1  **Review**

2  **Strengthening Causal Inference for Complex Disease Using Molecular**
3  **Quantitative Trait Loci**

4  **Sonja Neumeyer[1], Gibran Hemani[2], Eleftheria Zeggini[1*]**

5

6  [1]Institute of Translational Genomics, Helmholtz Zentrum München, German
7  Research Center for Environmental Health, Neuherberg, Germany.

8  [2]MRC Integrative Epidemiology Unit (IEU), Population Health Sciences, Bristol
9  Medical School, University of Bristol, Bristol, United Kingdom.

10

11  **\*Correspondence:** eleftheria.zeggini@helmholtz-muenchen.de (E. Zeggini).

12

15

16  **Abstract**

17  Large genome-wide association studies have identified loci associated with complex
18  traits and diseases, but often index variants are not causal and reside in non-coding
19  regions of the genome. To gain a better understanding of the relevant biological
20  mechanisms, intermediate traits such as gene expression or protein levels are
21  increasingly being investigated, as these are likely mediators between genetic variants
22  and disease outcome. Genetic variants associated with intermediate traits, termed
23  molecular quantitative trait loci (molQTLs), can then be used as instrumental variables
24  in a Mendelian randomization approach to identify causal features and mechanisms of
25  complex traits. Challenges such as pleiotropy and non-specificity of molQTLs remain
26  and further approaches and methods need to be developed.

27

28

29

## Genome-Wide Association Studies

Genome-wide association studies **(GWAS, Box 1)** have identified thousands of sequence variants that contribute to the genetic architecture of complex diseases and medically-relevant quantitative traits. This endeavour has been fuelled by two major ambitions: creating genetic predictors for disease; and identifying the genomic regions responsible for the disease to gain a better understanding of the relevant biological mechanisms [1, 2]. The latter objective is the focus of this review.

Typically, associated variants individually account for a very small proportion of phenotypic variation. This is common for quantitative or "complex" traits which are usually influenced by a large number of genes with small effects on the trait [3]. There is no simple Mendelian inheritance pattern but random sampling of alleles at each associated gene results in a normally distributed phenotype in the population [4]. Functional information on the underlying mechanisms of genetic variants identified by GWAS is often unclear, i.e. it is challenging to identify effector genes based on the observed association summary statistics only [3, 5]. The majority of complex trait variants reside in noncoding regions of the genome [6, 7] and it is possible that they confer their effect through modulating gene expression levels [8]. In their second decade of existence, GWAS are showing signs of maturity, with increasing diversity in populations studied [9], inclusion of low frequency and rare variants, and finer definition of phenotypic traits examined.

In this review we will describe how molecular traits are also being assayed and analysed for genetic associations, and how the understanding of complex disease aetiology is improving through combining genetic analysis of both the disease and molecular traits. The presiding manner in which these relationships are constructed is using a causal inference method known as Mendelian randomization (MR) which capitalizes on the abundance of GWAS results now available. We will describe MR in terms of both its current implementation and the future developments that are needed to address known limitations.

## Molecular Quantitative Trait Loci

60  The influence of a genetic variant associated with a disease is likely to be mediated
61  via molecular traits (Figure 1), which themselves are often complex. Quantitative
62  molecular traits, such as gene expression or protein abundance, are frequently
63  dysregulated in disease and can act as intermediate phenotypes, affording greater
64  power to detect association compared to the dichotomous definition of a disease
65  endpoint, which is the culmination of multiple biological processes being perturbed
66  [10].

67  Multiple studies have investigated mRNA levels combined with genome-wide genotype
68  information to identify expression quantitative trait loci (eQTLs), i.e. genetic variants
69  associated with gene expression levels [11]. The first studies to investigate molecular
70  quantitative trait loci (molQTLs) started out with small sample sizes. Due to challenges
71  associated with collecting human biospecimens using invasive procedures, analyses
72  initially focussed on using the most accessible tissues [12]. Today, sample sizes used
73  for molQTL investigation in blood have grown very large [13]. MolQTLs are generally
74  classified into cis-acting, which is typically defined as regulation of genes within 1Mb,
75  or trans-acting, defined as molQTLs affecting genes further away or on different
76  chromosomes [14]. Whereas detected cis-effects have generally been large and easily
77  found using small sample sizes, trans effects tend to be much smaller and larger
78  sample sizes are required. Large studies such as the eQTLGen Consortium [13] or
79  GoDMC (http://www.godmc.org.uk/) are emerging to identify these small effects that
80  might play central roles in disease etiology. Molecular trait loci seem to be highly tissue
81  dependent [15, 16]. However, tissue-sharing of cis-eQTLs seems to be bimodal. Either
82  cis-eQTLs seem to be shared across many tissues or they are very specific to only a
83  small subset of tissues [17].To provide a resource which enables the systematic study
84  of genetic variation on regulation of gene expression in multiple human tissues, the
85  Genotype-Tissue Expression (GTEx) project was initiated a decade ago [18]. The
86  current GTEx release provides a total of 11688 samples and 53 tissues across 714
87  donors (current release V7, dbGaP accession phs000424.v7.p2). Sample sizes of
88  other studies have also largely increased [19-21] and a variety of tissues have been
89  studied. The picture is far from complete, but has been massively enhanced since the
90  inception of these studies.

