

## Type 2 diabetes: genetic data sharing to advance complex disease research

Jason Flannick<sup>1,2</sup> and Jose C. Florez<sup>1–4</sup>

**Abstract** | As with other complex diseases, unbiased association studies followed by physiological and experimental characterization have for years formed a paradigm for identifying genes or processes of relevance to type 2 diabetes mellitus (T2D). Recent large-scale common and rare variant genome-wide association studies (GWAS) suggest that substantially larger association studies are needed to identify most T2D loci in the population. To hasten clinical translation of genetic discoveries, new paradigms are also required to aid specialized investigation of nascent hypotheses. We argue for an integrated T2D knowledgebase, designed for a worldwide community to access aggregated large-scale genetic data sets, as one paradigm to catalyse convergence of these efforts.

### Heritability

The proportion of phenotypic variance in a population owing to genetic differences, as opposed to environmental differences.

Genetic studies of human disease seek insight into disease processes or novel therapies from naturally occurring heritable variation in a population. Characterizing the degree to which disease-informative variants exist, identifying them and their effects, and understanding how they might guide preventive or therapeutic interventions are thus related but distinct goals of genetic association studies.

Type 2 diabetes mellitus (T2D) is a complex, heritable disease with a sibling relative risk of approximately 2 (REF. 1) and a heritability estimated at 30–70%<sup>2</sup>. For decades, analyses of natural genetic variation have been used to understand T2D mechanisms or to improve prognoses for the approximately 415 million people who suffer from its effects; these include prolonged hyperglycaemia from insulin resistance and relative insulin deficiency<sup>3</sup>, numerous micro- and macrovascular complications<sup>4</sup> and an increased risk of early death<sup>5</sup>.

Before common variant genome-wide association studies (GWAS), it was not known to what extent such ‘experiments of nature’ could lend insight into T2D. Early genetic mapping studies for T2D involved comparatively small-scale candidate gene and linkage studies. Although these approaches localized several disease genes (for example, *PPARG*<sup>6</sup>, *KCNJ11* (REF. 7) and *TCF7L2* (REF. 8)), they were on the whole unsuccessful<sup>9</sup>. By contrast, GWAS yielded — for the first time — a substantial number of genomic loci reproducibly linked to T2D risk, suggesting previously unsuspected biological pathways<sup>10,11</sup>. And yet, the modest fraction of heritability attributable to GWAS associations, together with insufficient resolution to identify specific causal variants or effector transcripts, led to suggestions that the endeavour produced minimal insights<sup>12,13</sup>.

In the era that has followed these early GWAS, studies have become much larger and more diverse, technologies assay rarer variants across the entire exome or genome, and biological experiments have begun to translate reported associations to insights into causal variants and disease processes. Here, we review the history and recent progress towards mapping genes for T2D and consequent implications for disease genetic architecture. We then review studies to understand pathogenic mechanisms that may point to new therapeutic modalities. We argue that future research into each area will be most powerful and synergistic with a new model for open and interactive sharing of data and results between previously loosely coupled communities.

### Mapping disease loci

**Common variant association studies.** Motivated by the increased power of association studies over prior family studies<sup>14</sup>, as well as the common disease common variant hypothesis (CDCV hypothesis)<sup>15</sup>, GWAS exploited inexpensive single-nucleotide polymorphism (SNP) microarrays and reference catalogues of population-level variation<sup>16</sup> to analyse through linkage disequilibrium nearly all common variants in the genome. The first T2D GWAS, which studied a few thousand samples each<sup>17–21</sup>, collectively identified 10 genomic loci, and many loci were reported by multiple studies. These findings validated the GWAS design but suggested that common variants of large effects were mostly absent from the population: all associated variants affected risk by less than 40%, and most affected risk by closer to 15%.

In order to identify common variants of weaker effects, collaboration and data sharing were necessary to increase power (FIG. 1). The Diabetes Genetics Replication

<sup>1</sup>Programs in Metabolism and Medical and Population Genetics, Broad Institute of MIT and Harvard, 415 Main Street, Cambridge, Massachusetts 02142, USA.

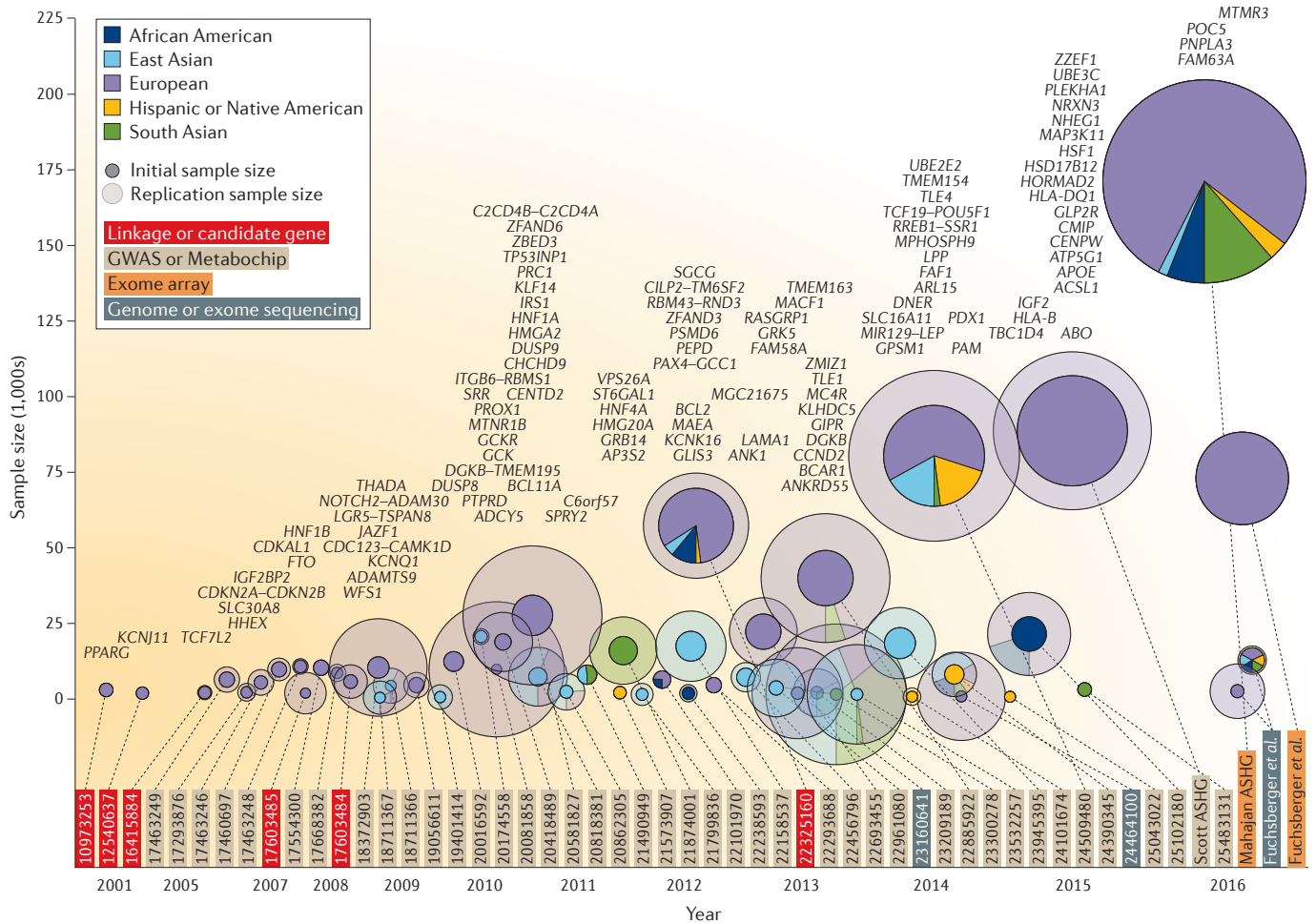
<sup>2</sup>Center for Human Genetic Research, Department of Medicine, Massachusetts General Hospital, 185 Cambridge Street, Boston, Massachusetts 02062, USA.

<sup>3</sup>Diabetes Research Center, Diabetes Unit, Department of Medicine, Massachusetts General Hospital, 185 Cambridge Street, Boston, Massachusetts 02062, USA.

<sup>4</sup>Department of Medicine, Harvard Medical School, 25 Shattuck Street, Boston, Massachusetts 02115, USA.

Correspondence to J.C.F. [JCFLOREZ@mgm.harvard.edu](mailto:JCFLOREZ@mgm.harvard.edu)

doi:10.1038/nrg.2016.56  
Published online 11 Jul 2016



**Figure 1 | The history of T2D GWAS.** Over the years, type 2 diabetes mellitus (T2D) genome-wide association studies (GWAS), which typically consist of a two-stage discovery and replication study design, have increased in size and diversity. Plotted are circles representing T2D GWAS, as specified by the US National Human Genome Research Institute–European Bioinformatics Institute (NHGRI–EBI) GWAS catalogue<sup>163</sup>, as well as additional candidate gene or sequencing studies of note. The x-axis shows the year of publication, whereas the y-axis shows discovery sample size. The inner (darker) circles are also scaled in proportion to discovery sample size, whereas the outer (lighter) circles are scaled in proportion to total (discovery + replication) sample size. Circles are coloured according to ethnic composition of the sample set: African American (dark blue), East Asian (light blue), European (purple), Hispanic or Native American (yellow) or South Asian (green). PubMed identifiers (or first author name) for each study are shown at the base of the figure and connected to the corresponding circle with a dotted line. Identifiers are coloured according to the technology used in the study: linkage or candidate gene studies (red), GWAS or Metabochip (beige), exome array (orange) or sequencing (dark grey). Additionally, T2D loci are listed directly above the first reporting study. With some exceptions based on retrospective knowledge, a  $P$  value of  $5 \times 10^{-8}$  is used as a threshold for significance, and loci are named as originally reported.

**Genome-wide association studies (GWAS).** An approach for genetic mapping that compares frequencies of variants across the genome between disease cases and matched controls. This is a paradigm for identifying genes or biological processes that are relevant to a phenotype by identifying correlations between polymorphic genetic markers and the presence of the phenotype.

and Meta-analysis (DIAGRAM) Consortium increased sample size above 10,000 (with approximately 4,500 T2D cases), identifying a further six loci<sup>22</sup>. Enabled by a willingness of researchers to collaborate rather than to compete, an availability of resources and tools to harmonize data sets<sup>23</sup> and a means to share statistics while protecting study participants<sup>24</sup>, GWAS meta-analysis consortia quickly became the standard for T2D and other complex traits<sup>25,26</sup>. The second DIAGRAM analysis included 45,000 samples (approximately 8,100 cases), identifying 12 additional T2D loci<sup>27</sup>, whereas other consortia analysed quantitative glycaemic traits<sup>28</sup> and other T2D-relevant phenotypes<sup>29–31</sup>

in either overlapping samples or more diverse ethnic groups<sup>32–34</sup>. The similar goals, organization and methodology of these consortia, as well as of those investigating other metabolic traits, led to the joint design of a custom genotyping array<sup>35</sup> that facilitated analysis of still larger sample sizes. A T2D association study used this array to analyse almost 150,000 individuals (approximately 34,800 cases) and identified ten further loci<sup>36</sup>, whereas a glycaemic trait association study with over 133,000 individuals led to a similar number of discoveries<sup>37</sup>. GWAS thus catalysed a clear breakthrough in the identification of genomic loci reproducibly associated with disease.

**Causal variants**

Specific mutations underlying the molecular cascade that produces a phenotypic trait; by design, in most genetic mapping studies, the associated genetic marker is merely correlated with the underlying causal variant.

**Effector transcripts**

The specific RNA transcript (for example, mRNA transcribed from a gene) for which the function or expression is altered by the causal variant, leading to a phenotypic difference.

**Genetic architecture**

The number, frequencies and effects on disease of genetic variants in a population.

**Common disease common variant hypothesis**

(CDCV hypothesis). The hypothesis that, owing to historical human population growth, some disease loci for common diseases may harbour alleles common in the population.

**Linkage disequilibrium**

Correlations among nearby variants, owing to historical patterns of demography and recombination, exploited by genome-wide association studies to map common variant associations.

**Glycaemic**

Traits pertaining to the physiology of blood glucose regulation, usually involving measures of glucose, insulin or other related hormones.

**Rare variant models**

A model of genetic architecture in which rare variants (for example, those with a frequency < 1%) explain most of the heritability.

**Synthetic associations**

A hypothesis based on simulations that multiple causal rare variants of strong effects might cause a common variant statistical association.

**Common variant models**

Models of genetic architecture in which common variants (for example, those with a frequency > 1%) explain most of the heritability.

**Models for the genetic architecture of T2D.** Despite the successes described above, the use of GWAS findings was debated, as the identified associations explained less than 10% of T2D heritability<sup>38</sup>. Whereas GWAS proponents cited previously unsuspected disease processes implicated by newly identified associations<sup>39</sup>, counter-arguments stated that modest-effect common variants offered little value towards understanding or predicting disease<sup>13</sup> and might even be misleading if caused, in fact, by distant rare variants<sup>40</sup> or population artefacts<sup>12</sup>. Using qualitative arguments<sup>12,41</sup> or simulations<sup>40,42</sup>, many argued for future studies of rare variants, which were expected from evolutionary arguments to have much larger effects on disease risk<sup>12,41</sup> and potentially to explain much of disease heritability<sup>38</sup>.

