two p-values
6 using Fisher's method yields a p-value of 0.0046. On the other hand, the pathway
7 meta-analysis produces a much more significant result ($p = 6.15 \times 10^{-7}$). This
8 demonstrates the power advantage of the pathway meta-analysis over the approach
9 that simply combines the pathways p-values from individual studies. The
10 aforementioned simulation studies also confirmed this observation (Table 2).

11

12 In Figure 3, we illustrate the connection between the 43 significant pathways and a
13 group of genes showing association evidence. For the purpose of illustration, in the
14 figure we only focus on 46 genes that are covered by the 43 pathways and have their
15 gene-level meta-analysis p-values less than 0.001. It is evident from Figure 3 that a
16 cluster of 4 genes, *UBE2Z*, *SNF8*, *GIP*, and *ATP5G1*, has the most significant gene-level
17 p-values (Table S3), and contribute association signals to 20 out of 43 significant
18 pathways (Figures S6-S25). These 4 genes overlap each other at chromosome 17q21.
19 This region contains a previously unidentified genome-wide significant synonymous
20 SNP rs1058018 (meta-analysis $p = 3.06 \times 10^{-8}$) after two rounds of inflation
21 adjustments. More detailed information on SNP rs1058018 and SNPs in that region
22 are given in Table S4, Figure S4, and Figure S5. By conditioning on rs1058018, none
23 of the other SNPs in this region are significant based on the conditional association

1 analysis using the GERA individual-level GWAS data. Based on GTEx data v6,
 2 rs1058018 is a *cis* eQTL for *UBE2Z* in blood ($p = 7.9 \times 10^{-15}$). *UBE2Z* is involved in
 3 Class I MHC antigen processing and presentation (GeneCards). The region at 17q21
 4 was previously implicated to be associated with T2D through a candidate gene/loci
 5 approach [36]. Although genes at the 17q21 region carry the strongest association
 6 signal, 11 out of those 20 pathways remain to be globally significant ($p < 1.06 \times 10^{-5}$)
 7 after excluding those genes from the pathway definition (Table 3).

8

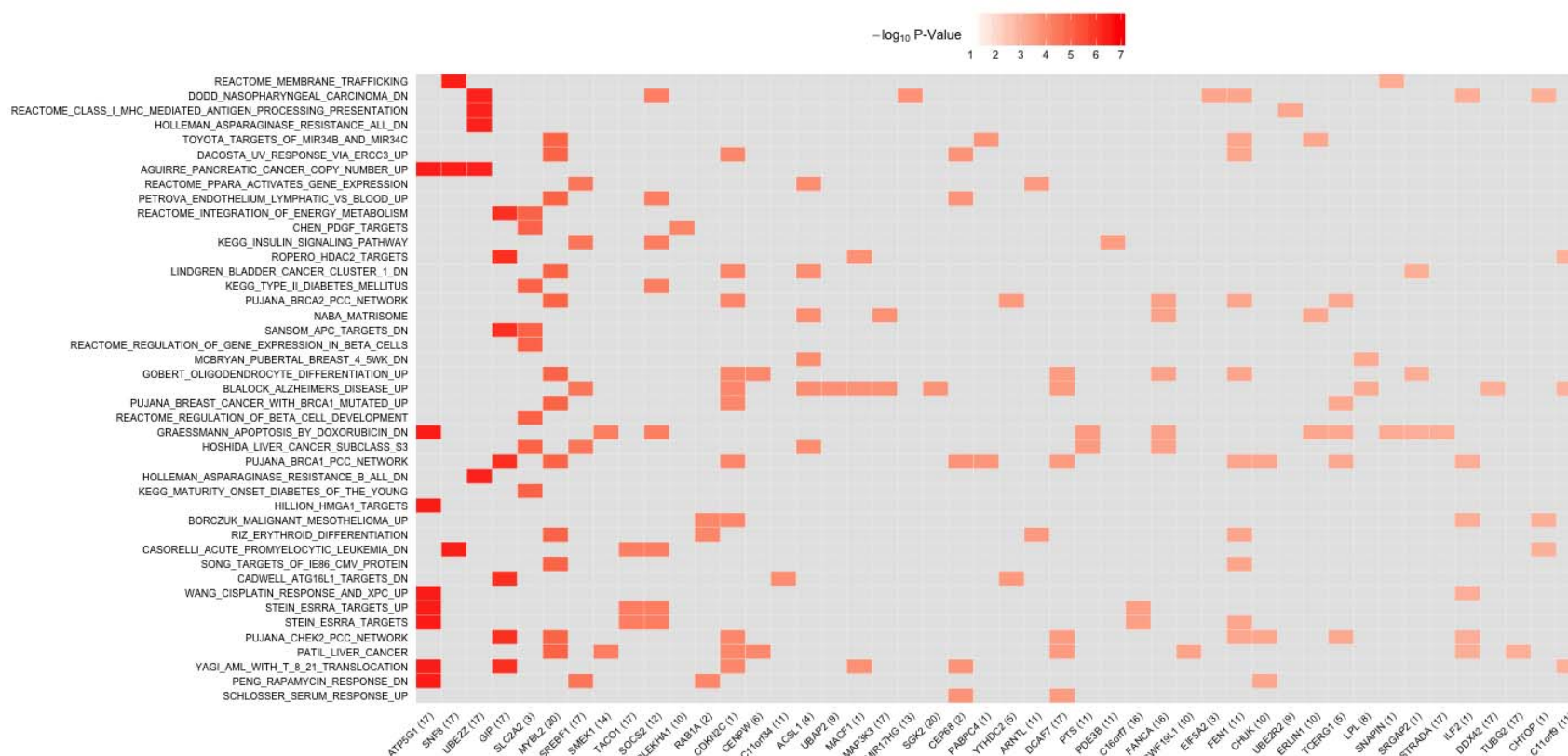


Figure 3. Heat map of gene-level p-values on selected genes within 43 significant pathways based on the DIAGRAM and GERA studies. There are 46 unique genes in the 43 significant pathways that have their gene-level meta-analysis p-values less than 0.001. Each row in the plot represents one of 43 significant pathways. Each column represents one of the 46 unique genes. The chromosome IDs of 46 unique genes are given in parentheses. The color of each cell represents the gene-level p-value (in the $-\log_{10}$ scale). A cell for a gene that is not included in a pathway is colored gray in the corresponding entry. The orders of genes (x-axis) and pathways (y-axis) are arranged according to their gene and pathway meta-analysis p-values.