The first expression phenotypes to be studied were gene transcript levels. They are highly heritable [22]. It is estimated that around 88% of all genes have at least one eQTL [13]. To date, many different molecular traits with a potential influence on gene regulation have been investigated [23]. They range from influencing the epigenome such as DNA methylation (meQTL), histone modification (hQTL) or chromatin accessibility (caQTL) to alternative splicing (sQTL), protein levels (pQTL), microRNA expression (mirQTL) or ribosome occupancy (rQTL) [23]. In addition, higher level intermediate phenotypes such as metabolites have been investigated and QTLs for metabolites such as carbohydrates, amino acids or fatty acids identified [24].

In an effort to find the molecular pathways that connect genetic variants to complex traits, overlapping/colocalisation methods between GWAS and molQTL signals have been developed. Colocalisation of an eQTL with a GWAS signal suggests that the eQTL target gene could be involved in the molecular pathway of the complex disease under investigation [25]. Several studies already discovered GWAS signals enriched for molQTLs in a tissue dependent-manner [26]. For example, the myocardial infarction and high LDL cholesterol-associated **1p13 locus (see Glossary)** had been fine mapped to the *CELSR2* gene. Using eQTL analyses, it was discovered that actually the expression of *SORT1* was influenced by this variant [27].

MolQTLs are being used as instrumental variables for molecular traits in a variety of ways: to infer the relative importance of different classes of molecular features on variation in complex traits; to identify the causal gene for a particular complex trait [23]; to identify the causal tissue for a complex trait [28] and to estimate causal relationships between different molecular traits [29]. In this review, we will focus on their use for identifying causal features of complex traits.

**Mendelian Randomization Studies Strengthen Causal Inference**

Mendelian randomization (**MR, Box 2**) studies use genetic variants as proxies for modifiable risk factors to test whether the risk factor is causally relevant to an outcome of interest [30, 31]. The advantage of such an approach is that unmeasured confounding, an issue of observational studies, and reverse causation can be minimized. It is, therefore, possible to use genetic information to draw causal inferences.

122

123    Early MR studies mainly used one-sample approaches, where the exposure and
124    outcome phenotypes along with the genetic variants that were being used to
125    instrument the exposure were available for all samples in a single dataset. Nowadays,
126    when many large-scale GWASs are conducted, it is much more powerful to use
127    published **SNP (single nucleotide polymorphism)** -trait associations from large
128    consortia. It is, therefore, common to use two-sample MR approaches where SNP-
129    exposure and SNP-outcome associations are estimated in different studies and
130    subsequently combined [32]. When using genome-wide significant SNPs as
131    instrumental variable for an exposure, the first MR assumption should be verified.

132

133    For two-sample MR methods, only summary statistics are required (per allele
134    regression coefficients, standard errors and effect allele) which are typically obtained
135    from published GWAS of the largest possible datasets [33]. The causal effect can be
136    estimated using the Wald ratio estimate, which is the ratio of SNP-outcome association
137    and SNP-exposure association.

138

139    SNP-exposure and SNP-outcome association statistics should ideally be obtained
140    from studies of non-overlapping individuals (two-sample MR). When using summary
141    statistics from only one sample or from partially overlapping samples, results might be
142    biased in the direction of the observational estimate, especially if the genetic effects
143    on the exposure are weak [34]. When several independent genetic variants are known
144    to be associated with the exposure of interest, these can be combined into a single MR
145    estimate using inverse variance weighted meta-analysis of the single Wald ratio
146    estimates [32]. In doing so, the MR framework can then be viewed as a meta-analysis
147    problem which itself has a rich set of tools to evaluate and correct for bias [35]. One
148    issue that has been of particular concern in MR is in proving that violation of the third
149    assumption, i.e. that the genetic instrument influences the outcome only through the
150    exposure, does not induce bias [36]. A suite of sensitivity analyses [37-41] are now
151    routinely implemented in MR studies that use multiple independent instruments to
152    model pleiotropy [42].

**Mendelian Randomization Studies Using Molecular QTLs as Instrumental Variables**

Whole genome approaches have indicated that the causal variants influencing complex traits are overrepresented by those that are also associated with eQTLs [43, 44]. This supports the notion that disease biology could be unravelled by mapping the causal path from genetic variant through the use of intermediate molQTLs [45]. At its most basic implementation, a Mendelian randomization framework for evaluating the causal influence of a molecular trait on a complex trait would be to test if a known molQTL is also associated with the complex trait (Key Figure, Figure 2). The Wald ratio of SNP-complex trait and SNP-molecular trait effects can then be obtained as an estimate of the causal effect. This simple method suffers from a number of potential pitfalls and is often performed as an initial screen to find, from amongst many molecular phenotypes (e.g. hundreds of thousands of DNA methylation levels), a few putative causal molecular phenotypes for more detailed follow up and sensitivity analysis [46-48]. Here we describe some of these approaches.