However, closer analysis of GWAS results revealed that evidence for rare variant models was far from conclusive. Simulations suggesting that rare variants (or synthetic associations) might explain many GWAS signals<sup>40,43</sup> were countered by empirical examples inconsistent with that model. These included an absence of linkage signals<sup>44</sup>, a common-shifted distribution of GWAS variant frequencies<sup>45</sup> and replicated associations across populations<sup>46</sup>. Furthermore, new methods for estimating heritability from all variants analysed in a GWAS, rather than only those reaching significance, suggested that common variants might explain more than half of T2D heritability, leaving hundreds or thousands of smaller-effect associations to be detected<sup>36,47</sup>.

This collection of analyses, examining similar data, thus offered contrasting models for the genetic architecture of T2D. In fact, a comprehensive simulation-based analysis suggested that, before 2012, T2D genetic studies and epidemiological parameters were consistent with either rare variant or common variant models<sup>48</sup>. Adjudicating between them required higher resolution studies to interrogate variants across the entire frequency spectrum.

**Rare variant association studies.** The first study to investigate the contribution of lower-frequency (<5%) variants to T2D or related traits appeared in 2012, when coding variant analysis in 8,229 individuals (using the Illumina exome array) identified five low-frequency variants associated with glycaemic traits (three at previously unidentified loci)<sup>49</sup>. Shortly thereafter, an exome sequencing study in 2,000 Danes (with 1,000 cases) identified two T2D-associated common coding variants at previously unidentified loci, although no low-frequency variants were identified<sup>50</sup>. In 2014, a genome sequencing study of 2,630 Icelanders, followed by statistical imputation of variants into 278,554 additional Icelanders, identified three more T2D-associated low-frequency variants<sup>51</sup>, all within loci that had previously been reported for T2D or related traits. Together with several studies that followed for T2D<sup>52</sup>, glucose<sup>53,54</sup> or insulin levels<sup>55</sup>, these reports demonstrated the ability of exome- or genome-wide analyses to identify lower-frequency variants relevant to T2D.

In parallel, candidate gene experiments were also successful, in limited contexts, at identifying broader allelic series within genes initially implicated through

common variants in T2D. An analysis of rare variants in *MTNR1B* demonstrated a strong effect of functional variants (according to a series of molecular or cellular assays) on T2D risk; non-functional variants had comparatively weaker effects<sup>56</sup>. A similar pattern was observed for rare variants in *PPARG*<sup>57</sup>. In *SLC30A8*, rare variants predicted to cause protein truncation were associated with protection from T2D<sup>58</sup>, which contrasted with expectations from previous animal<sup>59</sup> or cellular work<sup>60</sup>. The contribution of these variants to T2D heritability, however, was small.

**Deciphering genetic architecture through larger-scale sequencing studies.**

Although these newly reported associations were valuable, the extent to which they supported rare or common variant genetic architectural models for T2D was left unanswered. Analysis of the 2,000 Danish exomes rejected extreme models, with T2D explained by a small number (<20) of large-effect, low-frequency, non-synonymous variants<sup>61</sup>, but the power was insufficient for further conclusions. To characterize T2D genetic architecture to a much higher resolution, a more comprehensive set of sequencing studies was designed, spanning 12,940 multi-ethnic exome sequences (6,504 cases) with low-frequency variants genotyped in 79,854 additional individuals (28,305 cases), as well as 2,657 whole-genome sequences (1,326 cases) with all variants imputed into 44,414 additional individuals (12,971 cases)<sup>62</sup>. Strikingly, only a single low-frequency variant, which had been previously reported in the Icelandic sequencing study, was associated with T2D at genome-wide significance, despite high power to detect low-frequency variants of even moderate effect. These results suggested that low-frequency variation — and coding variation of any frequency — has, at most, limited roles in the genetic architecture of T2D. Consequently, larger common-variant GWAS would probably continue to identify additional T2D loci, whereas rare variant association studies would require comparatively much larger sample sizes<sup>48,63,64</sup>.

Indeed, larger GWAS meta-analyses have continued to expand the number of T2D-associated loci. A 2014 *trans*-ethnic meta-analysis of more than 110,000 individuals (approximately 26,500 cases), which was motivated by a consistency of common variant associations observed across different populations<sup>46</sup>, identified seven new loci<sup>65</sup>, whereas analysis of almost 160,000 European individuals (approximately 26,700 cases) subsequently identified 18 more<sup>66</sup>. GWAS of previously unstudied populations, such as Mexicans<sup>67</sup>, African Americans<sup>68</sup> or the Inuit<sup>69</sup>, have also yielded new associations, in many cases because different population histories have led to different allelic spectra and consequent gains in power. Conversely, even in larger samples (that is, >232,000 individuals, approximately 56,600 cases), coding variant analyses have identified few associations with low-frequency variants<sup>70</sup>. These findings, together with new methods to assess polygenicity from GWAS<sup>71</sup> and excess concordance observed for even nonsignificant common variant associations across ancestries<sup>65</sup>, continue to support a common variant model for the larger fraction of T2D genetic architecture.

## Box 1 | Population risk versus biologically relevant variants and genes

Genetic studies use ‘experiments of nature’ — that is, naturally occurring variation — to characterize both the genetic basis of disease in the population (on which stratification for intervention might be predicated) as well as genes or biological pathways that might guide insight into disease processes. These two goals are in many senses complementary: identifying most variants responsible for disease provides a catalogue of disease-relevant genes, and variants of disproportionate impact in the population may point to pathways of disproportionate relevance to disease. However, insight into disease biology is not predicated on a complete description of the genetic basis of disease, nor does explanation of heritability in the population necessarily lead to improved biological understanding.

These distinct goals thus lie at the heart of one notable debate in recent years: the degree to which genome-wide association studies (GWAS) have been ‘successful’. One interpretation of the first GWAS findings for type 2 diabetes mellitus (T2D) and other complex diseases was a focus on ‘missing heritability’: that is, the inability of GWAS findings to explain most of the genetic risk of disease in the population or to provide sufficient prognostic or stratifying information<sup>13,38</sup>. The modest effects of common variants, responsible for the limited fraction of heritability explained, were claimed to provide limited insight into disease biology, and rarer variants were claimed to offer a more profitable explanation for the remaining genetic basis of disease<sup>13</sup>.

GWAS, however, were not designed to explain the entirety of disease heritability but rather to identify the subset of variants that by chance may have reached common frequencies in the population<sup>15</sup>. The identification of some — let alone many — such variants offered key insights into T2D biology, such as the confirmation that most risk variants act through reduced  $\beta$ -cell function or mass, rather than many other previously hypothesized pathways<sup>27</sup>, and the demonstration that most molecular mechanisms of disease risk are regulatory<sup>110,113,162</sup>. Furthermore, the effects of genetic variants are constrained by population history and natural selection, suggesting that stronger perturbations of the same or different genes in a common pathway might lead to larger phenotypic effects. In this sense, rare variants can be of great value, even if they in fact explain less heritability than common variants: for example, the finding of protective rare loss-of-function variants in *SLC30A8* (REF. 58), which explained less heritability than the common variant by several orders of magnitude<sup>36</sup>, nonetheless provided an important suggestion about the direction of disease risk from protein inactivation. Similarly, *KCNJ11-ABCC8* and *PPARG* exemplify viable drug targets that could, in principle, have initially been identified by modest-effect common-variant associations<sup>6,7</sup>.

Conversely, recent studies that have succeeded in explaining missing heritability have not necessarily led to increased biological insight. Polygenic score<sup>47</sup> and mixed linear modelling<sup>36</sup> analyses have shown that GWAS variants in fact ‘tag’ most T2D heritability, but these analyses by design do not identify specific variants from which biological insight might be derived. Furthermore, an explanation of disease heritability does not imply an accurate model for disease risk prediction<sup>97,98</sup>, as the environmental component of T2D risk introduces substantial and inherent noise. Similarly, characterizing the genetic architecture of T2D — or the number, frequencies and effect sizes of variants that contribute to heritability — is, in some senses, easier than identifying specific causal variants and can greatly affect genetic study designs or genetic risk prediction.

Thus, the two goals of explaining disease heritability and identifying biologically relevant genes require, in many cases, distinct approaches by distinct communities. Diverse analyses of shared genetic data sets are one future route to empower both types of study.

**Imputation**

A technique to infer the unknown genotype of a variant in an individual based on correlations with nearby genotyped variants.

**Allelic series**

A number of alleles of a gene or locus with a range of phenotypic and/or molecular effects that are of use to infer a genetic–phenotypic dose–response curve.

**Polygenicity**

An idealized model in which a phenotype is caused by a large number of variants, each with small and normally distributed phenotypic effects.

**Transcriptomic**

The study of the expression levels of all transcripts in a cell.

**Epigenomic**

The study of all epigenetic modifications of a cell, including DNA methylation and histone modifications, which are largely responsible for the genes expressed in a specific tissue at a given developmental stage or metabolic state.

**Homeostasis model assessments**

A method based on fasting measures of glucose and insulin levels that is used to estimate  $\beta$ -cell function or insulin resistance.

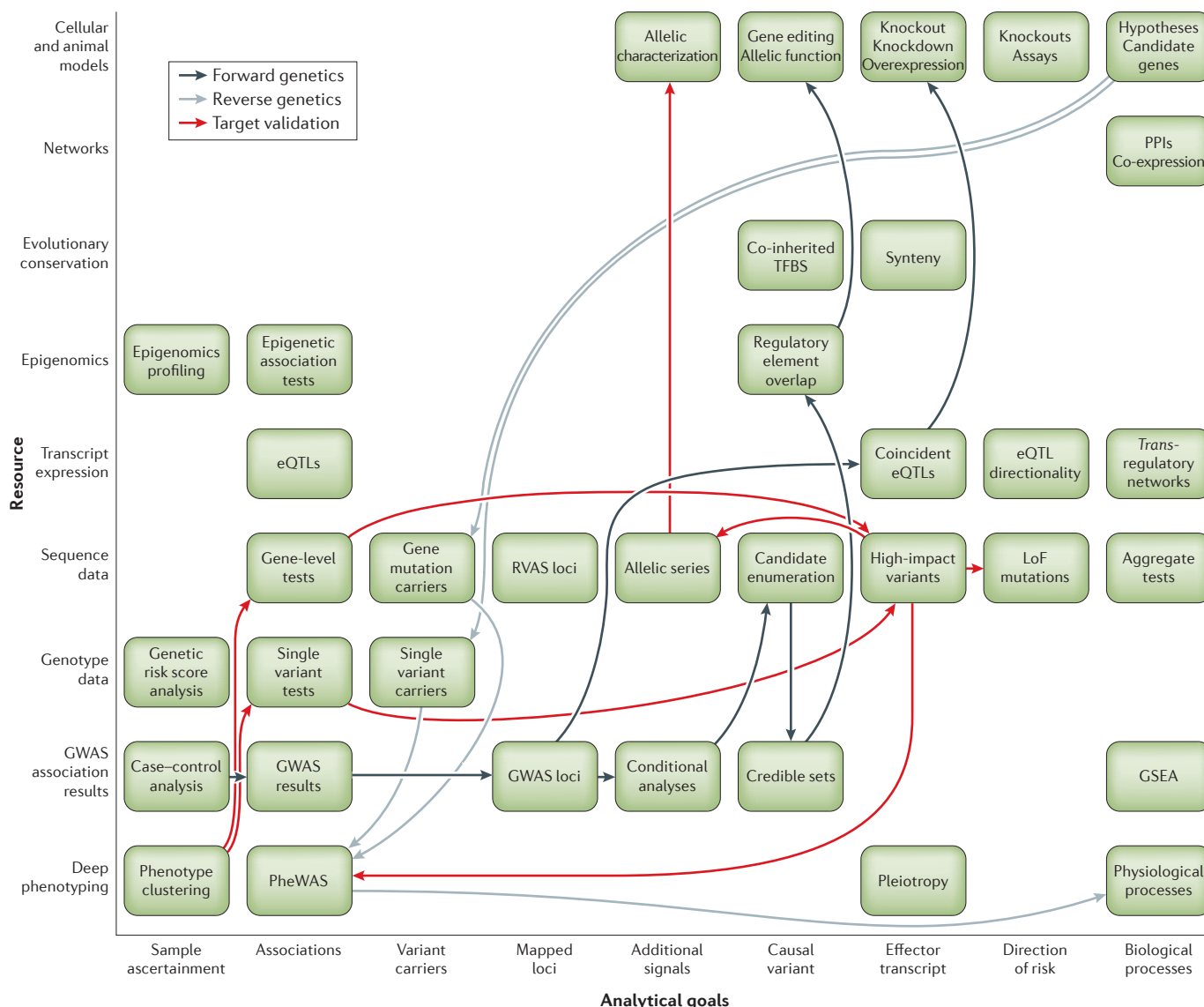
Thus, additional T2D loci will probably be discovered mostly by continued increases in the size and diversity of common-variant GWAS, incorporating new populations, additional phenotypes<sup>72,73</sup> and new resources for genotype imputation<sup>74,75</sup>. The collaborative networks and centralized analysis plans of GWAS consortia, refined over the past half-decade, will probably play a key part in the organization of these efforts.

