

1 **Developing a network view of type 2 diabetes risk pathways through integration of**
2 **genetic, genomic and functional data**

3

4 **Juan Fernández-Tajes^{1¶}, Kyle J Gaulton^{2¶*}, Martijn van de Bunt^{1,3,*}, Jason Torres^{1,3},**
5 **Matthias Thurner^{1,3}, Anubha Mahajan¹, Anna L Gloyn^{1,3,4}, Kasper Lage^{5,6}, Mark I**
6 **McCarthy^{1,3,4}.**

7

8 Affiliations

9

10 1. Wellcome Centre for Human Genetics, University of Oxford, Oxford, United Kingdom

11 2. University of California, Department of Pediatrics, San Diego, California, USA.

12 3. Oxford Centre for Diabetes, Endocrinology and Metabolism, University of Oxford, Oxford,
13 United Kingdom

14 4. Oxford NIHR Biomedical Research Centre, Churchill Hospital, Oxford, United Kingdom

15 5. Department of Surgery, Massachusetts General Hospital, Boston, Massachusetts, USA.

16 6. Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA.

17 7. Harvard Medical School, Boston, Massachusetts, USA.

18

19 [¶] These authors equally contributed to this work.

20 Corresponding author: Mark I McCarthy; mark.mccarthy@drl.ox.ac.uk.

21

22

23 *Current address: Department of Bioinformatics and Data Mining, Novo Nordisk A/S,

24 Maaloev, Denmark

25

26

27

28

29

30 **Abstract**

31

32 Genome wide association studies (GWAS) have identified several hundred susceptibility loci
33 for Type 2 Diabetes (T2D). One critical, but unresolved, issue concerns the extent to which
34 the mechanisms through which these diverse signals influencing T2D predisposition
35 converge on a limited set of biological processes. However, the causal variants identified by
36 GWAS mostly fall into non-coding sequence, complicating the task of defining the effector
37 transcripts through which they operate. Here, we describe implementation of an analytical
38 pipeline to address this question. First, we integrate multiple sources of genetic, genomic,
39 and biological data to assign positional candidacy scores to the genes that map to T2D
40 GWAS signals. Second, we introduce genes with high scores as seeds within a network
41 optimization algorithm (the asymmetric prize-collecting Steiner Tree approach) which uses
42 external, experimentally-confirmed protein-protein interaction (PPI) data to generate high
43 confidence subnetworks. Third, we use GWAS data to test the T2D-association enrichment
44 of the “non-seed” proteins introduced into the network, as a measure of the overall
45 functional connectivity of the network. We find: (a) non-seed proteins in the T2D protein-
46 interaction network so generated (comprising 705 nodes) are enriched for association to
47 T2D ($p=0.0014$) but not control traits; (b) stronger T2D-enrichment for islets than other
48 tissues when we use RNA expression data to generate tissue-specific PPI networks ; and (c)
49 enhanced enrichment ($p=3.9 \times 10^{-5}$) when we combine analysis of the islet-specific PPI
50 network with a focus on the subset of T2D GWAS loci which act through defective insulin
51 secretion. These analyses reveal a pattern of non-random functional connectivity between
52 causal candidate genes at T2D GWAS loci, and highlight the products of genes including
53 *YWHAG*, *SMAD4* or *CDK2* as contributors to T2D-relevant islet dysfunction. The approach we

54 describe can be applied to other complex genetic and genomic data sets, facilitating
55 integration of diverse data types into disease-associated networks.

56

57 **Author summary**

58

59 We were interested in the following question: as we discover more and more genetic
60 variants associated with a complex disease, such as type 2 diabetes, will the biological
61 pathways implicated by those variants proliferate, or will the biology converge onto a more
62 limited set of aetiological processes? To address this, we first took the 1895 genes that map
63 to ~100 type 2 diabetes association signals, and pruned these to a set of 451 for which
64 combined genetic, genomic and biological evidence assigned the strongest candidacy with
65 respect to type 2 diabetes pathogenesis. We then sought to maximally connect these genes
66 within a curated protein-protein interaction network. We found that proteins brought into
67 the resulting diabetes interaction network were themselves enriched for diabetes
68 association signals as compared to appropriate control proteins. Furthermore, when we
69 used tissue-specific RNA abundance data to filter the generic protein-protein network, we
70 found that the enrichment for type 2 diabetes association signals was enhanced within a
71 network filtered for pancreatic islet expression, particularly when we selected the subset of
72 diabetes association signals acting through reduced insulin secretion. Our data demonstrate
73 convergence of the biological processes involved in type 2 diabetes pathogenesis and
74 highlight novel contributors.

75

76

77

78

79 **Introduction**

80

81 The rising prevalence of type 2 diabetes (T2D) represents a major challenge to global health
82 [1]. Current strategies for both prevention and treatment of T2D are suboptimal, and
83 greater insight into the mechanisms responsible for the development of this condition is a
84 prerequisite for further advances in disease management [2].

85

86 The identification of human DNA sequence variants which influence predisposition to T2D
87 provides one of the most direct approaches for deriving mechanistic insight. However,
88 current understanding of the genetic architecture of T2D indicates that the genetic
89 component of T2D predisposition likely involves variation across many thousands of loci [3,
90 4]. Close to 500 independent genetic signals for which there is robust evidence of a
91 contribution to T2D predisposition have been identified, largely through genome-wide
92 association studies, supplemented by analysis of exome- and genome-sequence data [4-6].

93 This profusion of genetic signals has raised questions concerning the extent to which the
94 inherited susceptibility to complex traits such as T2D can be considered to occupy finite
95 biological space [7]. In other words, as the number of loci influencing T2D risk increases, will
96 the mechanisms through which these are found to mediate the development of this
97 condition continue to proliferate, or will they start to converge around a limited set of
98 pathways?

99

100 There are two main challenges in addressing this key question. First, whilst a minority of the
101 causal variants underlying these association signals are coding (and therefore provide direct
102 inference regarding the genes and proteins through which they act), most lie in regulatory

103 sequence. This makes assignment of their effector transcripts a non-trivial exercise, and
104 obscures the downstream mechanisms through which these variants impact T2D-risk [8-10].
105 This challenge can increasingly be addressed through the integration of diverse sources of
106 relevant data including (a) experimental data (e.g. from studies of cis-expression or
107 conformational capture) which link regulatory risk-variants to their likely effectors [11, 12];
108 and (b) evaluations of the biological evidence connecting each of the genes within a GWAS-
109 associated region to the disease of interest. In the present study, focussing on a set of
110 approximately 100 T2D-risk loci with the largest effects on T2D predisposition, we use a
111 range of information to derive “positional candidacy” scores for each of the coding genes
112 mapping to T2D-associated GWAS intervals.

113

114 The second challenge lies in the requirement to define functional relationships between sets
115 of candidate effector transcripts in ways that are robust, and, in particular, orthogonal to
116 the data used to assign candidacy in the first place [13, 14]. Solutions for the second
117 challenge are less well-developed but generally involve some type of network analysis (e.g.
118 weighted gene correlation network analysis [WGCNA]) and application of the “guilt-by-
119 association” framework to infer function [15-17]. However, recourse to co-expression
120 information, or functional pathway enrichment methods to generate and evaluate such
121 networks runs the risk of introducing circularity, given that information on expression and
122 function typically contributes (whether explicitly or not) to assignments of effector
123 transcript candidacy. The use of protein-protein interaction data provides one possible
124 solution to this conundrum [18]. In the present study, we make use of external protein-
125 protein interaction data from the InWeb3 dataset [19, 20] to evaluate and characterise the

126 connectivity of the T2D candidate effector transcripts in terms of their ability to nucleate
127 empirically-confirmed interactions between their encoded proteins.

128

129 **Materials and methods**

130

131 **Positional candidacy score derivation**

132

133 We developed a framework to score the candidacy of genes mapping to GWAS association
134 signals which aggregated data from multiple sources. The information collected fell into two
135 categories. First, we used regression-based approaches to link disease-associated variants
136 (most of which map into non-coding sequence and are therefore presumed to act through
137 transcriptional regulation of nearby genes) to their likely effector transcripts, using a
138 combination of variant-based annotations and expression QTL data. Second, we scored
139 each of the genes in these GWAS regions for disease-relevant biological function. We
140 combined the two measures to generate a “positional candidacy score” (PCS) for each gene.
141 We applied this framework to 1895 genes located within a 1Mb interval around the lead
142 variants from 101 T2D GWAS regions. These represent the loci with the largest effect sizes
143 for T2D, as identified in European subjects as of early 2017 [4, 6, 21]: Supplementary Table
144 1). The 1Mb intervals contained 1895 genes.

145

146 ***Mapping effector transcripts to GWAS signals***

147

148 At each of the 101 loci, we collected summary T2D case-control association data ($-\log_{10}p$
149 values) for all 1000 Genomes variants in the 1Mb interval surrounding the lead variants [6].
150 We then annotated variants in each interval using gene-based annotations for all genes in
151 the interval from several sources. First, we collected relevant discrete annotations for all

152 protein coding genes in GENCODE (version 19) [22] within the interval including (a) coding
153 exon location; (b) promoter location (defined as 1kb region upstream of the transcription
154 start site [TSS]); (c) distal regulatory elements correlated with gene activity from DNaseI
155 hypersensitivity (DHS) data (ENCODE version 3) [23]. We assigned each variant a binary
156 value based on whether it overlapped one of the discrete annotations for a gene in the
157 interval (exon, promoter, distal element). Second, we collected summary statistic
158 expression QTL (eQTL) data from liver, skeletal muscle, whole blood, subcutaneous adipose
159 and visceral adipose (GTEx version 6) [24] and pancreatic islets [11]. We assigned each
160 variant the $-\log_{10}p$ value of eQTL association for each cell type for each gene in the interval.
161 Third, we calculated the distance of each variant to the TSS of each gene in the interval, and
162 assigned each variant the inverse TSS distance for each gene (i.e. variants closer to the TSS
163 have higher values). Variants without values in the eQTL datasets were removed from the
164 analysis.

165
166 We then performed feature selection for each T2D locus separately using elastic net
167 regression (R package glmnet) with the T2D p-values as the outcome variable and binary
168 genomic annotations (exon, promoter, distal element), distance to TSS, and cell type cis-
169 eQTL p-values for each gene in the interval as the predictor variables. We also included
170 minor allele frequency and imputation quality of each variant at the locus as predictor
171 variables. We obtained the effects of features selected from the resulting model. We
172 applied a 10-fold scaling factor to coding exon features, based on known enrichment of T2D
173 variants in coding exons [25, 26]. Where multiple features were selected for the same gene
174 (e.g. distal DHS site and tissue eQTL) we summed the effects for that gene. We considered

175 the summed effects of features for each gene as the ‘variant link score’ in subsequent
176 analyses.