*Linkage disequilibrium links a causal variant for one trait with a different causal variant for another trait.*

A major lesson from GWAS is that complex traits follow a polygenic architecture [49, 50]. As a consequence, finding that a chosen SNP happens to show an association with a complex trait might not be surprising because many non-causal common variants are likely to be in **linkage disequilibrium (LD)** with a causal variant for a complex trait (Figure 3a). Colocalisation techniques seek to analyse specific genomic regions, determining whether the pattern of test statistics for one trait are concordant with the pattern from another, often with respect to the underlying LD structure. Evidence for shared causal variants at a locus is determined by the extent to which the test statistic patterns are shared between the two traits. An important recent finding is that the majority of genes that colocalise with a trait are not the genes that are closest to the biggest signal for the trait [11].

Typically, the proportion of overlapping signals between molecular and complex traits that appear to be due to LD is high. For example in [29] it was shown that two thirds of putative expression-trait MR relationships were due to LD, with a similar proportion

184  being found for DNA methylation-trait MR relationships. Nevertheless, when assessed

185  across hundreds of complex traits, there are now tens of thousands of examples of

186  colocalisation between gene expression levels and complex traits [51]. It remains

187  important to note that there are many colocalisation techniques [11, 52-54] and there

188  is not always strong agreement between them [54].

189  *The association is reverse causal*

190  One of the purported advantages of MR is that it protects against reverse causation.

191  This is true to the extent that the instrument is known to primarily influence the

192  hypothesised exposure. However it is conceivable that a molQTL arises because a

193  complex trait influences it. Mediation-based methods exist that require individual-level

194  data to orient the causal direction [55-57], but are susceptible to making the wrong

195  orientation under specific patterns of confounding or measurement error [58]. An

196  alternative approach is to perform MR in the reverse direction [47], identifying SNPs

197  that instrument the complex trait and testing for its association on the molecular trait.

198  Typically however, one would not expect reverse causal relationships to explain a

199  molQTL associated with a complex trait because in order for the molQTL to have been

200  detected in a small sample size it will necessarily be a large effect, which is impossible

201  if it were mediated through a polygenic trait [29].

202  *The instrumenting SNP is non-specific to the hypothesised exposure*

203  Often a single SNP is detected as an instrument for multiple molecular phenotypes.

204  For example, a SNP could be strongly associated with more than one gene expression

205  level, or the same gene expression level in different tissues or time points, or both a

206  gene expression level and a DNA methylation level (Figure 3). This is not necessarily

207  a problem, as all the molecular phenotypes that are associated with the trait could be

208  on the same causal pathway to the disease, and indeed it could be advantageous as

209  it presents us with multiple points of intervention. Non-specificity of genetic

210  associations is classically known as pleiotropy though care should be taken in using

211  the term. MR assumes a 'vertically' pleiotropic relationship, where the genetic

212  instrument is associated with the outcome because it is mediated by the exposure. By

213  contrast, 'horizontal' pleiotropy is a source of problems in MR, inducing bias or false

214  causal inference if the SNP influences the outcome through a pathway other than the

215 hypothesised exposure [59]. Proving that a putative MR finding is due to vertical and
216 not horizontal pleiotropy is far from trivial [36].

217 There are vastly more molecular phenotypes than independent genetic regions,
218 especially when temporal- and tissue-specific measurements are possible [60]. By
219 definition it is expected that many molQTL will not be specific to a particular molecular
220 trait. Therefore, it is difficult to prove which, from amongst the set of molecular traits
221 that are influenced by the molQTL, is the causal factor [51].

222 One approach is to focus on the use of cis-acting molQTLs, with the rationale that they
223 are biologically 'closer' to the intended molecular trait. Trans-acting QTLs are likely to
224 only influence the molecular trait because they are mediated by other molecular traits,
225 opening up a greater possibility that the instrument is non-specific to the intended
226 target (Figure 3b). Testing explicitly if the molQTL is associated with other molecular
227 traits is also sensible, as this can be used to (de-)prioritise a putative association
228 depending on how much evidence there is for (non-)specificity [2]. Methods are now
229 arising that attempt to model the MR estimates of multiple molecular exposures
230 simultaneously, thereby adjusting for potential horizontal pleiotropy [61]. While a useful
231 tool, interpretation remains difficult as the use of multivariable MR [62] requires that
232 there are marked differences in the genetic signatures across the exposures [63]. It
233 also requires measurement of all possible exposures that could be inducing the
234 pleiotropy, which is a similar assumption to observational study designs that prompted
235 the development of MR in the first place.

236 There are more standard MR sensitivity analyses that can be applied in the event that
237 multiple independent causal variants are available [42]. However, this typically requires
238 introducing trans-QTLs into the analysis which may not bring clarity, as they could have
239 systematically different properties to cis-QTLs. At this stage, if a molecular trait
240 colocalises with a complex trait, and doesn't appear to be reverse causal, it is still
241 extremely difficult to prove that it is causal and not simply one of many traits that are
242 all influenced by the same molQTL.

