

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

4

5

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

8

9

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

12

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

14

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

16

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

19

20 * Corresponding Authors

21

22 Kai Yu (yuka@mail.nih.gov)

23 Nilanjan Chatterjee (nchatte2@jhu.edu)

24

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 meta-analysis with sARTP on type 2 diabetes (T2D) by
13 integrating SNP-level summary statistics from two large studies consisting of 19,809
14 T2D cases and 111,181 controls with European ancestry. Among 4,713 candidate
15 pathways from which genes in neighborhoods of 170 GWAS established T2D loci
16 were excluded, we detected 43 T2D globally significant pathways (with Bonferroni
17 corrected p-values < 0.05), which included the insulin signaling pathway and T2D
18 pathway defined by KEGG, as well as the pathways defined according to specific
19 gene expression patterns on various tumor types, including pancreatic
20 adenocarcinoma, hepatocellular carcinoma, and bladder carcinoma. Using summary
21 data from eight eastern Asian T2D GWAS with 6,952 cases and 11,865 controls, we
22 showed 7 out of the 43 pathways identified in European populations remained to be
23 significant in eastern Asians at the false discovery rate of 0.1. We created a R

1 package and a web-based tool for sARTP with the capability to analyze pathways
2 with thousands of genes and tens of 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 eight Asian studies showed that some pathways were trans-
19 ethnically associated with T2D. This analysis clearly highlights novel T2D-associated
20 pathways beyond what has been known from single-variant association analysis
21 reported 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 sets) 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 in practice can
23 rarely identify new findings beyond what has already been discovered by the single-

1 locus tests. To maximize the chance of discovering novel outcome-associated
2 variants by increasing sample size, a number of consortia have been formed to
3 conduct single-locus meta-analysis on data across multiple GWAS [10-14]. The
4 single-locus meta-analysis aggregates easily accessible SNP-level summary statistics
5 from multiple studies. Similarly, the pathway-based meta-analysis [15-21] that
6 integrates the same type of summary data across participating studies could provide
7 us a greater opportunity for detecting novel pathway associations.

8

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

1 different. Ignoring the difference in sample sizes across SNPs in a pathway-based
2 meta-analysis would generate biased testing results.

3

4 Pathway analysis generally targets two types of null hypotheses [4], including the
5 competitive null hypothesis [15,16,18-20], i.e., the genes in a pathway of interest are
6 no more associated with the outcome than any other genes outside this pathway,
7 and the self-contained null hypothesis [17,21], i.e., none of the genes in a pathway is
8 associated with the outcome. The sARTP procedure focuses on the self-contained
9 null hypothesis, as our main goal is to identify outcome-associated genes or loci. One
10 may refer to Goeman and Buhlmann [27] and Wang et al. [4] for more discussions
11 and comparisons of these two types of hypotheses.

12

13 The pathways defined in many public databases can consist of thousands of genes
14 and tens of thousands of SNPs. To make the procedure applicable to large pathways,
15 or pathways with high statistical significance, we implement sARTP with very
16 efficient and parallelizable algorithms, and adopt the direct simulation approach [28]
17 to evaluate the significance of the pathway association.

18

19 We demonstrated the validity and power advantage of sARTP through simulated
20 and empirical data. We applied sARTP to conduct a pathway-based meta-analysis on
21 the association between type 2 diabetes (T2D) and 4,713 candidate pathways
22 defined in the Molecular Signatures Database (MSigDB) v5.0. The analysis used SNP-
23 level summary statistics from two sources with European ancestry. One is generated

1 from the Diabetes Genetics Replication and Meta-analysis (DIAGRAM) consortium
2 [13], which consists of 12,171 T2D cases and 56,862 controls across 12 GWAS. The
3 other one is based on a T2D GWAS with 7,638 T2D cases and 54,319 controls that
4 were extracted from the Genetic Epidemiology Research on Aging (GERA) study
5 [29,30]. The T2D-associated pathways detected in the European population were
6 further examined in Asians using summary data generated by the Asian Genetic
7 Epidemiology Network (AGEN) consortium meta-analysis, which combined eight
8 GWAS of T2D with a total of 6,952 and 11,865 controls from eastern Asian
9 populations [10].

10

11 **Material and Methods**

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

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

17

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

19 We assume we have data from L GWA studies, with each consisting of $n^{(l)}$ subjects,
20 $l = 1, \dots, L$. Each gene in that pathway can contain one or multiple SNP(s), while any
21 two genes may have some overlapped SNPs. For simplicity, we use superscript l to
22 represent an individual study. For subject i in study l , $i = 1, \dots, n^{(l)}$, let $y_i^{(l)}$ be the

1 dichotomous outcome (e.g., disease condition, case/control status) taking values
 2 from $\{0,1\}$, and let $X_i^{(l)}$ be the vector of covariates to be adjusted for. The
 3 centralized genotypes of q SNPs within a pathway are presented as a vector
 4 $G_i^{(l)} = (g_{i1}^{(l)}, \dots, g_{iq}^{(l)})^T$ for subject i . We assume the following logistic regression model
 5 as the risk model

$$6 \quad \text{logit } P(y_i^{(l)} = 1 | X_i^{(l)}, G_i^{(l)}) = (X_i^{(l)})^T \alpha^{(l)} + (G_i^{(l)})^T \gamma, \quad i = 1, \dots, n,$$

7 Under the self-contained null hypothesis $H_0: \gamma = 0$, we denote the maximum
 8 likelihood estimate of $\alpha^{(l)}$ as $\hat{\alpha}^{(l)}$. Let $\hat{y}_i^{(l)} = 1 / (1 + \exp(-X_i^{(l)} \hat{\alpha}^{(l)}))$ and
 9 $u_i^{(l)} = \hat{y}_i^{(l)}(1 - \hat{y}_i^{(l)})$. The Rao's score statistic vector on γ , which is the sum of score
 10 vectors from L participating studies, follows the asymptotic multivariate normal
 11 distribution $N(0, V)$, where

$$12 \quad S = (S_t)_{q \times 1} = \sum_{l=1}^L \sum_{i=1}^{n^{(l)}} G_i^{(l)} (y_i^{(l)} - \hat{y}_i^{(l)}) \quad (1)$$

13 and

$$14 \quad V = \sum_{l=1}^L \left(\sum_{i=1}^{n^{(l)}} u_i^{(l)} G_i^{(l)} (G_i^{(l)})^T \right. \\ \left. - \sum_{i=1}^{n^{(l)}} u_i^{(l)} G_i^{(l)} (X_i^{(l)})^T \left(\sum_{i=1}^{n^{(l)}} u_i^{(l)} X_i^{(l)} (X_i^{(l)})^T \right)^{-1} \sum_{i=1}^{n^{(l)}} u_i^{(l)} X_i^{(l)} (G_i^{(l)})^T \right). \quad (2)$$

15 For study l , let $n_t^{(l)}$ be the number of subjects having their genotypes measured as
 16 $H_t^{(l)}$ (or imputed) at SNP t , where $H_t^{(l)} = (g_{1t}^{(l)}, \dots, g_{n_t^{(l)}}^{(l)})^T$. As pointed out by Hu et al.

1 [24] if the covariates and genotypes are uncorrelated or weakly correlated, the
 2 covariance between scores at SNPs t and s can be approximated as

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

4 where $n_{ts}^{(l)}$ is the number of samples that have their genotypes available at both
 5 SNPs in study l , $\overline{u}^{(l)} = (n^{(l)})^{-1} \sum_{i=1}^{n^{(l)}} u_i^{(l)}$, and $\widehat{\text{Cov}}(H_t^{(l)}, H_s^{(l)}) = (n^{(l)})^{-1} \sum_{i=1}^{n^{(l)}} g_{it}^{(l)} g_{is}^{(l)}$. Here, we
 6 assume that the Pearson's correlation coefficient ρ_{ts} between two SNPs is the same
 7 among all participating studies. This assumption is valid as long as subjects from all
 8 studies are sampled from the same source population, or the population under
 9 study is relatively homogeneous, such as a study of subjects with European ancestry
 10 in the United States.

11
 12 When only the summary statistics, i.e., the estimated marginal log odds ratios $\hat{\beta}_t^{(l)}$
 13 and their standard errors $\tau_t^{(l)}$ are available for each of the L studies, the score
 14 statistic at SNP t , defined by (1) can be approximated as

$$15 \quad S_t \approx \sum_{l=1}^L (\tau_t^{(l)})^{-2} \hat{\beta}_t^{(l)}; \quad t = 1, \dots, q. \tag{4}$$

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

$$17 \quad V_{ts} \approx \sum_{l=1}^L \frac{n_{ts}^{(l)}}{\sqrt{n_t^{(l)} n_s^{(l)}}} \rho_{ts} \tau_t^{(l)} \tau_s^{(l)}. \tag{5}$$

1 Assume that ρ_{ts} can be estimated from a public dataset (e.g., 1000 Genomes Project)
 2 and the sample sizes $n_t^{(l)}$ and $n_s^{(l)}$ are known, we can approximately recover the
 3 variance-covariance matrix $V = (V_{ts})_{q \times q}$ of score statistics $S = (S_t)_{q \times 1}$. In cases when
 4 we only have the SNP p-value p and its marginal log odds ratio $\hat{\beta}$, we can compute
 5 its standard error as $\tau = |\hat{\beta}| / \sqrt{\chi_{1,p}^2}$, where $\chi_{1,p}^2$ is the quantile satisfying
 6 $P(\chi_1^2 \geq \chi_{1,p}^2) = p$, with χ_1^2 representing a 1-df chi-squared random variable.

7
 8 Many GWAS consortia usually publish their meta-analysis results by providing only
 9 the combined results from the fixed effect model, rather than the summary statistics
 10 from each participating study. The reported marginal log odds ratios for each SNP
 11 by using the fixed-effect inverse-variance weighting method is given by

$$12 \quad \hat{\beta}_t = \frac{\sum_{l=1}^L (\tau_t^{(l)})^{-2} \hat{\beta}_t^{(l)}}{\sum_{l=1}^L (\tau_t^{(l)})^{-2}},$$

13 with its standard error given by

$$14 \quad \tau_t = \left(\sum_{l=1}^L (\tau_t^{(l)})^{-2} \right)^{-1/2}. \quad (6)$$

15 Based on (4), we have $S_t \approx \tau_t^{-2} \hat{\beta}_t$. Assuming large sample sizes and certain
 16 conditions (see Appendix A), we can also approximate the covariance between S_t
 17 and S_s as

$$18 \quad V_{ts} \approx \frac{n_{ts}}{\sqrt{n_t n_s}} \frac{\rho_{ts}}{\tau_t \tau_s}, \quad (7)$$

1 where $n_t = \sum_{l=1}^L n_t^{(l)}$, and $n_{ts} = \sum_{l=1}^L n_{ts}^{(l)}$. Thus, given the summary results from the meta-
2 analysis, we can still approximately recover S_t and V_{ts} .