177

178 ***Semantic mapping of gene functional annotations***

179

180 We also derived a second score of the T2D-relevance for each gene within the 101 GWAS
181 intervals based on the annotations for each within data from gene ontology (GOA, version
182 157), the mouse genome database (MGD, version 6.08), and biological pathways (KEGG
183 (version 83.1), compiling these annotations into a single document per gene. We also
184 created a query document of empirically-compiled terms we considered relevant to T2D
185 pathophysiology (listed here: [https://github.com/kjgaulton/gene-](https://github.com/kjgaulton/gene-pred/blob/master/res/T2D.query.manual.txt)
186 [pred/blob/master/res/T2D.query.manual.txt](https://github.com/kjgaulton/gene-pred/blob/master/res/T2D.query.manual.txt)). Both gene documents were converted into a
187 word matrix. We calculated the total number of unique words *across* all documents N , after
188 removing a list of commonly used “stop” words from PubMed
189 (<https://www.ncbi.nlm.nih.gov/pubmed/>) and stemming the remaining words. We
190 weighted each word w for each gene document g using “term frequency (TF)” minus
191 “inverse document frequency” defined as:

$$TF = f_{g,w}$$
$$IDF = \log\left(\frac{N}{n_w}\right)$$

192

193 where n_w is the number of documents containing word w . We defined the value (g_w) of
194 word w in gene document g as:

$$g_w = TF * IDF$$

195

196 and applied latent semantic analysis (LSA) using singular value decomposition of the
197 weighted matrix M

$$M = TSD^T$$

198

199 where T is the left singular vector matrix of terms, D is the right singular vector matrix of
200 documents, S is the diagonal matrix of singular values, and the number of dimensions was
201 determined by the function `dimcalc_share` from the `lsa` package [27]. We used the resulting
202 matrices to identify genes with functional attributes that indicated relevance to T2D
203 pathogenesis. For each gene document vector g , we calculated similarity scores $S_{i,q}$ using
204 the dot product between the gene vector and the T2D query vector q

$$S_{g,q} = g \cdot q$$

205
206 From these data, we extracted similarity scores for the 1895 genes of interest, which we
207 considered the ‘semantic score’ in subsequent analyses.

208 209 **Combining gene scores**

210
211 For each of the 1895 genes, we scaled scores from these two analyses to the sum of scores
212 for each of the x genes at each locus resulting in a semantic score s_g and variant link score v_g .
213 To calculate a positional candidacy score (PCS), we averaged the two scores and rescaled
214 across all x genes at each locus.

$$CS_g = \frac{s_g + v_g}{\sum_{i=1}^x s_i + v_i}$$

215 216 217 **Network modelling**

218
219 **Selection of the “seed node set”.** At each GWAS locus, we defined the sets of genes that,
220 after ranking the genes for each locus by decreasing PCS, generated a cumulative PCS
221 exceeding 70%. This reduced the set of 1895 genes of interest to 451 “seed” nodes for
222 subsequent network analysis. We performed network analyses using an updated version of
223 InWeb3, a previously-described comprehensive map of protein-protein interactions,
224 containing 169,736 high-confidence interactions between 12,687 gene products compiled

225 from a variety of sources [19, 20]. We updated the version used in [20], by updating
226 outdated gene symbols and restricting interactions to those deemed “high-confidence”
227 (score >0.124).
228
229 **Prize-collecting Steiner Tree formulation.** We formulated the task of examining the
230 connectivity of GWAS positional candidates (the set of 451 “seed” genes) within protein-
231 protein interaction space as an asymmetric prize-collecting Steiner tree (APCST) problem.
232 APCST-like approaches have been widely used to solve network-design problems [28-30].
233 The APCST seeks to connect “seed” nodes (in formal nomenclature, “terminals”) to collect
234 “prizes”, using confirmed protein-protein interactions as edges. Prizes are weights added to
235 seed nodes: in our analysis, these correspond to the PCS values for each “seed” gene,
236 derived from the -omic integration approach. “Linking” (formally, “Steiner”) nodes (that is,
237 proteins/genes not included in the seed set) can be introduced into the network, where
238 necessary. Network expansion is controlled by the balance between the benefits of adding a
239 particular node (increased connectivity between seed genes, driven by the collection of
240 prizes) vs. the costs of adding additional edges (based on a function which penalises
241 expansion of the network). In mathematical terms, we defined the APCST as follows: given a
242 directed graph $G = (V, A)$, arc costs $c: A \mapsto \mathbb{R} \geq 0$, node prizes $p: V \mapsto \mathbb{R} \geq 0$ and a set
243 of fixed terminals T_f the goal is to find an arborescence $S = (V_s, A_s) \subseteq G$ that spans T_f such
244 that the following function is maximized:

$$P(S) = \beta \sum_{i \in V_s} p_i - \sum_{(i,j) \in A_s} c_{i,j}$$

245 In this formulation, we reward the inclusion of nodes $i \in V_s$ with higher prizes (that is,
246 higher PCS values) (first term of equation) while paying costs for including edges (second

247 term of equation). The parameter, β , scales the importance of node prizes versus edge costs
248 in the optimization and can be used to titrate the size of the generated network. We tested
249 different values of β (between 4 and 30) and selected $\beta = 8$ that produced a manageable
250 network size (~130 genes) and included >25% of the seed node set (**S1 Fig**).

251

252 Although this problem is NP-hard (nondeterministic polynomial time-hard) [31], the APCST
253 algorithm is found to be efficient in calculating exact and proximal solutions (DIMACs 11th
254 challenge, <http://dimacs11.zib.de/>). The branch-and-bound algorithm, implemented in
255 *dapstp* algorithm (<https://github.com/mluipersbeck/dapcstp>) and using default parameters,
256 was used to find the optimal (or near optimal) APCST solution.

257

258 **Generation of networks using dapcst algorithm**

259

260 We used a particular variation of the ACPST (“root-ACPST”) where the search for the
261 optimal solution starts in a specific node. This allowed us to force each seed node in turn to
262 be included in the network, in contrast to the default APCST method which initialises
263 network construction from the nodes with higher weights. For the main T2D analysis,
264 therefore, the algorithm was run 451 times, once for each “seed” node. Runs generating a
265 network of >10 nodes (353 networks, median 155 nodes) were combined to form an
266 ensemble network from the union of all n networks. This was reprojected onto the InWeb3
267 interactome to recover missing connections across nodes. As this final network represents a
268 superposition of many different networks, linking nodes may sometimes appear at the
269 periphery.

270

271 We assessed the specificity of each node in the final network by running the algorithm 100
272 times with the same parameter settings, but with random input data. We define specificity
273 in this context as the complement of the percentage with which a given seed or linking node
274 from the final network appears in runs generated from random input data. For each random
275 run, we selected, from the InWeb3 interactome, random seed nodes matching the binding
276 degree distribution of the observed set of seeds, and assigned them the same prize value as
277 the original. Using the final parameter settings, we found that the included linking nodes
278 were highly specific to our particular data, with 80% of them having a specificity higher than
279 75% (**S2 Fig**).

280

281 **Testing network for Enrichment in GWAS signal.**

282

283 To evaluate the extent to which the PPI network provided functional connectivity between
284 positional candidates across loci, we measured the enrichment of the linking nodes for T2D
285 association signals. This avoided the circularity of using co-expression or functional data to
286 evaluate connectivity (as both contributed to the PCS determination). We generated gene-
287 wise p values using the PASCAL method [32] from large-scale GWAS studies across a set of
288 33 traits (using data extracted from public repositories) including a recent meta-analysis of
289 T2D GWAS data from ~150,000 Europeans [6]. We mapped these gene-wise association p-
290 values to linking nodes, and converted them to Z-scores using the standard normal
291 cumulative distribution, $Z_i = \phi^{-1}(1 - p_i)$. We then quantified GWAS enrichment by
292 aggregating the Z-scores using Stouffer's method:

$$Z \sim \frac{\sum_{i=1}^k Z_i}{\sqrt{k}}$$

293 where Z_i is the Z-score for the gene-wise p value for linking node i and k is the number of
294 linking nodes in the network. Then, by permuting the InWeb3 network using a node
295 permutation scheme, we compared the observed enrichment in GWAS signals to a random
296 expectation, allowing us to calculate a nominal p value as:

$$(Nominal\ p\ value)\ P_n = \frac{\#(Z > Z_m)}{\#(total\ permtaions)}$$

297 where p_n is the permuted p value generated in the permutation scheme. In this last step,
298 the binding degree of all genes in the network is taken into full consideration (i.e. they all
299 have the same binding degree as provided by the APCST network). To minimise bias arising
300 from the co-localisation of genes with related functions (which is a feature of some parts of
301 the genome), in each of these permutations we only considered proteins whose genes
302 mapped outside a 1Mb window around the lead SNP for any significant GWAS association
303 for that trait.

304

305 **APCST model clustering**

306

307 To aid interpretation of the PPI networks, we used a community clustering algorithm that
308 maximizes network modularity and which breaks the full APCST model into smaller sub-
309 networks [33].

310

311 **GTEx and Islet RNAseq datasets**

312

313 The InWeb3 PPI network we used is generated from empirically-confirmed interactions, but
314 nevertheless includes many interactions that, owing to restricted tissue-specific expression,
315 are unlikely to be biologically relevant. We used tissue-specific RNA expression data to filter
316 the overall InWeb3 network and thereby generate *in silico* “tissue-specific” PPI networks,
317 using TPM counts from GTEx (version 7: <https://www.gtexportal.org/home/>, last accessed

318 21 Oct 2017), complemented by human pancreatic islet data from [11]. Proteins with mRNA
319 TPM counts <1 in over 50% of samples for that tissue were removed from the InWeb
320 network, allowing us to generate in silico PPI networks for 46 tissues.

321

322 ***Functional Enrichment Analysis***

323

324 Gene Set Enrichment (GSE) of networks and sub-networks were assessed with ClueGO [34]
325 using GO terms and REACTOME gene sets [35]. The enrichment results were grouped using
326 a Cohen's Kappa score of 0.4 and terms were considered significant when Bonferroni
327 adjusted p-value <0.05 and at least 3% of the genes contained in the tested gene set were
328 included in the network. Cohen's Kappa statistic measures the gene-set similarity of GO
329 terms and REACTOME pathways and allowed us to group enriched terms into functional
330 groups that improve visualization of enriched pathways.

331

332 **Results**

333

334 **Prioritizing positional candidates at T2D risk loci**

335

336 We implemented a framework to derive positional candidacy scores (PCS) for genes within
337 T2D GWAS loci through the aggregation of two main types of data (**Fig. 1; Methods**). First,
338 we used regression-based approaches to link disease-associated variants (most of which
339 map to non-coding sequence and are therefore presumed to act through transcriptional
340 regulation of nearby genes) to their likely effector transcripts, using a combination of
341 variant-based annotations and expression QTL data. Second, we scored each of the genes in
342 these GWAS regions for disease-relevant biological function using semantic mapping of
343 gene functional annotations from Gene Ontology, Mouse Genome Database, and KEGG. We

344 combined the evidence from both approaches, normalized across all genes at each GWAS
345 locus, to generate the PCS for each gene.