243 In the GoDMC study, which used 30k samples to discover instruments for DNA
244 methylation levels, multiple cis and trans instruments were used to model causal
245 relationships between DNA methylation levels and complex traits. It was found that,

8

246 while there were many putative colocalising signals with complex traits, there was
247 almost no agreement between the causal effect estimated using primary and
248 secondary molQTLs, implying that the majority of colocalising signals were due to
249 horizontal pleiotropy.

**Current Challenges and Issues**

251 The prospects of finding new drug targets has propelled forwards the data acquisition
252 and methodological development for mapping the pathways between molecular and
253 complex traits.

254 Genetic variation is finite, and though molecular traits are often polygenic the use of
255 more than the cis-region for instrumentation is currently not fully understood. This
256 incurs a limit on the extent to which current tools designed to protect against incorrect
257 causal inference due horizontal pleiotropy can be used. Conceptually, here we use
258 genetic instruments as a proxy for molecular phenotypes. However, molecular
259 phenotypic variation dwarfs the cis-genetic resource that is available for
260 instrumentation. Hence, the ubiquitous non-specificity of any molQTL makes it very
261 difficult to determine which molecular feature is actually mediating the genetic effect
262 on a trait. This could be because inference is for the wrong developmental time point
263 (e.g. genetic effects are very consistent over time [64] for DNA methylation) or the
264 wrong tissue (cis-QTLs are strongly shared across tissues [17]). Alternatively, it could
265 be that it was an entirely different molecular feature (e.g. gene expression, DNA
266 methylation and histone variation often share similar cis-regulatory features [65]).

267 Coupled with this problem of non-specificity, is the emerging evidence supporting a
268 model of ubiquitous horizontal pleiotropy [40, 66], in which any particular genetic
269 variant potentially influences a particular complex trait through multiple independent
270 pathways. The omnigenic model offers an extreme viewpoint on this problem, in which
271 polygenic architecture arises because every gene is related to every trait through an
272 underlying dense gene regulatory network [3].

273 Making meaningful inference from such an under-specified model requires a departure
274 from current practices of treating molecular features singly, and reliably incorporating
275 trans-instruments, which may exhibit tissue specificity [17]. Though any one instrument

276 might be non-specific, it is seldom the case that the genetic correlation of complex

277 traits is 1 [67], meaning that there are potentially combinations of instruments that

278 together provide some specificity. Large-scale pleiotropy maps are beginning to be

279 produced [40, 68, 69], and may provide an avenue into constructing instrument

280 combinations conditional on a background of complex pleiotropy.

**Concluding Remarks**

281 

282 Many genetic variants associated with complex traits and diseases have been

283 discovered, but often there is a lack of knowledge about mechanisms involved **(see**

284 **Clinician`s Corner)**. Investigation of intermediate traits and associated molQTLs has

285 been very helpful, as these better explain how genetic variants influence complex

286 traits. Using molQTLs combined with an MR approach, causal features of a complex

287 trait can be revealed. Challenges, such as the model of ubiquitous horizontal

288 pleiotropy and, therefore, a non-specificity of molQTLs to a particular molecular trait,

289 remain **(see Outstanding Questions)**. Therefore, new methods need to be

290 developed, including for example those that reliably incorporate trans-molQTLs,

291 which have a greater possibility for non-specificity of the instrument.

292 

293 Despite our growing understanding of the limitations of MR, the current data resources

294 and statistical frameworks for MR can be viewed as a resource with tremendous utility.

295 Most directly, using MR to support a negative association could be less prone to some

296 of the issues described. Of growing importance in causal inference is the concept of

297 triangulation, where information from orthogonal experimental designs are integrated

298 together to obtain a more reliable conclusion [70]. There are now open source data

299 and software repositories (including those that can be used in web browsers [42]) that

300 automate MR analyses. The inclusion of genetic evidence through MR should be a

301 natural part of any causal inquiry [71].

302

304      **BOX 1: Genome-wide association studies**

305      *Genome-wide association studies (GWAS) compare large numbers*
306      *of affected with unaffected individuals to identify sequence variants that*
307      *are associated with risk of complex diseases, or at the population-level*
308      *to identify associations with quantitative traits. The foundation for GWAS*
309      *was laid by the sequencing of the human genome [72], characterization*
310      *of the correlation patterns between pairs of variants genome-wide [73],*
311      *development of high-throughput genotyping platforms, and the*
312      *availability of large-scale sample sizes. Millions of single nucleotide*
313      *polymorphisms (SNPs) have been mapped [74]. For several reasons it*
314      *has been difficult to elucidate the underlying mechanism between*
315      *associated genetic variant and disease trait. One reason is the co-*
316      *inheritance of many genetic variants with the disease-associated variant*
317      *(linkage disequilibrium (LD)) [75]. Due to this complicated correlation*
318      *structure of human genome, the most strongly associated GWAS signal*
319      *(index variant) is often not causal [76]. Similarly, compounded by*
320      *complex regulatory mechanisms, the nearest gene to the top GWAS*
321      *signal is not necessarily the causal gene [11].*