The majority of 43 identified pathways are enriched with signals from multiple chromosomal regions as demonstrated by the Q-Q plots of their SNP-level and gene-level p-values (Figures S6-S48). For example, the strongest T2D-associated pathway, SCHLOSSER_SERUM_RESPONSE_UP, consists of 103 genes, which includes two genes with p-values < 0.001, and has 20 genes with p-values between 0.001 and 0.05 (Figure S26, and Supplemental data). We conducted the ingenuity pathway analysis on those 22 genes with p-values less than 0.05, and found enrichment of these genes in caveolae-mediated cytolysis (important for removal of low/high density lipoproteins), and lipid metabolism pathways, and in functions/diseases related to differentiation of phagocytes and transport of proteins.

It is assuring that our pathway analysis detected several pathways that are natural candidates underlying the development of T2D, including the pathways KEGG_MATURITY_ONSET_DIABETES_OF_THE_YOUNG (Figure S31), KEGG_TYPE_II_DIABETES_MELLITUS (Figure S41), KEGG_INSULIN_SIGNALING_PATHWAY (Figure S43), and REACTOME_REGULATION_OF_BETA_CELL_DEVELOPMENT (Figure S33). It is worth emphasizing that these pathways were analyzed after excluding genes in the neighborhood of 170 GWAS established T2D loci and 195 SNPs with p-values < 10^{-7} on either DIAGRAM or GERA data, which suggests that these well-defined T2D-related pathways are enriched with additional unidentified and contributory T2D-associated genes.

1 Among the 43 globally significant pathways, there are multiple ones that are defined
 2 according to specific gene expression patterns on various tumor types, including
 3 pancreatic adenocarcinoma (Figure S21), hepatocellular carcinoma (HCC) (Figure
 4 S27, and S32), bladder carcinoma (Figure S42), nasopharyngeal carcinoma (Figure
 5 S24), and familial breast cancer (Figures S34, and S37). It is well recognized that
 6 T2D patients have elevated risk of cancer at multiple cancer sites, such as the liver
 7 and pancreas [37, 38]. These findings can provide valuable insights into the genetic
 8 basis underlying the connection between T2D and a host of different cancers.

9

10 In the above analysis, we used the sARTP method with the gene-level association
 11 evidence summarized by one or two most significant SNPs within each gene, under
 12 the assumption that there are at most two independent association signals within a
 13 given gene. We also applied sARTP by using 3 SNP-level cut points (i.e.,
 14 $(c_1, c_2, c_3) = (1, 2, 3)$) to reanalyze the 4,713 pathways based on the combined data of
 15 DIAGRAM and GERA. It appears that results obtained by sARTP with 3 SNP-level cut
 16 points are very consistent with those with 2 cut points (Figure S49).

17

18 ***Findings from Eastern Asian Populations***

19 We reanalyzed the 43 significant pathways identified from the European
 20 populations using summary-level data generated by the AGEN-T2D study. An
 21 inflation factor $\lambda_{GC} = 1.03$ calculated from the AGEN-T2D meta-analysis was
 22 adjusted in the pathway meta-analysis. The genetic regions excluded from analyzing
 23 the DIAGRAM and GERA studies were also excluded from the AGEN study. The

1 results were summarized in Table 4. There are 10 out of 43 pathways with the
2 unadjusted p-value less than 0.05, suggesting that many pathways identified from
3 the European populations were also enriched with T2D-associated genes in the
4 eastern Asian populations. The Supplemental Data provides more details on results
5 of those 10 pathways. Among the 43 pathways, we were able to identify 4 significant
6 T2D-associated pathways at the false discovery rate (FDR [39]) of 0.05 (Figures S27,
7 S12, S32 and S36), and 3 additional T2D-associated pathways at the FDR of 0.1
8 (Figures S47, S44, and S25). All the pathway p-values remain basically the same
9 level if we further excluded genes within ± 500 kb regions surrounding the GWAS
10 T2D loci established in eastern Asian populations. These results support the
11 presence of trans-ethnic pathway effect on T2D in European and eastern Asian
12 populations [11, 12].

13

14 Table 4: Pathway p-values and FDR adjusted p-values based on the AGEN-T2D study.

Pathway	P-value ^a	FDR ^b
PATIL_LIVER_CANCER ^c	0.0014	0.029
CADWELL_ATG16L1_TARGETS_DN	0.0023	0.029
HOSHIDA_LIVER_CANCER_SUBCLASS_S3 ^c	0.0025	0.029
GOBERT_OLIGODENDROCYTE_DIFFERENTIATION_UP ^c	0.0027	0.029
DACOSTA_UV_RESPONSE_VIA_ERCC3_UP ^c	0.011	0.074
CHEN_PDGF_TARGETS ^c	0.011	0.074
REACTOME_MEMBRANE_TRAFFICKING	0.012	0.074
AGUIRRE_PANCREATIC_CANCER_COPY_NUMBER_UP	0.026	0.14
MCBRYAN_PUBERTAL_BREAST_4_5WK_DN ^c	0.041	0.19
LINDGREN_BLADDER_CANCER_CLUSTER_1_DN ^c	0.043	0.19
PUJANA_CHEK2_PCC_NETWORK	0.057	0.21
BLALOCK_ALZHEIMERS_DISEASE_UP ^c	0.059	0.21
SCHLOSSER_SERUM_RESPONSE_UP ^c	0.085	0.27
RIZ_ERYTHROID_DIFFERENTIATION ^c	0.089	0.27
DODD_NASOPHARYNGEAL_CARCINOMA_DN	0.097	0.28
TOYOTA_TARGETS_OF_MIR34B_AND_MIR34C ^c	0.10	0.28