**Understanding disease biology**

Characterizing the alleles that contribute to population disease risk is only one goal of genetic studies (BOX 1). Arguably more consequential is using the mapped associations to develop and test hypotheses about disease biology. Progress towards this latter goal has been variable, in part because experiments to understand the molecular, cellular and physiological mechanisms behind an association are typically highly domain-specific and thus are less systematized than the analysis protocols now deployed in GWAS. Moreover, the opportunity to use genetic findings to validate or disprove hypotheses from animal or cellular models is often impeded by the cumbersome access to results and opaque language of genetic studies.

Recent biological and clinical studies increasingly suggest one possible means to increase the translational use of human genetic findings: through convergence on common resources and workflows (FIG. 2; TABLE 1). Paradigms have begun to emerge for investigations of physiological mechanisms through deep patient phenotyping and stratification; for the analysis of potential causal variants and effector transcripts through the use of genetic, transcriptomic and epigenomic maps; and for the elucidation of molecular or cellular mechanisms through functional experiments, new techniques and genomic resources.

**Physiological investigations.** To gain physiological insight into T2D associations, many studies have investigated the many T2D-related phenotypes that are available in large numbers of genotyped samples. Analysing 31 T2D-associated variants for associations with homeostasis model assessments of  $\beta$ -cell function (HOMA-B) or insulin resistance (HOMA-IR) produced early evidence that most T2D variants increase risk through pancreatic  $\beta$ -cell dysfunction<sup>27</sup>. Ensuing similar approaches have defined finer-grained clusters of T2D-associated variants<sup>76</sup>,



**Figure 2 | Common resources, analyses and workflows for understanding T2D biology.** Although investigations of biological mechanisms for type 2 diabetes mellitus (T2D) are typically highly context-dependent, in recent years common approaches have emerged. Example goals or lines of investigation (x-axis) include ascertaining samples for genetic analysis (sample ascertainment), computing genotype–phenotype associations (associations), identifying genetic variant carriers (variant carriers), mapping T2D loci (mapped loci), identifying additional signals at mapped loci (additional signals), identifying the causal variant for an association signal (causal variant), identifying the transcript that mediates an association (effector transcript), elucidating the directional relationship between molecular activity and disease risk (direction of risk) and understanding or identifying T2D-relevant biological processes (biological processes). Example resources (y-axis) include deep human phenotypic measurements (deep phenotyping), catalogues of associations reported from genome-wide association studies (GWAS association results), genotypes at select variants (genotype data), full sequence data (sequence data), catalogues of transcript expression (transcript expression), catalogues of epigenomic marks (epigenomics), measurements of cross-species conservation (evolutionary conservation), networks of gene–gene or protein–protein interactions (PPIs) or associations (networks) and experimental results from model systems (cellular and animal models). Analyses that use a class of resource to address a specific question are shown as boxes; examples of these analyses are described in TABLE 1. Three common workflows, or series of analyses, are indicated by sets of arrows between boxes. An example ‘forward genetics’ workflow (dark grey) carries out a GWAS of T2D cases and matched controls, identifies new loci and then predicts causal variants (by fine mapping and mechanistic prediction from epigenomic annotations) or effector transcripts (by analysis of expression quantitative trait loci (eQTLs)). An example ‘reverse genetics’ workflow (light grey) investigates the function of a candidate gene by identifying carriers of high-impact variants and analysing their phenotypes. An example ‘target validation’ workflow (red) identifies a potential target through association analysis for a specific disease subtype (by stratifying individuals based on a range of phenotypic criteria) and then uses additional genetic and experimental analyses to inform the directional relationship between molecular and phenotypic effect, a therapeutic ‘dose–response’ curve and effects of target perturbation on a range of deeper phenotypes. GSEA, genome set enrichment analysis; LoF, loss of function; PheWAS, phenome-wide association study; RVAS, rare variant association study; TFBS, transcription factor binding site.

Table 1 | Example analyses for understanding T2D biology

Goal	Analysis	Description/rationale	Examples
Sample ascertainment	Phenotype clustering	Stratification of samples to identify variants for disease subtypes	Ascertainment based on T2D-related phenotypes <sup>89</sup> or high-dimensional EMRs <sup>91</sup>
	Case-control analysis	Classic association analysis	Most GWAS <sup>65</sup> or sequencing studies <sup>62</sup>
	Genetic risk score analysis	Measure aggregate effects of all variants carried by a patient	Weighted scores from GWAS variants for T2D <sup>97</sup> or related traits <sup>78,98</sup>
	Epigenomics profiling	Use epigenomics rather than genetics to stratify patients	Prediction based on blood methylation levels <sup>96</sup>
Associations	PheWAS	Test for association with all available phenotypes for a given variant	Associations with T2D-related traits <sup>27,104</sup> ; glucose tolerance tests <sup>164</sup> ; fasting versus 2-hour post-OGTT <sup>69</sup> ; glucose monitoring <sup>165</sup> ; insulin sensitivity and signalling <sup>87</sup> ; body mass distribution <sup>78</sup>
	GWAS results	Standard output of a GWAS	Variants reaching genome-wide significance in T2D GWAS <sup>65</sup>
	Single-variant tests	More nuanced association tests not typical of a standard GWAS	Replication of preliminary single-variant associations <sup>51,58</sup> ; custom analyses <sup>89</sup> ; custom conditional analyses <sup>62,90</sup> ; recessive models <sup>69</sup>
	Gene-level tests	Aggregate association analysis, commonly used in RVAS	Aggregate LoF tests <sup>58</sup> ; tests of functional variants <sup>56,57</sup>
	eQTLs	Associations with transcript levels	eQTL analysis of islet RNA-seq data <sup>119,120</sup>
	Epigenetic association tests	Associations with epigenomic differences	Methylation differences between cases and controls <sup>96,122,124,125</sup> or monozygotic twins <sup>124,125</sup>
Variant carriers	Single-variant carriers	Carriers of a specific variant	Identification of low-frequency variant <sup>51,104</sup> or homozygous variant carriers <sup>69</sup>
	Gene mutation carriers	Carriers of any variant meeting specified criteria within a gene	Diagnostic mutation screening <sup>87,165</sup>
Mapped loci	GWAS loci	The loci identified by GWAS	Reported from common variant GWAS <sup>27,36,65</sup>
	RVAS loci	The genes identified by RVAS	Reported from RVAS <sup>51</sup>
Additional signals	Conditional analysis	Used by GWAS or fine mapping to identify multiple signals at a locus	Conditional analysis on GWAS signals <sup>62,70,102</sup>
	Allelic series	Identify variants of a range of effects	Aggregate rare variant analysis for known genes <sup>56–58</sup> or gene sets <sup>62,99</sup>
	Allelic characterization	Functionally validate and calibrate the effects of an allelic series	Missense variant screening with functional assays <sup>56,57</sup>
Causal variant	Credible sets	Localize common variant association to individual variant	Credible set analysis of European <sup>105</sup> or multi-ethnic samples <sup>65</sup> , using dense genotyping <sup>102</sup> or sequencing <sup>62</sup>
	Candidate enumeration	Use sequence data to ensure all variants analysed for credible set	Candidates from 1000 Genomes Project <sup>74</sup> or large-scale sequencing <sup>62</sup>
	Regulatory element overlap	Prioritize variants based on predicted regulatory effect	Enrichment of epigenomic annotations <sup>102,110,111,113–116</sup> or causal variant prioritization <sup>102,115,127,129,135,136</sup>
	Co-inherited TFBS	Novel method to use evolutionary patterns to predict causal variant	Phylogenetic conservation of co-occurring TFBS <sup>138</sup>
	Gene editing, allelic function	Verify variant effect on regulatory landscape or expression	Allele-specific reporter assays <sup>102,113,115,116,129,135,136</sup> or EMSA <sup>102,113,129,135,136</sup> ; CRISPR <sup>127,138</sup>
Effector transcript	Pleiotropy	Effects on multiple traits can hypothesize effector transcript	Multiple variants near <i>RREB1</i> associated with related traits <sup>62</sup>
	High-impact variants	Elevated prior likelihood of causality	Missense variant associations in <i>TM6SF2</i> (REF. 62), <i>PPARG</i> <sup>6</sup> or <i>SLC30A8</i> (REF. 19)
	Coincident eQTLs	Technique used to suggest regulatory effect of association	Association coincident with eQTL from public data sets <sup>19</sup> , custom data sets <sup>51,126,127,131</sup> or islet RNA-seq <sup>120</sup> ; allelic expression profiling <sup>128,129</sup>
	Syntenic	Co-inheritance of variant and gene	Conserved genomic regulatory blocks <sup>130</sup>
	Knockout, knockdown, overexpression	Experimentally measure effects of gene perturbation	Gene overexpression or knockdown in human cell <sup>126,127</sup> or mouse models <sup>127</sup>
Direction of risk	LoF mutations	Observe 'human knockouts'	Effect of gene inactivation <sup>58</sup>
	eQTL directionality	Altered regulation in humans	Direction of expression change <sup>51,120</sup>
	Knockouts, assays	Effect of knockout in model systems	Inference from cellular <sup>126,127</sup> or animal phenotypes <sup>127</sup>

Table 1 cont. | Example analyses for understanding T2D biology

Goal	Analysis	Description/rationale	Examples
Biological processes	Physiological processes	Understanding of physiology from patient investigation	Insight from carrier phenotypes <sup>69,78,87,164,165</sup>
	GSEA	Pathways enriched for GWAS variants	Commonly reported in GWAS <sup>27,36</sup>
	Aggregate tests	Tests of variants across gene set	Aggregate variant effects in monogenic genes <sup>62,99</sup>
	Regulatory networks	Infer networks from variant perturbations of expression patterns	<i>Trans</i> -eQTL analyses <sup>133</sup> or enrichment <sup>132</sup>
	Co-expression, PPIs	Use interactome to organize variants	Networks to rank genes <sup>146</sup> or test hypotheses <sup>145</sup>
	Hypotheses, candidate genes	Classic functional studies and starting point for 'reverse genetics'	Experiments based on biological hypotheses <sup>142–144</sup>

Listed are examples of commonplace analyses, depicted in FIG. 1, to gain insight into T2D biology. For each line of investigation (first column) and analysis (second column) depicted in FIG. 1, listed is a more specific description of the analysis (third column) and examples, not intended as comprehensive, from the literature (fourth column). EMR, electronic medical record; EMSA, electrophoretic mobility shift assay; eQTL, expression quantitative trait loci; GSEA, gene set enrichment analysis; GWAS, genome-wide association study; LoF, loss of function; OGTT, oral glucose tolerance test; PheWAS, phenome-wide association study; PPI, protein–protein interaction; RNA-seq, RNA sequencing; RVAS, rare variant association study; TFBS, transcription factor binding site; T2D, type 2 diabetes mellitus.

classified glycaemic trait associations<sup>77,78</sup> and assessed cellular phenotypes<sup>79</sup>. Recently, Mendelian randomization<sup>80</sup> has gained favour as an approach to evaluate causal relationships between T2D-relevant endophenotypes or biomarkers, supporting causality between body mass index (BMI)<sup>81</sup> or bilirubin levels<sup>82</sup> and T2D, but rejecting causal relationships between triglyceride<sup>83</sup>, high-density lipoprotein (HDL) cholesterol<sup>84</sup>, adiponectin<sup>85</sup> or uric acid<sup>86</sup> levels and T2D.

For some genes, deeper phenotyping of variant carriers has led to quite specific mechanistic insights. For example, insulin resistance measures and muscle and fat biopsy samples of *PTEN* mutation carriers demonstrated a link with T2D via enhanced phosphatidylinositol 3-kinase (PI3K)–protein kinase B (PKB; also known as AKT) pathway signalling and hence increased insulin sensitivity<sup>87</sup>. Similarly, in the study that identified a common but large-effect Inuit-specific T2D variant in *TBC1D4*, access to measures of fasting and 2-hour post-oral glucose tolerance test (OGTT) glucose and insulin showed an association with glucose-stimulated insulin secretion but not with fasting levels of glucose<sup>69</sup>. This finding is consistent with mouse models and a proposed mechanism of action related to glucose uptake by skeletal muscle in the post-prandial setting<sup>88</sup>.

These distinct phenotypic spectra support a perspective widely held by clinicians, who contend with the heterogeneity of T2D on a daily basis, of the existence of distinct T2D subtypes. Indeed, stratification of patients into more homogeneous subgroups has produced novel associations, such as near *LAMA1*, through analysis of lean T2D cases<sup>89</sup> or additional associations with insulin after controlling for BMI<sup>90</sup>. The striking difference in association strength between homozygous *TBC1D4* variants and T2D as defined by glycated haemoglobin levels ( $P = 0.008$ ) versus World Health Organization standards ( $P = 1.6 \times 10^{-24}$ ) further emphasizes how phenotype definition can affect association results<sup>69</sup>. In addition, attempts to formalize specific disease clusters have begun to yield fruit<sup>91</sup>, while intriguing recent studies have suggested novel markers, such as somatic clonal mosaic events (CMEs)<sup>92</sup>, the gut microbiome<sup>93</sup>,

metabolite levels<sup>94,95</sup> or epigenomic marks<sup>96</sup>, to predict or classify T2D subtypes. The clinical benefit of genetic scores to predict diabetes risk, however, remains modest at best, not only for T2D<sup>97,98</sup>, but also for monogenic forms of diabetes mellitus<sup>99</sup>. By contrast, genetic risk scores for type 1 diabetes, in which the *HLA* region explains a substantial portion of heritability, have proved useful to distinguish type 1 from monogenic diabetes<sup>100</sup> or even T2D<sup>101</sup> in specific clinical scenarios.