346

347 **Figure 1 Overview of the Data Integration pipeline.** We collected variants in the 1Mb
348 interval surrounding index variants at each of the 101 T2D GWAS loci along with relevant
349 annotations for all protein coding genes in GENCODE including coding exon location,
350 promoter location, distal regulatory elements correlated with gene activity from DNaseI
351 hypersensitivity (DHS) data and summary statistic expression QTL (eQTL) data from T2D-
352 relevant tissues. This, combined with information at gene level from a semantic similarity
353 metric, allowed us to define positional candidacy scores for each gene in the GWAS
354 intervals. These genes were projected into the InWeb3 data set using a Steiner tree
355 algorithm to define a PPI network that maximises candidate gene connectivity. This network
356 was further analysed to find processes, pathways, and genes implicated in T2D
357 pathogenesis.

358

359 We applied this method to score 1,895 genes mapping within a 1Mb interval around the
360 lead variant at 101 T2D GWAS regions. This list of 101 T2D loci was assembled from a series
361 of recent large-scale T2D GWAS studies [4, 6, 21] and represents the largest-effect T2D
362 GWAS loci identified as of early 2017. The 1Mb interval was selected to capture the majority
363 of cis-acting regulatory effects (95% of cis-eQTLs map within 445kb of the lead SNP [24]),
364 and is therefore also likely to encompass most potential effector genes [36]. We observed
365 only weak correlation between the semantic and risk variant link scores for the 1,895
366 positional candidates ($r^2=0.05$, $p=0.01$), indicating that these provide distinct information
367 **(S3 Fig).**

368

369 Most (71%) of the 1,895 genes had minimal evidence linking them to a causal role in T2D
370 pathogenesis ($PCS<0.05$) **(S3 Fig)**. However, 95% of T2D loci included at least one gene
371 (median, 3) with $PCS>0.10$, and at 70% of loci, there was at least one gene with $PCS>0.20$
372 **(S3 Fig)**. The top-scoring genes across the 101 loci (such as *IRS1* [$PCS=0.69$], *SLC30A8*

373 [PCS=0.77], *HNF1B* [PCS=0.54]) include many of the genes with the strongest prior claims for
374 involvement in T2D risk, prior claims which arise in part from data used to generate the
375 PCSs. For example, these genes each contain rare coding variants directly implicated in
376 development of T2D (or related conditions): these variants are independent of the common
377 variant GWAS signals, but their relationship to diabetes is likely to have been captured
378 through the semantic mapping. The PCS also highlighted several other highly-scoring
379 candidates with known causal roles in relation to diabetes and obesity such as *MC4R*
380 (PCS=0.43), *WFS1* (0.41), *ABCC8* (0.37), *LEP* (0.27), *GCK* (0.24), and *HNF1A* (0.23). At other
381 loci, these analyses highlighted candidates that have received scant attention to date: for
382 example, *CENPW* (PCS=0.83) scored highly both in terms of semantic links to T2D-relevant
383 processes and an adipose cis-eQTL linking the T2D GWAS SNP to *CENPW* expression.

384

385 To define the seed-genes for subsequent PPI analyses, we gathered the sets of genes that,
386 after ranking the transcripts for each locus by decreasing PCS, cumulatively accounted for at
387 least 70% of the candidacy score for each locus. For example, at the *TP53/INP1* locus, where
388 the gene-specific PCSs range from 0.01 to 0.16 across a total of 17 mapped genes, the seed-
389 gene set includes the first six (**S4 Fig**). This filter identified a total of 451 positional
390 candidates across the loci, reducing the median number of genes per locus from 19 to 6 (**S4**
391 **Fig**). This filtering mostly removes genes with low PCS values: the proportion of genes with
392 $PCS < 0.05$ falls from 71% to 12%, while most genes with $PCS > 0.1$ or > 0.2 are retained (**S3**
393 **Fig**).

394

395 This prioritisation process ensures that genes with the strongest combined causal evidence
396 are favoured for network modelling, resulting in sets of seed genes that are more extensive

397 than selection based on proximity alone (such as “nearest gene” approaches that seek to
398 generate networks from only the genes mapping closest to the lead variants) but smaller
399 than those which consider all regional genes of equal weight (“all gene” approaches). Note
400 that our strategy does not require complete ascertainment of all true causal genes within
401 this set of 451 genes: true effector genes excluded from the prioritised set of 451 genes (e.g.
402 because they map more distal to the lead variant than 500kb) remain available for
403 “discovery” through the network modelling described below.

404

405 **Building a T2D-relevant protein-protein interactome**

406

407 We set out to test whether this list of prioritised candidates could be used to characterise
408 the functional relationships between genes (and proteins) implicated in T2D pathogenesis.
409 Because the PCS scores used to prioritise the genes already incorporated (explicitly or
410 otherwise) diverse types of functional and expression data, biasing any assessment of
411 connectivity in these domains, we focused the network analysis around protein-protein
412 interaction (PPI) data. To do so, we projected these 451 genes onto externally-derived,
413 empirically-driven PPI resources (InWeb3) [19, 20] using an established network modelling
414 strategy (the Asymmetric Prize-Collecting Steiner Tree (APCST)) (**Fig1; Methods**). In this
415 analysis, the 451 positional candidates represent “seed” nodes which are used by the APCST
416 algorithm to generate PPI networks which seek (with appropriate penalties to prevent
417 frivolous propagation) to connect as many seed nodes as possible to each other, either
418 directly, or using other (non-seed) proteins as links (“linking” nodes). The network topology
419 is dependent only on the PCS values of the “seed” genes which are carried forward as
420 weights into the APCST analysis, the confidence scores for each of the empirical PPI

421 interactions in InWeb3, and the beta value used to tune the overall size of the PPI network
422 generated (see **Methods**).

423

424 We operationalised the PPI network as follows (see **Methods**). Using each “seed” gene in
425 turn, we used InWeb3 data to generate a PPI network that maximised the connectivity to
426 other seed genes within the constraints of the APCST model. Of the 451 seed genes, 98
427 failed to produce a network exceeding 10 nodes. The remaining 353 networks had a median
428 of 110 seed and 45 linking nodes and were combined into an ensemble network, which was
429 again projected into the InWeb3 interactome to recover missing connections between
430 nodes. The final network contained 705 nodes (431 seed nodes, 274 linking nodes) and
431 2678 interactions (**Fig 2**). Based on random networks generated with the same algorithm
432 (see **Methods**), 80% of the linking nodes have a specificity for membership of the final
433 network exceeding 75%, indicating that these linking nodes do not simply reflect generic
434 hubs in PPI space (**S2 Fig**).

435

436 **Figure 2 APCST final network.** The final PPI network generated from the T2D GWAS interval
437 genes includes 431 seed nodes and 274 linking nodes connected by 2,678 interactions. We
438 divided this network into 20 sub-networks (communities) using a community clustering
439 algorithm that maximizes network modularity (33), and highlighted enrichment of specific
440 biological processes contained within these based on Gene Ontology terms and REACTOME
441 pathways. Coloured nodes represent seed nodes, whereas grey nodes represent linking
442 nodes.

443

444

445 **The T2D PPI network is enriched for T2D associations**

446

447 If the final network truly provides novel insights into the functional relationships between
448 genes thought to be mediating T2D predisposition, we reasoned that the “linking” genes
449 (those brought into the network purely on the basis of external data indicating their

450 protein-level interaction with seed genes) should be enriched for other seed gene
451 characteristics. To avoid circularity arising from validation using data types that had
452 contributed to the generation of the original PCS weights, including measures of gene
453 function (eg GO, KEGG) or RNA expression data, we turned to T2D GWAS data, looking for
454 evidence that the genes encoding the linking proteins were themselves enriched for T2D-
455 association signals. For this, we used T2D-association data from a set of ~150,000 European
456 T2D case-control subjects imputed to 1000 Genomes [6]. Briefly, the linking nodes were
457 mapped to gene-wise association p-values generated from the GWAS results using PASCAL
458 [32]. The significance of the collective enrichment of these gene-wise p-values was obtained
459 by permuting the observed set of linking nodes with equivalent sets of “random” nodes
460 from the InWeb3 database, matched for binding degrees (see **Methods**). To minimise the
461 prospects of picking up false signals arising from the combination of local LD and the non-
462 random genomic location of functionally-related genes, we excluded all genes from the
463 1Mb-window around the 101 lead variants from these analyses.

464

465 Compared to the distribution of scores in the permuted background, the gene-wise p-values
466 for linking genes in the empirical reconstructed network demonstrated significant
467 enrichment of T2D association ($p=0.0014$). To confirm that this enrichment was specific to
468 T2D, we repeated the analysis, retaining the same PPI final network, but instead using
469 GWAS data (and PASCAL-derived gene-wise p-values) from 33 different traits across a wide
470 range of disease areas. The only other traits displaying evidence of GWAS enrichment within
471 the linking nodes of the T2D PPI network were those for anthropometric traits with known
472 relevance to T2D pathophysiology (**S5 Fig**).

473

474 To gain insights into how the linking nodes of our final network contribute to T2D biology,
475 we used the DisGeNET database [37], which collates gene-disease information from public
476 data as well as from literature via natural language processing tools. We focused on the 274
477 linking nodes included in our model to avoid circularity arising from using the seeds, and
478 identified 92 (~33%) with known links to T2D (**S2 Table**). Examples include: (a) *NEUROD1*
479 which encodes a transcription factor that is involved in the development of the endocrine
480 cell lineage and has been implicated in monogenic diabetes [38]; (b) *PRKCB* involved in
481 insulin resistance [39], and (c) *GNAS*, implicated in beta-cell proliferation [40]. For this last
482 gene, mice knockouts have been shown to produce phenotypes concordant with diabetes
483 [41]. These examples demonstrate the potential of these analyses to draw in “linking” nodes
484 as related to T2D even when they are not located within genome-wide association signals.

485

486 **The T2D PPI network captures biological processes relevant to disease pathogenesis**

487

488 To increase biological interpretability, we next sought to split the large final PPI network of
489 705 nodes into smaller sub-networks of closely-interacting proteins (“communities”). Using
490 the algorithm proposed by [33], we identified 18 such communities (each containing
491 between 2 and 186 nodes) (**Fig 2**). We performed enrichment analyses on each community
492 using GO and REACTOME datasets, this time including both seed and linking nodes. We
493 observed that the individual sub-networks were enriched for processes including “glucose
494 homeostasis” and “insulin receptor signalling cascade” (sub-network 1), “Wnt” and “NIK/NF-
495 kappaB signalling pathways” and “cellular response to stress” (subnetwork 2), “COPII vesicle
496 coating” and “Wnt ligand biogenesis and trafficking” (sub-network 3), “regulation of insulin
497 secretion” (sub-network 8), and “glucagon signalling in metabolic regulation” (sub-network
498 12) (**Fig 2, S3 Table**). This pattern of functional enrichment is broadly consistent with

499 existing knowledge regarding aspects of T2D pathogenesis [42-44]. We saw no evidence in
500 support of certain processes that have been proposed as contributors to T2D pathogenesis
501 such as mitochondrial function, or oxidative phosphorylation [45, 46], in line with the
502 paucity of evidence linking these processes to T2D risk in standard gene-set enrichment
503 analyses [4, 21].