3

4 ***Combining Score Statistics for Pathway Analysis***

5 With recovered score statistics vector S and its variance-covariance matrix V , we
6 can conduct a pathway association test using the framework of the ARTP method.
7 The ARTP method first combines p-values of individual SNPs within a gene to form a
8 gene-based association statistic (i.e., the gene-level p-value), and then combines the
9 gene-level p-values into a final testing statistic for the pathway-outcome association.
10 In the original ARTP method, Yu et al. [9] proposed the use of a resampling-based
11 method to evaluate the significance level of the pathway association test. Here we
12 adopt the use of score statistics into the ARTP framework and propose the use of
13 DSA [28] to evaluate the significance level, which would be much faster than the
14 original ARTP algorithm [31]. Below is a brief summary of the improved ARTP
15 algorithm.

16

17 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

18 $p_t^{(0)} = P(\chi_1^2 \geq S_t^2 / V_{tt})$. Let $p_{j(1)}^{(0)}, \dots, p_{j(q_j)}^{(0)}$ be their order statistics such that

19 $p_{j(1)}^{(0)} \leq \dots \leq p_{j(q_j)}^{(0)}$. For any predefined integer K and SNP-level cut points $c_1 < \dots < c_K$,

20 we define the observed negative log product statistics for that gene at cut point c_k

21 as

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

2 We sample M copies of vectors of the score statistic from the null distribution

3 $N(0, V)$ and convert each of them to be the tail probability of χ_1^2 as $p_{i_1}^{(m)}, \dots, p_{i_{q_j}}^{(m)}$,

4 $m = 1, \dots, M$, which are then used to calculate $w_{jk}^{(m)}$, $m = 1, \dots, M$. The significance of

5 $w_{jk}^{(0)}$ can be estimated as

$$6 \quad \xi_{jk}^{(0)} = \frac{\#\{w_{jk}^{(m)} \geq w_{jk}^{(0)}; m = 1, \dots, M\}}{M + 1}.$$

7

8 The ARTP statistic for testing association between gene j and the outcome is

9 defined as $T_j^{(0)} = \min_{k=1, \dots, K} \xi_{jk}^{(0)}$. Note that for any $w_{jk}^{(m)}$, the set

10 $\{w_{jk}^{(m')} : m' \in \{0, \dots, M\} \text{ and } m' \neq m\}$ forms its empirical null distribution. The

11 significance of $w_{jk}^{(m)}$ therefore can be estimated as

$$12 \quad \xi_{jk}^{(m)} = \frac{\#\{w_{jk}^{(m')} \geq w_{jk}^{(m)}; m' \neq m \text{ and } m' = 1, \dots, M\}}{M + 1}, \quad m = 1, \dots, M.$$

13 This idea, which was given by Ge et al. [32] can be used to avoid the computationally

14 challenging nested two-layer resampling procedure for evaluating p-values. The p-

15 value of $T_j^{(0)}$ can be readily calculated as

$$16 \quad z_j^{(0)} = \frac{\#\{T_j^{(m)} \leq T_j^{(0)} : m = 1, \dots, M\}}{M + 1}, \quad j = 1, \dots, J.$$

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

10

11 In the following discussion, we will use the term sARTP to represent the proposed
12 pathway analysis procedure using the SNP-level summary statistics as input, and
13 reserve the term ARTP to represent the original ARTP procedure that requires the
14 individual-level genetic data. Both procedures adopt the DSA algorithm to accelerate
15 evaluating the significance level. When using sARTP and ARTP for the pathway
16 analysis, we set SNP-level cut points as $(c_1, c_2) = (1, 2)$, and gene-level cut points as
17 $d_k = k \max(1, \lceil J/20 \rceil)$, $k = 1, \dots, 10$, where J is the number of genes in a pathway,
18 and $\lceil J/20 \rceil$ is the largest integer that is less or equal to $J/20$. We used $M = 10^5$
19 DSA steps to assess the significance level of each pathway in the initial screening.
20 For pathways with estimated p-values $< 10^{-4}$, we further refined their p-value
21 estimates with $M = 10^7$ or 10^8 DSA steps.

1

2 **Study Materials**

3 *Pathway and Gene Definition*

4 We downloaded definitions for 4,716 human and murine (mammalian) pathways
5 (gene sets) from the MSigDB v5.0 (C2: curated gene sets). Genomic definitions for all
6 genes were downloaded from Homo sapiens genes NCBI36 and reference genome
7 GRCh37.p13 using the Ensemble BioMart tool.

8

9 *DIAGRAM Study*

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

22

1 We sorted 2.5 million autosomal SNPs by their corresponding meta-analysis sample
2 sizes in Figure S1, which shows that there are two major groups of SNPs with equal
3 sample sizes. One group of 469,985 SNPs (19.0%) had 12,171 cases and 56,862
4 controls, which included all the available samples in the meta-analysis; another
5 group of 1,431,361 SNPs (57.9%) had 9,580 cases and 53,810 controls. Since the
6 calculation of covariance V_{ts} in (7) relies on n_{ts} , the number of samples having
7 genotypes available at both SNP s and SNP t , in order to have an accurate estimate
8 of n_{ts} , we focused on these two groups of SNPs, which in combination had a total of
9 1,901,346 SNPs. For any two SNPs in this reduced set, it is certain $n_{ts} = \min(n_t, n_s)$.
10 The Pearson's correlation coefficients ρ_{ts} were estimated using an external
11 reference panel consisting of genotypes on 503 European subjects from the 1000
12 Genomes Project (Phase 3, v5, 2013/05/02).

13

14 ***GERA Study***

15 We assembled a GWAS on T2D from the Genetic Epidemiology Research on Adult
16 Health and Aging (GERA, dbGaP Study Accession: phs000674.v1.p1). The GERA
17 project includes a cohort of over 100,000 adults who are members of the Kaiser
18 Permanente Medical Care Plan, Northern California Region, and participating in the
19 Kaiser Permanente Research Program on Genes, Environment, and Health (RPGEH).
20 From the GERA data, we compiled a GWAS with 7,638 T2D cases and 54,319
21 controls (subjects without T2D) who self-reported to be non-Hispanic White
22 Europeans in the RPGEH survey. We performed the genotype imputation with
23 IMPUTE2 [34] using CEU reference samples from Phase II of the International

1 HapMap Project. After removing SNPs with low imputation quality ($r^2 < 0.3$), we
2 ended up with 2.4 million SNPs for further analysis. In the single-locus analysis, we
3 adjusted for the categorized body mass index (BMI) provided in the downloaded
4 dataset (adding a category for missing BMI), gender, year of birth (in five-year
5 categories), a binary indicator on whether or not a participant was diagnosed with
6 cancer (includes malignant tumors, neoplasms, lymphoma and sarcoma), and the
7 top five eigenvectors for the adjustment of population stratification. In the following
8 discussion, we refer this assembled T2D GWAS as the GERA study.

9

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

13

14 ***AGEN-T2D Study***

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

1 1000 Genomes Project for the variance-covariance matrix estimation in the pathway
2 analysis.

3

4 **Results**

5 **Simulation Studies**

6 Firstly, we conducted a simulation study to evaluate the empirical size of sARTP and
7 ARTP. Secondly, we compared empirical powers of different strategies for carrying
8 out pathway-based meta-analysis that integrated summary statistics from multiple
9 studies. We also evaluated whether results from sARTP were consistent with the
10 ones from ARTP. Thirdly, we compared our method to the recently developed
11 method aSPUpath [8] in the scenarios where there was individual-level genotype
12 data from one study. The aSPUpath method has been shown through simulation
13 studies to be more powerful than several existing pathway association tests [8]. We
14 used the R package, aSPU (version 1.39), with the default settings suggested in their
15 papers to conduct the aSPUpath test.

16

17 ***Empirical Size of sARTP***

18 To evaluate the empirical size of sARTP, we conducted a simulation study by using
19 individual-level GWAS data of the pathway
20 PUJANA_BREAST_CANCER_WITH_BRCA1_MUTATED_UP (including 777 SNPs in 45
21 genes) from the GERA study. We picked 10,000 samples randomly, on which we
22 simulated 10,000 case-control study datasets. For each dataset, we randomly
23 assigned 5,000 samples as cases and the other 5,000 samples as controls, while

1 keeping their original genotypes in the pathway unchanged. For each generated
2 dataset, we first calculated the SNP-level summary statistics, and then applied the
3 sARTP procedure with the variance-covariance matrix estimated by either an
4 external reference panel (i.e., 503 European reference samples from the 1000
5 Genomes Project), or an internal reference panel (i.e., 500 randomly selected GERA
6 samples). We also analyzed each simulated study by applying the ARTP procedure
7 to the individual-level genetic data. Based on results from the 10,000 generated
8 datasets, the simulation study showed that at all three approaches can properly
9 control their sizes (Table 1).

10

11 Table 1: Empirical sizes of the sARTP and ARTP procedure

Method	Size		
	0.05	0.01	0.005
ARTP ^a	0.0518	0.0112	0.0058
sARTP (Internal reference) ^b	0.0518	0.0103	0.0044
sARTP (External reference) ^c	0.0506	0.0089	0.0047

12 Empirical sizes are estimated based on 10,000 datasets simulated from the GERA
13 study, with each dataset consisting of 5,000 cases and 5,000 controls.

14 ^aARTP method using individual-level GWAS data as input;

15 ^bsARTP method using summary statistics calculated from the simulated GWAS
16 dataset as input, in combination with an internal reference panel;

17 ^csARTP method using summary statistics calculated from the simulated GWAS data
18 as input, in combination with an external reference panel.