323      **Box 2: Mendelian randomization studies**
324      *Due to the laws of Mendelian inheritance, alleles are assigned at*
325      *conception to individuals independent of environmental risk factors and*
326      *confounders. To obtain valid estimates using Mendelian randomization*
327      *(MR), three assumptions have to be met: firstly, the genetic variants need*
328      *to be sufficiently associated with the exposure of interest; secondly, the*
329      *genetic variants should not be associated to any confounder of the risk*
330      *factor – outcome relationship; finally there should not be any other*
331      *pathway from genetic variants to outcome except through the exposure*
332      *of interest. Except for the first assumption, which can be tested, the other*
333      *two assumptions can only be addressed by sensitivity analyses [77].*

11

*Clinician`s Corner*

- *Poor efficacy and poor safety are the two major reasons for the very high failure rate of drug trials, ultimately driving up the cost of drugs and their development times. This can be partly framed as a causal inference problem, where the objective is to identify which molecular targets are causal for the disease of interest and filter out those that are likely to fail prior to initiating trials.*

- *Randomized controlled trials (RCTs) are ideal for making causal inference but are expensive, slow and often impracticable for a particular causal enquiry. The Mendelian randomization statistical framework leverages genetic associations to mimic randomized control trials. The potential of this strategy is increasingly being exploited due to the ready availability of data to quickly and cheaply evaluate the causal importance for thousands of molecular features on complex diseases.*

- *To interrogate the causal influence of a particular molecular trait on a particular disease, knowledge of robust genetic factors for the molecular trait, and the corresponding effect of those factors on the disease, are both required. Thanks to over a decade of genome-wide association studies and the recent emergence of national genetic biobanks, most complex diseases have genome-wide genetic associations from large sample sizes made publicly available. In addition, the genetic influences on a range of molecular features such as protein levels, gene expression levels, DNA methylation levels, metabolites etc are being mapped and made publicly available.*

- *Though it is impossible to mimic an RCT perfectly using such observational data, statistical techniques and data continue to improve, and Mendelian randomization is poised to further help make causal claims about a molecular trait on complex disease.*

**Glossary**

**1p13 locus:** GWAS analysis in humans demonstrated that this locus on chromosome 1 is strongly associated with plasma low-density lipoprotein cholesterol (LDL-C) levels, which in turn is a major risk factor for myocardial infarction. SNPs (see below) in this locus have also been linked to coronary artery disease. This locus alters the expression of SORT1 (see below) in the liver.

*CELSR2:* Cadherin EGF LAG seven-pass G-type receptor 2, a receptor with possible role in cell/cell signaling during nervous system formation. *CELSR2* is physically linked to the 1p13 locus. Because of this, CELSR2 expression was thought to be controlled by the 1p13 locus until eQTL analysis showed that this was not the case.

**LD:** linkage disequilibrium, the non-random association of alleles at different loci. Based on the assumption that over time recombination events will result in a random association of alleles at two loci, linkage disequilibrium is defined as the difference between the observed frequency of a particular combination of alleles at two loci compared to the frequency expected at random. When analyzing causal SNPs in GWAS analysis, special care must be taken to not wrongly interpret a non-causal SNP that is in LD with a causal SNP.

**SNP:** single nucleotide polymorphism, a DNA sequence variant within a population. SNPs can be linked to disease development and response to pathogens or medication in humans, which makes them invaluable in personalized medicine. Comparison of SNP composition in genomic regions between different cohorts (e.g. with and without disease) is of great importance in biomedical research on a larger scale (e.g. GWAS).

*SORT1:* Sortilin, which is localized in intracellular compartments, notably the Golgi apparatus. It is involved in endocytosis and functions as a sorting receptor in the Golgi compartment and clearance receptor on the cell surface. SORT1 expression is modulated by the 1p13 locus (see above). In liver cells of mouse models, LDL-C levels are significantly decreased by SORT1 overexpression whereas SORT1 knockdown resulted in an increase of LDL-C levels.

## Acknowledgements

**References**

1. Nelson, M.R. et al. (2015) The support of human genetic evidence for approved drug indications. Nat. Genet. 47, 856-860.

2. Zheng, J. et al. (2019) Phenome-wide Mendelian randomization mapping the influence of the plasma proteome on complex diseases. bioRxiv, 627398.

3. Boyle, E.A. et al. (2017) An expanded view of complex traits: From polygenic to omnigenic. Cell 169, 1177-1186.

4. Fisher, R.A. (1919) XV.—The correlation between relatives on the supposition of Mendelian inheritance. Trans. R. Soc. Edinburgh 52, 399-433.

5. Ruiz-Narváez, E.A. (2011) What is a functional locus? Understanding the genetic basis of complex phenotypic traits. Med. Hypotheses 76, 638-642.