504

505 **Information on tissue-specificity enhances the model**

506

507 The APCST model described above was constructed from a generic, tissue-agnostic PPI
508 network. As a result, it features edges that, whilst they may be supported by the empirical
509 data used to generate the InWeb3 database, are unlikely to be pathophysiologically
510 relevant, due to mutually-exclusive tissue-specific expression patterns. We hypothesised
511 that the use of tissue-specific interactomes, focused on T2D-relevant tissues, would allow us
512 to refine the reconstructed PPI network, and might enhance the GWAS enrichment signal. In
513 the absence of empirical PPI data for all relevant tissues, we generated these tissue-specific
514 PPI networks by filtering on RNA transcript abundance. Starting from the generic final APCST
515 network, we removed, for each tissue, all nodes (and their corresponding edges) with little
516 or no transcriptional activity (see **Methods**). In all, we generated tissue-specific PPI
517 networks, using RNA-Seq data sourced from 46 different tissues, 45 (including fat, liver and
518 skeletal muscle) from GTEx (v7) [24][www.gtexportal.org] (median number of individuals =
519 235) together with a set of human islet RNA-seq data (n=118) [11], which had been
520 reprocessed through a GTEx-aligned pipeline.

521

522 We then repeated the T2D GWAS signal enrichment analysis (“linking” nodes only; 100,000
523 permutations) across each of these 46 tissue-specific PPI networks. We detected broad

524 enrichment for T2D association in linking nodes across many of these tissue-specific
525 networks: this likely reflects the fact that these tissue-specific networks remain highly
526 overlapping (**S6 Fig**). Nonetheless, with the exception of whole blood, the strongest
527 enrichment signal for T2D GWAS data was observed in the islet-specific PPI network (**Fig 3**).
528 This enrichment was less significant ($p=0.019$) than that observed in the full network
529 ($p=0.0014$), but this, at least in part, reflects the reduction in the number of linking nodes in
530 the islet-specific network (from 274 to 229). Other tissues implicated in T2D pathogenesis
531 such as adipose, skeletal muscle or liver generated more limited evidence of enrichment (**Fig**
532 **3**). This pattern of enrichment (favoring islets, and to a lesser degree, adipose) mirrors
533 equivalent observations for other tissue-specific annotations (including cis-eQTL signals and
534 active enhancers) with respect to T2D association data [10, 11].

535

536 **Figure 3. GWAS signal enrichment in tissue-specific interactomes.** RNA-Seq data was used
537 to filter the overall InWeb3 network and generate in silico tissue-specific networks that
538 maximise connectivity between GWAS interval genes. Linking nodes within these networks
539 were then tested for enrichment for GWAS signals using a permutation scheme. Each dot in
540 the figure depicts the $-\log_{10}$ p-value for enrichment for signals in a given GWAS dataset, for
541 each of the 46 tissues. Dot colors reflect the GWAS phenotypes with T2D in the larger red
542 color. The dotted red line represents the nominal value of significance ($p=0.05$). Islet
543 showed the second strongest enrichment signal for T2D.

544

545 **Further enhancement of model using GWAS locus subsets**

546

547 To further refine the analysis, we took account of the multi-organ nature of T2D and,
548 specifically, of evidence that it is possible, using patterns of association across T2D-related
549 quantitative traits such as BMI, lipids and insulin levels, to define subsets of T2D GWAS loci
550 which impact primarily on insulin secretion and those that perturb insulin action [47-49].
551 We reasoned that the former would be expected to show preferential enrichment within
552 the islet-filtered PPI network. Accordingly, we built APCST networks (both generic and

553 filtered for expression in islets exactly as above) formed from the sets of high-PCS seed
554 genes mapping to each of seven T2D GWAS locus subsets defined in two recent publications
555 [48], [49].

556

557 In both the islet-specific (**Fig 4**) and the generic network (**S7 Fig**), the strongest signals for
558 GWAS enrichment were seen for loci in the three subsets (beta-cell [BC] in [48]; acute
559 insulin response [AIR] and peak insulin response in [49]) comprised of T2D GWAS loci which
560 influence T2D risk primarily through a detrimental effect on insulin secretion (**Fig 4; S6 Fig**).

561 In particular, there was striking enrichment in the islet-specific PPI network for linking nodes
562 in analyses of the BC ($p=3.9 \times 10^{-5}$) and AIR ($p=1.9 \times 10^{-4}$) T2D GWAS locus subsets.

563

564 **Figure 4. GWAS signal enrichment in islet-specific network derived from T2D GWAS**
565 **subsets.** We built APCST networks filtered for islet RNA-expression for each of the subsets
566 of T2D GWAS loci defined by shared mechanistic mediation (refs[48], [49]). Enrichment in
567 GWAS signals for linking nodes only was tested using a permutation scheme. Each dot in the
568 figure depicts the $-\log_{10}$ p-value of enrichment for association signals in a particular GWAS
569 analysis. The results for T2D GWAS enrichment for the APCST networks built around the
570 different T2D GWAS subsets are also represented (large red dots). The dotted red line
571 represents nominal significance ($p=0.05$). The strongest enrichment for T2D GWAS data in
572 islet-filtered PPI data is observed for subsets of loci acting through reduced insulin
573 secretion. In the cluster hairballs for the seven T2D GWAS locus subset categories, nodes
574 are coloured according to their PCS with grey nodes representing linking nodes.

575

576 As before, we were interested to see whether this marked convergence of PPI signal (as
577 assessed by the enrichment of T2D association signals in linking nodes) was T2D-specific.

578 We therefore repeated the enrichment analysis using GWAS data from 33 additional traits.

579 For each trait, we took the APCST networks generated using the seven T2D locus subsets

580 and assessed the “linking” nodes in those networks with respect to enrichment for

581 respective gene-wise association p-values. We found broad levels of enrichment for

582 association signals for T2D-related phenotypes including (quantitative) glycemc traits, lipid
583 levels, anthropometric and cardiovascular traits, which are consistent with known GWAS
584 signal overlap. However, we saw very limited enrichment for other (non-diabetes related)
585 traits. Furthermore, the patterns of enrichment were consistent with underlying
586 physiological expectation: GWAS enrichment for anthropometric and lipid phenotypes was
587 most marked in the APCST networks generated from the insulin-resistant subset of T2D loci
588 (category “insulin response” in [48]), whilst T2D remained the most enriched phenotype for
589 the subsets related to insulin secretion (**Fig 4**).

590

591 These analyses demonstrated that parallel efforts to refine the phenotypic impact of T2D
592 GWAS loci, and the tissue-specificity of the underlying PPI dataset used to generate the
593 APCST network, resulted in progressive, biologically-appropriate, improvement of the
594 enrichment signal observed at the “non-seed” proteins represented within the network.

595

596 **Biological insights**

597

598 To better understand the biological function of the highly-enriched PPI network generated
599 by the intersection of islet-specific expression, and the subset of T2D GWAS loci acting
600 through reduced islet function (henceforth, the “islet network”), we performed a Gene Set
601 enrichment analysis using GO and REACTOME terms (**S4 Table**). Captured pathways
602 included well-known biological processes of “glucose homeostasis” ($p = 1.5 \times 10^{-4}$),
603 “regulation of WNT signalling pathway” ($p = 8.9 \times 10^{-3}$), “response to insulin” ($p = 6.2 \times 10^{-4}$), and
604 “pancreas development” ($p = 3.0 \times 10^{-5}$).

605

606 This islet network included many “seed” genes with a high T2D PCS score (**Fig 4**) including
607 *SOX4* ([PCS=0.62] at the locus usually named for *CDKAL1*), and *ATXN7* ([PCS=0.57] at the
608 locus named for *ADAMTS9*). Some of the loci (e.g. *TLE4*, *CAGE1* and *GCK*) are represented by
609 a single “seed” because the PCS for the highest-ranking gene exceeded 0.70. At other loci,
610 this islet network does not include the gene with the highest PCS score for the respective
611 GWAS signal, but instead features an alternative gene from the same locus on the basis of
612 its better connectivity within the network. Examples such as the gene *TBS* [PCS=0.21] at the
613 *ZBED3* locus, and *THRB* [PCS=0.43] at the *UBE2E2* locus, demonstrate how the PPI data
614 provides information additional to that used to derive the PCS.

615

616 In addition, several of the linking nodes introduced into this islet network through their PPI
617 connections represent interesting candidates for a role in T2D pathogenesis. Cyclin-
618 dependent kinase 2 (*CDK2*), for example, has been shown to influence beta-cell mass in a
619 compensatory mechanism related to age and diet-induced stress, connecting beta-cell
620 dysfunction and progressive beta-cell mass deterioration [50]; *YHWAG* is a member of the
621 14-3-3 family, known to be signalling hubs for beta-cell survival [51]; and disruption of
622 *SMAD4* drives islet hypertrophy [52].

623

624 **Discussion**

625 In this study, we set out to overcome two challenges that have impeded efforts to
626 synthesise the biological information that is captured in the growing number of association
627 signals emerging from GWAS studies. In the case of type 2 diabetes, for example, there are
628 now well over a hundred independent common variant signals [6, 21], but most of these

629 map to regulatory sequence, and the molecular mechanisms whereby these, individually
630 and/or collectively, contribute to differences in T2D predisposition remain largely
631 unresolved. A key question, of direct relevance to the opportunities for translational use of
632 this information, is the extent to which, as the number of loci expands, there will be
633 “saturation” or “convergence” of the biological mechanisms through which they operate, or
634 whether, on the contrary, the range of networks and pathways implicated will continue to
635 proliferate.

636

637 The first challenge concerns the identification of the effector transcripts through which the
638 T2D predisposition effects at each of the GWAS signals (most obviously those that are
639 regulatory) are mediated. We approached this challenge by integrating, for each of the
640 genes within each of the GWAS signals, two types of data, one based around the fine-
641 mapping of the causal variant, and the use of cis-eQTL data (in the case of regulatory
642 variants) or direct coding variant inference to highlight the most likely effectors, the other
643 making use of diverse sources of biological information concerning the candidate effector
644 genes and their protein products. Using this framework, we were able to assign candidacy
645 scores to each regional gene, and then to deploy these scores as summaries of diverse
646 sources of data that could be propagated into subsequent network analyses. We recognise
647 that, given the sparse nature of the data used, not all such candidacy assignments will be
648 accurate. However, these scores provide a principled and objective way of synthesising
649 current knowledge, and the framework allows for iterative improvements in candidacy
650 assignments as additional sources of relevant data become available. These are likely for
651 example, to include further refinements in fine-mapping, additional links from associated
652 variants to their effectors arising from chromatin conformation analyses, detection of rare

653 coding variant signals through exome sequencing, and genome-wide screens of transcript
654 function.

655

656 The second challenge relates to the objective evaluation of the extent to which the
657 strongest positional candidates at these GWAS loci occupy overlapping biological space.

