1    **Expression based polygenic scores - A gene network perspective to capture**

2    **individual differences in biological processes.**

3    Barbara Barth[1,2,3], Euclides José de Mendonça Filho[2,3], Danusa Mar Arcego[2,3], Irina Pokhvisneva

4    [2,3], Michael J. Meaney[2,3,4,5], Patrícia Pelufo Silveira[2,3,4]

5    1.  Integrated Program in Neurosciences, McGill University, Montreal, QC, Canada.
6    2.  Douglas Mental Health University Institute, McGill University, Montreal, QC, Canada.
7    3.  Ludmer Centre for Neuroinformatics and Mental Health, Douglas Research Centre,
8        McGill University, Montreal, QC, Canada
9    4.  Department of Psychiatry, Faculty of Medicine, McGill University, Montreal, QC,
10       Canada.
11   5.  Translational Neuroscience Program, Singapore Institute for Clinical Sciences, Agency
12       for Science, Technology and Research (A*STAR), Singapore, Singapore
13
14
15

16   *Corresponding author:

17   Patricia Pelufo Silveira, MD, PhD

18   Department of Psychiatry, Faculty of Medicine, McGill University

19   Douglas Research Centre, 6875 Boulevard LaSalle, Montreal, QC, H4H 1R3, Canada.

20   Phone: 514-761-6131 (ext.2776)

21   Fax: 514-761-6131

22   patricia.silveira@mcgill.ca

23

0

24    **Incorporating functional aspects into polygenic scores may accelerate early diagnosis and**

25    **the discovery of therapeutic targets. Yet, existing polygenic scores summarize information**

26    **from genome wide statistical associations between SNPs and phenotypes. We developed the**

27    **novel biologically informed, expression-based polygenic scores (ePRS or ePGS). The**

28    **method characterizes tissue specific gene co-expression networks from genome-wide RNA**

29    **sequencing data and incorporates this information into polygenic scores. Performance and**

30    **characteristics of the ePGS were compared to traditional polygenic risk score (PRS). We**

31    **observed that ePGS differs from PRS for aggregating information on; i. the relation**

32    **between different genes (co-expression); ii. the levels of tissue-specific gene expression; iii.**

33    **the genetic variation of the target sample; iv. the tissue-specific effect size of the association**

34    **between genotyping and gene expression; v. the portability across different ancestries.**

35    **Variations in the ePGS represent individual variations in the expression of a tissue-specific**

36    **gene co-expression network, and this methodology may profoundly influence the way we**

37    **study human disease biology.**

38

39

40

41

42

43

44

45

# Main

47     Genome wide association studies (GWAS) are used to identify genetic variants statistically

48     associated with a disease or trait[1] by comparing single nucleotide polymorphisms (SNPs) across

49     the genome in cases and controls. An initial objective was to identify individual common

50     variants closely linked to phenotype that might account for a substantial portion of inter-

51     individual variation. However, it is now clear that common disorders and complex traits are

52     instead highly polygenic, reflecting the influence of thousands of polymorphisms, each with

53     relatively small effects. Polygenicity led to the development of polygenic risk scores (PRSs) that

54     are calculated from GWAS results in target samples to reflect a cumulative influence of risk

55     alleles. PRS aggregates the GWAS information by summing the risk alleles count weighted by

56     the effect size for each SNP presented in the GWAS [2,3]. PRS combines the isolated small effects

57     of multiple genetic variants in a single score that represents the genetic risk for a disease or

58     variation in the expression of a trait. The use of PRSs has proven effective in defining main

59     effects of heritable genetic variations in relation to a wide range of outcomes. Moreover, PRSs

60     are a continuous measure that offer a plausible alternative to candidate gene approaches.

61         Polygenicity involves the function of diverse genes and molecules that interact with each

62     other in cellular networks[4]. Genes do not operate in isolation but conjointly in tissue-specific

63     networks that regulate molecular events and precise biological functions[5]. A gene network

64     involves a number of genes co-expressed within a specific tissue or brain region that exert a

65     concerted effect on a target biological process. Since they rely solely on DNA sequence

66     variation, existing PRS methods do not capture these biological intricacies and functional

67     relations of tissue-specific gene networks. The challenge was to create a genomic metric that

68     would reflect the influence of genetic variation, as does the PRS method, but do so within the

69    context of a tissue-specific gene network. To meet this challenge, we created an innovative

70    approach to genomic profiling that characterizes gene networks based on the levels of co-

71    expression within a specific tissue[6-17]. The co-expression based polygenic score (ePRS or ePGS)

72    method integrates information from both GWAS and tissue-specific RNA sequencing (RNAseq)

73    data sets.