#### **Prediction of causal variants and effector transcripts.**

To link predicted physiological mechanisms to causal variants or effector transcripts that mediate a T2D association, the first step is typically to fine-map the locus in large sample sizes. Fine mapping starts with conditional analysis to determine the number of independent signals at the locus. Some loci, in fact, harbour a substantial number of signals (for example, five at the *KCNQ1* locus)<sup>102</sup>, which can collectively discern independent risk and protective signals (at the 9p21 locus)<sup>103</sup> or resolve directional relationships between different phenotypes (such as T2D and glucose at the *G6PC2* locus)<sup>53</sup>. Additional signals can also implicate effector transcripts<sup>51</sup> or highlight new phenotypic associations<sup>51,104</sup>, both of which are exemplified by a second association at the *CCND2* locus<sup>51</sup>.

For each signal, fine-mapping approaches then construct credible sets<sup>105</sup> that quantify the likelihood of causality for each variant. Large-scale analyses have, in some cases, identified a single causal candidate with near certainty, such as at the *CDKN2A–CDKN2B*<sup>65</sup> or *MTNR1B*<sup>102</sup> loci. In addition, the transferability of associations between ethnic groups<sup>46,65,68</sup> can enable multiethnic fine mapping studies to increase credible set resolution through differential linkage disequilibrium patterns: for example, at the *JAZF1* locus<sup>65</sup>. Large sample sizes, as well as comprehensive catalogues of variation from efforts such as the 1000 Genomes Project<sup>74</sup>, are crucial to the success of fine mapping<sup>62</sup>.

Among the variants that are identified by fine mapping, perhaps the most valuable are coding variants that have a clear molecular impact on protein function.

#### **Mendelian randomization**

A technique that uses genetic variation to infer causal relationships between correlated phenotypes.

#### **Fine mapping**

An approach to localize common variant association signals to potentially causal variants, using exhaustive candidate enumeration and genotyping in large case–control samples.

Common coding variants in linkage disequilibrium with T2D associations have an elevated prior likelihood of causality and facilitate hypotheses about effector transcripts, such as for *RREB1* at the *RREB1*–*SSR1* locus and *TM6SF2* at the *CILP2*–*TM6SF2* locus<sup>62</sup>. An allelic series of coding variants can provide further insights into gene function. For *SLC30A8*, the identification of protective protein-truncating variants hypothesized an unsuspected relationship between decreased, rather than increased, activity and T2D protection<sup>58</sup>. In genes for monogenic forms of diabetes<sup>62,99</sup>, rare risk-increasing coding variants suggested mechanistic links between T2D and comparatively well-understood Mendelian disease processes<sup>62,101</sup>. For *MTNR1B*<sup>56</sup> or *PPARG*<sup>57</sup>, the ability to functionally discriminate T2D-associated from neutral variants provided a potential foundation to develop disease-relevant assays for therapeutic development<sup>106</sup>.

**Resources for prioritization of causal non-coding variants.** Credible sets for most T2D associations do not include coding variants. Instead, they span non-coding regions of the genome and suggest causal variant regulatory effects. Although early fine-mapping studies<sup>103,105</sup> had limited access to noncoding functional annotations, efforts such as the ENCODE<sup>107</sup>, Epigenomics Roadmap<sup>108</sup> and GTEx<sup>109</sup> projects have since provided reference regulatory landscape maps of numerous human tissues. Indeed, GWAS variants have since been shown to cluster within regulatory elements<sup>110,111</sup>, and T2D-associated variants show the strongest localization to islet enhancers<sup>112,113</sup>. More recent studies have also shown enrichment of T2D-associated variants within transcription-factor binding sites, most notably for forkhead box A2 (*FOXA2*) in islets<sup>102</sup> and *PPAR $\gamma$*  in fat<sup>114</sup>.

To further capture the range of developmental stages and metabolic states of specific relevance to T2D pathophysiology, numerous smaller-scale efforts have increasingly characterized the pancreatic islet regulatory landscape along multiple axes. Early formaldehyde-assisted isolation of regulatory elements sequencing (FAIRE-seq), DNase sequencing (DNase-seq) and chromatin immunoprecipitation followed by sequencing (ChIP-seq) studies identified thousands of islet-selective open chromatin regions<sup>115</sup> and regulatory elements<sup>116</sup>. RNA sequencing (RNA-seq) transcriptome maps have since catalogued long non-coding RNAs (lncRNAs)<sup>117</sup>, microRNAs<sup>118</sup>, and expression quantitative trait loci (eQTLs)<sup>119,120</sup> that are specific to pancreatic islets or  $\beta$ -cells<sup>121</sup>. Islet methylation profiling has been carried out in patients with T2D<sup>122,123</sup> and monozygotic twin pairs<sup>124,125</sup>. Progress has also been made towards integrated catalogues of these resources, identifying regulatory element clusters that are of high relevance to islet function<sup>112,113</sup>.

**Experimental investigation of molecular and cellular risk mechanisms.** These new epigenomic and transcriptomic resources have helped to guide functional experiments for many T2D GWAS loci. To support hypothesized effector transcripts, genotype-dependent expression can provide key evidence, as exemplified by

studies of the *ADCY5* (REF. 126) and *FTO*<sup>127</sup> loci. T2D GWAS thus routinely investigate whether new signals coincide with *cis*-eQTLs<sup>27,36,65</sup>; these analyses are limited by the tissues available in public data sets but can be successful, for example at the *KLF14* (REF. 27) locus. Analyses of eQTLs in custom data sets, such as pancreatic islets, have shown further success, predicting and experimentally validating *ZMIZ1* as a gene involved in glucose homeostasis<sup>120</sup>.

Alternatives have been used to address limitations of *cis*-eQTL mapping, such as potential environmental confounders or restrictions to transcripts near the original signal. Allelic expression profiling is less frequently subject to confounders and has been used in some studies<sup>128</sup>, for example, to implicate *ARAP1* (REF. 129) as an effector transcript (although more recent studies of the locus provide evidence in favour of *STARD10* (REF. 120)). Mapping of chromatin interactions has been used to identify longer-range eQTL relationships, exemplified by implication of the megabase-distant *IRX3* transcript at the *FTO* locus<sup>127,130,131</sup>. Finally, *trans*-eQTLs can be a valuable resource to investigate global regulatory effects of T2D-associated variants<sup>132</sup>, showing, for example, that *KLF14* affects multiple metabolic traits as a master *trans*-regulator of adipose gene expression<sup>133</sup>.

In analogy to transcriptomic maps hypothesizing effector transcripts, epigenomic annotations have proven valuable in hypothesizing causal variants. An early example was a variant near *TCF7L2*, which was shown to lie within islet-selective open chromatin<sup>115</sup> and later an islet-specific enhancer<sup>113</sup>. This variant suggested that allele-specific reporter assays could be used in  $\beta$ -cell lines<sup>115,116</sup> to confirm a previously suggested<sup>134</sup>  $\beta$ -cell-specific risk mechanism. Variants within predicted islet enhancers at the *ZFAND3* (REF. 113) and *MTNR1B*<sup>102</sup> loci have also been verified to have an impact on enhancer activity; in both cases, islet transcription factor binding site (TFBS) maps motivated further experiments to verify a molecular mechanism of decreased neurogenic differentiation factor 1 (*NEUROD1*) binding<sup>102,113</sup>. Similarly, at the *JAZF1* (REF. 135), *CDC123*–*CAMK1D*<sup>136</sup>, *ARAP1* (REF. 129) and *FTO*<sup>127</sup> loci, maps of open chromatin prioritized variants subsequently tested for enhancer activity, allele-specific expression and transcription factor binding.

Other functional experiments show the increasingly diverse genomic resources used in mechanistic characterization. Imprinting patterns led to predictions of effector transcripts at the *KCNQ1* locus<sup>137</sup> and potential explanations for seemingly conflicting phenotype associations at the *GRB10* locus<sup>73</sup>. Patterns of cross-species conservation predicted *IRX3* as the effector transcript at the *FTO* locus<sup>130</sup> and a regulatory variant as causal at the *PPARG* locus<sup>138</sup>. Elucidation of a molecular, cellular and physiological mechanism that underlies obesity risk at the *FTO* locus<sup>127</sup> demonstrates to date perhaps the most comprehensive use of genomic information in a functional setting: epigenomic annotations guided enhancer assays to suggest adipocyte precursors as the affected cell type; long-range chromatin interactions guided eQTL experiments to predict *IRX3* as the effector transcript;

#### Protein-truncating variants

Variants, such as nonsense, frameshift, readthrough or splice site mutations, that lead to incomplete protein sequences and possibly non-functional proteins.

#### Expression quantitative trait loci

(eQTLs). Associations between a genetic marker and expression levels of a transcript.

#### *cis*-eQTLs

Expression quantitative trait loci (eQTLs) on the same chromosome and typically near the location of the gene that encodes the associated transcript.

#### *trans*-eQTLs

Expression quantitative trait loci (eQTLs) in a different chromosome from the gene encoding the associated transcript.



Table 2 | Potential users of an integrated T2D knowledge base

	Statistical geneticists	Biologists (informed by genetics)	Biologists (uninformed by genetics)	Pharmaceutical researchers	Clinicians
<b>Major goals</b>	<ul style="list-style-type: none"> <li>Identify novel associations</li> <li>Explain heritability</li> </ul>	Elucidate molecular, cellular and physiological mechanisms of association	Examine human genetic support for hypothesis from model system or pathway	<ul style="list-style-type: none"> <li>Identify potential novel targets</li> <li>Obtain support for targets from human genetics</li> </ul>	Estimate phenotypic effect for variant of interest
<b>Relevant data sets</b>	<ul style="list-style-type: none"> <li>Large collections of genotype and sequence data</li> <li>Basic phenotypes</li> </ul>	<ul style="list-style-type: none"> <li>Association catalogue</li> <li>Credible set results</li> <li>Epigenomic and transcriptomic data sets</li> </ul>	<ul style="list-style-type: none"> <li>Coding variation</li> <li>Annotation of functional impact</li> <li>Rich phenotypes</li> </ul>	<ul style="list-style-type: none"> <li>GWAS associations</li> <li>Coding variation</li> <li>Annotations</li> <li>Rich phenotypes</li> </ul>	<ul style="list-style-type: none"> <li>Large sequence datasets</li> <li>Rich phenotypes</li> </ul>
<b>Typical workflow</b>	Carry out novel association analysis on entire body of data	Examine overlap of functional annotations with variants in credible set	Identify high-impact variants in gene of interest and test for association with range of phenotypes	<ul style="list-style-type: none"> <li>Predict causal gene, directionality from GWAS signal</li> <li>Assess phenotypic associations for LoF variants in gene of interest</li> </ul>	Obtain variant frequencies in individuals with various phenotypes
<b>Preferred interface</b>	<ul style="list-style-type: none"> <li>Programmatic APIs</li> <li>Exploratory ‘sandbox’</li> </ul>	Visualizations or genome browser within a portal	Query builders and analysis modules within a portal	<ul style="list-style-type: none"> <li>Portal with support for workflows</li> <li>Programmatic APIs</li> </ul>	<ul style="list-style-type: none"> <li>Simple portal</li> <li>Programmatic APIs</li> </ul>
<b>Knowledge of genetics</b>	Expert	Intermediate	Basic or none	Basic or intermediate	Basic
<b>Needed curated content</b>	Description of methods applied to data and results	<ul style="list-style-type: none"> <li>Qualitative interpretation of statistics</li> <li>Potential caveats</li> <li>Knowledge relevant to gene or pathway</li> </ul>	<ul style="list-style-type: none"> <li>Guide for workflow to run</li> <li>Relevant statistics and their interpretation</li> <li>Potential caveats</li> </ul>	<ul style="list-style-type: none"> <li>Summary of validated genetic results</li> <li>Qualitative interpretation of statistics</li> </ul>	Documentation of ascertainment protocol and phenotypic measurements
<b>Type of finding made possible</b>	LoF variants in <i>SLC30A8</i> (REF. 58)	Causal variants, effector transcripts or mechanisms at <i>MTNR1B</i> <sup>102</sup> or <i>ZMIZ1</i> (REF. 120) loci	Associations with T2D for mutations in <i>LIN28-let-7</i> (REF. 143) pathway or PPAR $\gamma$ binding sites <sup>114</sup>	Protective LoF variants in <i>SLC30A8</i> (REF. 58), assay calibration with PPAR $\gamma$ variants <sup>57</sup> , deep phenotyping of <i>PTEN</i> variant carriers <sup>87</sup>	Investigation of phenotypic effects of variants in monogenic diabetes genes <sup>99</sup>

Many communities might benefit from analyses that can be carried out on integrated genomic data sets. The table shows five potential user groups that might find use from a type 2 diabetes mellitus (T2D) knowledge base, together with their major analytical goals, data sets of high relevance, an example workflow that might be supported, the preferred interface to the knowledge base, the expected level of genetics expertise, the curated content that might be necessary to augment the analysis and the type of finding that a knowledge base might facilitate in the future. The groups listed are statistical geneticists who typically carry out genome-wide association studies (GWAS) or complex genetic association analyses, biologists who aim to translate genetic associations into mechanistic insight, biologists who investigate hypotheses from animal or cellular models, pharmaceutical researchers who seek to identify or prioritize targets based on human genetics, and clinicians who seek to interpret genetic variants identified in patient genome sequences. Both the examples and classification of users are simplistic and not intended to be comprehensive. Furthermore, a typical researcher may in reality move between different user groups depending on the line of research. API, application programming interface; LoF, loss of function; PPAR $\gamma$ , peroxisome proliferator activated receptor- $\gamma$ .

co-expression and *trans*-eQTL analysis guided cellular phenotyping, gene perturbation and mouse models to predict adipocyte browning as a cellular mechanism; and cross-species conservation of TFBS patterns guided CRISPR–Cas9 editing to predict the causal variant.