658 Standard approaches to network analysis applied to GWAS data – such as gene-set
659 enrichment [32], or co-expression analyses [53] – were not an option for this study since
660 source data relevant to these had already been factored into the assessments of positional
661 candidacy. Instead, we focused on the relationships between positional candidates as
662 revealed by protein-protein interaction data, which we considered to be independent of the
663 data in the earlier stages. We used the enrichment of T2D association signals in linking
664 nodes (i.e. proteins included in the network which did not map to known GWAS loci) as our
665 principal metric of network convergence.

666

667 This strategy uncovered a highly-interconnected network associated with T2D, which was
668 built around proteins involved in processes such as autophagy, lipid transport, cell growth,
669 and insulin receptor signalling pathways. We were able to show that this signal of
670 enrichment was enhanced when we constrained the generic PPI network to reflect only
671 genes expressed in pancreatic islets, and, concomitantly, limited the set of GWAS loci to
672 those at which the T2D predisposition was mediated by defective islet function. These
673 analyses reinforce the importance of the pancreatic islet as a critical tissue for the
674 development of T2D, and highlight multiple proteins (both those that map within GWAS
675 loci, and those that fall outside) that are represented within this core islet network. These
676 findings provide compelling hypotheses that can be explored further through direct

677 experimental study, and also highlight the need to generate tissue-specific protein-protein
678 interaction data. They also provide evidence to support a convergence of the mechanisms
679 mediating predisposition across diverse T2D association signals.

680

681 Finally, these analyses demonstrate a valuable approach for the interrogation of large-scale
682 GWAS data to capture biologically-plausible disease-specific processes, one which can
683 readily be applied to other complex diseases.

684

685 **Acknowledgments**

686 We would like to acknowledge Dr. Heiko Horn and Dr Loukas Moutsianas for their valuable
687 comments and useful discussion about methods during the elaboration of this manuscript.

688 **References**

689

- 690 1. Federation ID. IDF Diabetes Atlas, 8th edn. International Diabetes Federation; 2017
691 2017.
- 692 2. McCarthy MI. Genomics, type 2 diabetes, and obesity. *N Engl J Med.*
693 2010;363(24):2339-50.
- 694 3. Agarwala V, Flannick J, Sunyaev S, Go TDC, Altshuler D. Evaluating empirical
695 bounds on complex disease genetic architecture. *Nat Genet.* 2013;45(12):1418-27.
- 696 4. Fuchsberger C, Flannick J, Teslovich TM, Mahajan A, Agarwala V, Gaulton KJ, et al.
697 The genetic architecture of type 2 diabetes. *Nature.* 2016;536(7614):41-7.
- 698 5. Mahajan A, Wessel J, Willems SM, Zhao W, Robertson NR, Chu AY, et al. Refining
699 the accuracy of validated target identification through coding variant fine-mapping in type 2
700 diabetes. *Nat Genet.* 2018;50(4):559-71.
- 701 6. Scott RA, Scott LJ, Magi R, Marullo L, Gaulton KJ, Kaakinen M, et al. An Expanded
702 Genome-Wide Association Study of Type 2 Diabetes in Europeans. *Diabetes.*
703 2017;66(11):2888-902.
- 704 7. Boyle EA, Li YI, Pritchard JK. An Expanded View of Complex Traits: From
705 Polygenic to Omnigenic. *Cell.* 2017;169(7):1177-86.
- 706 8. Gaulton KJ, Ferreira T, Lee Y, Raimondo A, Magi R, Reschen ME, et al. Genetic fine
707 mapping and genomic annotation defines causal mechanisms at type 2 diabetes susceptibility
708 loci. *Nat Genet.* 2015;47(12):1415-25.

- 709 9. Varshney A, Scott LJ, Welch RP, Erdos MR, Chines PS, Narisu N, et al. Genetic
710 regulatory signatures underlying islet gene expression and type 2 diabetes. *Proc Natl Acad*
711 *Sci U S A*. 2017;114(9):2301-6.
- 712 10. Thurner M, van de Bunt M, Torres JM, Mahajan A, Nylander V, Bennett AJ, et al.
713 Integration of human pancreatic islet genomic data refines regulatory mechanisms at Type 2
714 Diabetes susceptibility loci. *Elife*. 2018;7.
- 715 11. van de Bunt M, Manning Fox JE, Dai X, Barrett A, Grey C, Li L, et al. Transcript
716 Expression Data from Human Islets Links Regulatory Signals from Genome-Wide
717 Association Studies for Type 2 Diabetes and Glycemic Traits to Their Downstream Effectors.
718 *PLoS Genet*. 2015;11(12):e1005694.
- 719 12. Hughes JR, Roberts N, McGowan S, Hay D, Giannoulatou E, Lynch M, et al.
720 Analysis of hundreds of cis-regulatory landscapes at high resolution in a single, high-
721 throughput experiment. *Nat Genet*. 2014;46(2):205-12.
- 722 13. Thomsen SK, Ceroni A, van de Bunt M, Burrows C, Barrett A, Scharfmann R, et al.
723 Systematic Functional Characterization of Candidate Causal Genes for Type 2 Diabetes Risk
724 Variants. *Diabetes*. 2016;65(12):3805-11.
- 725 14. Thomsen SK, Gloyn AL. The pancreatic beta cell: recent insights from human
726 genetics. *Trends Endocrinol Metab*. 2014;25(8):425-34.
- 727 15. Kohler S, Bauer S, Horn D, Robinson PN. Walking the interactome for prioritization
728 of candidate disease genes. *Am J Hum Genet*. 2008;82(4):949-58.
- 729 16. Lee SA, Tsao TT, Yang KC, Lin H, Kuo YL, Hsu CH, et al. Construction and
730 analysis of the protein-protein interaction networks for schizophrenia, bipolar disorder, and
731 major depression. *BMC Bioinformatics*. 2011;12 Suppl 13:S20.
- 732 17. Hou L, Chen M, Zhang CK, Cho J, Zhao H. Guilt by rewiring: gene prioritization
733 through network rewiring in genome wide association studies. *Hum Mol Genet*.
734 2014;23(10):2780-90.
- 735 18. Lundby A, Rossin EJ, Steffensen AB, Acha MR, Newton-Cheh C, Pfeufer A, et al.
736 Annotation of loci from genome-wide association studies using tissue-specific quantitative
737 interaction proteomics. *Nat Methods*. 2014;11(8):868-74.
- 738 19. Lage K, Karlberg EO, Storling ZM, Olason PI, Pedersen AG, Rigina O, et al. A
739 human phenome-interactome network of protein complexes implicated in genetic disorders.
740 *Nat Biotechnol*. 2007;25(3):309-16.
- 741 20. Rossin EJ, Lage K, Raychaudhuri S, Xavier RJ, Tatar D, Benita Y, et al. Proteins
742 encoded in genomic regions associated with immune-mediated disease physically interact
743 and suggest underlying biology. *PLoS Genet*. 2011;7(1):e1001273.
- 744 21. Morris AP, Voight BF, Teslovich TM, Ferreira T, Segre AV, Steinthorsdottir V, et al.
745 Large-scale association analysis provides insights into the genetic architecture and
746 pathophysiology of type 2 diabetes. *Nat Genet*. 2012;44(9):981-90.
- 747 22. Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, et al.
748 GENCODE: the reference human genome annotation for The ENCODE Project. *Genome*
749 *Res*. 2012;22(9):1760-74.
- 750 23. Consortium EP. An integrated encyclopedia of DNA elements in the human genome.
751 *Nature*. 2012;489(7414):57-74.
- 752 24. Consortium G. Genetic effects on gene expression across human tissues. *Nature*.
753 2017;550:204.
- 754 25. Mahajan A, Go MJ, Zhang W, Below JE, Gaulton KJ, Ferreira T, et al. Genome-wide
755 trans-ancestry meta-analysis provides insight into the genetic architecture of type 2 diabetes
756 susceptibility. *Nat Genet*. 2014;46(3):234-44.

- 757 26. Steinhorsdottir V, Thorleifsson G, Sulem P, Helgason H, Grarup N, Sigurdsson A, et
758 al. Identification of low-frequency and rare sequence variants associated with elevated or
759 reduced risk of type 2 diabetes. *Nat Genet.* 2014;46(3):294-8.
- 760 27. Wild F. *lsa: Latent Semantic Analysis*. R Package Version 0.57 ed2005.
- 761 28. Dittrich MT, Klau GW, Rosenwald A, Dandekar T, Muller T. Identifying functional
762 modules in protein-protein interaction networks: an integrated exact approach.
763 *Bioinformatics.* 2008;24(13):i223-31.
- 764 29. Tuncbag N, McCallum S, Huang SS, Fraenkel E. SteinerNet: a web server for
765 integrating 'omic' data to discover hidden components of response pathways. *Nucleic Acids*
766 *Res.* 2012;40(Web Server issue):W505-9.
- 767 30. Balbin OA, Prensner JR, Sahu A, Yocum A, Shankar S, Malik R, et al.
768 Reconstructing targetable pathways in lung cancer by integrating diverse omics data. *Nat*
769 *Commun.* 2013;4:2617.
- 770 31. Garey MR, Johnson DS. *Computers and intractability : a guide to the theory of NP-*
771 *completeness.* New York: W.H. Freeman; 1979. x, 340 p. p.
- 772 32. Lamparter D, Marbach D, Rueedi R, Kutalik Z, Bergmann S. Fast and Rigorous
773 Computation of Gene and Pathway Scores from SNP-Based Summary Statistics. *PLoS*
774 *Comput Biol.* 2016;12(1):e1004714.
- 775 33. Clauset A, Newman ME, Moore C. Finding community structure in very large
776 networks. *Phys Rev E Stat Nonlin Soft Matter Phys.* 2004;70(6 Pt 2):066111.
- 777 34. Bindea G, Mlecnik B, Hackl H, Charoentong P, Tosolini M, Kirilovsky A, et al.
778 ClueGO: a Cytoscape plug-in to decipher functionally grouped gene ontology and pathway
779 annotation networks. *Bioinformatics.* 2009;25(8):1091-3.
- 780 35. Croft D, Mundo AF, Haw R, Milacic M, Weiser J, Wu G, et al. The Reactome
781 pathway knowledgebase. *Nucleic Acids Res.* 2014;42(Database issue):D472-7.
- 782 36. Shin SY, Fauman EB, Petersen AK, Krumsiek J, Santos R, Huang J, et al. An atlas of
783 genetic influences on human blood metabolites. *Nat Genet.* 2014;46(6):543-50.
- 784 37. Pinero J, Bravo A, Queralt-Rosinach N, Gutierrez-Sacristan A, Deu-Pons J, Centeno
785 E, et al. DisGeNET: a comprehensive platform integrating information on human disease-
786 associated genes and variants. *Nucleic Acids Res.* 2017;45(D1):D833-D9.
- 787 38. Rubio-Cabezas O, Minton JA, Kantor I, Williams D, Ellard S, Hattersley AT.
788 Homozygous mutations in NEUROD1 are responsible for a novel syndrome of permanent
789 neonatal diabetes and neurological abnormalities. *Diabetes.* 2010;59(9):2326-31.
- 790 39. Yuan W, Xia Y, Bell CG, Yet I, Ferreira T, Ward KJ, et al. An integrated epigenomic
791 analysis for type 2 diabetes susceptibility loci in monozygotic twins. *Nat Commun.*
792 2014;5:5719.
- 793 40. Kimple ME, Moss JB, Brar HK, Rosa TC, Truchan NA, Pasker RL, et al. Deletion of
794 GalphaZ protein protects against diet-induced glucose intolerance via expansion of beta-cell
795 mass. *J Biol Chem.* 2012;287(24):20344-55.
- 796 41. Weinstein MM, Goulbourne CN, Davies BS, Tu Y, Barnes RH, 2nd, Watkins SM, et
797 al. Reciprocal metabolic perturbations in the adipose tissue and liver of GPIHBP1-deficient
798 mice. *Arterioscler Thromb Vasc Biol.* 2012;32(2):230-5.
- 799 42. Bergman BC, Cornier MA, Horton TJ, Bessesen DH. Effects of fasting on insulin
800 action and glucose kinetics in lean and obese men and women. *Am J Physiol Endocrinol*
801 *Metab.* 2007;293(4):E1103-11.
- 802 43. Bano G. Glucose homeostasis, obesity and diabetes. *Best Pract Res Clin Obstet*
803 *Gynaecol.* 2013;27(5):715-26.
- 804 44. Arnold AC, Robertson D. Defective Wnt Signaling: A Potential Contributor to
805 Cardiometabolic Disease? *Diabetes.* 2015;64(10):3342-4.