74         In the examples presented here using the ePGS technique, we focus on specific brain

75    regions, but the method can be applied to any tissue. There are two approaches to the definition

76    of the co-expression networks that depend upon the research objective. One approach is designed

77    to test specific hypothesis regarding the function of a specific gene network in a specific brain

78    region. In this instance (see **Figure 1**) a gene network is constructed by focusing on a target gene

79    in a specific brain region.  In a series of studies, we focused on dopamine signaling in the

80    prefrontal cortex and thus created a co-expression network comprised of genes in which the

81    expression is significantly (i.e., r>0.5) correlated with that of *SLC6A3*, which encodes the

82    dopamine transporter. For the sake of comparison, we created *SLC6A3*-based co-expression

83    networks from RNAseq data sets in an alternative brain region. This approach allows the

84    researcher to define the region-specificity for any outcomes. A virtue of this approach is the

85    ability to test hypotheses, often derived from studies with model systems, using human data sets.

86    The second approach is aligned to discovery and employs Whole-Genome Co-Expression

87    Network Analysis (WGCNA)[18] to identify co-expression modules from RNAseq data. The

88    resulting modules can then be statistically tested for associations with treatment or traits of

89    interest. The module statistically related to the trait of interest then serves as the gene network

90    for the calculation of the ePGS.

91         The gene network of interest then serves as the basis for the selection of genes used in the

92    formulation of the ePGS. SNPs from these genes are functionally annotated and subjected to

3

93    linkage disequilibrium clumping for removal of highly correlated SNPs. A count function of the

94    number of effect alleles at a given SNP is established and weighted by the effect size of the

95    association between the individual SNP and the expression of the related gene in a specified

96    tissue using the Gene Tissue Expression (GTEx [19]) human RNAseq data sets. The sum of these

97    values from the total number of SNPs defines the ePGS at the level of the individual subject

98    (**Figure 1, Supplemental Figure 2**) (**Supplemental Methods**).

99         The ePGS combines information on: i. the relation between different genes (co-

100   expression); ii. the levels of tissue-specific gene expression (bulk or single-cell genome wide

101   RNAseq); iii. the genetic variation of the target sample (genotyping data); iv. the tissue-specific

102   effect size of the association between variants and gene expression (GTEx). Therefore, variations

103   in the ePGS represent individual variations in the genetically-determined capacity for the

104   expression of the genes that comprise the tissue-specific gene co-expression network. In this

105   paper we present the ePGS technique, its method of calculation and compare its features and

106   score content with a traditional PRS.

107

## Results

*Expression-based polygenic scores (ePGS) calculation:*

110   The steps by which an ePGS is constructed are summarized in **Figure 1**. We first describe the

111   methods for the identification of tissue-specific gene networks, which are the essential feature of

112   the ePGS approach. Researchers can use both co-expression[6-16] and differential expression[17]

113   data, from publicly available or their own datasets, see **Supplementary Figure 2**. Publicly

114   available data sets include RNAseq databases for both rodents (e.g. GeneNetwork [20]) and

115    humans (e.g. BrainEAC [21]) that can be used to identify gene networks. In the examples presented

116    here we focus on a specific brain region, but the method can be applied to any tissue.

117        Since expression of gene networks vary from region to region, obtaining gene networks

118    that are tissue specific informs on the relevance of both the gene network and the brain region or

119    tissue. A formidable advantage of the ePGS approach is the ability to create a genomic metric by

120    which to test hypothesis concerning tissue-specific gene expression profiles in any human data

121    sets for which there is both genotyping and the target phenotypic measure. In this instance a gene

122    network is constructed by focusing on a target gene in a specific brain region. For the examples

123    that will be discussed here, we have focused on dopamine signaling in the mesocortical pathway,

124    more specifically, the prefrontal cortex (PFC), the final target of this pathway. To achieve this,

125    we constructed a co-expression network comprising genes whose expression is notably

126    correlated (i.e., $r \geq 0.5$) in the PFC with either SLC6A3, responsible for encoding the dopamine

127    transporter, or with the dopamine receptor D2 gene (DRD2), two important regulators of

128    dopamine neurotransmission in the brain (See **Supplemental Table 1**) (see **Figure 1** for

129    schematic representation and **Supplemental Figure 2** for gene co-expression rationale). The

130    calculations were performed separately for each gene network of interest using the GeneNetwork

131    (http://genenetwork.org) database from RNAseq data from mice. Note, the cut-off for the

132    correlation coefficient is arbitrary, based on conventionally regarded as moderate to high

133    correlation. For the sake of comparison and to establish tissue specificity, we create a co-

134    expression network from RNAseq data sets in an alternative brain region. This feature allows the

135    researcher to statistically establish associations that are tissue or brain region specific.