These recent experiments thus augment more traditional functional investigations of genetic associations, such as using gene knockdown to validate *ATM* in metformin response<sup>29</sup>, mouse models to investigate *CDKAL1* (REF. 139), mRNA expression analysis to study the tissue of action for *SLC2A2* (REF. 140), or biochemical assays to study the function of *GCKR*<sup>141</sup>. Conversely, genetic associations and genomic resources can also be used to test predictions that were originally made by functional studies; for example, a prediction from mouse adipose tissue that variants in PPAR $\gamma$ -binding DNA elements account for strain-specific expression patterns and drug responses was tested through enrichment of metabolic trait associations within human PPAR $\gamma$ -binding elements<sup>114</sup>. Other predictions from mice, such as the relevance of *Rfx6* (REF. 142), the *Lin28-let-7* pathway<sup>143</sup> and *Imp2* (REF. 144) to glycaemic traits or diabetes, could also

be investigated through phenotypic analysis of human mutation carriers. The same can be said for predictions from functional genomic approaches, such as interactome analysis<sup>145,146</sup> or knockout screens from new assays<sup>147</sup>.

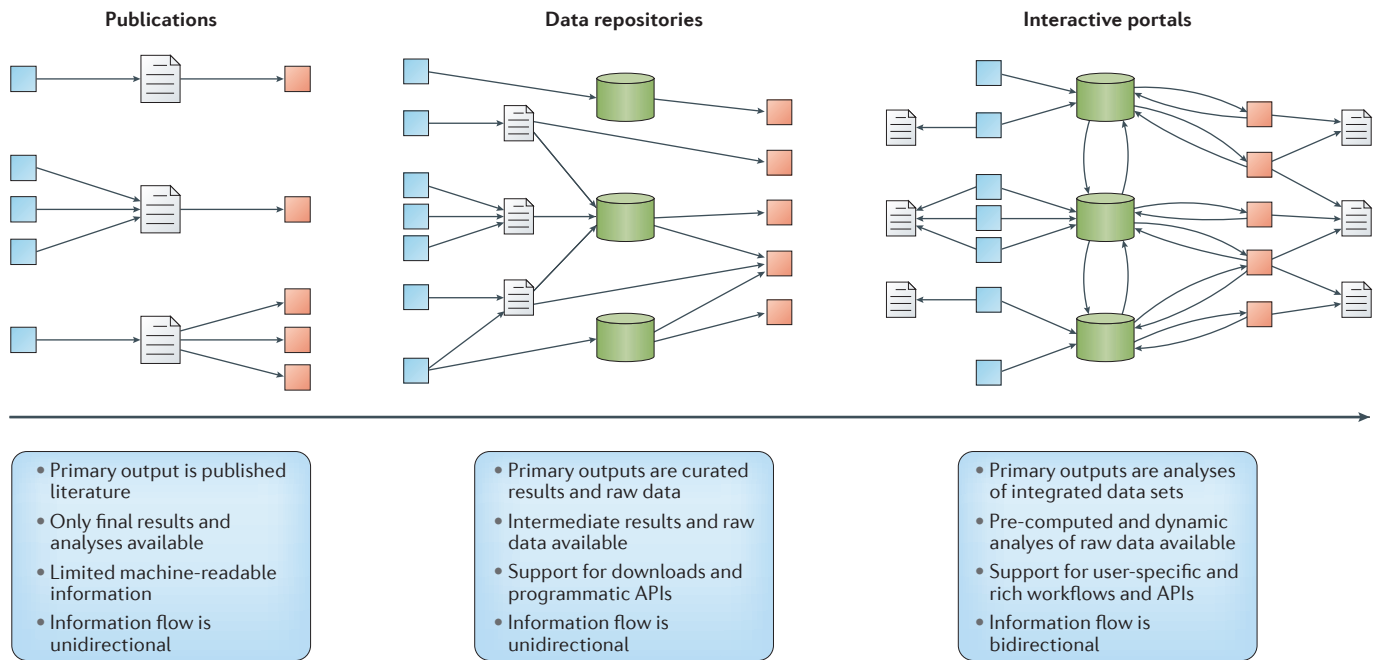
### A need for a public T2D genetics knowledge base

The state of genetic mapping and functional studies for T2D highlight two different but complementary trends. The need for larger consortium-led GWAS and sequencing studies will produce centralized analyses of large genomic data sets. The need for increasingly diverse and specialized approaches to investigate these findings will produce a need for common genomic resources. And yet, today, the value of genetic data sets is only partially realized. For example, protective loss-of-function mutations in *SLC30A8* were undetected in a consortium-led genome-wide analysis of 13,000 exomes<sup>62</sup> and only identified after years of focused genetic and functional follow up<sup>58</sup>.

For the biological community to make fuller and more intelligent use of the data produced by genomics consortia, barriers must be removed. Current consortia

**CRISPR–Cas9 editing**  
A technique for precise and efficient editing of genetic information within a cell.

**Interactome**  
The study of all protein–protein interactions in a system.



**Figure 3 | The evolution of human genetic knowledge bases.** For experimentalists to make wider use of the analyses and data produced by genetics consortia (FIG. 1), new access mechanisms are necessary. Traditionally, the results produced by genome-wide association study (GWAS) consortia have been primarily accessible through publications, which list loci that meet stringent genome-wide significance. Catalogues of associations require manual curation from the literature. The left panel shows the flow of information from data producers (blue boxes on the left) to publications (white papers), which are then read by data consumers (pink boxes on the right). Over time, many studies have begun to make available a wider range of intermediate results alongside publications, such as files that contain associations for every analysed variant. Genomic resources commonly used in experiments to understand GWAS associations (FIG. 2) have also been made accessible through databases and portals. The middle panel shows a revised flow of information, in which not only published results but also intermediate analyses or raw data are made available through multiple knowledge bases (green cylinders). In the future, we argue that much more information could be extracted from genetic analyses if portals made available results as well as facilitated novel analyses on data sets integrated across multiple studies. The right panel shows a future flow of information, in which consumers of data specify complex analytical workflows, which are carried out on data within interconnected knowledge bases. By enabling a broad community of users to carry out custom analyses, progress towards understanding type 2 diabetes mellitus (T2D) biology may accelerate (shown by publications on the right). API, application programming interface.

settings are seldom the ideal context for producing the analyses required for any specific experiment; although custom analyses can be requested, disparate data set locations and multiple analysts can make it challenging to compute results in the most rigorous and expedited fashion. Furthermore, experimentalists or other ‘consumers’ of genomic data have varied levels of expertise and may be unable to interpret current analytical practices (TABLE 2).

The increased convergence of experimental approaches and resources, however, suggests that many investigators could benefit from access to genomic data sets through a tractable number of user-friendly workflows. A logical mechanism, recently embraced by the wider T2D community<sup>148</sup>, is a web-based portal to an integrated T2D genetics knowledge base, designed explicitly to aid translation of association results to biological insights. Other genotype–phenotype databases have become commonplace, most notably for Mendelian diseases<sup>149,150</sup>, and the inherent challenges in building and supporting them are well reviewed elsewhere<sup>151</sup>.

Additional reference epigenomic<sup>107,108</sup>, transcriptomic<sup>109</sup> and genetic<sup>152</sup> data sets commonly used in functional experiments are also increasingly available through public web portals: in some cases through integrated analyses<sup>153,154</sup>.

To facilitate functional studies of T2D loci, a portal must be designed to support experimentalists, rather than to serve as a repository for consortium-produced results (FIG. 3). This will require increased interactivity, data drawn from multiple interconnected sources and access to on-demand analyses of sensitive genotype and phenotype data. Such a portal would forge new connections among GWAS consortia and experimentalists in the academic, government or private sectors, diversifying the downstream uses of genetic data sets.

### Building a knowledge base

To realize the vision of a shared T2D knowledge base and portal (FIG. 4), several efforts offer insight. Modern portals, such as those for the [GTEx](#) or [WashU Epigenome Browser](#) projects, are increasingly interactive

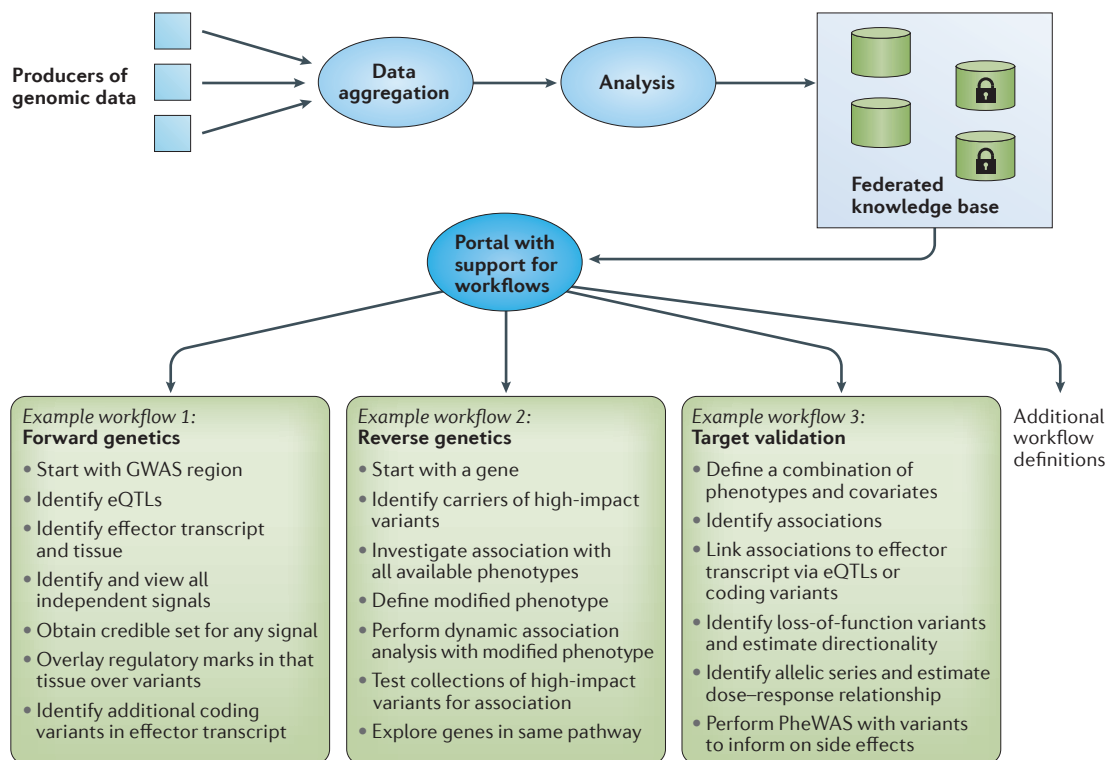


Figure 4 | **An integrated knowledge base of T2D genetics.** An integrated knowledge base of type 2 diabetes mellitus (T2D) genetics requires a new approach to aggregate, analyse and democratize genomic data sets. The most foundational need is to encourage cohort investigators to share genotypic and phenotypic information. Where regulations permit, these data could be aggregated at a central location, whereas federated knowledge bases will be necessary to integrate data sets that must be kept at their original locations. To enable a broad community of users to access information within these data, the methods and approaches developed by genome-wide association study (GWAS) consortia must be systematically applied to compute the needed analyses, either ‘offline’ in advance or ‘online’ on demand. Interactive web portals are then necessary for users to access these results, probably through workflows that span multiple questions or analyses (FIG. 2). eQTL, expression quantitative trait locus; PheWAS, phenome-wide association study.

and responsive. Collaborations such as the exome aggregation consortium<sup>155</sup> have pioneered approaches to aggregate raw sequence data from independent projects, whereas more aspirational efforts, such as the [Global Alliance for Genomics and Health \(GA4GH\)](#), seek to create standards and incentives for responsible genomic and clinical data sharing. Web platforms such as [GenomeSpace](#) and [Galaxy](#)<sup>156</sup> allow a wide community of (even non-expert) researchers to reproducibly carry out analyses, whereas tools such as the [Michigan](#) or [Sanger](#) imputation services and [xBrowse](#) allow secure upload and automated expert analysis of user data. Increasingly, new initiatives recognize the value in collaboration across prior boundaries, such as between industry and academia in the case of the Accelerating Medicines Partnership (AMP)<sup>157</sup>, Center for Therapeutic Target Validation (CTTV)<sup>158</sup> and Innovative Medicines Initiative (IMI)<sup>159</sup>.