Pathway	P-value ^a	FDR ^b
REACTOME_CLASS_I_MHC_MEDIATED_ANTIGEN_PROCESSING_PRESENTATION	0.14	0.32
PUJANA_BRCA2_PCC_NETWORK ^c	0.14	0.32
WANG_CISPLATIN_RESPONSE_AND_XPC_UP	0.14	0.32
STEIN_ESRRA_TARGETS	0.16	0.35
PUJANA_BREAST_CANCER_WITH_BRCA1_MUTATED_UP ^c	0.17	0.35
GRAESSMANN_APOPTOSIS_BY_DOXORUBICIN_DN	0.18	0.36
PUJANA_BRCA1_PCC_NETWORK	0.23	0.43
HOLLEMAN_ASPARAGINASE_RESISTANCE_B_ALL_DN	0.35	0.62
PENG_RAPAMYCIN_RESPONSE_DN	0.38	0.62
ROPERO_HDAC2_TARGETS	0.38	0.62
KEGG_MATURITY_ONSET_DIABETES_OF_THE_YOUNG ^c	0.39	0.62
HILLION_HMGA1_TARGETS	0.43	0.65
HOLLEMAN_ASPARAGINASE_RESISTANCE_ALL_DN	0.44	0.65
PETROVA_ENDOTHELIUM_LYMPHATIC_VS_BLOOD_UP ^c	0.45	0.65
REACTOME_PPARA_ACTIVATES_GENE_EXPRESSION ^c	0.52	0.72
REACTOME_REGULATION_OF_GENE_EXPRESSION_IN_BETA_CELLS ^c	0.55	0.74
KEGG_INSULIN_SIGNALING_PATHWAY ^c	0.58	0.74
STEIN_ESRRA_TARGETS_UP	0.59	0.74
YAGI_AML_WITH_T_8_21_TRANSLOCATION	0.64	0.76
NABA_MATRISOME ^c	0.65	0.76
KEGG_TYPE_II_DIABETES_MELLITUS ^c	0.67	0.76
SANSOM_APC_TARGETS_DN	0.69	0.76
REACTOME_INTEGRATION_OF_ENERGY_METABOLISM	0.70	0.76
REACTOME_REGULATION_OF_BETA_CELL_DEVELOPMENT ^c	0.70	0.76
BORCZUK_MALIGNANT_MESOTHELIOMA_UP ^c	0.75	0.79
SONG_TARGETS_OF_IE86_CMV_PROTEIN ^c	0.86	0.88
CASORELLI_ACUTE_PROMYELOCYTIC_LEUKEMIA_DN	0.88	0.88

1 These 43 pathways are nominated through the pathway meta-analysis on DIAGRAM
2 and GERA studies. The analysis is carried out on the summary data from the AGEN-
3 T2D study.

4 ^a P-values based on summary statistics from the AGEN-T2D study;

5 ^b FDR adjusted p-values;

6 ^c Pathways that do not contain genes in the 17q21 region.

7

8 Given the existing epidemiologic evidence on the close connection between T2D and
9 the liver cancer, it is noteworthy that the two HCC related pathways (Figures S27,

1 S32) identified in European populations remain to be significant in eastern Asian
2 populations at the FDR of 0.05 (Table 4). The pathway PATIL_LIVER_CANCER
3 consists of 653 genes (after data preprocessing) that are highly expressed in HCC
4 and are enriched with genes having functions related to cell growth, cell cycle,
5 metabolism, and cell proliferation [40]. The other pathway,
6 HOSHIDA_LIVER_CANCER_SUBCLASS_S3 consists of 240 genes that show similar
7 gene expression variation patterns and together define a HCC subtype with its
8 unique histologic, molecular and clinical characteristics [41]. These two pathways
9 have only 6 genes in common, and none of the 6 genes has a gene-level p -value $<$
10 0.05 in either European or eastern Asian data. More in depth investigations of these
11 two complementary pathways could lead to further understanding the connection
12 between T2D and the liver cancer.

13

14 The genome-wide significant SNP rs1058018 at the 17q21 region identified through
15 the combined analysis of DIAGRAM and GERA studies turned out to be null in the
16 AGEN-T2D study ($p = 0.29$). This could be due to the relatively small sample size of
17 the AGEN-T2D study, or the genetic risk heterogeneity at the 17q21 locus among
18 different ethnic populations. Nevertheless, 2 out of the 20 pathways (Figures S12
19 and S25) that contain genes within the 17q21 region are still significant at the FDR
20 of 0.1. Among the 23 pathways that do not contain any gene within the 17q21 region,
21 5 pathways remain significant at the FDR of 0.1 (Figure S27, S32, S36, S47, and S44).

22

23 Discussion

1 We developed a general statistical procedure sARTP for pathway analysis using
 2 SNP-level summary statistics generated from multiple GWAS. By applying sARTP to
 3 summary statistics from two large studies with a total of 19,809 T2D cases and
 4 111,181 controls with European ancestry, we were able to identify 43 globally
 5 significant T2D-associated pathways after excluding genes in neighborhoods of
 6 GWAS established T2D loci. Using summary data generated from 8 T2D GWAS with
 7 6,952 cases and 11,865 controls from eastern Asian populations, we further showed
 8 that 7 out of 43 pathways identified in the European populations were also
 9 significant in the eastern Asian populations at the FDR of 0.1. The analysis clearly
 10 highlights novel T2D-associated genes and pathways beyond what has been known
 11 from single-SNP association analysis reported from largest GWAS to date. Since the
 12 new procedure requires only SNP-level summary statistics, it provides a flexible
 13 way for conducting pathway analysis, alleviating the burden of handling large
 14 volumes of individual-level GWAS data.

15

16 We have developed a computationally efficient R package called ARTP2
 17 implementing the ARTP and sARTP procedures, so that it can be used for conducting
 18 pathway analysis based on individual-level genetic data, as well as SNP-level
 19 summary data from one or multiple GWAS. The R package also supports the
 20 parallelization on Unix-like OS, which can substantially accelerate the computation
 21 of small p-values when a large number of resampling steps are needed. The ARTP2
 22 package has a user-friendly interface and provides a comprehensive set of data
 23 preprocessing procedures to ensure that all the input information (e.g., allele

1 information of SNP-level summary statistics and genotype reference panel) can be
 2 processed coherently. To make the sARTP method accessible to a wider research
 3 community, we have also developed a web-based tool that allows investigators to
 4 conduct their pathway analyses using the computing resource at the National
 5 Cancer Institutes through simple on-line inputs of summary data.

6

7 Single-locus analysis of GWAS usually has its genomic control inflation factor larger
 8 than 1.0. Some proportion of the inflation can be attributed to various confounding
 9 biases, such as the one caused by population stratification, while the other part can
 10 be due to the real polygenic effect. In the pathway analysis it is important to
 11 minimize the confounding bias at the SNP-level summary statistic. Otherwise a
 12 small bias at the SNP level can be accumulated in the pathway analysis, and lead to
 13 an elevated false discovery rate. Here we try to remove the confounding bias by
 14 adjusting for the genomic control inflation factor observed at the GWAS study. This
 15 approach is conservative because part of the inflation can be caused by the real
 16 polygenic effect. Recently, [42] developed the LD score regression method to
 17 quantify the level of inflation caused solely by the confounding bias. Adjusting for
 18 the inflation factor estimated by this method, instead of the genomic control
 19 inflation factor, can potentially increase the power of the pathway analysis.
 20 However, the LD score regression method relies on a specific polygenic risk model,
 21 and its estimate might not be robust for this model assumption. More investigations
 22 are needed to evaluate the impact of this new inflation adjustment on the pathway
 23 analysis.