- 806 45. Wang CH, Wang CC, Wei YH. Mitochondrial dysfunction in insulin insensitivity:
807 implication of mitochondrial role in type 2 diabetes. *Ann N Y Acad Sci.* 2010;1201:157-65.
808 46. Antoun G, McMurray F, Thrush AB, Patten DA, Peixoto AC, Slack RS, et al.
809 Impaired mitochondrial oxidative phosphorylation and supercomplex assembly in rectus
810 abdominis muscle of diabetic obese individuals. *Diabetologia.* 2015;58(12):2861-6.
811 47. Voight BF, Scott LJ, Steinthorsdottir V, Morris AP, Dina C, Welch RP, et al. Twelve
812 type 2 diabetes susceptibility loci identified through large-scale association analysis. *Nat*
813 *Genet.* 2010;42(7):579-89.
814 48. Dimas AS, Lagou V, Barker A, Knowles JW, Magi R, Hivert MF, et al. Impact of
815 type 2 diabetes susceptibility variants on quantitative glycemic traits reveals mechanistic
816 heterogeneity. *Diabetes.* 2014;63(6):2158-71.
817 49. Wood AR, Jonsson A, Jackson AU, Wang N, van Leewen N, Palmer ND, et al. A
818 Genome-Wide Association Study of IVGTT-Based Measures of First-Phase Insulin Secretion
819 Refines the Underlying Physiology of Type 2 Diabetes Variants. *Diabetes.* 2017;66(8):2296-
820 309.
821 50. Kim SY, Lee JH, Merrins MJ, Gavrilova O, Bisteau X, Kaldis P, et al. Loss of
822 Cyclin-dependent Kinase 2 in the Pancreas Links Primary beta-Cell Dysfunction to
823 Progressive Depletion of beta-Cell Mass and Diabetes. *J Biol Chem.* 2017;292(9):3841-53.
824 51. Lim GE, Piske M, Johnson JD. 14-3-3 proteins are essential signalling hubs for beta
825 cell survival. *Diabetologia.* 2013;56(4):825-37.
826 52. Simeone DM, Zhang L, Treutelaar MK, Zhang L, Graziano K, Logsdon CD, et al.
827 Islet hypertrophy following pancreatic disruption of Smad4 signaling. *Am J Physiol*
828 *Endocrinol Metab.* 2006;291(6):E1305-16.
829 53. Calabrese GM, Mesner LD, Stains JP, Tommasini SM, Horowitz MC, Rosen CJ, et
830 al. Integrating GWAS and Co-expression Network Data Identifies Bone Mineral Density
831 Genes SPTBN1 and MARK3 and an Osteoblast Functional Module. *Cell Syst.* 2017;4(1):46-
832 59 e4.
833

834 **Supporting information**

835

836 **S1 Fig. Correlation between β values and PPI network size**

837 We tested different values of β to characterise the impact on network size and the
838 percentage of seed genes represented in the network. The figure depicts the relationship
839 between β values and both the number of nodes and the percentage of seed genes in the
840 optimal solution generated with the Steiner tree approach. We selected as optimal a β value
841 of 10 which produces a network of ~150 nodes which contains at least 25 % of the seed
842 genes.

843

844 **S2 Fig. Specificity of linking nodes in the final network.**

845 We assessed the specificity of each node in the final network solution by running the
846 algorithm 100 times with the same parameter settings, but with random input data. We
847 define specificity in this context as the complement of the percentage with which a given
848 linking node from the final network appears in runs generated from random input data. 80%
849 of linking nodes have a specificity exceeding 0.75, indicating that these linking nodes do not
850 simply reflect generic hubs in PPI space.

851

852 **S3 Fig. Distribution of PCS and correlation of semantic and risk variant link scores. a)**

853 Distribution of PCS values for the 1,895 candidate genes (left histogram) and for the 451
854 prioritised candidate genes (those that, for each locus contribute collectively to at least 70%
855 of the total PCS) (right histogram); b) the number of genes per locus stratified in terms of
856 PCS ranges for the 1,895 candidate genes (left boxplot), and the 451 prioritised candidate
857 genes (right boxplot); c) the correlation between the semantic and risk variant link scores
858 for the 1,895 positional candidates.

859

860 **S4 Fig. Summary of characteristics of PCS values.**

861 a) Distribution of gene number per locus for the 1,895 candidate genes (top histogram); and
862 for the 451 prioritised candidate genes (bottom histogram); b) example of the distribution
863 of PCS scores for the *TP53INP1* locus under “nearest-gene” selection (top figure), “all gene”
864 selection (median figure), and under our prioritisation method (bottom figure); c) Scatter
865 plot displaying the correlation between maximum PCS values for each locus under “nearest
866 gene” our prioritisation approach; d) distribution of the maximum PCS per locus with our
867 prioritisation strategy.

868

869 **S5 Fig. Enrichment of GWAS signals in the final PPI network.**

870 Using the generic PPI network generated from optimisation of seed node connectivity, we
871 performed GWAS signal enrichment analyses (linking nodes only) for 33 GWAS datasets
872 including T2D. Each point represents the $-\log_{10}$ p-value of the enrichment signal for the
873 specified GWAS dataset. The strongest signal of enrichment was observed between the
874 generic network and T2D ($p=0.0014$), with other significant associations for related
875 phenotypes such as BMI. The dotted red line represents nominal significance ($p=0.05$).

876

877 **S6 Fig. Correlations between tissue-specific PPI networks.**

878 The composition of 46 tissue-specific networks generated by filtering the generic PPI
879 network using tissue-specific RNA-Seq abundance data, was compared by applying a Jaccard
880 similarity index to network nodes. Many tissue-specific networks showed high similarity
881 with grouping by higher-level tissue of origin (e.g. brain, artery).

882

883 **S7 Fig. GWAS signal enrichment in the PPI-generic network derived from T2D GWAS
884 subsets.**

885 We built APCST networks using the generic PPI network for each of the subsets of T2D
886 GWAS loci defined by shared mechanistic mediation (refs[48], [49]). Enrichment in GWAS
887 signals for linking nodes only was tested using a permutation scheme. Each point in the
888 figure depicts the $-\log_{10}$ p-value of enrichment for association signals derived from a
889 particular GWAS analysis. The results for T2D GWAS enrichment for APCST networks built
890 around the different T2D GWAS subsets are also represented (large red dots). The dotted
891 red line represents nominal significance ($p=0.05$). The strongest enrichment for T2D GWAS

892 data in the generic PPI network is observed for subsets of loci acting through reduced
893 insulin secretion. In the cluster hairballs for the seven T2D GWAS locus subset categories,
894 nodes are coloured according to their PCS with grey nodes representing linking nodes.
895 Comparison with Figure 4 (the equivalent figure generated using the islet-filtered PPI
896 network) demonstrates increased levels of enrichment for the subset of T2D loci influencing
897 insulin secretion when filtering for nodes reflecting pancreatic islet expression.

898

899 **S1 Table. 101 loci and candidate genes by loci used to calculate the Positional Candidacy**
900 **Score (PCS).** Note: We developed a framework to score the candidacy of genes mapping to
901 GWAS association signals which aggregated data from multiple sources. The information
902 collected fell into two categories. First, we used regression-based approaches to link
903 disease-associated variants to their likely effector transcripts, using a combination of
904 variant-based annotations and expression QTL data (Link score). Second, we scored each of
905 the genes in these GWAS regions for disease-relevant biological function (Semantic score).
906 We combined the two measures to generate a “positional candidacy score” (PCS) for each
907 gene. Cumulative: cumulative frequency of PCS for each loci. References: Bibliographic
908 references describing the loci as associated to type II diabetes.

909

910 **S2 Table. DisGeNet results.** MeSH = Medical Subject Headings; DPI score = disease
911 pleiotropic index; DSI score = disease specific index; GDA Score = Gene-Disease Association
912 Score; EI = Evidence Score.

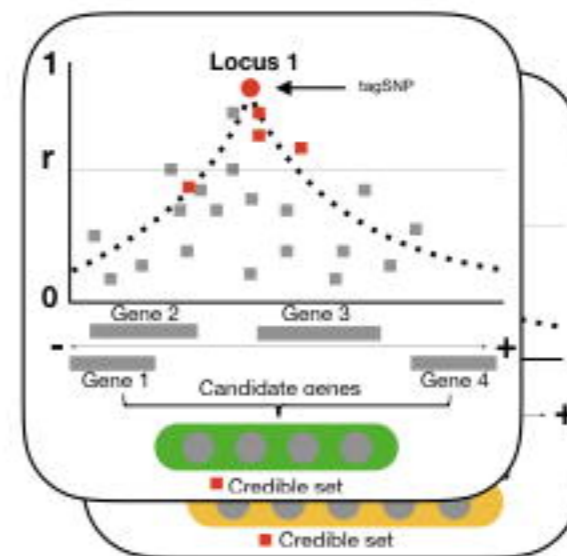
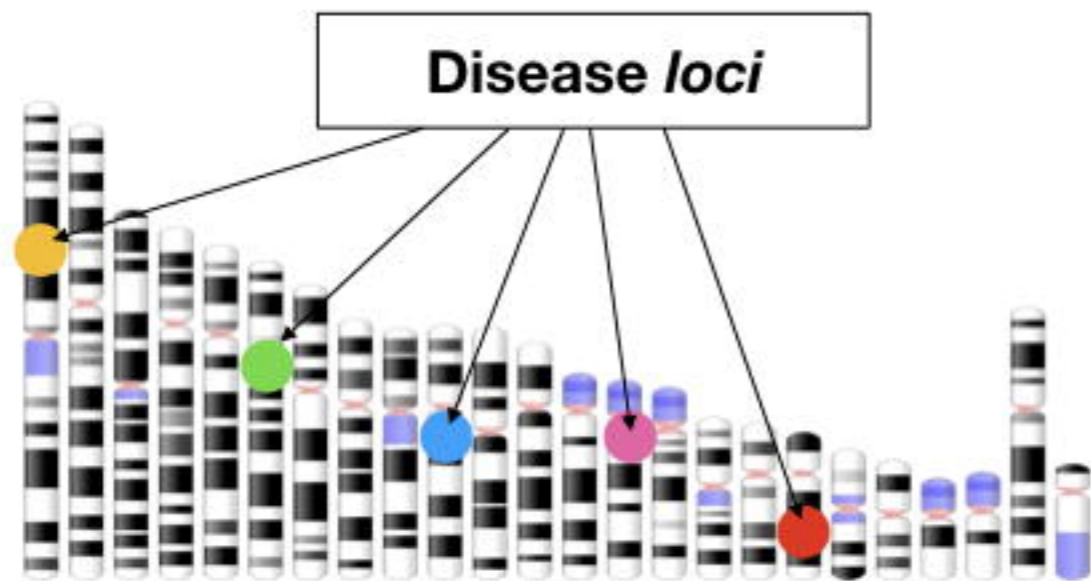
913