6. Edwards, Stacey L. et al. (2013) Beyond GWASs: Illuminating the dark road from association to function. Am. J. Hum. Genet. 93, 779-797.

7. Farh, K.K.-H. et al. (2014) Genetic and epigenetic fine mapping of causal autoimmune disease variants. Nature 518, 337-343.

8. Maurano, M.T. et al. (2012) Systematic localization of common disease-associated variation in regulatory DNA. Science 337, 1190-1195.

9. Gurdasani, D. et al. (2019) Genomics of disease risk in globally diverse populations. Nat. Rev. Genet., doi: 10.1038/s41576-019-0144-0.

10. Civelek, M. and Lusis, A.J. (2013) Systems genetics approaches to understand complex traits. Nat. Rev. Genet. 15, 34-48.

11. Zhu, Z. et al. (2016) Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. Nat. Genet. 48, 481-487.

12. Emilsson, V. et al. (2008) Genetics of gene expression and its effect on disease. Nature 452, 423-428.

13. Võsa, U. et al. (2018) Unraveling the polygenic architecture of complex traits using blood eQTL metaanalysis. bioRxiv, 447367.

14. Westra, H.-J. et al. (2013) Systematic identification of trans eQTLs as putative drivers of known disease associations. Nat. Genet. 45, 1238-1243.

15. Brown, C.D. et al. (2013) Integrative modeling of eQTLs and cis-regulatory elements suggests mechanisms underlying cell type specificity of eQTLs. PLoS Genet. 9, e1003649.

429 16. Sun, B.B. et al. (2018) Genomic atlas of the human plasma proteome. Nature

430 558, 73-79.

431 17. The GTEx Consortium et al. (2017) Genetic effects on gene expression across

432 human tissues. Nature 550, 204-213.

433 18. Lonsdale, J. et al. (2013) The Genotype-Tissue Expression (GTEx) project. Nat.

434 Genet. 45, 580-585.

435 19. Hannon, E. et al. (2018) Leveraging DNA-methylation quantitative-trait loci to

436 characterize the relationship between methylomic variation, gene expression, and

437 complex traits. Am. J. Hum. Genet. 103, 654-665.

438 20. Yao, C. et al. (2018) Genome-wide mapping of plasma protein QTLs identifies

439 putatively causal genes and pathways for cardiovascular disease. Nat. Commun. 9,

440 3268.

441 21. Taylor, K. et al. (2019) Prioritizing putative influential genes in cardiovascular

442 disease susceptibility by applying tissue-specific Mendelian randomization. Genome

443 Med. 11, 6.

444 22. Morley, M. et al. (2004) Genetic analysis of genome-wide variation in human

445 gene expression. Nature 430, 743-747.

446 23. Vandiedonck, C. (2018) Genetic association of molecular traits: A help to identify

447 causative variants in complex diseases. Clin. Genet. 93, 520-532.

448 24. Kastenmüller, G. et al. (2015) Genetics of human metabolism: an update. Hum.

449 Mol. Genet. 24, R93-R101.

450 25. Wen, X. et al. (2017) Integrating molecular QTL data into genome-wide genetic

451 association analysis: Probabilistic assessment of enrichment and colocalization.

452 PLoS Genet. 13, e1006646.

453 26. Nicolae, D.L. et al. (2010) Trait-associated SNPs are more likely to be eQTLs:

454 Annotation to enhance discovery from GWAS. PLoS Genet. 6, e1000888.

455 27. Musunuru, K. et al. (2010) From noncoding variant to phenotype via *SORT1* at

456 the 1p13 cholesterol locus. Nature 466, 714-719.

457 28. Ongen, H. et al. (2017) Estimating the causal tissues for complex traits and

458 diseases. Nat. Genet. 49, 1676-1683.

459 29. Taylor, D.L. et al. (2019) Integrative analysis of gene expression, DNA

460 methylation, physiological traits, and genetic variation in human skeletal muscle.

461 Proc. Natl. Acad. Sci. 116, 10883-10888.

462  30. Davey Smith, G. and Ebrahim, S. (2003) 'Mendelian randomization': can genetic

463  epidemiology contribute to understanding environmental determinants of disease?

464  Int. J. Epidemiol. 32, 1-22.

465  31. Ebrahim, S. and Davey Smith, G. (2008) Mendelian randomization: can genetic

466  epidemiology help redress the failures of observational epidemiology? Hum. Genet.

467  123, 15-33.

468  32. Pierce, B.L. and Burgess, S. (2013) Efficient design for Mendelian randomization

469  studies: subsample and 2-sample instrumental variable estimators. Am. J. Epidemiol

470  178, 1177-1184.

471  33. Hartwig, F.P. et al. (2017) Two-sample Mendelian randomization: avoiding the

472  downsides of a powerful, widely applicable but potentially fallible technique. Int. J.

473  Epidemiol. 45, 1717-1726.

474  34. Burgess, S. et al. (2016) Bias due to participant overlap in two-sample Mendelian

475  randomization. Genet. Epidemiol. 40, 597-608.