19

20 ***Empirical Power of sARTP for Pathway-based Meta-analysis***

21 We conducted a set of simulation studies to investigate the power of pathway meta-
22 analysis using sARTP when SNP-level summary statistics were available from two
23 GWA studies. We considered a hypothetical pathway consisting of 50 genes
24 randomly selected from chromosome 17, each with 20 randomly chosen SNPs. The
25 joint genotype distribution at the 20 SNPs within a gene was defined by the

1 observed genotypes in the GERA study. We further assumed that all genes in that
2 pathway are independent. This assumption is not required by sARTP and ARTP, but
3 it is introduced for simplifying the simulation. For the risk model, we assumed the
4 first \mathcal{M} ($\mathcal{M} = 5, 10, 15$) genes were associated with the outcome. Within each
5 outcome associated gene, we picked the SNP with MAF closest to the median MAF
6 level within the gene to be functional. We considered the following risk model

$$7 \quad \text{logit } P(y = 1 | g_1^*, \dots, g_{\mathcal{M}}^*) = \alpha + \sum_{l=1}^{\mathcal{M}} \gamma_l^* g_l^*, \quad (8)$$

8 where g_l^* is the genotype (encoded as 0, 1, or 2 according to counts of minor alleles)
9 at the functional SNP within gene l . Under this model, γ_l^* is also the marginal log
10 odds ratio for the l th functional SNP [9]. Given the sample sizes of cases and
11 controls, and the MAF of the l th functional SNP, γ_l^* was chosen such that the
12 theoretical power of the trend test to detect the l th functional SNP is equal to \mathcal{P}
13 ($\mathcal{P} = 0.3, 0.4$), with 0.05 as the targeted type I error rate. For every pair of $(\mathcal{M}, \mathcal{P})$,
14 we generated 1,000 datasets, each consisting of two case-control studies, with each
15 study consisting of 1,000 cases and 1,000 controls. Given the genotype distribution
16 in the general population, individual-level genetic data for a case-control study can
17 be generated according to the assumed risk model (8).

18

19 We analyzed each simulated dataset in several different ways. First, we applied
20 ARTP to the individual-level genotype data combined from two studies, and
21 considered its results as the gold standard. A study indicator was added to the null
22 model when computing the score statistics. Next, we assumed that only SNP-level

1 summary statistics from each of the two studies were available. We applied sARTP
2 to conduct the pathway meta-analysis by combining summary statistics from two
3 studies with formulas (4) and (5), using either the internal or external reference
4 panel for the variance-covariance matrix estimation. The sARTP approach integrates
5 association evidence across SNP-level summary statistics, which are obtained by
6 pooling information from all participating studies on individual SNPs. As a
7 comparison, we also considered a naïve meta-analysis approach, in which we first
8 applied sARTP to analyze the summary statistics from each study separately, and
9 then combined the two pathway p-values with Fisher's method. The empirical
10 powers are computed at the type I error level of 0.05, and are summarized in Table
11 2. It is obvious that the pathway meta-analysis using sARTP, with either the internal
12 or external reference panel, have almost the same level of power as the ARTP
13 method using the individual-level genetic data. It is also evident that the sARTP
14 approach is much more powerful than the naïve approach using the Fisher's method
15 to combine pathway association p-values from two studies.

16

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

\mathcal{P}^a	\mathcal{M}^b	Study 1 + Study 2					Study 1		Study 2	
		ARTP ^c	sARTP		Fisher's method		sARTP		sARTP	
			Internal ref ^d	External ref ^e	Internal ref ^f	External ref ^g	Internal ref ^h	External ref ⁱ	Internal ref ^j	External ref ^k
0.3	5	0.163	0.141	0.141	0.076	0.07	0.072	0.067	0.064	0.068
	10	0.335	0.322	0.322	0.155	0.144	0.114	0.114	0.118	0.115
	15	0.555	0.523	0.523	0.247	0.236	0.177	0.173	0.173	0.168
0.4	5	0.285	0.252	0.252	0.116	0.111	0.100	0.095	0.101	0.106
	10	0.579	0.554	0.554	0.257	0.259	0.170	0.173	0.189	0.186
	15	0.826	0.804	0.804	0.428	0.438	0.281	0.284	0.292	0.283

2 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. Every
 3 dataset contains two studies, with each consisting of 1,000 cases and 1,000 controls. The pathway has 50 independent genes,
 4 each with 20 SNPs. Fisher's method is used to combine the two pathway p-values obtained by applying sARTP to the SNP-level
 5 summary data from study 1 and study 2 separately.

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

8 ^bThe number of genes including the functional SNPs;

9 ^cARTP method using individual-level genotype data from two studies as input;

10 ^dsARTP method using summary data calculated from two studies as input, in combination with the internal reference panel;

11 ^esARTP method using summary data calculated from two studies as input, in combination with the external reference panel;

12 ^fFisher's method combining the p-values from sARTP using the internal reference panel;

13 ^gFisher's method combining the p-values from sARTP using the external reference panel;

14 ^hsARTP method using summary data calculated from the first study as input, in combination with the internal reference panel;

15 ⁱsARTP method using summary data calculated from the first study as input, in combination with the external reference panel;

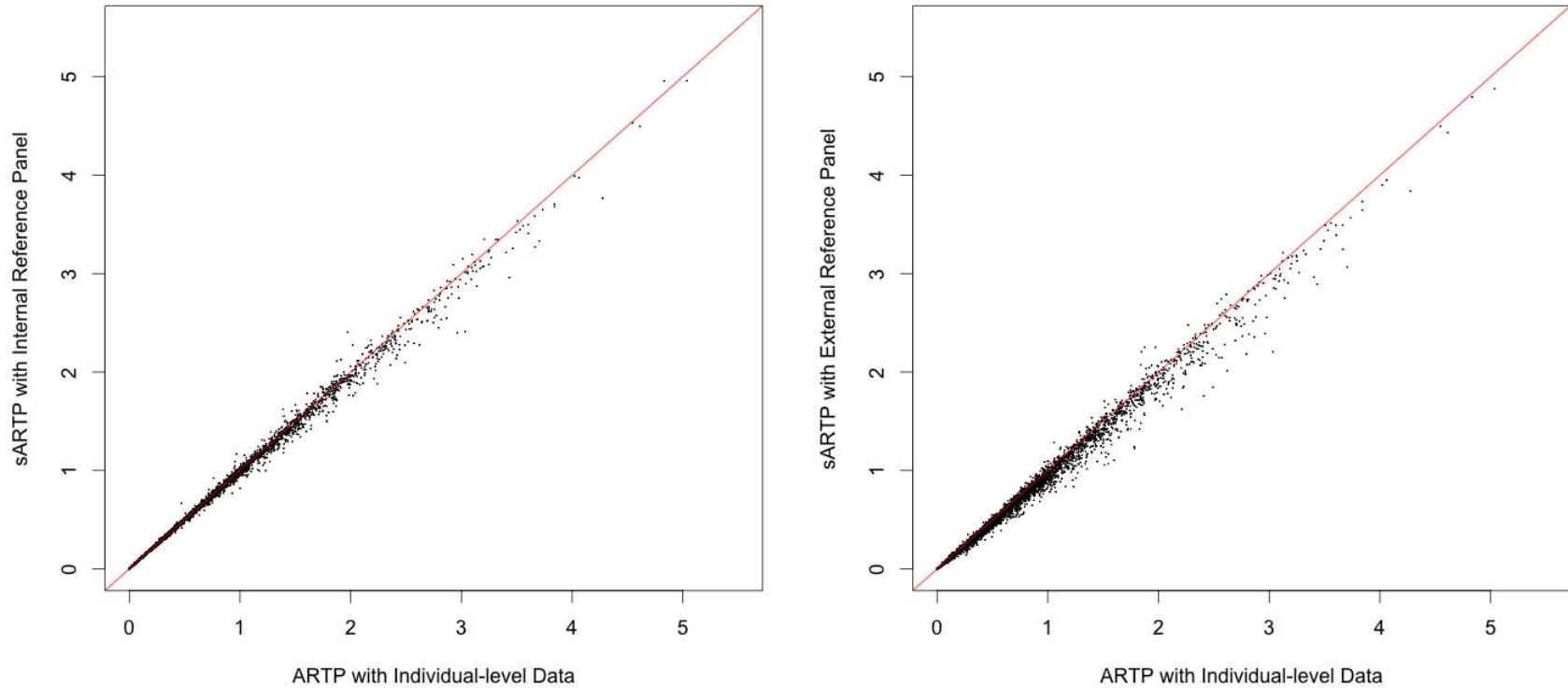
16 ^jsARTP method using summary data calculated from the second study as input, in combination with the internal reference
 17 panel;

18 ^ksARTP method using summary data calculated from the second study as input, in combination with the external reference
 19 panel.

1 To further demonstrate the consistency between results obtained by sARTP and the
2 ones by ARTP, we compared pathway analysis results from three different
3 procedures on the 4,713 candidate pathways using the GERA GWAS data. Details on
4 how those 4,713 pathways were pre-processed are given in the Results of T2D
5 Pathway Analysis Section. We applied sARTP to the SNP-level summary statistics
6 generated from the GERA study, using either an internal or an external reference
7 panel for the variance-covariance matrix estimation. We also obtained the pathway
8 p-values by directly applying the ARTP method to the individual-level GERA GWAS
9 data. Figure 1 shows the comparison among p-values from the three analyses, and
10 demonstrates that all three approaches can generate very consistent results.

11

12



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

3
4 Based on the GERA data, 4,713 pathways are analyzed in three different ways. Pathway p-values obtained by ARTP using the
5 GERA individual-level genetic data (x-axis) are compared with the ones obtained by sARTP using summary statistics in
6 combination with the internal reference panel that consists of 500 randomly selected GERA samples (left), and the ones using
7 the summary statistics in combination with the external reference panel that consists of 503 European subjects from the 1000
8 Genomes Project (right).

1 ***Comparison with aSPUpath method***

2 As we have shown in the previous section, ARTP and sARTP methods generate very
3 similar results (with Pearson's correlation coefficient of their p-values > 0.98). The
4 aSPUpath method also has a version, called aSPUsPath [17], for using SNP-level
5 summary statistics. However, based on our numerical experiments, p-values from
6 aSPUpath and aSPUsPath are less correlated, with their Pearson's correlation
7 coefficient being at the range of 0.75-0.85. To reduce the variability caused by the
8 use of reference panel, here we focused on the comparison between ARTP and
9 aSPUpath using individual-level genetic data. We adopted the similar simulation
10 strategy as the one used by Pan et al. [8].

11

12 We simulated haplotypes on a set of SNPs within a gene in the general population
13 using the algorithm of Wang and Elston [35]. Then the joint genotypes on a subject
14 can be formed by randomly pairing two haplotypes. In brief, we first chose the MAF
15 for each SNP by randomly sampling a value from the uniform distribution
16 $U(0.1,0.4)$. Then for the set of SNPs in a gene we sampled a latent vector
17 $Z = (z_1, \dots, z_q)^T$ from a multivariate normal distribution with a covariance matrix
18 $\text{Cov}(z_i, z_j) = \rho^{|i-j|}, 1 \leq i, j \leq q$, where ρ was sampled from the uniform distribution
19 $U(0,0.8)$ for a given gene. We randomly picked 50% of the SNPs and converted their
20 simulated z_i into minor and major alleles (coded as 0, 1), with the cuts chosen for
21 each z_i such that the resultant minor allele has its frequency defined by the
22 specified MAF. For the remaining SNPs, we used the same algorithm to dichotomize

1 their $-z_i$ into minor and major alleles. This created a more realistic haplotype
2 structure such that a haplotype can consist of a mixture of minor and major alleles.
3 Genotypes on SNPs from different genes were generated independently.