1

2 There are several other strategies to increase the power of pathway analysis besides
 3 increasing sample size [4]. One area of active research is to find better ways to
 4 define the gene-level summary statistic using observed genotypes on multiple SNPs,
 5 so that it can accurately characterize the impact of the gene on the outcome [43-46].
 6 In our proposed procedure, we adopt a data driven approach to select a subset of
 7 SNPs within a gene that collectively show the strongest association evidence.
 8 Because of this, we have to pay the penalty of multiple-comparison in the final
 9 pathway significance assessment. However, it is well recognized that SNPs at
 10 different loci can have varied levels of functional implications. We can potentially
 11 reduce the burden of multiple-comparisons and thus improve the power of the
 12 pathway analysis, by prioritizing SNPs according to existing genomic knowledge and
 13 other data resources. For example, [47] recently proposed a new gene-level
 14 summary statistic based on a prediction model that was trained with external
 15 transcriptome data. The gene-level summary statistic is defined as the predicted
 16 value that estimates the component of gene expression regulated by a subject's
 17 genotypes within the neighborhood of the considered gene. Pathway analysis
 18 procedures using this kind of biologically informed gene-level summary statistic can
 19 be easily incorporated into the ARTP2 framework.

20

21 The sARTP method can be easily expanded to adopt other multi-locus statistics in
 22 accumulating association within a gene, as long as they can be written in terms of
 23 SNP-level score statistics and their variance-covariance matrix. For example, the

1 current ARTP2 package provides the option for conducting the pathway meta-
2 analysis using the joint test statistics proposed by [31].

3

4 When conducting pathway analysis with individual-level genetic data, we could run
5 into a computing memory issue if the study has a large sample size and the pathway
6 consists of a large number of genes and SNPs (Figure S4). The ability of performing
7 pathway analysis using summary data provides a convenient and efficient solution
8 in those situations. We can first calculate the SNP-level summary statistics based on
9 the individual-level genetic data, and then randomly sample a small proportion of
10 the original data as an internal reference to estimate the variance-covariant matrix
11 for score statistics at considered SNPs. Based on our experiments, using 500 or
12 more subjects to form a reference panel would be good enough to generate accurate
13 pathway p-values. As shown in Figure 1, the testing results using this approach are
14 very consistent with those based on individual-level genotype data.

15

16 The sARTP approach can be applied directly to SNP-level meta-analysis results. This
17 is very convenient as meta-analysis results are in general easily accessible. But we
18 want to emphasize that it is important to know the set of the SNPs studied by each
19 participating study in order to apply sARTP properly, as the SNP coverage
20 information is essential for accurately estimating the variance-covariance matrix of
21 SNP-level score statistics. GWAS consortia usually do not post the SNP coverage
22 information when releasing their meta-analysis results. Many statistical packages
23 designed for conducting multi-locus analysis based on meta-analysis results often

1 assume the uniform coverage [15-18, 24, 25, 48]. As we already have demonstrated
2 in the context of pathway analysis, this type of over-simplification could lead to
3 inflated false positive rate.

4

5 The proposed procedure assumes that all participating studies are conducted with
6 subjects with the same ancestry background. If this is not the case, a simple
7 approach is to use the Fisher's method to combine pathway p-values estimated on
8 different ethnic populations. However, if there were no evidence for the existence of
9 cross ethnic risk heterogeneity, it would be more powerful to assume a fixed effects
10 model on the SNP-level association when performing the pathway analysis. In that
11 case, since the LD structures in different ethnic populations are different, we need a
12 separate reference panel for each ethnic group to derive the corresponding variance-
13 covariance matrix of the score statistics. The current ARTP2 package needs to be
14 modified to accommodate such a more complicated case.

15

16 As already demonstrated by many successful GWAS meta-analysis, increasing the
17 sample size through combining results from multiple studies is a very effective way
18 to improve our chance for new findings. For the same reason, pathway-based meta-
19 analysis can provide us with new opportunities to uncover biological pathways that
20 are previously undetectable due to the limitation on the sample size. With more
21 summary data from meta-analysis becoming increasingly available, we expect the
22 ARTP2 package would be a valuable tool for further exploring the genome in search
23 for the hidden heritability.

1

2 **Appendix A: Recovering Score Statistics and Its Variance-** 3 **Covariance Matrix Using Summary Results from the Fixed** 4 **Effect Model**

5 Here we derive the approximated score statistic S and its variance-covariance
6 matrix V using summary statistics from the fixed effect model. Based on (4), it is
7 straightforward to see that $S_t \approx \tau_t^{-2} \hat{\beta}_t$. Note that V_{ts} in equation (5) depends on $\tau_t^{(l)}$
8 estimated from individual studies, which cannot be derived from τ_t . However,
9 assume that $\sqrt{n_t^{(l)}} \tau_t^{(l)}$ can be approximated as an unknown but common constant
10 value ν_t across all studies, and if $\bar{y}^{(l)}(1 - \bar{y}^{(l)}) \approx \bar{y}(1 - \bar{y})$, we have $\nu_t \approx \sqrt{n_t} \tau_t$, and
11 $V_{ts} \approx \frac{n_{ts}}{\sqrt{n_t n_s}} \frac{\rho_{ts}}{\tau_t \tau_s}$. The similar argument has been used in Lin and Zeng (49) to
12 demonstrate that the meta-analysis is as efficient as the pooled analysis under those
13 conditions.