914 **S3 Table. Gene Set Enrichment Analysis by community.** GOID: Gene Ontology ID; GOTerm:
915 Gene Ontology Term; Note: Gene Set Enrichment (GSE) of networks and sub-networks was

916 performed with ClueGO using GO terms and REACTOME gene sets. The enrichment results
917 were considered significant when bonferroni adjusted p-value < 0.05 and at least 3% of the
918 genes contained in the tested gene set is included in the network. Gene sets were also
919 grouped using kappa score into functional groups to improve visualization of enriched
920 pathways.

921

922 **S4 Table. Gene Set Enrichment Analysis in Beta-cell Islet-specific network.** GOID: Gene
923 Ontology ID; GOTerm: Gene Ontology Term. Note: Gene Set Enrichment (GSE) of networks
924 and sub-networkswas performed with ClueGO using GO terms and REACTOME gene sets.
925 The enrichment results were considered significant when bonferroni adjusted p-value < 0.05
926 and at least 3% of the genes contained in the tested gene set is included in the network.
927 Gene sets were also grouped using kappa score into functional groups to improve
928 visualization of enriched pathways.

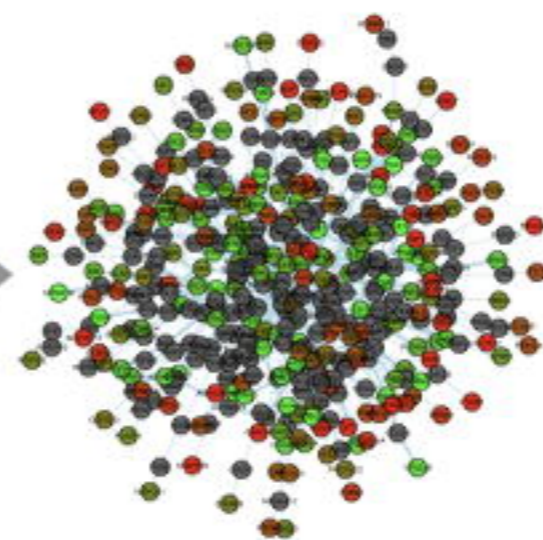
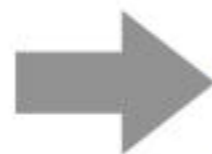
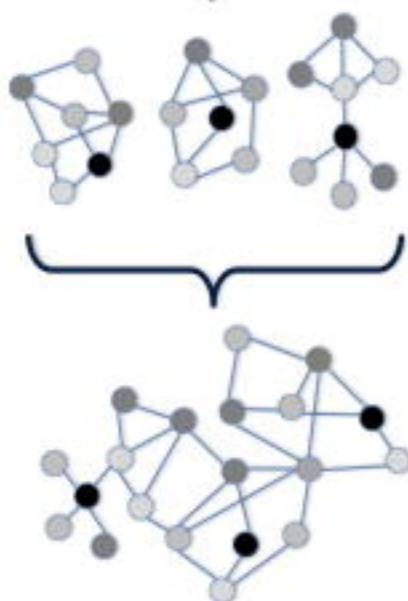


Locus	Gene	Score	Nearest	Equal
CDC123	<i>CAMK1D</i>	0.309	0	0.14
CDC123	<i>CDC123</i>	0.235	1	0.14
CDC123	<i>PROSER2</i>	0.167	0	0.14
CDC123	<i>UPF2</i>	0.138	0	0.14
CDC123	<i>SEC61A2</i>	0.101	0	0.14
CDC123	<i>NUDT5</i>	0.035	0	0.14
CDC123	<i>DHTKD1</i>	0.013	0	0.14

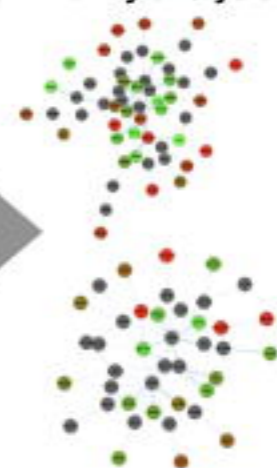
Candidacy scores



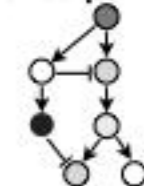
Steiner Tree approach



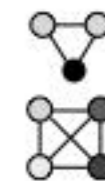
community analysis



enriched pathways



cofunction modules

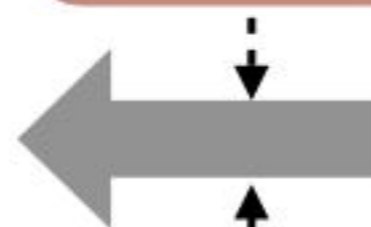


Variant level

Information from eQTL
Enhancer links, and coding
variants info

Variant level

ENCODE regulatory annotations
T2D relevant tissues: pancreatic
islets, adipose, ...