4

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

20

21 All simulation results are given in Table S1 and Table S2. It is clear that ARTP are
22 generally more powerful than aSPUpath, especially when the signal-to-noise ratio
23 (the proportion of genes including a functional SNP) is relatively low. The two types

1 of tests tend to have comparable performance when the signal-to-noise ratio
2 increases to 30%, although it is uncommon for a candidate pathway to have such a
3 high signal-to-noise ratio in real applications. For example, among the 4,713
4 candidate pathways analyzed in the next section, only 4.2% and 0.9% of the
5 pathways have over 20% and 30% of their genes that are likely to contain
6 association signals (i.e., with gene-level p-values < 0.05).

7

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

17

18 **Results of T2D Pathway Analysis**

19 *Findings from the European Populations*

20 Since our goal was to identify new susceptibility loci for T2D through the pathway
21 analysis, we excluded 170 high evidence T2D associated SNPs that were either listed
22 in Morris et al. [13] or found from the GWAS Catalog satisfying the following three
23 conditions simultaneously: (1) were investigated by GWAS of samples with

1 European ancestry; (2) had reported p-values $<10^{-7}$ on the initial study; and (3)
2 were replicated on independent studies. We further excluded 195 SNPs that has
3 their single-locus testing p-values less than 10^{-7} in either DIAGRAM or GERA data to
4 ensure that the pathway analysis result was not driven by a single SNP. In addition,
5 we further excluded genes within a ± 500 kb region from each of the removed SNPs
6 to eliminate potential association signals that could be caused by linkage
7 disequilibrium (LD) with the index SNPs.

8

9 We conducted three types of pathway analyses using sARTP, including the one using
10 the DIAGRAM SNP-level summary statistics, the one using the GERA SNP-level
11 summary statistics, and the pathway meta-analysis combining SNP-level summary
12 statistics from both DIAGRAM and GERA studies. When applying the pathway meta-
13 analysis to a single gene, we refer to this as the gene-level meta-analysis. We used
14 the external reference panel of 503 Europeans from the 1000 Genomes Project for
15 the variance-covariance matrix estimate in the analysis.

16

17 Before performing a pathway analysis, we applied LD filtering to remove redundant
18 SNPs. For any two SNPs with their pairwise squared Pearson's correlation
19 coefficient > 0.9 estimated from the external reference panel from the 1000
20 Genomes Project, we removed the one with a smaller value defined as,
21 $2f(1-f)n_0n_1n^{-1}$, where n_0 and n_1 are numbers of controls and cases, $n = n_0 + n_1$,
22 and f is the MAF based on the reference panel. This value is proportional to the
23 non-centrality parameter of the trend test statistic at a given SNP. We also excluded

1 SNPs with MAF < 1%. After all SNP filtering steps, we had a total of 4,713 pathways
2 for the analysis. The summary of the number of genes and SNPs used in each
3 pathway analysis is given in Figure S2.

4

5 The DIAGRAM study had a genomic control inflation factor $\lambda_{GC} = 1.10$ based on the
6 published meta-analysis result. The assembled GERA T2D GWAS had $\lambda_{GC} = 1.08$.

7 When conducting the pathway analysis on each of two studies, we adjusted the

8 inflation by using the corresponding $\sqrt{\lambda_{GC}}$ to rescale the standard error of

9 estimated log odds ratio at each SNP. The single-locus meta-analysis combining

10 results from DIAGRAM and GERA datasets had an inflation factor $\lambda_{GC} = 1.067$ after

11 each study had adjusted for its own inflation factor. We further adjusted this

12 inflation in the pathway and gene-level meta-analysis when combining SNP-level

13 summary statistics from both studies using formulas (4) and (5).

14

15 The Q-Q plots of gene-level and pathway p-values are given in Figure 2. Gene-level

16 p-value Q-Q plots based on the three analyses show no sign of inflation with their

17 λ_{GC} close to 1.0, but suggest that there are enriched gene-level association signals at

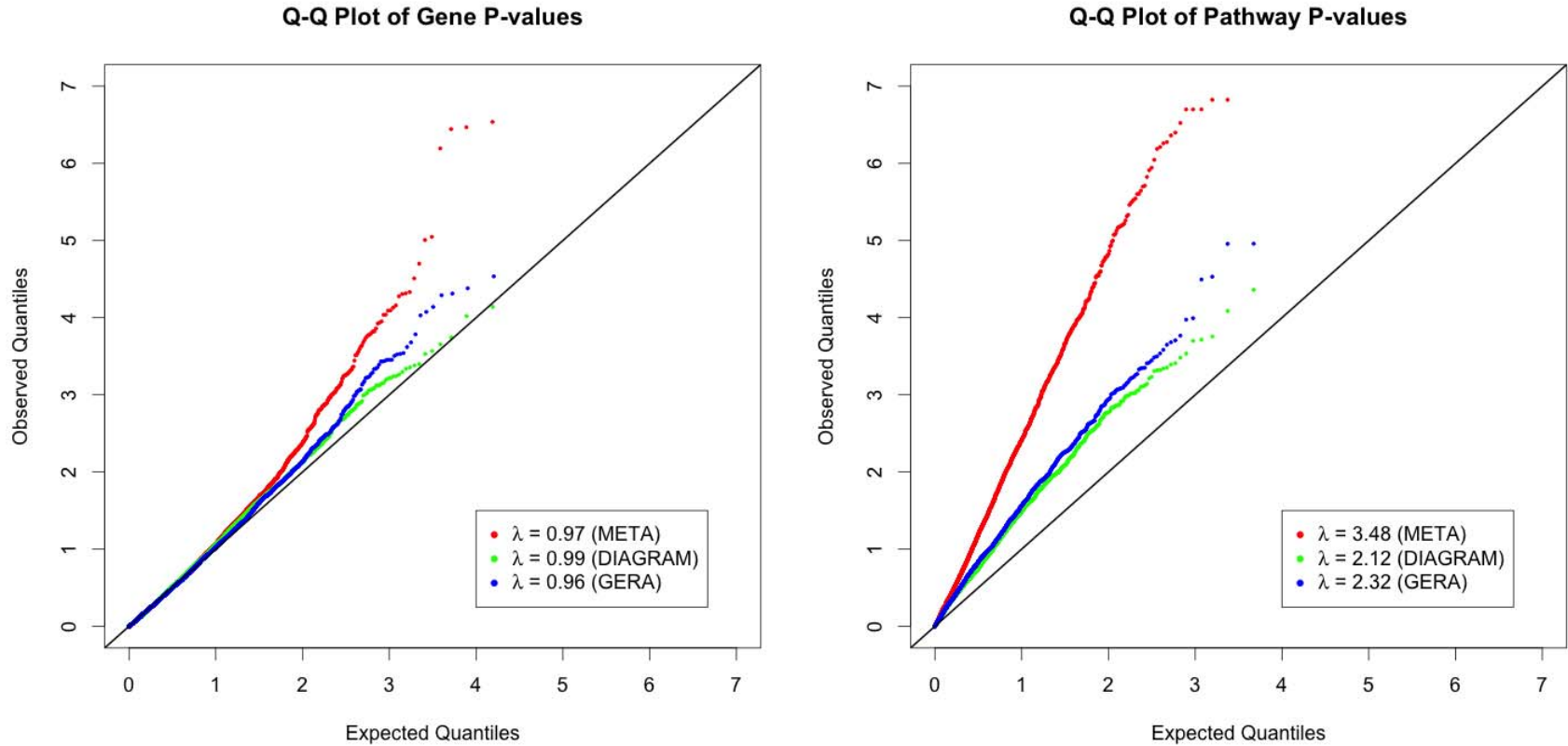
18 the tail end. The pathway p-value Q-Q plots, on the other hand, shift away from the

19 diagonal identify line and have much higher λ_{GC} , which suggests that T2D

20 associated genes are preferably included in pathways under study. In fact, it can be

21 seen from Figure S3 that a gene with a smaller gene-level meta-analysis p-value

- 1 tends to be included in more pathways, even though the 4,713 pathways collected
- 2 from MSigDB v5.0 are not specifically defined for the study of T2D.



1
2 Figure 2. Q-Q plots of gene-level and pathway p-values based on the sARTP procedure on the DIAGRAM study, the GERA study,
3 and the two studies combined
4
5 Left: Q-Q plots of gene-level p-values on 15,946 genes based on the sARTP gene-based analysis of the DIAGRAM study
6 (DIAGRAM), the GERA study (GERA), and the two studies combined (META)
7 Right: Q-Q plots of pathway p-values on 4,713 pathways based on the sARTP pathway analysis of the DIAGRAM study
8 (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 S4-S46. There are
18 a total of 15,946 unique genes in all 4,713 pathways. The top 50 genes with smallest
19 gene-level p-values based on the gene meta-analysis are listed in Table 4.

20

21

22

23

24

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

Pathway ^a	META ^b	DIAGRAM ^c	GERA ^d
SCHLOSSER_SERUM_RESPONSE_UP ^e	2.50E-08	2.92E-04	1.77E-03
PENG_RAPAMYCIN_RESPONSE_DN	1.50E-07	1.68E-03	2.08E-04
YAGI_AML_WITH_T_8_21_TRANSLOCATION	1.50E-07	4.33E-03	4.46E-04
PATIL_LIVER_CANCER ^e	2.00E-07	4.35E-05	3.89E-03
PUJANA_CHEK2_PCC_NETWORK	2.00E-07	1.18E-02	3.39E-03
STEIN_ESRRA_TARGETS	2.00E-07	9.37E-04	8.39E-04
STEIN_ESRRA_TARGETS_UP	3.00E-07	6.38E-03	1.02E-04
WANG_CISPLATIN_RESPONSE_AND_XPC_UP	4.00E-07	6.75E-03	1.18E-01
CADWELL_ATG16L1_TARGETS_DN	4.35E-07	1.45E-03	9.59E-03
SONG_TARGETS_OF_IE86_CMV_PROTEIN ^e	5.30E-07	7.88E-04	3.20E-05
CASORELLI_ACUTE_PROMYELOCYTIC_LEUKEMIA_DN	5.50E-07	2.48E-02	5.70E-03
RIZ_ERYTHROID_DIFFERENTIATION ^e	6.15E-07	2.33E-02	2.31E-02
BORCZUK_MALIGNANT_MESOTHELIOMA_UP ^e	6.50E-07	7.61E-03	2.52E-02
HILLION_HMGA1_TARGETS	9.00E-07	3.39E-01	1.10E-05
KEGG_MATURITY_ONSET_DIABETES_OF_THE_YOUNG ^e	1.14E-06	1.68E-02	3.58E-04
HOLLEMAN_ASPARAGINASE_RESISTANCE_BALL_DN	1.22E-06	3.42E-02	1.67E-03
PUJANA_BRCA1_PCC_NETWORK	1.50E-06	1.69E-02	5.11E-03
HOSHIDA_LIVER_CANCER_SUBCLASS_S3 ^e	1.95E-06	6.14E-03	5.83E-03
GRAESSMANN_APOPTOSIS_BY_DOXORUBICIN_DN	2.00E-06	9.57E-04	6.41E-04
REACTOME_REGULATION_OF_BETA_CELL_DEVELOPMENT ^e	2.26E-06	4.81E-02	1.15E-03
PUJANA_BREAST_CANCER_WITH_BRCA1_MUTATED_UP ^e	2.48E-06	7.61E-03	2.39E-02
BLALOCK_ALZHEIMERS_DISEASE_UP ^e	2.50E-06	3.52E-02	3.26E-02
GOBERT_OLIGODENDROCYTE_DIFFERENTIATION_UP ^e	2.85E-06	3.68E-02	5.37E-04
MCBRYAN_PUBERTAL_BREAST_4_5WK_DN ^e	2.95E-06	2.20E-01	1.10E-05
REACTOME_REGULATION_OF_GENE_EXPRESSION_IN_BETA_CELLS ^e	3.11E-06	2.81E-02	2.33E-03
SANSOM_APC_TARGETS_DN	3.25E-06	3.08E-01	2.84E-03
NABA_MATRISOME ^e	3.45E-06	4.46E-02	1.30E-02
PUJANA_BRCA2_PCC_NETWORK ^e	4.65E-06	1.25E-02	6.32E-02
KEGG_TYPE_II_DIABETES_MELLITUS ^e	4.85E-06	6.38E-02	9.93E-04
LINDGREN_BLADDER_CANCER_CLUSTER_1_DN ^e	5.50E-06	5.20E-02	4.36E-03