14

15 **Appendix B: Further Evaluation of sARTP Under the Null**

16 We conducted additional experiments to evaluate the empirical size of sARTP. Based
17 on the GERA T2D GWAS, we created 20 GWAS data under the null by randomly
18 permuting the outcome, while keeping individual genotypes unchanged. On each
19 null data, we excluded 274 pathways with over 10,000 SNPs for the sake of reducing
20 computational burden, and conducted a pathway-based meta-analysis with sARTP

on the remaining 4,439 pathways defined in MSigDB v5.0. The Q-Q plots of the pathway p-values of these 20 experiments are shown in Figure S50. Since there are extensive overlaps between pathways, their pathway p-values in each experiment are correlated. As a result, the Q-Q plot has a large variation around the diagonal line. But on average, there is no apparent genomic control inflation across 20 experiments. Based on those 20 experiments, Table S5 shows the genomic control inflation factors, Spearman's rank correlation coefficient between the pathway size (in terms of the number of unique SNPs, or genes in a pathway) and its pathway p-value. By inspecting those correlation coefficients, we did not see any evidence suggesting that the association significance level of a pathway is influenced by its size under the null.

Supporting Information

Table S1. Power comparison between sARTP and aSPUsPath under the scenario where each outcome-associated gene contains one functional SNP

Table S2. Power comparison between sARTP and aSPUsPath under the scenario where each outcome-associated gene contains one or two functional SNP(s) with equal probability

Table S3. Summary of top 50 genes with smallest gene-level p-values from the gene-level meta-analysis based on the DIAGRAM and GERA studies

1 Table S4. Effect of SNP rs1058018 on type 2 diabetes

2

3 Table S5. The genomic control inflation factors, Spearman's rank correlation
4 coefficients between the pathway size and its p-value based on results obtained by
5 applying sARTP to 20 simulated GWAS under the null

6

7 Figure S1. The sample size used for the study of each SNP in the DIAGRAM meta-
8 analysis

9

10 Figure S2: Histograms of numbers of SNPs and genes after SNP filtering within each
11 of 4,718 pathways in pathway analyses of the DIAGRAM study, the GERA study, and
12 the two studies combined (META)

13

14 Figure S3: Boxplot of the number of pathways containing genes with p-values in a
15 given range

16

17 Figure S4. The LocusZoom plot showing $\pm 100\text{kb}$ region of rs1058018 in European
18 populations

19

20 Figure S5: The LocusZoom plot showing $\pm 100\text{kb}$ region of rs1058018 in eastern
21 Asian populations

22

1 Figure S6: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
2 PENG_RAPAMYCIN_RESPONSE_DN.

3

4 Figure S7: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
5 YAGI_AML_WITH_T_8_21_TRANSLOCATION

6

7 Figure S8: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
8 PUJANA_CHEK2_PCC_NETWORK

9

10 Figure S9: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
11 STEIN_ESRRA_TARGETS

12

13 Figure S10: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
14 STEIN_ESRRA_TARGETS_UP

15

16 Figure S11: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
17 WANG_CISPLATIN_RESPONSE_AND_XPC_UP

18

19 Figure S12: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
20 CADWELL_ATG16L1_TARGETS_DN

21

22 Figure S13: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
23 CASORELLI_ACUTE_PROMYELOCYTIC_LEUKEMIA_DN

1

2 Figure S14: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

3 HILLION_HMGA1_TARGETS

4

5 Figure S15: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

6 HOLLEMAN_ASPARAGINASE_RESISTANCE_B_ALL_DN

7

8 Figure S16: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

9 PUJANA_BRCA1_PCC_NETWORK

10

11 Figure S17: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

12 GRAESSMANN_APOPTOSIS_BY_DOXORUBICIN_DN

13

14 Figure S18: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

15 SANSOM_APC_TARGETS_DN

16

17 Figure S19: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

18 ROPERO_HDAC2_TARGETS

19

20 Figure S20: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

21 REACTOME_INTEGRATION_OF_ENERGY_METABOLISM

22

1 Figure S21: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

2 AGUIRRE_PANCREATIC_CANCER_COPY_NUMBER_UP

3

4 Figure S22: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

5 HOLLEMAN_ASPARAGINASE_RESISTANCE_ALL_DN

6

7 Figure S23: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

8 REACTOME_CLASS_I_MHC_MEDIATED_ANTIGEN_PROCESSING_PRESENTATION

9

10 Figure S24: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

11 DODD_NASOPHARYNGEAL_CARCINOMA_DN

12

13 Figure S25: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

14 REACTOME_MEMBRANE_TRAFFICKING

15

16 Figure S26: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

17 SCHLOSSER_SERUM_RESPONSE_UP

18

19 Figure S27: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

20 PATIL_LIVER_CANCER

21

22 Figure S28: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

23 SONG_TARGETS_OF_IE86_CMV_PROTEIN

1

2 Figure S29: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

3 RIZ_ERYTHROID_DIFFERENTIATION

4

5 Figure S30: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

6 BORCZUK_MALIGNANT_MESOTHELIOMA_UP

7

8 Figure S31: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

9 KEGG_MATURITY_ONSET_DIABETES_OF_THE_YOUNG

10

11 Figure S32: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

12 HOSHIDA_LIVER_CANCER_SUBCLASS_S3

13

14 Figure S33: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

15 REACTOME_REGULATION_OF_BETA_CELL_DEVELOPMENT

16

17 Figure S34: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

18 PUJANA_BREAST_CANCER_WITH_BRCA1_MUTATED_UP

19

20 Figure S35: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

21 BLALOCK_ALZHEIMERS_DISEASE_UP

22

1 Figure S36: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
2 GOBERT_OLIGODENDROCYTE_DIFFERENTIATION_UP

3

4 Figure S37: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
5 MCBRYAN_PUBERTAL_BREAST_4_5WK_DN

6

7 Figure S38: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
8 REACTOME_REGULATION_OF_GENE_EXPRESSION_IN_BETA_CELLS

9

10 Figure S39: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
11 NABA_MATRISOME

12

13 Figure S40: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
14 PUJANA_BRCA2_PCC_NETWORK

15

16 Figure S41: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
17 KEGG_TYPE_II_DIABETES_MELLITUS

18

19 Figure S42: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
20 LINDGREN_BLADDER_CANCER_CLUSTER_1_DN

21

22 Figure S43: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
23 KEGG_INSULIN_SIGNALING_PATHWAY

1

2 Figure S44: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

3 CHEN_PDGF_TARGETS

4

5 Figure S45: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

6 PETROVA_ENDOTHELIUM_LYMPHATIC_VS_BLOOD_UP

7

8 Figure S46: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

9 REACTOME_PPARA_ACTIVATES_GENE_EXPRESSION

10

11 Figure S47: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

12 DACOSTA_UV_RESPONSE_VIA_ERCC3_UP

13

14 Figure S48: Q-Q plots for SNP p-values and sARTP gene p-values of pathway

15 TOYOTA_TARGETS_OF_MIR34B_AND_MIR34C

16

17 Figure S49. Comparison of sARTP p-values obtained with 2 or 3 SNP-level cut points

18 based on combined data of DIAGRAM and GERA

19

20 Figure S50. Q-Q plots of pathway p-values based on 20 GWAS datasets generated

21 under the null

22

23 **Acknowledgments**

1 This study utilized the computational resources of the NIH HPC Biowulf cluster
 2 (<https://hpc.nih.gov/>). The authors thank the AGEN-T2D consortium and DIAGRAM
 3 consortium for sharing the meta-analysis summary data. The authors acknowledge
 4 Sue Pan, JJ Pan, Wesley Obenshain, Tony Hall, Jim Zhou, Ye Wu, Cuong Nguyen for
 5 their help in developing the web-based tool for ARTP2.