476  35. Bowden, J. and Holmes, M.V. (2019) Meta-analysis and Mendelian

477  randomization: A review. Res. Synth. Methods, doi: 10.1002/jrsm.1346.

478  36. Hemani, G. et al. (2018) Evaluating the potential role of pleiotropy in Mendelian

479  randomization studies. Hum. Mol. Genet. 27, R195-R208.

480  37. Bowden, J. et al. (2015) Mendelian randomization with invalid instruments: effect

481  estimation and bias detection through Egger regression. Int. J. Epidemiol. 44, 512-

482  525.

483  38. Bowden, J. et al. (2016) Consistent estimation in Mendelian randomization with

484  some invalid instruments using a weighted median estimator. Genet. Epidemiol. 40,

485  304-314.

486  39. Hartwig, F.P. et al. (2017) Robust inference in summary data Mendelian

487  randomization via the zero modal pleiotropy assumption. Int. J. Epidemiol. 46, 1985-

488  1998.

489  40. Hemani, G. et al. (2017) Automating Mendelian randomization through machine

490  learning to construct a putative causal map of the human phenome. bioRxiv, 173682.

491  41. Zhu, Z. et al. (2018) Causal associations between risk factors and common

492  diseases inferred from GWAS summary data. Nat. Commun. 9, 224.

493  42. Hemani, G. et al. (2018) The MR-Base platform supports systematic causal

494  inference across the human phenome. eLife 7, e34408.

495 43. Gusev, A. et al. (2014) Partitioning heritability of regulatory and cell-type-specific

496 variants across 11 common diseases. Am. J. Hum. Genet. 95, 535-552.

497 44. Banos, D.T. et al. (2018) Bayesian reassessment of the epigenetic architecture of

498 complex traits. bioRxiv, 450288.

499 45. Schadt, E.E. et al. (2003) Genetics of gene expression surveyed in maize, mouse

500 and man. Nature 422, 297-302.

501 46. Hannon, E. et al. (2015) Methylation QTLs in the developing brain and their

502 enrichment in schizophrenia risk loci. Nat. Neurosci. 19, 48-54.

503 47. Richardson, T.G. et al. (2017) Mendelian randomization analysis identifies CpG

504 sites as putative mediators for genetic influences on cardiovascular disease risk. Am.

505 J. Hum. Genet. 101, 590-602.

506 48. Richardson, T.G. et al. (2018) Systematic Mendelian randomization framework

507 elucidates hundreds of CpG sites which may mediate the influence of genetic

508 variants on disease. Hum. Mol. Genet. 27, 3293-3304.

509 49. The International Schizophrenia Consortium et al. (2009) Common polygenic

510 variation contributes to risk of schizophrenia and bipolar disorder. Nature 460, 748-

511 752.

512 50. Yang, J. et al. (2010) Common SNPs explain a large proportion of the heritability

513 for human height. Nat. Genet. 42, 565-569.

514 51. Richardson, T.G. et al. (2019) A transcriptome-wide Mendelian randomization

515 study to uncover tissue-dependent regulatory mechanisms across the human

516 phenome. bioRxiv, 563379.

517 52. Giambartolomei, C. et al. (2014) Bayesian test for colocalisation between pairs of

518 genetic association studies using summary statistics. PLoS Genet. 10, e1004383.

519 53. Hormozdiari, F. et al. (2016) Colocalization of GWAS and eQTL signals detects

520 target genes. Am. J. Hum. Genet. 99, 1245-1260.

521 54. Barbeira, A.N. et al. (2018) Exploring the phenotypic consequences of tissue

522 specific gene expression variation inferred from GWAS summary statistics. Nat.

523 Commun. 9, 1825.

524 55. Aten, J.E. et al. (2008) Using genetic markers to orient the edges in quantitative

525 trait networks: The NEO software. BMC Syst. Biol. 2, 34.

526 56. Millstein, J. et al. (2009) Disentangling molecular relationships with a causal

527 inference test. BMC Genet. 10, 23.

528    57. Waszak, Sebastian M. et al. (2015) Population variation and genetic control of

529    modular chromatin architecture in humans. Cell 162, 1039-1050.

530    58. Hemani, G. et al. (2017) Orienting the causal relationship between imprecisely

531    measured traits using GWAS summary data. PLoS Genet. 13, e1007081.

532    59. Davey Smith, G. and Hemani, G. (2014) Mendelian randomization: genetic

533    anchors for causal inference in epidemiological studies. Hum. Mol. Genet. 23, R89-

534    R98.

535    60. Houle, D. et al. (2010) Phenomics: the next challenge. Nat. Rev. Genet. 11, 855-

536    866.

537    61. Porcu, E. et al. (2019) Mendelian Randomization integrating GWAS and eQTL

538    data reveals genetic determinants of complex and clinical traits. bioRxiv, 377267.

539    62. Burgess, S. and Thompson, S.G. (2015) Multivariable Mendelian randomization:

540    The use of pleiotropic genetic variants to estimate causal effects. Am. J. Epidemiol

541    181, 251-260.