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

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

4 ^aThe name of the pathway given by the Molecular Signatures Database;

5 ^bP-values from the sARTP pathway meta-analysis combining summary statistics
6 from the DIAGRAM and GERA studies;

7 ^cP-values from the sARTP pathway analysis using summary statistics from the
8 DIAGRAM study;

9 ^dP-values from the sARTP pathway analysis using summary statistics from the
10 GERA study;

11 ^eThe 23 pathways that do not have any gene belonging to the cluster of genes at
12 chromosome 17q21 that consist of genes *ATP5G1*, *SNF8*, *UBE2Z*, and *GIP*.

13

14

15

16

17

18

19

20

21

1 Table 4: Summary of top 50 genes with smallest gene-level p-values from the gene-
 2 level meta-analysis based on the DIAGRAM and GERA studies

Gene	Chromosome	META ^a	DIAGRAM ^b	GERA ^c
ATP5G1	17	2.90E-07	2.08E-03	2.91E-05
SNF8	17	3.40E-07	1.00E-03	4.87E-05
UBE2Z	17	3.60E-07	1.59E-03	5.13E-05
GIP	17	6.40E-07	1.02E-03	2.40E-04
SLC2A2	3	8.97E-06	9.97E-03	2.95E-04
MYBL2	20	9.84E-06	4.41E-04	7.07E-03
IFT52	20	2.00E-05	2.71E-04	3.26E-03
SREBF1	17	3.09E-05	1.45E-02	1.62E-03
SMEK1	14	4.65E-05	9.53E-05	7.41E-02
TACO1	17	4.84E-05	1.86E-03	4.77E-03
SOCS2	12	4.94E-05	7.76E-02	2.88E-04
IL20RA	6	5.29E-05	5.59E-03	1.61E-02
PLEKHA1	10	6.94E-05	7.23E-05	3.82E-02
YPEL2	17	7.41E-05	3.60E-03	1.73E-02
RAB1A	2	8.05E-05	3.97E-03	1.89E-02
CDKN2C	1	8.11E-05	1.20E-02	3.61E-03
CENPW	6	9.15E-05	1.96E-03	2.70E-02
RAI1	17	9.20E-05	8.82E-03	6.64E-03
C11orf34	11	1.12E-04	1.14E-02	4.59E-03
KIAA0754	1	1.16E-04	8.05E-02	7.38E-04
ACSL1	4	1.19E-04	3.87E-03	4.52E-02
UBAP2	9	1.38E-04	9.94E-03	2.40E-02
MACF1	1	1.51E-04	1.14E-01	2.09E-04
MAP3K3	17	1.51E-04	4.88E-03	2.01E-02
GTSF1L	20	1.59E-04	4.67E-03	1.91E-02
MIR17HG	13	1.66E-04	1.39E-02	5.80E-03
SGK2	20	1.68E-04	5.05E-04	5.56E-02
CEP68	2	1.75E-04	9.38E-03	2.00E-02
BLOC1S2	10	1.84E-04	5.75E-04	1.07E-01
PABPC4	1	1.97E-04	1.33E-01	9.85E-04
PPIEL	1	2.10E-04	1.29E-01	1.03E-03
CALCOCO2	17	2.29E-04	5.50E-02	1.41E-03
YTHDC2	5	2.37E-04	6.21E-03	4.33E-02
ARNTL	11	2.54E-04	4.30E-02	2.03E-02
DCAF7	17	2.73E-04	6.11E-03	1.61E-02
PTS	11	2.89E-04	7.46E-03	2.33E-03
PDE3B	11	3.06E-04	4.44E-01	2.99E-04
ZNF276	16	3.08E-04	9.09E-02	7.26E-05
C16orf7	16	3.60E-04	1.11E-01	8.42E-05
FANCA	16	4.31E-04	1.03E-01	9.35E-05

Gene	Chromosome	META ^a	DIAGRAM ^b	GERA ^c
CWF19L1	10	4.43E-04	7.50E-04	8.81E-02
GPR151	5	4.68E-04	3.22E-02	1.04E-03
KCNH6	17	4.82E-04	1.41E-02	2.55E-02
EIF5A2	3	4.93E-04	6.43E-02	1.47E-03
FEN1	11	5.11E-04	1.06E-02	2.27E-02
CPPED1	16	5.21E-04	5.05E-03	7.08E-02
ZNF664	12	5.35E-04	2.61E-02	2.29E-02
C11orf10	11	5.46E-04	1.06E-02	2.47E-02
LIMD2	17	5.50E-04	1.52E-02	2.85E-02
CHUK	10	5.59E-04	8.33E-04	1.23E-01

1 The top 50 genes are chosen from 15,946 unique genes included in 4,713 candidate
2 pathways for having smallest gene-level p-values based on the gene-level meta-
3 analysis.

4 ^aP-values from the sARTP gene-level meta-analysis combining summary statistics
5 from the DIAGRAM and GERA studies;

6 ^bP-values from the sARTP gene-level analysis using summary statistics from the
7 DIAGRAM study;

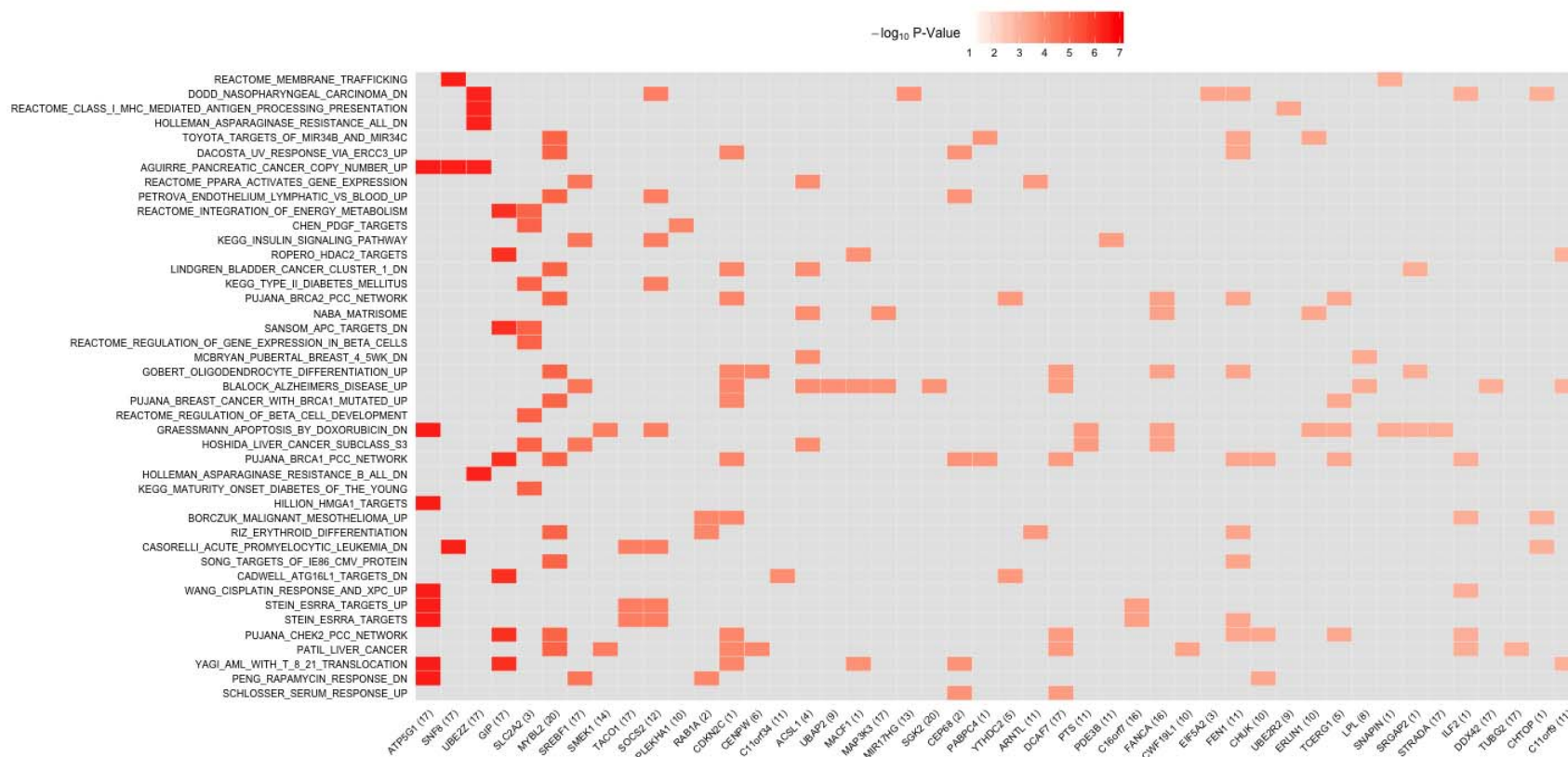
8 ^cP-values from the sARTP gene-level analysis using summary statistics from the
9 GERA study;