What is needed in order to exploit these trends to advance genetic or biological understanding of T2D? On the basis of recent T2D genetic studies, extremely large sample sizes will be necessary to map or to characterize substantially more loci in the population<sup>62</sup>. Rich phenotype information will be needed for physiological

characterizations, and additional epigenomic and transcriptomic data sets will be necessary to support molecular or cellular experiments. The foundational task for a T2D knowledge base is thus to support large-scale genetic data aggregation, harmonization and integration with other resources. Although ‘business intelligence’ approaches commonly address such challenges through ‘big data’ analytics or ‘data warehouses’ (REF. 160), applying analogous computational techniques to human genetic data sets will require substantial research and investment.

As a T2D knowledge base will contain sensitive patient data, secure mechanisms will be necessary to enforce access to only consented data sets. Where regulations prevent transfer to a central location, federated databases — in which data reside locally but support central analyses — will need development. Incentives such as increased access to cutting-edge analytics or fair publication and data embargo policies must be developed to inspire investigators, who typically spend years or decades collecting patient data, to contribute to the knowledge base for the sake of global discovery. Potentially, a shared knowledge base might enable entirely new categories of data sets, such as those from clinical trials

**Business intelligence**

A term, commonly used in business, that denotes a set of techniques for transforming raw data into meaningful insights.

**Big data**

A term for data sets that are so large or complex that new paradigms are needed to extract meaningful insights from them.

**Data warehouses**

A system for carrying out integrated analyses across multiple initially disparate data sources.

with rich information on patient drug responses<sup>161</sup>, to be integrated with genetic association analyses for the first time.

Data sets in a knowledge base must be further analysed to produce broadly applicable results. Consortium-developed methods and workflows for genetic association analysis must be made sufficiently robust and efficient to analyse large and diverse data set collections. ‘Offline’ analyses will need regular (and automatic, where possible) application to update pre-computed statistics, while ‘online’ analyses will be needed for queries that cannot be pre-computed, such as novel patient stratifications to define homogeneous disease subtypes or aggregate association tests<sup>63</sup> of custom variant sets from functional assays<sup>56,57</sup>. Some online analyses may be complex, offering an opportunity for users to create and share custom workflows with the community. Where necessary, standards or new approaches will be needed to extend the analyses offered by a knowledge base to a federated setting.

Association results will also require additional processing to produce interpretable answers to biological queries. Statistical significance must be carefully conveyed to non-technical investigators to limit both type 1 errors, which arise from the multiple testing burden of ‘online’ analyses, and type 2 errors, which arise from false equation of absence of evidence with evidence of absence. Analyses from the knowledge base should also be packaged to suggest actionable hypotheses or worthwhile experiments, such as affected cell types or ranked lists of causal variants. Ideally, researchers who investigate these hypotheses should be encouraged to contribute their experimental results back to the knowledge base to empower future analyses by others.

Finally, these data and analyses must be made available through a public portal. Analytical results will need accompanying summaries and visualizations that answer simple but specific biological questions. Users will need the capabilities to carry out interactive workflows and have the ability to save work or rerun analyses reproducibly even if the underlying data are updated. Clear documentation and educational content must accompany the data and methods and must be tailored to a broad community with varied levels of expertise.

In summary, an ideal T2D knowledge base must be comprehensive, secure, compliant, automated and rigorous, yet also interpretable and inviting. If successful, it would greatly enhance the value of genetic association studies by generating a synergistic link between a network of data contributors and a community of experimentalists. It can therefore not only serve as a central repository of a large trove of genomic information but also aspire to ‘democratize genetics’ to the global community, thus becoming a transformative engine for discovery and paradigm for other disease fields.

### Conclusions

Two trends have emerged from the study of T2D genetics in recent years. First, evidence has mounted that T2D genetic architecture is likely to be polygenic and characterized by many loci that are detectable only in hundreds of thousands of samples, arguing that larger and larger collections of genetic data will be necessary to discover disease-relevant variants in the population. Second, the increasing number of genes or processes that are linked to T2D will only increase the diversity of approaches necessary to translate these associations to mechanistic insight, although common resources and workflows have begun to emerge.

We have argued that a T2D knowledge base and portal could exploit both of these trends and create additional value from the consortium-led analyses that will probably continue to be the main paradigm for genetic mapping. Through an interactive portal that connects researchers around the world, the current — mostly unidirectional — flow of information from genetic discovery to functional characterization may one day be augmented by a cycle in which functional experiments inspire additional genetic experiments as well. One such portal, the [Type 2 Diabetes Knowledge Portal](#), represents the efforts of over 100 investigators to increase access to genetic analyses of tens to hundreds of thousands of samples. Whatever the means to facilitate collaboration, a diverse collection of communities using increasingly common resources will be necessary to continue to advance our understanding of T2D biology and help one day to improve patient treatment and outcomes.

- Hemminki, K., Li, X., Sundquist, K. & Sundquist, J. Familial risks for type 2 diabetes in Sweden. *Diabetes Care* **33**, 293–297 (2010).
- Almgren, P. *et al.* Heritability and familiarity of type 2 diabetes and related quantitative traits in the Botnia Study. *Diabetologia* **54**, 2811–2819 (2011).
- Kahn, S. E., Cooper, M. E. & Del Prato, S. Pathophysiology and treatment of type 2 diabetes: perspectives on the past, present, and future. *Lancet* **383**, 1068–1083 (2014).
- Fowler, M. J. Microvascular and macrovascular complications of diabetes. *Clin. Diabetes* **26**, 77–82 (2008).
- Tancredi, M. *et al.* Excess mortality among persons with type 2 diabetes. *N. Engl. J. Med.* **373**, 1720–1732 (2015).
- Altshuler, D. *et al.* The common PPAR $\gamma$  Pro12Ala polymorphism is associated with decreased risk of type 2 diabetes. *Nat. Genet.* **26**, 76–80 (2000).
- Gloyn, A. L. *et al.* Large-scale association studies of variants in genes encoding the pancreatic  $\beta$ -cell K<sub>ATP</sub> channel subunits Kir6.2 (*KCNJ11*) and SUR1 (*ABCC8*) confirm that the *KCNJ11* E23K variant is associated with type 2 diabetes. *Diabetes* **52**, 568–572 (2003).
- Grant, S. F. *et al.* Variant of transcription factor 7-like 2 (*TCF7L2*) gene confers risk of type 2 diabetes. *Nat. Genet.* **38**, 320–323 (2006).
- Guan, W., Pluzhnikov, A., Cox, N. J. & Boehnke, M. Meta-analysis of 23 type 2 diabetes linkage studies from the international type 2 diabetes linkage analysis consortium. *Hum. Hered.* **66**, 35–49 (2008).
- Altshuler, D., Daly, M. J. & Lander, E. S. Genetic mapping in human disease. *Science* **322**, 881–888 (2008).
- McCarthy, M. I. Genomics, type 2 diabetes, and obesity. *N. Engl. J. Med.* **363**, 2339–2350 (2010).
- McClellan, J. & King, M. C. Genetic heterogeneity in human disease. *Cell* **141**, 210–217 (2010).
- Goldstein, D. B. Common genetic variation and human traits. *N. Engl. J. Med.* **360**, 1696–1698 (2009).
- Risch, N. & Merikangas, K. The future of genetic studies of complex human diseases. *Science* **273**, 1516–1517 (1996).
- Reich, D. E. & Lander, E. S. On the allelic spectrum of human disease. *Trends Genet.* **17**, 502–510 (2001).
- International HapMap Consortium. A haplotype map of the human genome. *Nature* **437**, 1299–1320 (2005).
- Diabetes Genetics Initiative. Genome-wide association analysis identifies loci for type 2 diabetes and triglyceride levels. *Science* **316**, 1331–1336 (2007).
- Scott, L. J. *et al.* A genome-wide association study of type 2 diabetes in Finns detects multiple susceptibility variants. *Science* **316**, 1341–1345 (2007).
- Sladek, R. *et al.* A genome-wide association study identifies novel risk loci for type 2 diabetes. *Nature* **445**, 881–885 (2007).
- Steinthorsdottir, V. *et al.* A variant in *CDKAL1* influences insulin response and risk of type 2 diabetes. *Nat. Genet.* **39**, 770–775 (2007).
- Zeggini, E. *et al.* Replication of genome-wide association signals in UK samples reveals risk loci for type 2 diabetes. *Science* **316**, 1336–1341 (2007).
- Zeggini, E. *et al.* Meta-analysis of genome-wide association data and large-scale replication identifies additional susceptibility loci for type 2 diabetes. *Nat. Genet.* **40**, 638–645 (2008).