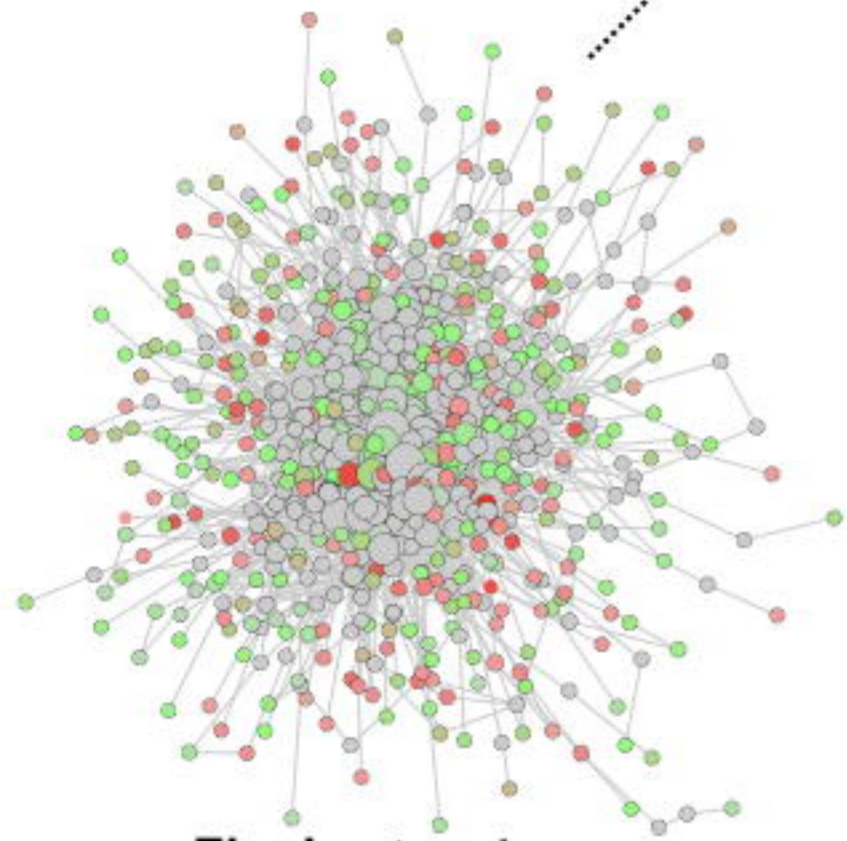
Gene level

Biological annotations from
MGD, KEGG, and GO

**Fine-mapped
credible sets**

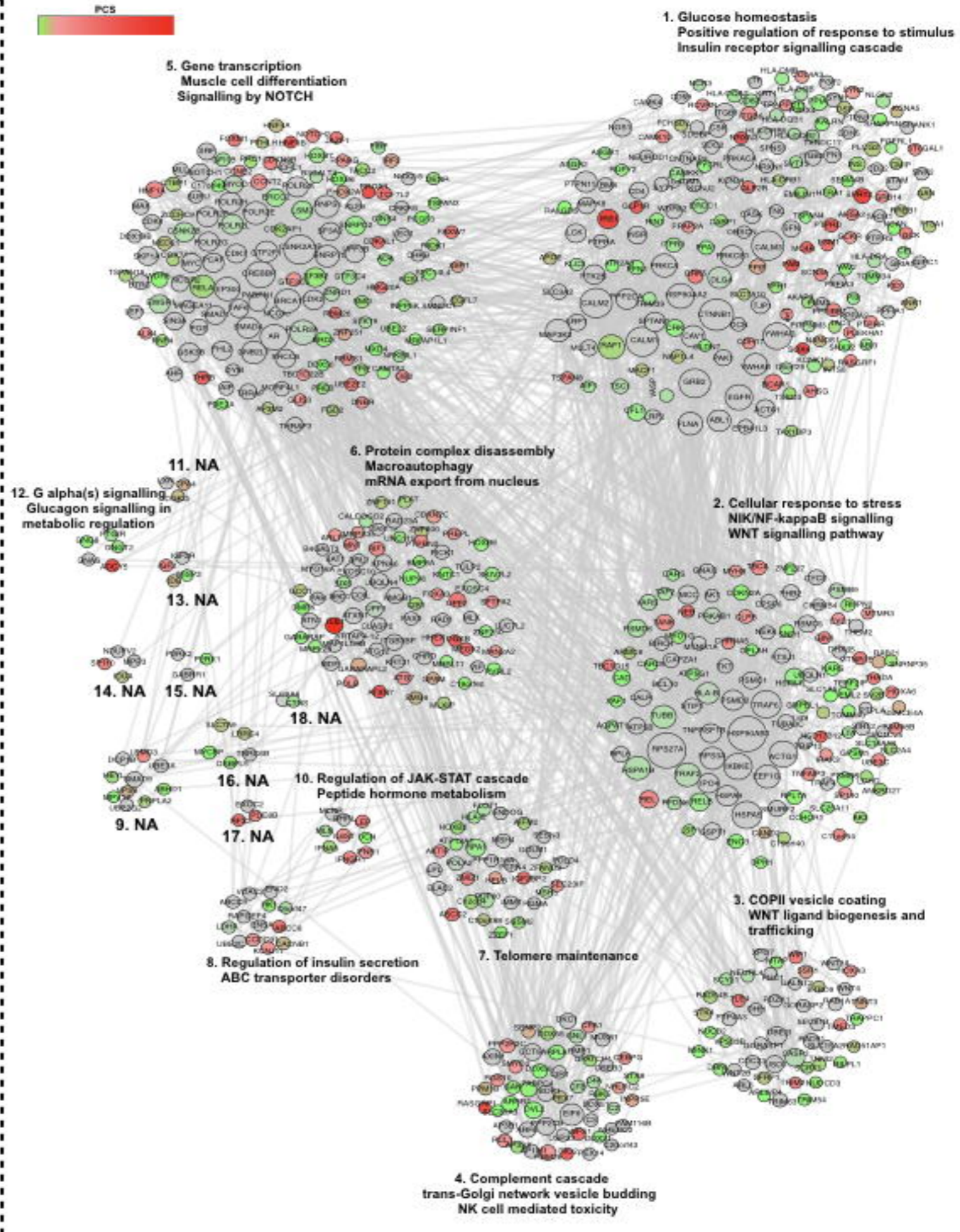
n Loci



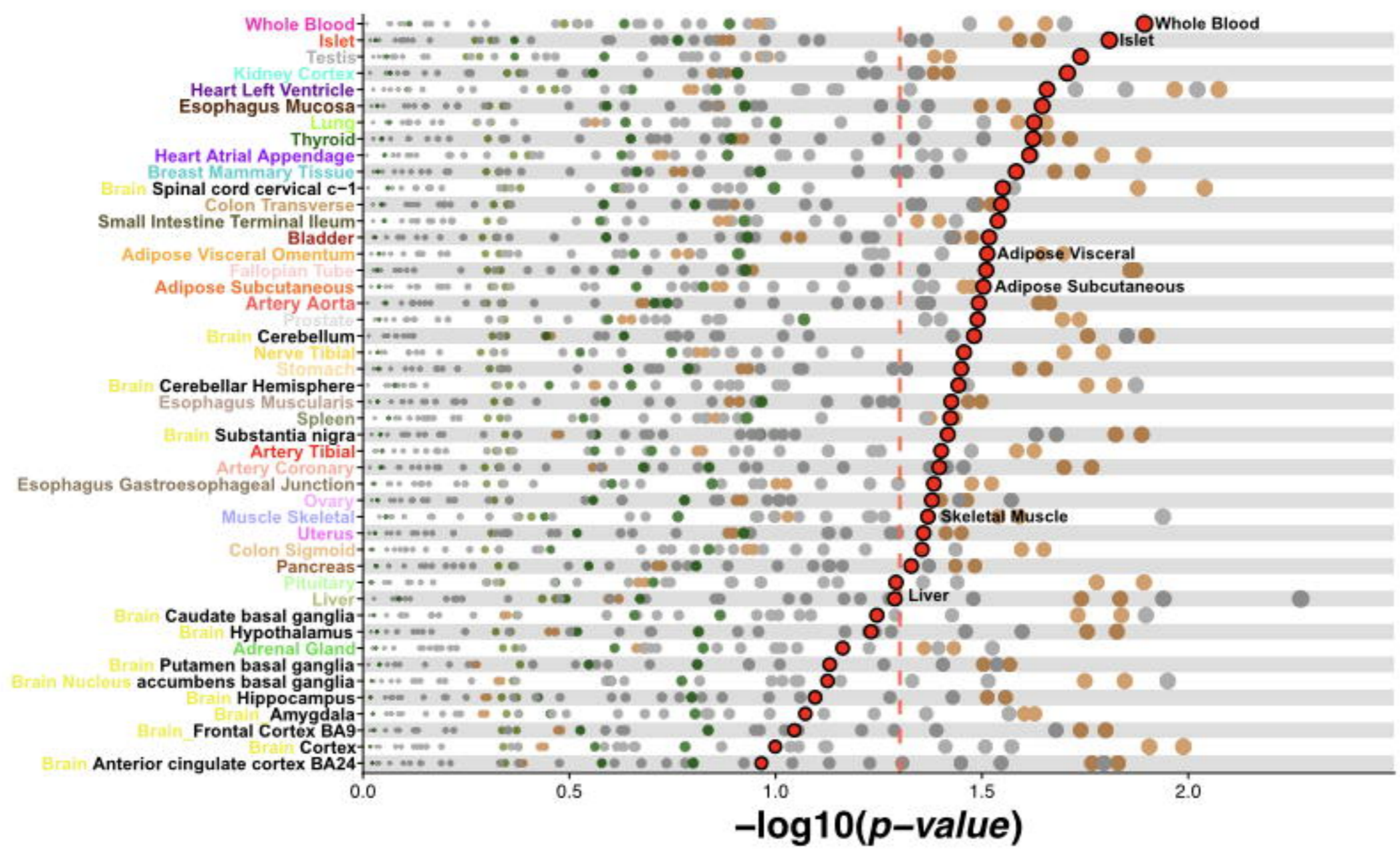


Final network
(705 nodes, 2678 edges)

Communities were discovered using the algorithm of Clauset et al 2014)

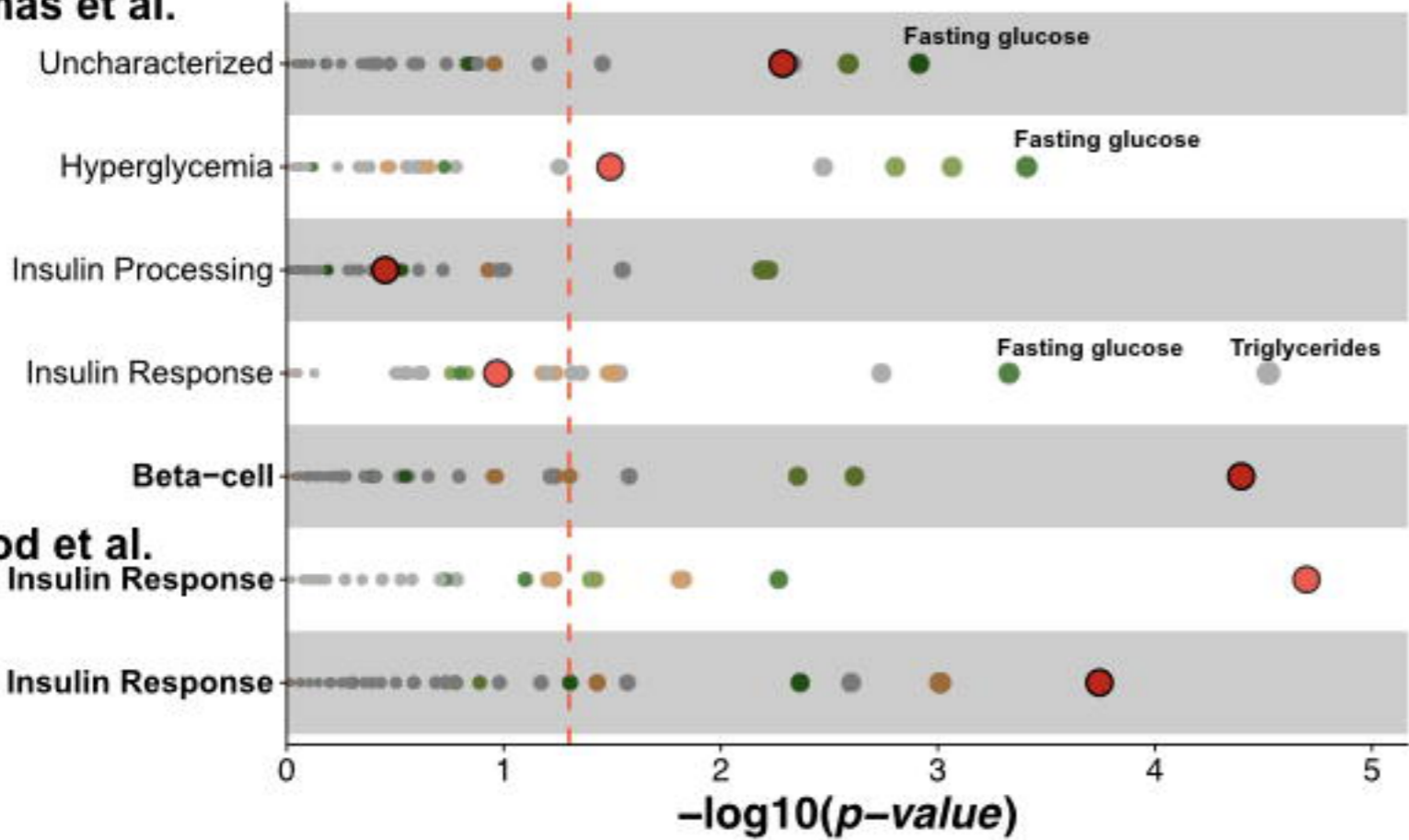


PPI intersected with GTEx tissues

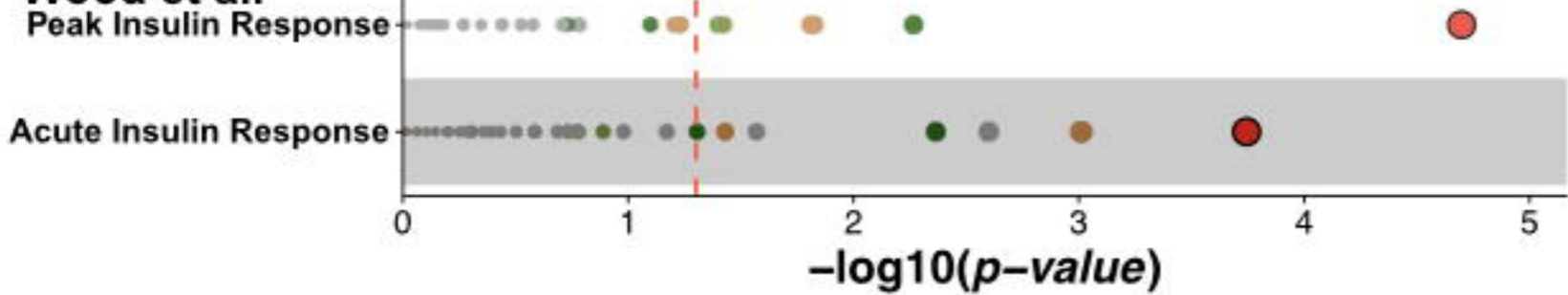


- GWAS
- T2D (DIAGRAM)
 - Glycemic traits (MAGIC)
 - Anthropometric traits (GIANT)
 - Birth weight (EGG)
 - OTHERS

Dimas et al.



Wood et al.



- GWAS
- DIAGRAM (T2D)
 - MAGIC (Glycemic traits)
 - GIANT (Anthropometric traits)
 - EGG (Birth weight)
 - OTHERS