10

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

21

1 In Figure 3, we illustrate the connection between the 43 significant pathways and a
2 group of genes showing association evidence. For the purpose of illustration, in the
3 figure we only focus on 46 genes that are covered by the 43 pathways and have their
4 gene-level meta-analysis p-values less than 0.001. It is evident from Figure 3 that a
5 cluster of 4 genes, *UBE2Z*, *SNF8*, *GIP*, and *ATP5G1*, have the most significant gene-
6 level p-values (Table 4), and contribute association signals to 20 out of 43
7 significant pathways (Figures S4-S23). These 4 genes overlap each other at
8 chromosome 17q21. This region contains a previously unidentified genome-wide
9 significant synonymous SNP rs1058018 (meta-analysis $p = 3.06 \times 10^{-8}$) after two
10 rounds of inflation adjustments. By conditioning on rs1058018, none of the other
11 SNPs in this region are significant based on the conditional association analysis
12 using the GERA individual-level GWAS data. Based on GTEx data v6, rs1058018 is a
13 *cis* eQTL for *UBE2Z* in blood ($p = 7.9 \times 10^{-15}$). *UBE2Z* is involved in Class I MHC
14 antigen processing and presentation (GeneCards). The region at 17q21 was
15 previously implicated to be associated with T2D through a candidate gene/loci
16 approach [36]. Although genes at the 17q21 region carry the strongest association
17 signal, 10 out of those 20 pathways remain to be globally significant ($p < 1.06 \times 10^{-5}$)
18 after excluding those genes from the pathway definition.
19



1
 2 Figure 3. Heat map of gene-level p-values on selected genes within 43 significant pathways based on the DIAGRAM and GERA
 3 studies
 4 There are 46 unique genes in the 43 significant pathways that have their gene-level meta-analysis p-values less than 0.001.
 5 Each row in the plot represents one of 43 significant pathways. Each column represents one of the 46 unique genes. The
 6 chromosome IDs of 46 unique genes are given in parentheses. The color of each cell represents the gene-level p-value (in the
 7 $-\log_{10}$ scale). A cell for a gene that is not included in a pathway is colored gray in the corresponding entry. The orders of genes
 8 (x-axis) and pathways (y-axis) are arranged according to their gene and pathway meta-analysis p-values.

1 The majority of 43 identified pathways are enriched with signals from multiple
2 chromosomal regions as demonstrated by the Q-Q plots of their SNP-level and gene-
3 level p-values (Figures S4-S46). For example, the strongest T2D-associated pathway,
4 SCHLOSSER_SERUM_RESPONSE_UP, consists of 103 genes, which includes two
5 genes with p-values < 0.001, and has 20 genes with p-values between 0.001 and
6 0.05 (Figure S24). We conducted the ingenuity pathway analysis on the 22 genes
7 with p-values less than 0.05, and found enrichment of these genes in caveolae-
8 mediated cytolysis (important for removal of low/high density lipoproteins), and lipid
9 metabolism pathways, and in functions/diseases related to differentiation of
10 phagocytes and transport of proteins.

11

12 It is assuring that our pathway analysis detected several pathways that are natural
13 candidates underlying the development of T2D, including the pathways
14 KEGG_MATURITY_ONSET_DIABETES_OF_THE_YOUNG (Figure S29),
15 KEGG_TYPE_II_DIABETES_MELLITUS (Figure S39),
16 KEGG_INSULIN_SIGNALING_PATHWAY (Figure S41), and
17 REACTOME_REGULATION_OF_BETA_CELL_DEVELOPMENT (Figure S31). It is worth
18 emphasizing that these pathways were analyzed after excluding genes in the
19 neighborhood of 170 GWAS established T2D loci and 195 SNPs with p-values < 10^{-7}
20 on either DIAGRAM or GERA data, which suggests that these well-defined T2D-
21 related pathways are enriched with additional unidentified and contributory T2D-
22 associated genes.

23

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 S19), hepatocellular carcinoma (HCC) (Figure
4 S25, and S30), bladder carcinoma (Figure S40), nasopharyngeal carcinoma (Figure
5 S22), and familial breast cancer (Figures S32, and S35). It is well recognized that
6 T2D patients have elevated risk of cancer at multiple sites, such as the liver and
7 pancreas [37,38]. These findings can provide valuable insights into the genetic basis
8 underlying the connection between T2D and a host of different cancers.

9

10 ***Findings from Eastern Asian Populations***

11 We reanalyzed the 43 significant pathways identified from the European
12 populations using summary-level data generated by the AGEN-T2D study. An
13 inflation factor $\lambda_{GC}=1.03$ calculated from the AGEN-T2D meta-analysis was
14 adjusted in the pathway meta-analysis. The results were summarized in Table 5. We
15 had 10 out of 43 pathways with the unadjusted p-value less than 0.05, suggesting
16 that pathways identified from the European populations were also enriched with
17 T2D-associated genes in the eastern Asian populations. Among the 43 pathways, we
18 were able to identify 4 significant T2D-associated pathways at the false discovery
19 rate (FDR [39]) of 0.05 (Figures S25, S10, S30 and S34), and 3 additional T2D-
20 associated pathways at the FDR of 0.1 (Figures S45, S42, and S23). These results
21 support the presence of trans-ethnic pathway effect on T2D in European and
22 eastern Asian populations [11,12].

23

1

2 Table 5: Pathway p-values and FDR adjusted p-values based on the AGEN-T2D study.

Pathway	sARTP ^a	FDR ^b
PATIL_LIVER_CANCER	0.0013	0.029
CADWELL_ATG16L1_TARGETS_DN	0.0023	0.029
HOSHIDA_LIVER_CANCER_SUBCLASS_S3	0.0025	0.029
GOBERT_OLIGODENDROCYTE_DIFFERENTIATION_UP	0.0027	0.029
DACOSTA_UV_RESPONSE_VIA_ERCC3_UP	0.011	0.074
CHEN_PDGF_TARGETS	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	0.041	0.19
LINDGREN_BLADDER_CANCER_CLUSTER_1_DN	0.043	0.19
PUJANA_CHEK2_PCC_NETWORK	0.057	0.21
BLALOCK_ALZHEIMERS_DISEASE_UP	0.059	0.21
SCHLOSSER_SERUM_RESPONSE_UP	0.085	0.27
RIZ_ERYTHROID_DIFFERENTIATION	0.089	0.27
DODD_NASOPHARYNGEAL_CARCINOMA_DN	0.097	0.28
TOYOTA_TARGETS_OF_MIR34B_AND_MIR34C	0.10	0.28
REACTOME_CLASS_I_MHC_MEDIATED_ANTIGEN_PROCESSING_PRESENTATION	0.14	0.32
PUJANA_BRCA2_PCC_NETWORK	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	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	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	0.45	0.65
REACTOME_PPARA_ACTIVATES_GENE_EXPRESSION	0.52	0.72
REACTOME_REGULATION_OF_GENE_EXPRESSION_IN_BETA_CELLS	0.55	0.74
KEGG_INSULIN_SIGNALING_PATHWAY	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	0.65	0.76
KEGG_TYPE_II_DIABETES_MELLITUS	0.67	0.76
SANSOM_APC_TARGETS_DN	0.69	0.76

Pathway	sARTP ^a	FDR ^b
REACTOME_INTEGRATION_OF_ENERGY_METABOLISM	0.70	0.76
REACTOME_REGULATION_OF_BETA_CELL_DEVELOPMENT	0.70	0.76
BORCZUK_MALIGNANT_MESOTHELIOMA_UP	0.75	0.79
SONG_TARGETS_OF_IE86_CMV_PROTEIN	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 from the sARTP pathway meta-analysis;

5 ^b FDR adjusted p-values.

6

7 Given the existing epidemiologic evidence on the close connection between T2D and
8 the liver cancer, it is noteworthy that the two HCC related pathways (Table 5,
9 Figures S25, S30) identified in European populations remain to be significant in
10 eastern Asian populations at the FDR of 0.05. The pathway PATIL_LIVER_CANCER
11 consists of 653 genes (after data preprocessing) that are highly expressed in HCC
12 and are enriched with genes having functions related to cell growth, cell cycle,
13 metabolism, and cell proliferation [40]. The other pathway,
14 HOSHIDA_LIVER_CANCER_SUBCLASS_S3 consists of 240 genes that show similar
15 gene expression variation patterns and together define a HCC subtype with its
16 unique histologic, molecular and clinical characteristics [41]. These two pathways
17 have only 6 genes in common, and none of the 6 genes has a gene-level p-value <
18 0.05 in either European or eastern Asian data. More in depth investigations of these
19 two complementary pathways could lead to further understanding the connection
20 between T2D and the liver cancer.

21

1 The genome-wide significant SNP rs1058018 at the 17q21 region identified through
2 the combined analysis of DIAGRAM and GERA studies turned out to be null in the
3 AGEN-T2D study ($p = 0.29$). This could be due to the relatively small sample size of
4 the AGEN-T2D study, or the genetic risk heterogeneity at the 17q21 locus among
5 different ethnic populations. However, 2 out of the 20 pathways (Figures S10 and
6 S23) that contain genes within the 17q21 region can still be significant at the FDR of
7 0.1 because of enriched association signals throughout the other chromosomal
8 regions.

9

10 **Discussion**

11 We developed a general statistical procedure sARTP for pathway analysis using
12 SNP-level summary statistics generated from multiple GWAS. By applying sARTP to
13 summary statistics from two large studies with a total of 19,809 T2D cases and
14 111,181 controls with European ancestry, we were able to identify 43 globally
15 significant T2D-associated pathways after excluding genes in neighborhoods of
16 GWAS established T2D loci. Using summary data generated from eight T2D GWAS
17 with 6,952 cases and 11,865 controls from eastern Asian populations, we can
18 further showed that 7 out of 43 pathways identified in the European populations are
19 also significant in the eastern Asian populations at the FDR of 0.1. The analysis
20 clearly highlights novel T2D-associated genes and pathways beyond what has been
21 known from single-SNP association analysis reported from largest GWAS to date.
22 Since the new procedure requires only SNP-level summary statistics, it provides a

1 flexible way for conducting pathway analysis, alleviating the burden of handling
2 large volumes of individual-level GWAS data.

3

4 We have developed a computationally efficient R package called ARTP2
5 implementing both the ARTP and sARTP procedures, so that it can be used for
6 conducting pathway analysis based on individual-level genetic data, as well as SNP-
7 level summary data from one or multiple GWAS. The R package also supports the
8 parallelization on Unix-like OS, which can substantially accelerate the computation
9 of small p-values when a large number of resampling steps are needed. The ARTP2
10 package has a user-friendly interface and provides a comprehensive set of data
11 preprocessing procedures to ensure that all the input information (e.g., allele
12 information of SNP-level summary statistics and genotype reference panel) can be
13 processed coherently. To make the sARTP method accessible to a wider research
14 community, we have also developed a web-based tool that allows investigators to
15 conduct their pathway analyses using the computing resource at the National
16 Cancer Institutes through simple on-line inputs of summary data.