23. de Bakker, P. I. *et al.* Practical aspects of imputation-driven meta-analysis of genome-wide association studies. *Hum. Mol. Genet.* **17**, R122–R128 (2008).
24. Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of genome-wide association scans. *Bioinformatics* **26**, 2190–2191 (2010).
25. Speliotes, E. K. *et al.* Association analyses of 249,796 individuals reveal 18 new loci associated with body mass index. *Nat. Genet.* **42**, 937–948 (2010).
26. Teslovich, T. M. *et al.* Biological, clinical and population relevance of 95 loci for blood lipids. *Nature* **466**, 707–713 (2010).
27. Voight, B. F. *et al.* Twelve type 2 diabetes susceptibility loci identified through large-scale association analysis. *Nat. Genet.* **42**, 579–589 (2010).  
**This is an early illustration of the paradigm now established for GWAS analyses.**
28. Dupuis, J. *et al.* New genetic loci implicated in fasting glucose homeostasis and their impact on type 2 diabetes risk. *Nat. Genet.* **42**, 105–116 (2010).
29. Zhou, K. *et al.* Common variants near *ATM* are associated with glycaemic response to metformin in type 2 diabetes. *Nat. Genet.* **43**, 117–120 (2011).
30. Saxena, R. *et al.* Genetic variation in *GIPR* influences the glucose and insulin responses to an oral glucose challenge. *Nat. Genet.* **42**, 142–148 (2010).
31. Strawbridge, R. J. *et al.* Genome-wide association identifies nine common variants associated with fasting proinsulin levels and provides new insights into the pathophysiology of type 2 diabetes. *Diabetes* **60**, 2624–2634 (2011).
32. Cho, Y. S. *et al.* Meta-analysis of genome-wide association studies identifies eight new loci for type 2 diabetes in East Asians. *Nat. Genet.* **44**, 67–72 (2012).
33. Kooner, J. S. *et al.* Genome-wide association study in individuals of South Asian ancestry identifies six new type 2 diabetes susceptibility loci. *Nat. Genet.* **43**, 984–989 (2011).
34. Yamauchi, T. *et al.* A genome-wide association study in the Japanese population identifies susceptibility loci for type 2 diabetes at *UBE2E2* and *C2CD4A–C2CD4B*. *Nat. Genet.* **42**, 864–868 (2010).
35. Voight, B. F. *et al.* The metabochip, a custom genotyping array for genetic studies of metabolic, cardiovascular, and anthropometric traits. *PLoS Genet.* **8**, e1002793 (2012).
36. Morris, A. P. *et al.* Large-scale association analysis provides insights into the genetic architecture and pathophysiology of type 2 diabetes. *Nat. Genet.* **44**, 981–990 (2012).
37. Scott, R. A. *et al.* Large-scale association analyses identify new loci influencing glycaemic traits and provide insight into the underlying biological pathways. *Nat. Genet.* **44**, 991–1005 (2012).
38. Manolio, T. A. *et al.* Finding the missing heritability of complex diseases. *Nature* **461**, 747–753 (2009).  
**A thorough summary of the debate around different models of genetic architecture is presented here, following the results of early GWAS.**
39. Hirschhorn, J. N. Genomewide association studies — illuminating biologic pathways. *N. Engl. J. Med.* **360**, 1699–1701 (2009).
40. Dickson, S. P., Wang, K., Krantz, I., Hakonarson, H. & Goldstein, D. B. Rare variants create synthetic genome-wide associations. *PLoS Biol.* **8**, e1000294 (2010).
41. Bodmer, W. & Bonilla, C. Common and rare variants in multifactorial susceptibility to common diseases. *Nat. Genet.* **40**, 695–701 (2008).
42. Pritchard, J. K. Are rare variants responsible for susceptibility to complex diseases? *Am. J. Hum. Genet.* **69**, 124–137 (2001).
43. Goldstein, D. B. The importance of synthetic associations will only be resolved empirically. *PLoS Biol.* **9**, e1001008 (2011).
44. Anderson, C. A., Soranzo, N., Zeggini, E. & Barrett, J. C. Synthetic associations are unlikely to account for many common disease genome-wide association signals. *PLoS Biol.* **9**, e1000580 (2011).
45. Wray, N. R., Purcell, S. M. & Visscher, P. M. Synthetic associations created by rare variants do not explain most GWAS results. *PLoS Biol.* **9**, e1000579 (2011).
46. Waters, K. M. *et al.* Consistent association of type 2 diabetes risk variants found in Europeans in diverse racial and ethnic groups. *PLoS Genet.* **6**, e1001078 (2010).
47. Stahl, E. A. *et al.* Bayesian inference analyses of the polygenic architecture of rheumatoid arthritis. *Nat. Genet.* **44**, 483–489 (2012).
48. Agarwala, V., Flannick, J., Sunyaev, S., Go, T. D. C. & Altshuler, D. Evaluating empirical bounds on complex disease genetic architecture. *Nat. Genet.* **45**, 1418–1427 (2013).  
**This is a comprehensive study of the support from empirical data sets for different models of T2D genetic architecture.**
49. Huyghe, J. R. *et al.* Exome array analysis identifies new loci and low-frequency variants influencing insulin processing and secretion. *Nat. Genet.* **45**, 197–201 (2013).
50. Albrechtsen, A. *et al.* Exome sequencing-driven discovery of coding polymorphisms associated with common metabolic phenotypes. *Diabetologia* **56**, 298–310 (2013).
51. Steinthorsdottir, V. *et al.* Identification of low-frequency and rare sequence variants associated with elevated or reduced risk of type 2 diabetes. *Nat. Genet.* **46**, 294–298 (2014).
52. Estrada, K. *et al.* Association of a low-frequency variant in *HNF1A* with type 2 diabetes in a Latino population. *JAMA* **311**, 2305–2314 (2014).
53. Mahajan, A. *et al.* Identification and functional characterization of *G6PC2* coding variants influencing glycaemic traits define an effector transcript at the *G6PC2-ABCB11* locus. *PLoS Genet.* **11**, e1004876 (2015).
54. Wessel, J. *et al.* Low-frequency and rare exome chip variants associate with fasting glucose and type 2 diabetes susceptibility. *Nat. Commun.* **6**, 5897 (2015).
55. Manning, A. K. *et al.* A low frequency *AKT2* coding variant enriched in the Finnish population is associated with fasting insulin levels. (Abstract #56) *The 64th Annual Meeting of The American Society of Human Genetics, San Diego, California* [http://www.ashg.org/2014meeting/pdf/2014\\_ASHG\\_Meeting\\_Platform\\_Abstracts.pdf](http://www.ashg.org/2014meeting/pdf/2014_ASHG_Meeting_Platform_Abstracts.pdf) (18–22 Oct 2014).
56. Bonnefond, A. *et al.* Rare *MTNR1B* variants impairing melatonin receptor 1B function contribute to type 2 diabetes. *Nat. Genet.* **44**, 297–301 (2012).  
**This is an early example of the ability of functional assays to filter benign from deleterious alleles and improve the power of aggregate association tests.**
57. Majithia, A. R. *et al.* Rare variants in *PPARG* with decreased activity in adipocyte differentiation are associated with increased risk of type 2 diabetes. *Proc. Natl Acad. Sci. USA* **111**, 13127–13132 (2014).
58. Flannick, J. *et al.* Loss-of-function mutations in *SLC30A8* protect against type 2 diabetes. *Nat. Genet.* **46**, 357–363 (2014).  
**This is one of the first studies to identify protective rare variants associated with T2D, demonstrating the need for large-scale data aggregation.**
59. Rutter, G. A. Think zinc: new roles for zinc in the control of insulin secretion. *Islets* **2**, 49–50 (2010).
60. Nicolson, T. J. *et al.* Insulin storage and glucose homeostasis in mice null for the granule zinc transporter ZnT8 and studies of the type 2 diabetes-associated variants. *Diabetes* **58**, 2070–2083 (2009).
61. Lohmueller, K. E. *et al.* Whole-exome sequencing of 2,000 Danish individuals and the role of rare coding variants in type 2 diabetes. *Am. J. Hum. Genet.* **93**, 1072–1086 (2013).
62. Fuchsberger, C. *et al.* The genetic architecture of type 2 diabetes. *Nature* <http://dx.doi.org/10.1038/nature18642> (2016).  
**A comprehensive characterization is presented here of the genetic architecture of T2D using various sequencing approaches.**
63. Moutsianas, L. *et al.* The power of gene-based rare variant methods to detect disease-associated variation and test hypotheses about complex disease. *PLoS Genet.* **11**, e1005165 (2015).  
**This simulation study shows the need for potentially numerous aggregate association analyses to identify disease genes with different allelic spectra.**
64. Zuk, O. *et al.* Searching for missing heritability: designing rare variant association studies. *Proc. Natl Acad. Sci. USA* **111**, E455–E464 (2014).  
**Approaches for rare variant association studies are thoroughly explained and examined.**
65. Mahajan, A. *et al.* Genome-wide trans-ancestry meta-analysis provides insight into the genetic architecture of type 2 diabetes susceptibility. *Nat. Genet.* **46**, 234–244 (2014).
66. Scott, R. A. *et al.* Genome-wide association study imputed to 1000 Genomes reveals 18 novel associations with type 2 diabetes. (Abstract 53). *The 64th Annual Meeting of The American Society of Human Genetics, San Diego, California* [http://www.ashg.org/2014meeting/pdf/2014\\_ASHG\\_Meeting\\_Platform\\_Abstracts.pdf](http://www.ashg.org/2014meeting/pdf/2014_ASHG_Meeting_Platform_Abstracts.pdf) (18–22 Oct 2014).
67. Sigma Type 2 Diabetes Consortium. Sequence variants in *SLC16A11* are a common risk factor for type 2 diabetes in Mexico. *Nature* **506**, 97–101 (2014).
68. Ng, M. C. Y. *et al.* Meta-analysis of genome-wide association studies in African Americans provides insights into the genetic architecture of type 2 diabetes. *PLoS Genet.* **10**, e1004517 (2014).
69. Moltke, I. *et al.* A common Greenlandic *TBC1D4* variant confers muscle insulin resistance and type 2 diabetes. *Nature* **512**, 190–193 (2014).  
**The largest effect association observed for T2D to date is described here, demonstrating the power of studying population isolates as well as insights that can be learned from variant carrier phenotyping.**
70. Mahajan, A. *et al.* Large-scale exome chip association analysis identifies novel type 2 diabetes susceptibility loci and highlights candidate effector genes. (Abstract #299). *The 65th Annual Meeting of The American Society of Human Genetics, Baltimore, Maryland* [http://www.ashg.org/2015meeting/pdf/57175\\_Platform.pdf](http://www.ashg.org/2015meeting/pdf/57175_Platform.pdf) (6–10 Oct 2015).
71. Bulik-Sullivan, B. K. *et al.* LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).
72. Iyengar, S. K. *et al.* Genome-wide association and trans-ethnic meta-analysis for advanced diabetic kidney disease: Family Investigation of Nephropathy and Diabetes (FIND). *PLoS Genet.* **11**, e1005352 (2015).
73. Prokopenko, I. *et al.* A central role for *GRB10* in regulation of islet function in man. *PLoS Genet.* **10**, e1004235 (2014).
74. Horikoshi, M. *et al.* Discovery and fine-mapping of glycaemic and obesity-related trait loci using high-density imputation. *PLoS Genet.* **11**, e1005230 (2015).
75. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *bioRxiv* <http://dx.doi.org/10.1101/035170> (2015).
76. Dimas, A. S. *et al.* Impact of type 2 diabetes susceptibility variants on quantitative glycaemic traits reveals mechanistic heterogeneity. *Diabetes* **63**, 2158–2171 (2014).
77. Ingelsson, E. *et al.* Detailed physiologic characterization reveals diverse mechanisms for novel genetic loci regulating glucose and insulin metabolism in humans. *Diabetes* **59**, 1266–1275 (2010).
78. Scott, R. A. *et al.* Common genetic variants highlight the role of insulin resistance and body fat distribution in type 2 diabetes, independent of obesity. *Diabetes* **63**, 4378–4387 (2014).
79. Rosengren, A. H. *et al.* Reduced insulin exocytosis in human pancreatic  $\beta$ -cells with gene variants linked to type 2 diabetes. *Diabetes* **61**, 1726–1733 (2012).
80. Lawlor, D. A., Harbord, R. M., Sterne, J. A., Timpson, N. & Davey Smith, G. Mendelian randomization: using genes as instruments for making causal inferences in epidemiology. *Stat. Med.* **27**, 1133–1163 (2008).
81. Fall, T. *et al.* The role of adiposity in cardiometabolic traits: a Mendelian randomization analysis. *PLoS Med.* **10**, e1001474 (2013).
82. Abbasi, A. *et al.* Bilirubin as a potential causal factor in type 2 diabetes risk: a Mendelian randomization study. *Diabetes* **64**, 1459–1469 (2015).
83. De Silva, N. M. *et al.* Mendelian randomization studies do not support a role for raised circulating triglyceride levels influencing type 2 diabetes, glucose levels, or insulin resistance. *Diabetes* **60**, 1008–1018 (2011).
84. Haase, C. L., Tybjaerg-Hansen, A., Nordestgaard, B. G. & Frikke-Schmidt, R. H. D. L. Cholesterol and risk of type 2 diabetes: a Mendelian randomization study. *Diabetes* **64**, 3328–3333 (2015).
85. Yaghoobkar, H. *et al.* Mendelian randomization studies do not support a causal role for reduced circulating adiponectin levels in insulin resistance and type 2 diabetes. *Diabetes* **62**, 3589–3598 (2013).
86. Sluijs, I. *et al.* A mendelian randomization study of circulating uric acid and type 2 diabetes. *Diabetes* **64**, 3028–3036 (2015).