542    63. Sanderson, E. et al. (2018) An examination of multivariable Mendelian

543    randomization in the single-sample and two-sample summary data settings. Int. J.

544    Epidemiol., doi: 10.1093/ije/dyy262.

545    64. Gaunt, T.R. et al. (2016) Systematic identification of genetic influences on

546    methylation across the human life course. Genome Biol. 17, 61.

547    65. Grubert, F. et al. (2015) Genetic control of chromatin states in humans involves

548    local and distal chromosomal interactions. Cell 162, 1051-1065.

549    66. Jordan, D.M. et al. (2019) The landscape of pervasive horizontal pleiotropy in

550    human genetic variation is driven by extreme polygenicity of human traits and

551    diseases. bioRxiv, 311332.

552    67. Bulik-Sullivan, B. et al. (2015) An atlas of genetic correlations across human

553    diseases and traits. Nat. Genet. 47, 1236-1241.

554    68. Cortes, A. et al. (2017) Bayesian analysis of genetic association across tree-

555    structured routine healthcare data in the UK Biobank. Nat. Genet. 49, 1311-1318.

556    69. Cho, Y. et al. (2018) MR-TRYX: Exploiting horizontal pleiotropy to infer novel

557    causal pathways. bioRxiv, 476085.

558    70. Lawlor, D.A. et al. (2017) Triangulation in aetiological epidemiology. Int. J.

559    Epidemiol. 45, 1866-1886.

560    71. Walker, V. et al. (2019) Using the MR-Base platform to investigate risk factors

561    and drug targets for thousands of phenotypes [version 1; peer review: 2 approved].

562    Wellcome Open Research 4.

563    72. Lander, E.S. (2011) Initial impact of the sequencing of the human genome.

564    Nature 470, 187-197.

565    73. Thorisson, G.A. et al. (2005) The International HapMap Project Web site.

566    Genome Res. 15, 1592-1593.

567    74. Wang, D.G. et al. (1998) Large-scale identification, mapping, and genotyping of

568    single-nucleotide polymorphisms in the Human genome. Science 280, 1077-1082.

569    75. Altshuler, D. et al. (2005) A haplotype map of the human genome. Nature 437,

570    1299-1320.

571    76. Wu, Y. et al. (2017) Quantifying the mapping precision of genome-wide

572    association studies using whole-genome sequencing data. Genome Biol. 18, 86.

573    77. Burgess, S. et al. (2017) Sensitivity analyses for robust causal inference from

574    Mendelian randomization analyses with multiple genetic variants. Epidemiology 28,

575    30-42.

576

**Figure 1: Molecular quantitative trait loci influencing intermediate traits.** Left graph: Molecular quantitative trait loci (molQTL) are genetic variants associated to a molecular trait and have an influence on intermediate traits (genotypes AA, AG, GG). Right graph: The GG genotype (blue) is associated with higher expression levels of the molecular quantitative trait compared to the AG (yellow) and AA (rose) genotype. These molecular traits can modulate the expression of further target genes (green).

**Key Figure, Figure 2: Schematic representation of a Mendelian randomization study using quantitative trait loci as instrumental variables.** Due to random distribution of alleles at conception, genetic variants are unrelated to environmental confounders. If genetic variants are sufficiently associated with the modifiable exposure of interest (here: methylation levels, RNA expression levels or protein levels) and not associated to the outcome by a different pathway, then they can be used as instrumental variable for the exposure.

**Figure 3: Simplified directed acyclic graphs of possible systems that would lead to an apparent causal effect of gene expression on a trait.** Gene regulation may be regulated by several elements. In all the situations depicted, a naïve Mendelian randomization (MR) analysis would return a causal signal for any of the regulatory elements though most often they are not on the causal pathway. A) Three scenarios for cis molecular quantitative trait loci (molQTL) regulation are presented. Vertical: Both gene expression and DNA methylation (DNAm) are on the causal pathway, hence MR using the cis-genetic variant will give valid causal estimates whether it is used to instrument either of these elements. Horizontal: Using the instrument for DNAm will be invalid due to horizontal pleiotropy. Different causal variants: The molQTL is in linkage disequilibrium (LD) with another variant that influences the trait, hence neither regulatory element is causally influenced though naïve MR could indicate otherwise. B) Four scenarios for molQTL regulation are similar to A) except the molQTL for DNAm is on a different chromosome. There are now more opportunities for horizontal pleiotropy as there needs to be a longer path from the trans chromosome to the DNA methylation level.

610 **Figure 1**



Coding gene of a molecular quantitative trait    Protein coding gene

611

612

613 **Figure 2**



Genetic variants

Exposure
methylation status
RNA level
protein level

Outcome
(Human complex
disease / trait)

Confounders

614

615

**Figure 3**



A. Cis molQTL

Vertical  Horizontal  Different causal variants

DNAm

Gene expression

Genetic variant

B. Trans molQTL

Vertical  Horizontal

Horizontal  Different causal variants

617

618

619

620