17

18 Single-locus analysis of GWAS usually has its genomic control inflation factor larger
19 than 1.0. Some proportion of the inflation can be attributed to various confounding
20 biases, such as the one caused by population stratification, while the other part can
21 be due to the real polygenic effect. In the pathway analysis it is important to
22 minimize the confounding bias at the SNP-level summary statistic. Otherwise a
23 small bias at the SNP level can be accumulated in the pathway analysis, and lead to

1 an elevated false discovery rate. Here we try to remove the confounding bias by
2 adjusting for the genomic control inflation factor observed at the GWAS study. This
3 approach is conservative because part of the inflation can be caused by the real
4 polygenic effect. Bulik-Sullivan et al. [42] recently developed the LD score
5 regression method to quantify the level of inflation caused solely by the confounding
6 bias. Adjusting for the inflation factor estimated by this method, instead of the
7 genomic control inflation factor, can potentially increase the power of the pathway
8 analysis. However, the LD score regression method relies on a specific polygenic
9 risk model, and its estimate might not be robust for this model assumption. More
10 investigations are needed to evaluate the impact of this new inflation adjustment on
11 the pathway analysis.

12

13 There are several other strategies to increase the power of pathway analysis besides
14 increasing sample size [4]. One area of active research is to find better ways to
15 define the gene-level summary statistic using observed genotypes on multiple SNPs,
16 so that it can accurately characterize the impact of the gene on the outcome [43-46].
17 In our proposed procedure, we adopt a data driven approach to select a subset of
18 SNPs within a gene that collectively show the strongest association evidence.
19 Because of this, we have to pay the penalty of multiple-comparison in the final
20 pathway significance assessment. However, it is well recognized that SNPs at
21 different loci can have varied levels of functional implications. We can potentially
22 reduce the burden of multiple-comparisons and thus improve the power of the
23 pathway analysis, by prioritizing SNPs according to existing genomic knowledge and

1 other data resources. For example, Gamazon et al. [47] recently proposed a new
2 gene-level summary statistic based on a prediction model that was trained with
3 external transcriptome data. The gene-level summary statistic is defined as the
4 predicted value that estimates the component of gene expression regulated by a
5 subject's genotypes within the neighborhood of the considered gene. Pathway
6 analysis procedures using this kind of biologically informed gene-level summary
7 statistic can be easily incorporated into the ARTP2 framework.

8

9 When conducting pathway analysis with individual-level genetic data, we could run
10 into a computing memory issue if the study has a large sample size and the pathway
11 consists of a large number of genes and SNPs (Figure S2). The ability of performing
12 pathway analysis using summary data provides a convenient and efficient solution
13 in those situations. We can first calculate the SNP-level summary statistics based on
14 the individual-level genetic data, and then randomly sample a small proportion of
15 the original data as an internal reference to estimate the variance-covariant matrix
16 for score statistics at considered SNPs. As we have shown in the Result section, the
17 testing results using this approach are very consistent with those based on the
18 original dataset.

19

20 As already demonstrated by many successful GWAS meta-analysis, increasing the
21 sample size through combining results from multiple studies is a very effective way
22 to improve our chance for new findings. For the same reason, pathway-based meta-
23 analysis can provide us with new opportunities to uncover biological pathways that

1 are previously undetectable due to the limitation on the sample size. With more
2 summary data from meta-analysis becoming increasingly available, we expect the
3 ARTP2 package would be a valuable tool for further exploring the genome in search
4 for the hidden heritability.

5

6 **Appendix A: Recovering Score Statistics and Its Variance-** 7 **Covariance Matrix Using Summary Results from the Fixed** 8 **Effect Model**

9 Here we derive the approximated score statistic S and its variance-covariance
10 matrix V using summary statistics from the fixed effect model. Based on (4), it is
11 straightforward to see that $S_i \approx \tau_i^{-2} \hat{\beta}_i$. Note that V_{is} in equation (5) depends on $\tau_i^{(l)}$
12 from individual studies, and cannot be estimated by using τ_i only. However, assume
13 that $\sqrt{n_i^{(l)}} \tau_i^{(l)}$ can be approximated as an unknown but common constant value v_i ,
14 across all studies, and if $\bar{y}^{(l)}(1-\bar{y}^{(l)}) \approx \bar{y}(1-\bar{y})$, we have $v_i \approx \sqrt{n_i} \tau_i$, and

15 $V_{is} \approx \frac{n_{is}}{\sqrt{n_i n_s}} \frac{\rho_{is}}{\tau_i \tau_s}$. The similar argument has been used in Lin and Zeng [48] to

16 demonstrate that the meta-analysis is as efficient as the pooled analysis under those
17 conditions.

18

19 **Supporting Information**

1 Table S1. Power comparison between ARTP and aSPUpath under the scenrio where
2 each outcome-associated gene contains one functional SNP

3

4 Table S2. Power comparison between ARTP and aSPUpath under the scenrio where
5 each outcome-associated gene contains one or two functional SNP(s) with equal
6 probability

7

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

10

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

14

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

17

18 Figure S4: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
19 PENG_RAPAMYCIN_RESPONSE_DN.

20

21 Figure S5: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
22 YAGI_AML_WITH_T_8_21_TRANSLOCATION

23

- 1 Figure S6: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
- 2 PUJANA_CHEK2_PCC_NETWORK
- 3
- 4 Figure S7: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
- 5 STEIN_ESRRA_TARGETS
- 6
- 7 Figure S8: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
- 8 STEIN_ESRRA_TARGETS_UP
- 9
- 10 Figure S9: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
- 11 WANG_CISPLATIN_RESPONSE_AND_XPC_UP
- 12
- 13 Figure S10: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
- 14 CADWELL_ATG16L1_TARGETS_DN
- 15
- 16 Figure S11: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
- 17 CASORELLI_ACUTE_PROMYELOCYTIC_LEUKEMIA_DN
- 18
- 19 Figure S12: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
- 20 HILLION_HMGA1_TARGETS
- 21
- 22 Figure S13: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
- 23 HOLLEMAN_ASPARAGINASE_RESISTANCE_B_ALL_DN

1

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

3 PUJANA_BRCA1_PCC_NETWORK

4

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

6 GRAESSMANN_APOPTOSIS_BY_DOXORUBICIN_DN

7

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

9 SANSOM_APC_TARGETS_DN

10

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

12 ROPERO_HDAC2_TARGETS

13

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

15 REACTOME_INTEGRATION_OF_ENERGY_METABOLISM

16

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

18 AGUIRRE_PANCREATIC_CANCER_COPY_NUMBER_UP

19

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

21 HOLLEMAN_ASPARAGINASE_RESISTANCE_ALL_DN

22

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

3

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

6

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

9

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

12

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

15

16 Figure S26: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
17 SONG_TARGETS_OF_IE86_CMV_PROTEIN

18

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

21

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

1

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

3 KEGG_MATURITY_ONSET_DIABETES_OF_THE_YOUNG

4

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

6 HOSHIDA_LIVER_CANCER_SUBCLASS_S3

7

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

9 REACTOME_REGULATION_OF_BETA_CELL_DEVELOPMENT

10

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

12 PUJANA_BREAST_CANCER_WITH_BRCA1_MUTATED_UP

13

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

15 BLALOCK_ALZHEIMERS_DISEASE_UP

16

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

18 GOBERT_OLIGODENDROCYTE_DIFFERENTIATION_UP

19

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

21 MCBRYAN_PUBERTAL_BREAST_4_5WK_DN

22

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

3

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

6

7 Figure S38: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
8 PUJANA_BRCA2_PCC_NETWORK

9

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

12

13 Figure S40: Q-Q plots for SNP p-values and sARTP gene p-values of pathway
14 LINDGREN_BLADDER_CANCER_CLUSTER_1_DN

15

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

18

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

21

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

1

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

3 REACTOME_PPARA_ACTIVATES_GENE_EXPRESSION

4

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

6 DACOSTA_UV_RESPONSE_VIA_ERCC3_UP

7

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

9 TOYOTA_TARGETS_OF_MIR34B_AND_MIR34C

10

11

12

13 **Acknowledgments**

14 This study utilized the computational resources of the NIH HPC Biowulf cluster

15 (<https://hpc.nih.gov/>). The authors thank the AGEN-T2D consortium and DIAGRAM

16 consortium for sharing the meta-analysis summary data. The authors acknowledge

17 Sue Pan, JJ Pan, Wesley Obenshain, Tony Hall, Jim Zhou, Ye Wu, Cuong Nguyen for

18 their help in developing the web-based tool for ARTP2.

19

20 **Web Resources**

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

- 1 DIAbetes Genetics Replication And Meta-analysis (DIAGRAMv3), [http://diagram-](http://diagram-consortium.org/)
- 2 [consortium.org/](http://diagram-consortium.org/)
- 3 Genetic Epidemiology Research on Aging (GERA, dbGaP Study Accession:
- 4 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)
- 5 [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)
- 6 Molecular Signatures Database (C2: curated gene sets),
- 7 <http://software.broadinstitute.org/gsea/msigdb/collections.jsp#C2>
- 8 BioMart (Homo sapiens genes NCBI36 and GRCh37.p13),
- 9 <http://feb2014.archive.ensembl.org/>
- 10 IMPUTE2, https://mathgen.stats.ox.ac.uk/impute/impute_v2.html
- 11 GWAS Catalog, <http://www.ebi.ac.uk/gwas/>
- 12 1000 Genomes Project (Phase 3, v5, 2013/05/02),
- 13 <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>
- 14 GTEx Portal v6, <http://gtexportal.org/home/>
- 15 GeneCards Human Gene Database, <http://www.genecards.org/>
- 16 Ingenuity Pathway Analysis, <http://www.ingenuity.com/>
- 17 ARTP2 package, <https://cran.r-project.org/web/packages/ARTP2/>
- 18 Web-based tool of ARTP2, <http://analysistools.nci.nih.gov/pathway/>
- 19

20 **References**

- 21
- 22 1. Welter D, MacArthur J, Morales J, Burdett T, Hall P, et al. (2014) The NHGRI GWAS
- 23 Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res* 42:
- 24 D1001-1006.