87. Pal, A. *et al.* *PTEN* mutations as a cause of constitutive insulin sensitivity and obesity. *N. Engl. J. Med.* **367**, 1002–1011 (2012).  
**This paper provides a nice illustration of the power of deep phenotyping studies to gain molecular, cellular and physiological insights into disease pathways.**
88. Wang, L. *et al.* *PTEN* deletion in pancreatic  $\alpha$ -cells protects against high-fat diet–induced hyperglucagonemia and insulin resistance. *Diabetes* **64**, 147–157 (2015).
89. Perry, J. R. B. *et al.* Stratifying type 2 diabetes cases by BMI identifies genetic risk variants in *LAMA1* and enrichment for risk variants in lean compared to obese cases. *PLoS Genetics* **8**, e1002741 (2012).  
**This is an example in which phenotypic sample stratification identified a novel GWAS locus.**
90. Manning, A. K. *et al.* A genome-wide approach accounting for body mass index identifies genetic variants influencing fasting glycaemic traits and insulin resistance. *Nat. Genet.* **44**, 659–669 (2012).
91. Li, L. *et al.* Identification of type 2 diabetes subgroups through topological analysis of patient similarity. *Sci. Transl. Med.* **7**, 311ra174 (2015).
92. Bonnefond, A. *et al.* Association between large detectable clonal mosaicism and type 2 diabetes with vascular complications. *Nat. Genet.* **45**, 1040–1043 (2013).
93. Qin, J. *et al.* A metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature* **490**, 55–60 (2012).
94. Walford, G. A. *et al.* Metabolite traits and genetic risk provide complementary information for the prediction of future type 2 diabetes. *Diabetes Care* **37**, 2508–2514 (2014).
95. Wurtz, P. *et al.* Circulating metabolite predictors of glycemia in middle-aged men and women. *Diabetes Care* **35**, 1749–1756 (2012).
96. Chambers, J. C. *et al.* Epigenome-wide association of DNA methylation markers in peripheral blood from Indian Asians and Europeans with incident type 2 diabetes: a nested case-control study. *Lancet Diabetes Endocrinol.* **3**, 526–534 (2015).
97. Talmud, P. J. *et al.* Sixty-five common genetic variants and prediction of type 2 diabetes. *Diabetes* **64**, 1830–1840 (2015).
98. Vassy, J. L. *et al.* Polygenic type 2 diabetes prediction at the limit of common variant detection. *Diabetes* **63**, 2172–2182 (2014).
99. Flannick, J. *et al.* Assessing the phenotypic effects in the general population of rare variants in genes for a dominant Mendelian form of diabetes. *Nat. Genet.* **45**, 1380–1385 (2013).
100. Patel, K., Weedon, M. N., Ellard, S., Oram, R. A. & Hattersley, A. T. Type 1 diabetes genetic risk score — a novel tool to differentiate monogenic diabetes from T1D. (Abstract 1746-P) *Diabetes* **64** (Suppl. 1), A453 (2015).
101. Oram, R. A. *et al.* A type 1 diabetes genetic risk score can aid discrimination between type 1 and type 2 diabetes in young adults. *Diabetes Care* **39**, 337–344 (2016).
102. Gaulton, K. J. *et al.* Genetic fine mapping and genomic annotation defines causal mechanisms at type 2 diabetes susceptibility loci. *Nat. Genet.* **47**, 1415–1425 (2015).  
**This is a nice example of the common workflow to localize common variant associations to causal variants and molecular disease mechanisms.**
103. Shea, J. *et al.* Comparing strategies to fine-map the association of common SNPs at chromosome 9p21 with type 2 diabetes and myocardial infarction. *Nat. Genet.* **43**, 801–805 (2011).
104. Yaghootkar, H. *et al.* Association analysis of 29,956 individuals confirms that a low-frequency variant at *CCND2* halves the risk of type 2 diabetes by enhancing insulin secretion. *Diabetes* **64**, 2279–2285 (2015).
105. Maller, J. B. *et al.* Bayesian refinement of association signals for 14 loci in 3 common diseases. *Nat. Genet.* **44**, 1294–1301 (2012).  
**An early fine mapping study is presented here that established the now commonly used ‘credible set’ methodology.**
106. Plenge, R. M., Scolnick, E. M. & Altshuler, D. Validating therapeutic targets through human genetics. *Nat. Rev. Drug Discov.* **12**, 581–594 (2013).
107. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
108. Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–330 (2015).
109. GTEx Consortium. Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* **348**, 648–660 (2015).
110. Maurano, M. T. *et al.* Systematic localization of common disease-associated variation in regulatory DNA. *Science* **337**, 1190–1195 (2012).  
**This is the first prominent and clear demonstration that GWAS variants cluster within regulatory regions of the human genome.**
111. Trynka, G. *et al.* Chromatin marks identify critical cell types for fine mapping complex trait variants. *Nat. Genet.* **45**, 124–130 (2012).
112. Parker, S. C. J. *et al.* Chromatin stretch enhancer states drive cell-specific gene regulation and harbor human disease risk variants. *Proc. Natl Acad. Sci. USA* **110**, 17921–17926 (2013).
113. Pasquali, L. *et al.* Pancreatic islet enhancer clusters enriched in type 2 diabetes risk-associated variants. *Nat. Genet.* **46**, 136–143 (2014).
114. Soccio, Raymond, E. *et al.* Genetic variation determines PPAR $\gamma$  function and anti-diabetic drug response *in vivo*. *Cell* **162**, 33–44 (2015).
115. Gaulton, K. J. *et al.* A map of open chromatin in human pancreatic islets. *Nat. Genet.* **42**, 255–259 (2010).  
**This is one of the earliest studies that used epigenomic annotations to prioritize potentially causal variants at T2D GWAS loci.**
116. Stitzel, M. L. *et al.* Global epigenomic analysis of primary human pancreatic islets provides insights into type 2 diabetes susceptibility loci. *Cell Metab.* **12**, 443–455 (2010).
117. Morán, I. *et al.* Human  $\beta$  cell transcriptome analysis uncovers lncRNAs that are tissue-specific, dynamically regulated, and abnormally expressed in type 2 diabetes. *Cell Metab.* **16**, 435–448 (2012).
118. van de Bunt, M. *et al.* The miRNA profile of human pancreatic islets and Beta-Cells and relationship to type 2 diabetes pathogenesis. *PLoS ONE* **8**, e55272 (2013).
119. Fadista, J. *et al.* Global genomic and transcriptomic analysis of human pancreatic islets reveals novel genes influencing glucose metabolism. *Proc. Natl Acad. Sci. USA* **111**, 13924–13929 (2014).
120. van de Bunt, M. *et al.* Transcript expression data from human islets links regulatory signals from genome-wide association studies for type 2 diabetes and glycaemic traits to their downstream effectors. *PLoS Genet.* **11**, e1005694 (2015).  
**This paper shows how eQTL analysis can be used to hypothesize effector transcripts at GWAS loci.**
121. Nica, A. C. *et al.* Cell-type, allelic, and genetic signatures in the human pancreatic beta cell transcriptome. *Genome Res.* **23**, 1554–1562 (2013).
122. Dayeh, T. *et al.* Genome-wide DNA methylation analysis of human pancreatic islets from type 2 diabetic and non-diabetic donors identifies candidate genes that influence insulin secretion. *PLoS Genet.* **10**, e1004160 (2014).
123. Volkmar, M. *et al.* DNA methylation profiling identifies epigenetic dysregulation in pancreatic islets from type 2 diabetic patients. *EMBO J.* **31**, 1405–1426 (2012).
124. Nilsson, E. *et al.* Altered DNA methylation and differential expression of genes influencing metabolism and inflammation in adipose tissue from subjects with type 2 diabetes. *Diabetes* **63**, 2962–2976 (2014).
125. Yuan, W. *et al.* An integrated epigenomic analysis for type 2 diabetes susceptibility loci in monozygotic twins. *Nat. Commun.* **5**, 5719 (2014).
126. Hodson, D. J. *et al.* ADCY5 couples glucose to insulin secretion in human islets. *Diabetes* **65**, 3009–3021 (2014).
127. Claussnitzer, M. *et al.* FTO obesity variant circuitry and adipocyte browning in humans. *N. Engl. J. Med.* **373**, 895–907 (2015).  
**This is a comprehensive investigation of the disease mechanisms responsible for a GWAS association, using nearly all relevant modern experimental and computation techniques.**
128. Locke, J. M., Hysenaj, G., Wood, A. R., Weedon, M. N. & Harries, L. W. Targeted allelic expression profiling in human islets identifies *cis*-regulatory effects for multiple variants identified by type 2 diabetes genome-wide association studies. *Diabetes* **64**, 1484–1491 (2015).
129. Kulzer, Jennifer, R. *et al.* A common functional regulatory variant at a type 2 diabetes locus upregulates *ARAP1* expression in the pancreatic beta cell. *Am. J. Hum. Genet.* **94**, 186–197 (2014).
130. Ragvin, A. *et al.* Long-range gene regulation links genomic type 2 diabetes and obesity risk regions to *HHEX*, *SOX4*, and *IRX3*. *Proc. Natl Acad. Sci. USA* **107**, 775–780 (2010).
131. Smemo, S. *et al.* Obesity-associated variants within *FTO* form long-range functional connections with *IRX3*. *Nature* **507**, 371–375 (2014).
132. Elbein, Steven, C. *et al.* Genetic risk factors for type 2 diabetes: a *trans*-regulatory genetic architecture? *Am. J. Hum. Genet.* **91**, 466–477 (2012).
133. Small, K. S. *et al.* Identification of an imprinted master *trans* regulator at the *KLF14* locus related to multiple metabolic phenotypes. *Nat. Genet.* **43**, 561–564 (2011).  
**This is a nice example of the power of *trans*-eQTL analysis to identify disease mechanisms responsible for a GWAS association.**
134. Florez, J. C. *et al.* *TCF7L2* polymorphisms and progression to diabetes in the Diabetes Prevention Program. *N. Engl. J. Med.* **355**, 241–250 (2006).
135. Fogarty, M. P., Panhuis, T. M., Vadlamudi, S., Buchkovich, M. L. & Mohlke, K. L. Allele-specific transcriptional activity at type 2 diabetes — associated single nucleotide polymorphisms in regions of pancreatic islet open chromatin at the *JAZF1* locus. *Diabetes* **62**, 1756–1762 (2013).
136. Fogarty, M. P., Cannon, M. E., Vadlamudi, S., Gaulton, K. J. & Mohlke, K. L. Identification of a regulatory variant that binds, FOXA1 and FOXA2 at the *CDC123/CAMK1D* type 2 diabetes, GWAS locus. *PLoS Genet.* **10**, e1004633 (2014).
137. Travers, M. E. *et al.* Insights into the molecular mechanism for type 2 diabetes susceptibility at the *KCNQ1* locus from temporal changes in imprinting status in human islets. *Diabetes* **62**, 987–992 (2013).
138. Claussnitzer, M. *et al.* Leveraging cross-species transcription factor binding site patterns: from diabetes risk loci to disease mechanisms. *Cell* **156**, 343–358 (2014).
139. Wei, F.-Y. *et al.* Deficit of tRNA<sup>lys</sup> modification by Cdkal1 causes the development of type 2 diabetes in mice. *J. Clin. Invest.* **121**, 3598–3608 (2011).
140. McCulloch, L. J. *et al.* *GLUT2 (SLC2A2)* is not the principal glucose transporter in human pancreatic beta cells: implications for understanding genetic association signals at this locus. *Mol. Genet. Metab.* **104**, 648–653 (2011).
141. Beer, N. L. *et al.* The P446L variant in *GCKR* associated with fasting plasma glucose and triglyceride levels exerts its effect through increased glucokinase activity in liver. *Hum. Mol. Genet.* **18**, 4081–4088 (2009).
142. Smith, S. B. *et al.* Rfx6 directs islet formation and insulin production in mice and humans. *Nature* **463**, 775–780 (2010).
143. Zhu, H. *et al.* The *Lin28/let-7* axis regulates glucose metabolism. *Cell* **147**, 81–94 (2011).
144. Dai, N. *et al.* IGF2BP2/IMP2-Deficient mice resist obesity through enhanced translation of *Ucp1* mRNA and other mRNAs encoding mitochondrial proteins. *Cell Metab.* **21**, 609–621 (2015).
145. Mercader, J. M. *et al.* Identification of novel type 2 diabetes candidate genes involved in the crosstalk between the mitochondrial and the insulin signaling systems. *PLoS Genet.* **8**, e1003046 (2012).
146. Taneera, J. *et al.* Expression profiling of cell cycle genes in human pancreatic islets with and without type 2 diabetes. *Mol. Cell. Endocrinol.* **375**, 35–42 (2013).
147. Burns, S. M. *et al.* High-throughput luminescent reporter of insulin secretion for discovering regulators of pancreatic beta-cell function. *Cell Metab.* **21**, 126–137 (2015).
148. The American Diabetes Association. Genetics portal for type 2 diabetes debuts. *Diabetes Dispatch* <http://www.diabetesdispatchextra.org/genetics-portal-for-type-2-diabetes-debuts> (2015).
149. Landrum, M. J. *et al.* ClinVar: public archive of interpretations of clinically relevant variants. *Nucleic Acids Res.* **44**, D862–D868 (2016).
150. Stenson, P. D. *et al.* Human Gene Mutation Database (HGMD): 2003 update. *Hum. Mutat.* **21**, 577–581 (2003).
151. Brookes, A. J. & Robinson, P. N. Human genotype-phenotype databases: aims, challenges and opportunities. *Nat. Rev. Genet.* **16**, 702–715 (2015).

152. 1000 Genomes Project Consortium. An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56–65 (2012).
153. Rosenbloom, K. R. *et al.* The UCSC Genome Browser database: 2015 update. *Nucleic Acids Res.* **43**, D670–D681 (2015).
154. Ward, L. D. & Kellis, M. HaploReg: a resource for exploring chromatin states, conservation, and regulatory motif alterations within sets of genetically linked variants. *Nucleic Acids Res.* **40**, D930–D934 (2012).
155. Exome Aggregation Consortium. Analysis of protein-coding genetic variation in 60,706 humans. *bioRxiv* <http://dx.doi.org/10.1101/030338> (2015).
156. Blankenberg, D. *et al.* in *Current Protocols in Molecular Biology* (eds Ausubel, F. M. *et al.*) (Wiley, 2010).
157. Reardon, S. Pharma firms join NIH on drug development. *Nature* <http://dx.doi.org/10.1038/nature.2014.14672> (2014).
158. Barrett, J. C., Dunham, I. & Birney, E. Using human genetics to make new medicines. *Nat. Rev. Genet.* **16**, 561–562 (2015).
159. Goldman, M. The innovative medicines initiative: a European response to the innovation challenge. *Clin. Pharmacol. Ther.* **91**, 418–425 (2012).
160. Kimball, R., Ross, M., Thornthwaite, W., Mundy, J. & Becker, K. *The Data Warehouse Lifecycle Toolkit* 2nd edn (Wiley, 2008).
161. Espeland, M. A. *et al.* Consent for genetics studies among clinical trial participants: findings from Action for Health in Diabetes (Look AHEAD). *Clin. Trials* **3**, 443–456 (2006).
162. Gusev, A. *et al.* Partitioning heritability of regulatory and cell-type-specific variants across 11 common diseases. *Am. J. Hum. Genet.* **95**, 535–552 (2014).
163. Welter, D. *et al.* The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res.* **42**, D1001–D1006 (2014).
164. Pal, A. *et al.* Loss-of-function mutations in the cell-cycle control gene *CDKN2A* impact on glucose homeostasis in humans. *Diabetes* **65**, 527–533 (2015).
165. Torekov, S. S. *et al.* KCNQ1 long QT syndrome patients have hyperinsulinemia and symptomatic hypoglycemia. *Diabetes* **63**, 1315–1325 (2014).

#### Acknowledgements

The authors thank B. Alexander for help with figure design and creation, as well as helpful discussions.

#### Competing interests statement

The authors declare no competing interests.

#### FURTHER INFORMATION

GenomeSpace: <http://www.genomespace.org>

Global Alliance for Genomics and Health (GA4GH):

<http://genomicsandhealth.org>

GTEx Portal: <http://gtexportal.org/home>

Michigan Imputation Server: <https://imputationserver.sph.umich.edu>

Sanger Imputation Service: <https://imputation.sanger.ac.uk>

Type 2 Diabetes Knowledge Portal: <http://www.type2diabetesgenetics.org>

WashU Epigenome Browser: <http://epigenomegateway.wustl.edu>

xBrowse: <https://atgu.mgh.harvard.edu/xbrowse>

ALL LINKS ARE ACTIVE IN THE ONLINE PDF