6

7 **Web Resources**

8 The URLs for data and software presented herein are as follows:

9 DIAbetes Genetics Replication And Meta-analysis (DIAGRAMv3), [http://diagram-](http://diagram-consortium.org/)
 10 [consortium.org/](http://diagram-consortium.org/)

11 Genetic Epidemiology Research on Aging (GERA, dbGaP Study Accession:
 12 phs000674.v1.p1), [http://www.ncbi.nlm.nih.gov/projects/gap/cgi-](http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000674.v1.p1)
 13 [bin/study.cgi?study_id=phs000674.v1.p1](http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000674.v1.p1)

14 Molecular Signatures Database (C2: curated gene sets),
 15 <http://software.broadinstitute.org/gsea/msigdb/collections.jsp#C2>

16 BioMart (Homo sapiens genes NCBI36 and GRCh37.p13),
 17 <http://feb2014.archive.ensembl.org/>

18 IMPUTE2, https://mathgen.stats.ox.ac.uk/impute/impute_v2.html

19 GWAS Catalog, <http://www.ebi.ac.uk/gwas/>

20 1000 Genomes Project (Phase 3, v5, 2013/05/02),
 21 <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>

22 aSPU, <https://cran.r-project.org/web/packages/aSPU/index.html>

23 GTEx Portal v6, <http://gtexportal.org/home/>

- 1 GeneCards Human Gene Database, <http://www.genecards.org/>
- 2 Ingenuity Pathway Analysis, <http://www.ingenuity.com/>
- 3 LocusZoom , <http://locuszoom.sph.umich.edu/locuszoom/>
- 4 ARTP2 package, <https://cran.r-project.org/web/packages/ARTP2/>
- 5 Web-based tool of ARTP2, <http://analysis-tools.nci.nih.gov/pathway/>

6

7 **References**

8

- 9 1. Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, et al. The
10 NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic acids*
11 *research*. 2014;42(Database issue):D1001-6. doi: 10.1093/nar/gkt1229. PubMed
12 PMID: 24316577; PubMed Central PMCID: PMC3965119.
- 13 2. Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorff LA, Hunter DJ, et al.
14 Finding the missing heritability of complex diseases. *Nature*. 2009;461(7265):747-
15 53. doi: 10.1038/nature08494. PubMed PMID: 19812666; PubMed Central PMCID:
16 PMC2831613.
- 17 3. Moore JH, Asselbergs FW, Williams SM. Bioinformatics challenges for
18 genome-wide association studies. *Bioinformatics*. 2010;26(4):445-55. doi:
19 10.1093/bioinformatics/btp713. PubMed PMID: 20053841; PubMed Central
20 PMCID: PMC2820680.
- 21 4. Wang K, Li M, Hakonarson H. Analysing biological pathways in genome-wide
22 association studies. *Nature reviews Genetics*. 2010;11(12):843-54. doi:
23 10.1038/nrg2884. PubMed PMID: 21085203.
- 24 5. Chen LS, Hutter CM, Potter JD, Liu Y, Prentice RL, Peters U, et al. Insights into
25 colon cancer etiology via a regularized approach to gene set analysis of GWAS data.
26 *American journal of human genetics*. 2010;86(6):860-71. doi:
27 10.1016/j.ajhg.2010.04.014. PubMed PMID: 20560206; PubMed Central PMCID:
28 PMC3032068.
- 29 6. Evangelou M, Rendon A, Ouwehand WH, Wernisch L, Dudbridge F.
30 Comparison of methods for competitive tests of pathway analysis. *PloS one*.
31 2012;7(7):e41018. doi: 10.1371/journal.pone.0041018. PubMed PMID: 22859961;
32 PubMed Central PMCID: PMC3409204.
- 33 7. Li MX, Kwan JS, Sham PC. HYST: a hybrid set-based test for genome-wide
34 association studies, with application to protein-protein interaction-based
35 association analysis. *American journal of human genetics*. 2012;91(3):478-88. doi:
36 10.1016/j.ajhg.2012.08.004. PubMed PMID: 22958900; PubMed Central PMCID:
37 PMC3511992.