- 1 2. Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorff LA, et al. (2009) Finding the
2 missing heritability of complex diseases. *Nature* 461: 747-753.
- 3 3. Moore JH, Asselbergs FW, Williams SM (2010) Bioinformatics challenges for
4 genome-wide association studies. *Bioinformatics* 26: 445-455.
- 5 4. Wang K, Li M, Hakonarson H (2010) Analysing biological pathways in genome-
6 wide association studies. *Nat Rev Genet* 11: 843-854.
- 7 5. Chen LS, Hutter CM, Potter JD, Liu Y, Prentice RL, et al. (2010) Insights into colon
8 cancer etiology via a regularized approach to gene set analysis of GWAS data.
9 *Am J Hum Genet* 86: 860-871.
- 10 6. Evangelou M, Rendon A, Ouwehand WH, Wernisch L, Dudbridge F (2012)
11 Comparison of methods for competitive tests of pathway analysis. *PLoS One*
12 7: e41018.
- 13 7. Li MX, Kwan JS, Sham PC (2012) HYST: a hybrid set-based test for genome-wide
14 association studies, with application to protein-protein interaction-based
15 association analysis. *Am J Hum Genet* 91: 478-488.
- 16 8. Pan W, Kwak IY, Wei P (2015) A Powerful Pathway-Based Adaptive Test for
17 Genetic Association with Common or Rare Variants. *Am J Hum Genet* 97: 86-
18 98.
- 19 9. Yu K, Li Q, Bergen AW, Pfeiffer RM, Rosenberg PS, et al. (2009) Pathway analysis
20 by adaptive combination of P-values. *Genet Epidemiol* 33: 700-709.
- 21 10. Cho YS, Chen CH, Hu C, Long J, Ong RT, et al. (2012) Meta-analysis of genome-
22 wide association studies identifies eight new loci for type 2 diabetes in east
23 Asians. *Nat Genet* 44: 67-72.
- 24 11. DIAbetes Genetics Replication Meta-analysis C, Consortium AGENTD,
25 Consortium SATD, Consortium MATD, Consortium TDGEbN-gsim-ES, et al.
26 (2014) Genome-wide trans-ancestry meta-analysis provides insight into the
27 genetic architecture of type 2 diabetes susceptibility. *Nat Genet* 46: 234-244.
- 28 12. Imamura M, Takahashi A, Yamauchi T, Hara K, Yasuda K, et al. (2016) Genome-
29 wide association studies in the Japanese population identify seven novel loci
30 for type 2 diabetes. *Nat Commun* 7: 10531.
- 31 13. Morris AP, Voight BF, Teslovich TM, Ferreira T, Segre AV, et al. (2012) Large-
32 scale association analysis provides insights into the genetic architecture and
33 pathophysiology of type 2 diabetes. *Nat Genet* 44: 981-990.
- 34 14. Speliotes EK, Willer CJ, Berndt SI, Monda KL, Thorleifsson G, et al. (2010)
35 Association analyses of 249,796 individuals reveal 18 new loci associated
36 with body mass index. *Nat Genet* 42: 937-948.
- 37 15. Burren OS, Guo H, Wallace C (2014) VSEAMS: a pipeline for variant set
38 enrichment analysis using summary GWAS data identifies IKZF3, BATF and
39 ESRRA as key transcription factors in type 1 diabetes. *Bioinformatics* 30:
40 3342-3348.
- 41 16. Evangelou M, Smyth DJ, Fortune MD, Burren OS, Walker NM, et al. (2014) A
42 method for gene-based pathway analysis using genomewide association
43 study summary statistics reveals nine new type 1 diabetes associations.
44 *Genet Epidemiol* 38: 661-670.
- 45 17. Kwak IY, Pan W (2015) Adaptive gene- and pathway-trait association testing
46 with GWAS summary statistics. *Bioinformatics*.

- 1 18. Lamparter D, Marbach D, Rueedi R, Kutalik Z, Bergmann S (2016) Fast and
2 Rigorous Computation of Gene and Pathway Scores from SNP-Based
3 Summary Statistics. *PLoS Comput Biol* 12: e1004714.
- 4 19. Network and Pathway Analysis Subgroup of Psychiatric Genomics C (2015)
5 Psychiatric genome-wide association study analyses implicate neuronal,
6 immune and histone pathways. *Nat Neurosci* 18: 199-209.
- 7 20. Segre AV, Consortium D, investigators M, Groop L, Mootha VK, et al. (2010)
8 Common inherited variation in mitochondrial genes is not enriched for
9 associations with type 2 diabetes or related glycemic traits. *PLoS Genet* 6.
- 10 21. Swanson DM, Blacker D, Alchawa T, Ludwig KU, Mangold E, et al. (2013)
11 Properties of permutation-based gene tests and controlling type 1 error
12 using a summary statistic based gene test. *BMC Genet* 14: 108.
- 13 22. Genomes Project C, Abecasis GR, Altshuler D, Auton A, Brooks LD, et al. (2010) A
14 map of human genome variation from population-scale sequencing. *Nature*
15 467: 1061-1073.
- 16 23. Genomes Project C, Abecasis GR, Auton A, Brooks LD, DePristo MA, et al. (2012)
17 An integrated map of genetic variation from 1,092 human genomes. *Nature*
18 491: 56-65.
- 19 24. Hu YJ, Berndt SI, Gustafsson S, Ganna A, Genetic Investigation of ATC, et al.
20 (2013) Meta-analysis of gene-level associations for rare variants based on
21 single-variant statistics. *Am J Hum Genet* 93: 236-248.
- 22 25. Liu JZ, McRae AF, Nyholt DR, Medland SE, Wray NR, et al. (2010) A versatile
23 gene-based test for genome-wide association studies. *Am J Hum Genet* 87:
24 139-145.
- 25 26. Yang J, Ferreira T, Morris AP, Medland SE, Genetic Investigation of ATC, et al.
26 (2012) Conditional and joint multiple-SNP analysis of GWAS summary
27 statistics identifies additional variants influencing complex traits. *Nat Genet*
28 44: 369-375, S361-363.
- 29 27. Goeman JJ, Buhlmann P (2007) Analyzing gene expression data in terms of gene
30 sets: methodological issues. *Bioinformatics* 23: 980-987.
- 31 28. Seaman SR, Muller-Myhsok B (2005) Rapid simulation of P values for product
32 methods and multiple-testing adjustment in association studies. *Am J Hum*
33 *Genet* 76: 399-408.
- 34 29. Hoffmann TJ, Kvale MN, Hesselson SE, Zhan Y, Aquino C, et al. (2011) Next
35 generation genome-wide association tool: design and coverage of a high-
36 throughput European-optimized SNP array. *Genomics* 98: 79-89.
- 37 30. Hoffmann TJ, Zhan Y, Kvale MN, Hesselson SE, Gollub J, et al. (2011) Design and
38 coverage of high throughput genotyping arrays optimized for individuals of
39 East Asian, African American, and Latino race/ethnicity using imputation and
40 a novel hybrid SNP selection algorithm. *Genomics* 98: 422-430.
- 41 31. Zhang H, Shi J, Liang F, Wheeler W, Stolzenberg-Solomon R, et al. (2014) A fast
42 multilocus test with adaptive SNP selection for large-scale genetic-
43 association studies. *Eur J Hum Genet* 22: 696-702.
- 44 32. Ge Y, Dudoit S, Speed T (2003) Resampling-based multiple testing for
45 microarray data analysis. *Test* 12: 1-77.

- 1 33. Devlin B, Roeder K, Wasserman L (2001) Genomic control, a new approach to
2 genetic-based association studies. *Theor Popul Biol* 60: 155-166.
- 3 34. Marchini J, Howie B (2010) Genotype imputation for genome-wide association
4 studies. *Nat Rev Genet* 11: 499-511.
- 5 35. Wang T, Elston RC (2007) Improved power by use of a weighted score test for
6 linkage disequilibrium mapping. *Am J Hum Genet* 80: 353-360.
- 7 36. Johnson ME, Zhao J, Schug J, Deliard S, Xia Q, et al. (2014) Two novel type 2
8 diabetes loci revealed through integration of TCF7L2 DNA occupancy and
9 SNP association data. *BMJ Open Diabetes Res Care* 2: e000052.
- 10 37. Lin CC, Chiang JH, Li CI, Liu CS, Lin WY, et al. (2014) Cancer risks among patients
11 with type 2 diabetes: a 10-year follow-up study of a nationwide population-
12 based cohort in Taiwan. *BMC Cancer* 14: 381.
- 13 38. Tsilidis KK, Kasimis JC, Lopez DS, Ntzani EE, Ioannidis JP (2015) Type 2 diabetes
14 and cancer: umbrella review of meta-analyses of observational studies. *BMJ*
15 350: g7607.
- 16 39. Benjamini Y, Hochberg Y (1995) Controlling the False Discovery Rate - a
17 Practical and Powerful Approach to Multiple Testing. *Journal of the Royal*
18 *Statistical Society Series B-Methodological* 57: 289-300.
- 19 40. Patil MA, Chua MS, Pan KH, Lin R, Lih CJ, et al. (2005) An integrated data analysis
20 approach to characterize genes highly expressed in hepatocellular
21 carcinoma. *Oncogene* 24: 3737-3747.
- 22 41. Hoshida Y, Nijman SM, Kobayashi M, Chan JA, Brunet JP, et al. (2009) Integrative
23 transcriptome analysis reveals common molecular subclasses of human
24 hepatocellular carcinoma. *Cancer Res* 69: 7385-7392.
- 25 42. Bulik-Sullivan BK, Loh PR, Finucane HK, Ripke S, Yang J, et al. (2015) LD Score
26 regression distinguishes confounding from polygenicity in genome-wide
27 association studies. *Nat Genet* 47: 291-295.
- 28 43. Li M, Wang K, Grant SF, Hakonarson H, Li C (2009) ATOM: a powerful gene-
29 based association test by combining optimally weighted markers.
30 *Bioinformatics* 25: 497-503.
- 31 44. Wessel J, Schork NJ (2006) Generalized genomic distance-based regression
32 methodology for multilocus association analysis. *Am J Hum Genet* 79: 792-
33 806.
- 34 45. Wu MC, Kraft P, Epstein MP, Taylor DM, Chanock SJ, et al. (2010) Powerful SNP-
35 set analysis for case-control genome-wide association studies. *Am J Hum*
36 *Genet* 86: 929-942.
- 37 46. Zhang H, Wheeler W, Wang Z, Taylor PR, Yu K (2014) A fast and powerful tree-
38 based association test for detecting complex joint effects in case-control
39 studies. *Bioinformatics* 30: 2171-2178.
- 40 47. Gamazon ER, Wheeler HE, Shah KP, Mozaffari SV, Aquino-Michaels K, et al.
41 (2015) A gene-based association method for mapping traits using reference
42 transcriptome data. *Nat Genet* 47: 1091-1098.
- 43 48. Lin DY, Zeng D (2010) On the relative efficiency of using summary statistics
44 versus individual-level data in meta-analysis. *Biometrika* 97: 321-332.
- 45