8. Pan W, Kwak IY, Wei P. A Powerful Pathway-Based Adaptive Test for Genetic Association with Common or Rare Variants. *American journal of human genetics*. 2015;97(1):86-98. doi: 10.1016/j.ajhg.2015.05.018. PubMed PMID: 26119817.
9. Yu K, Li Q, Bergen AW, Pfeiffer RM, Rosenberg PS, Caporaso N, et al. Pathway analysis by adaptive combination of P-values. *Genetic epidemiology*. 2009;33(8):700-9. doi: 10.1002/gepi.20422. PubMed PMID: 19333968; PubMed Central PMCID: PMC2790032.
10. Cho YS, Chen CH, Hu C, Long J, Ong RT, Sim X, et al. Meta-analysis of genome-wide association studies identifies eight new loci for type 2 diabetes in east Asians. *Nature genetics*. 2012;44(1):67-72. doi: 10.1038/ng.1019. PubMed PMID: 22158537; PubMed Central PMCID: PMC3582398.
11. DIABetes Genetics Replication Meta-analysis C, Consortium AGENTD, Consortium SATD, Consortium MATD, Consortium TDGEbN-gsim-ES, Mahajan A, et al. Genome-wide trans-ancestry meta-analysis provides insight into the genetic architecture of type 2 diabetes susceptibility. *Nature genetics*. 2014;46(3):234-44. doi: 10.1038/ng.2897. PubMed PMID: 24509480; PubMed Central PMCID: PMC3969612.
12. Imamura M, Takahashi A, Yamauchi T, Hara K, Yasuda K, Grarup N, et al. Genome-wide association studies in the Japanese population identify seven novel loci for type 2 diabetes. *Nature communications*. 2016;7:10531. doi: 10.1038/ncomms10531. PubMed PMID: 26818947.
13. Morris AP, Voight BF, Teslovich TM, Ferreira T, Segre AV, Steinthorsdottir V, et al. Large-scale association analysis provides insights into the genetic architecture and pathophysiology of type 2 diabetes. *Nature genetics*. 2012;44(9):981-90. doi: 10.1038/ng.2383. PubMed PMID: 22885922; PubMed Central PMCID: PMC3442244.
14. Speliotes EK, Willer CJ, Berndt SI, Monda KL, Thorleifsson G, Jackson AU, et al. Association analyses of 249,796 individuals reveal 18 new loci associated with body mass index. *Nature genetics*. 2010;42(11):937-48. doi: 10.1038/ng.686. PubMed PMID: 20935630; PubMed Central PMCID: PMC3014648.
15. Burren OS, Guo H, Wallace C. VSEAMS: a pipeline for variant set enrichment analysis using summary GWAS data identifies IKZF3, BATF and ESRRA as key transcription factors in type 1 diabetes. *Bioinformatics*. 2014;30(23):3342-8. doi: 10.1093/bioinformatics/btu571. PubMed PMID: 25170024; PubMed Central PMCID: PMC4296156.
16. Evangelou M, Smyth DJ, Fortune MD, Burren OS, Walker NM, Guo H, et al. A method for gene-based pathway analysis using genomewide association study summary statistics reveals nine new type 1 diabetes associations. *Genetic epidemiology*. 2014;38(8):661-70. doi: 10.1002/gepi.21853. PubMed PMID: 25371288; PubMed Central PMCID: PMC4258092.
17. Kwak IY, Pan W. Adaptive gene- and pathway-trait association testing with GWAS summary statistics. *Bioinformatics*. 2015. doi: 10.1093/bioinformatics/btv719. PubMed PMID: 26656570.
18. Lamparter D, Marbach D, Rueedi R, Kutalik Z, Bergmann S. Fast and Rigorous Computation of Gene and Pathway Scores from SNP-Based Summary Statistics. *PLoS computational biology*. 2016;12(1):e1004714. doi: 10.1371/journal.pcbi.1004714. PubMed PMID: 26808494; PubMed Central PMCID: PMC4726509.

19. Network and Pathway Analysis Subgroup of Psychiatric Genomics C. Psychiatric genome-wide association study analyses implicate neuronal, immune and histone pathways. *Nature neuroscience*. 2015;18(2):199-209. doi: 10.1038/nn.3922. PubMed PMID: 25599223; PubMed Central PMCID: PMC4378867.
20. Segre AV, Consortium D, investigators M, Groop L, Mootha VK, Daly MJ, et al. Common inherited variation in mitochondrial genes is not enriched for associations with type 2 diabetes or related glycemic traits. *PLoS genetics*. 2010;6(8). doi: 10.1371/journal.pgen.1001058. PubMed PMID: 20714348; PubMed Central PMCID: PMC2920848.
21. Swanson DM, Blacker D, Alchawa T, Ludwig KU, Mangold E, Lange C. Properties of permutation-based gene tests and controlling type 1 error using a summary statistic based gene test. *BMC genetics*. 2013;14:108. doi: 10.1186/1471-2156-14-108. PubMed PMID: 24199751; PubMed Central PMCID: PMC3831057.
22. Genomes Project C, Abecasis GR, Altshuler D, Auton A, Brooks LD, Durbin RM, et al. A map of human genome variation from population-scale sequencing. *Nature*. 2010;467(7319):1061-73. doi: 10.1038/nature09534. PubMed PMID: 20981092; PubMed Central PMCID: PMC3042601.
23. Genomes Project C, Abecasis GR, Auton A, Brooks LD, DePristo MA, Durbin RM, et al. An integrated map of genetic variation from 1,092 human genomes. *Nature*. 2012;491(7422):56-65. doi: 10.1038/nature11632. PubMed PMID: 23128226; PubMed Central PMCID: PMC3498066.
24. Hu YJ, Berndt SI, Gustafsson S, Ganna A, Genetic Investigation of ATC, Hirschhorn J, et al. Meta-analysis of gene-level associations for rare variants based on single-variant statistics. *American journal of human genetics*. 2013;93(2):236-48. doi: 10.1016/j.ajhg.2013.06.011. PubMed PMID: 23891470; PubMed Central PMCID: PMC3738834.
25. Liu JZ, McRae AF, Nyholt DR, Medland SE, Wray NR, Brown KM, et al. A versatile gene-based test for genome-wide association studies. *American journal of human genetics*. 2010;87(1):139-45. doi: 10.1016/j.ajhg.2010.06.009. PubMed PMID: 20598278; PubMed Central PMCID: PMC2896770.
26. Yang J, Ferreira T, Morris AP, Medland SE, Genetic Investigation of ATC, Replication DIG, et al. Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies additional variants influencing complex traits. *Nature genetics*. 2012;44(4):369-75. S1-3. doi: 10.1038/ng.2213. PubMed PMID: 22426310; PubMed Central PMCID: PMC3593158.
27. Goeman JJ, Buhlmann P. Analyzing gene expression data in terms of gene sets: methodological issues. *Bioinformatics*. 2007;23(8):980-7. doi: 10.1093/bioinformatics/btm051. PubMed PMID: 17303618.
28. Seaman SR, Muller-Myhsok B. Rapid simulation of P values for product methods and multiple-testing adjustment in association studies. *American journal of human genetics*. 2005;76(3):399-408. doi: 10.1086/428140. PubMed PMID: 15645388; PubMed Central PMCID: PMC1196392.
29. Hoffmann TJ, Kvale MN, Hesselson SE, Zhan Y, Aquino C, Cao Y, et al. Next generation genome-wide association tool: design and coverage of a high-throughput European-optimized SNP array. *Genomics*. 2011;98(2):79-89. doi:

- 1 10.1016/j.ygeno.2011.04.005. PubMed PMID: 21565264; PubMed Central PMCID:
2 PMC3146553.
- 3 30. Hoffmann TJ, Zhan Y, Kvale MN, Hesselson SE, Gollub J, Iribarren C, et al.
4 Design and coverage of high throughput genotyping arrays optimized for individuals
5 of East Asian, African American, and Latino race/ethnicity using imputation and a
6 novel hybrid SNP selection algorithm. *Genomics*. 2011;98(6):422-30. doi:
7 10.1016/j.ygeno.2011.08.007. PubMed PMID: 21903159; PubMed Central PMCID:
8 PMC3502750.
- 9 31. Zhang H, Shi J, Liang F, Wheeler W, Stolzenberg-Solomon R, Yu K. A fast
10 multilocus test with adaptive SNP selection for large-scale genetic-association
11 studies. *European Journal of Human Genetics*. 2014;22(5):696-702. doi:
12 10.1038/ejhg.2013.201. PubMed PMID: 24022295; PubMed Central PMCID:
13 PMC3992564.
- 14 32. Ge Y, Dudoit S, Speed T. Resampling-based multiple testing for microarray
15 data analysis. *Test*. 2003;12:1-77.
- 16 33. Devlin B, Roeder K, Wasserman L. Genomic control, a new approach to
17 genetic-based association studies. *Theoretical population biology*. 2001;60(3):155-
18 66. doi: 10.1006/tpbi.2001.1542. PubMed PMID: 11855950.
- 19 34. Marchini J, Howie B. Genotype imputation for genome-wide association
20 studies. *Nature reviews Genetics*. 2010;11(7):499-511. doi: 10.1038/nrg2796.
21 PubMed PMID: 20517342.
- 22 35. Wang T, Elston RC. Improved power by use of a weighted score test for
23 linkage disequilibrium mapping. *American journal of human genetics*.
24 2007;80(2):353-60. doi: 10.1086/511312. PubMed PMID: 17236140; PubMed
25 Central PMCID: PMC1785334.
- 26 36. Johnson ME, Zhao J, Schug J, Deliard S, Xia Q, Guy VC, et al. Two novel type 2
27 diabetes loci revealed through integration of TCF7L2 DNA occupancy and SNP
28 association data. *BMJ open diabetes research & care*. 2014;2(1):e000052. doi:
29 10.1136/bmjdr-2014-000052. PubMed PMID: 25469308; PubMed Central PMCID:
30 PMC4250976.
- 31 37. Lin CC, Chiang JH, Li CI, Liu CS, Lin WY, Hsieh TF, et al. Cancer risks among
32 patients with type 2 diabetes: a 10-year follow-up study of a nationwide population-
33 based cohort in Taiwan. *BMC cancer*. 2014;14:381. doi: 10.1186/1471-2407-14-
34 381. PubMed PMID: 24884617; PubMed Central PMCID: PMC4057814.
- 35 38. Tsilidis KK, Kasimis JC, Lopez DS, Ntzani EE, Ioannidis JP. Type 2 diabetes and
36 cancer: umbrella review of meta-analyses of observational studies. *Bmj*.
37 2015;350:g7607. doi: 10.1136/bmj.g7607. PubMed PMID: 25555821.
- 38 39. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate - a Practical
39 and Powerful Approach to Multiple Testing. *J Roy Stat Soc B Met*. 1995;57(1):289-
40 300. PubMed PMID: WOS:A1995QE45300017.
- 41 40. Patil MA, Chua MS, Pan KH, Lin R, Lih CJ, Cheung ST, et al. An integrated data
42 analysis approach to characterize genes highly expressed in hepatocellular
43 carcinoma. *Oncogene*. 2005;24(23):3737-47. doi: 10.1038/sj.onc.1208479. PubMed
44 PMID: 15735714.
- 45 41. Hoshida Y, Nijman SM, Kobayashi M, Chan JA, Brunet JP, Chiang DY, et al.
46 Integrative transcriptome analysis reveals common molecular subclasses of human

hepatocellular carcinoma. *Cancer research*. 2009;69(18):7385-92. doi: 10.1158/0008-5472.CAN-09-1089. PubMed PMID: 19723656; PubMed Central PMCID: PMC3549578.

42. Bulik-Sullivan BK, Loh PR, Finucane HK, Ripke S, Yang J, Schizophrenia Working Group of the Psychiatric Genomics C, et al. LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nature genetics*. 2015;47(3):291-5. doi: 10.1038/ng.3211. PubMed PMID: 25642630; PubMed Central PMCID: PMC4495769.

43. Li M, Wang K, Grant SF, Hakonarson H, Li C. ATOM: a powerful gene-based association test by combining optimally weighted markers. *Bioinformatics*. 2009;25(4):497-503. doi: 10.1093/bioinformatics/btn641. PubMed PMID: 19074959; PubMed Central PMCID: PMC2642636.

44. Wessel J, Schork NJ. Generalized genomic distance-based regression methodology for multilocus association analysis. *American journal of human genetics*. 2006;79(5):792-806. doi: 10.1086/508346. PubMed PMID: 17033957; PubMed Central PMCID: PMC1698575.

45. Wu MC, Kraft P, Epstein MP, Taylor DM, Chanock SJ, Hunter DJ, et al. Powerful SNP-set analysis for case-control genome-wide association studies. *American journal of human genetics*. 2010;86(6):929-42. doi: 10.1016/j.ajhg.2010.05.002. PubMed PMID: 20560208; PubMed Central PMCID: PMC3032061.

46. Zhang H, Wheeler W, Wang Z, Taylor PR, Yu K. A fast and powerful tree-based association test for detecting complex joint effects in case-control studies. *Bioinformatics*. 2014;30(15):2171-8. doi: 10.1093/bioinformatics/btu186. PubMed PMID: 24794927; PubMed Central PMCID: PMC4103596.

47. Gamazon ER, Wheeler HE, Shah KP, Mozaffari SV, Aquino-Michaels K, Carroll RJ, et al. A gene-based association method for mapping traits using reference transcriptome data. *Nature genetics*. 2015;47(9):1091-8. doi: 10.1038/ng.3367. PubMed PMID: 26258848; PubMed Central PMCID: PMC4552594.

48. Mishra A, Macgregor S. VEGAS2: Software for More Flexible Gene-Based Testing. *Twin research and human genetics : the official journal of the International Society for Twin Studies*. 2015;18(1):86-91. doi: 10.1017/thg.2014.79. PubMed PMID: 25518859.

49. Lin DY, Zeng D. On the relative efficiency of using summary statistics versus individual-level data in meta-analysis. *Biometrika*. 2010;97(2):321-32. doi: 10.1093/biomet/asq006. PubMed PMID: 23049122; PubMed Central PMCID: PMC3412575.