136        When the identification of the gene network is anchored to a specific target gene, the

137    gene network is composed of the genes significantly co-expressed with that target gene in a

5

138    specific brain region or tissue (**Figure 1**). Using biomaRT R package[22,23] (Ensembl GRCh37) the

139    co-expressed genes are converted to human homologous genes, and all the existing SNPs from

140    these genes are gathered. Common SNPs were selected between the three sources (the SNPs

141    gathered from the gene networks of interest, the SNPs from the GTEx project[19] data in human

142    PFC and with the SNPs from the study sample (1000 Genomes Project[24])) and were subjected to

143    linkage disequilibrium clumping ($r^2<0.2$) within 500kb radius, to inform the removal of highly

144    correlated SNPs. The number of effect alleles at a given SNP is weighted using the estimated

145    effect of the tissue specific genotype-gene expression association from the GTEx project[19]. We

146    also accounted for the direction of the co-expression of each gene with *SLC6A3* or *DRD2* by

147    multiplying the weight by -1 in case the expression of a gene was negatively correlated with the

148    expression of the *SLC6A3* or *DRD2* genes. The sum of the weighted values from all SNPs,

149    divided by the number of SNPs, provided the region-specific ePGS scores.

150            The ePGS scores were calculated separately for each ancestry in the 1000 Genomes

151    Project, which includes African (N=661), American (N=347), East Asian (N=504), European

152    (N=503) and South Asian (N=489).  Since the majority of donors in the GTEx project were of

153    European ancestry[25] (see donor information at: https://gtexportal.org/home/tissueSummaryPage),

154    most of the comparisons demonstrated here used 1000 Genomes Project European sample, for

155    both ePGS and PRS (see Supplemental material, the exception being the analysis comparing the

156    scores across all ancestries). The *SLC6A3* network for European ancestry included 262 genes and

157    15387 SNPs. The *DRD2* network for European ancestry had 281 genes and 12595 SNPs (See

158    **Supplemental Table 1** for a description of genes and SNPs included in all scores described in

159    the study).

160

161    **ePGSs reflect cohesive, biologically meaningful gene networks**

162    We then compared the gene network structure represented by same size ePGS and PRS. To

163    achieve that, we mined gene co-expression information from GeneMANIA[26,27]

164    (http://genemania.org) to identify and quantify connections between the genes from each score.

165    GeneMANIA provides coexpression information between all genes from a queried gene list. We

166    also used the Centiscape tool[28] in Cytoscape®[29], to estimate two centrality measures of the

167    networks: degree, which is the number of connections between each node (each gene) and

168    betweenness, that estimates the number of times a node lies on the shortest path between other

169    nodes. **Figure 2a** depicts the gene network for *SLC6A3* PFC ePGS (number of genes = 262),

170    with a dense connection pattern between genes. Similar sized PRSs for broad depression resulted

171    in a network, depicted in **Figure 2b** (number of genes = 265). When comparing the total degree

172    between genes in the different scores using a one-way ANOVA, results show that the *SLC6A3*

173    PFC ePGS derived gene network has significantly more total connections than the broad

174    depression PRS (**Figure 2c**). The same results were found for the *DRD2* PFC ePGS (281 genes,

175    **Supplemental Figure 1a**) and its comparable size broad depression PRS (**Supplemental**

176    **Figures 1b and 1c**).

177    It is important to highlight main conceptual differences between ePGS and PRS that can

178    explain dissimilarities in total connectivity. PRSs are built selecting SNPs from a GWAS based

179    on their genome-wide significance level, and for that reason both intron and exon DNA

180    sequences are considered. Introns are non-coding DNA sequences within the genome, and

181    therefore are not mapped to genes. Introns embody 25% of the human genome and are 4 to 5

182    times the size of exons[30]. In fact, a large number of significant SNPs from GWAS are in intronic

183    and intergenic regions[31,32]. On the other hand, the ePGS is built from gene co-expression

184    information, and therefore considers only protein-coding DNA sequences, the exons, resulting in

185    every SNP being mapped to a gene. The ePGS maps into a dense group of genes (higher

186  connectivity) that interact with each other, possibly representing associated molecular functions

187  as described below.

188

**ePGS and PRS represent different biological mechanisms**

190  Because of the differences in SNP selection between ePGS (a gene co-expression network

191  identified in RNAseq data) and PRS (statistically significant SNPs from a GWAS), it is expected

192  that the two scores will differ in the biological mechanisms that they represent. We compared

193  PRS and ePGS enrichment analyses using MetaCore™ (Clarivate Analytics, version 21.4)

194  (https://portal.genego.com) and the function "compare experiments". We identified a significant

195  common gene ontology (GO) term and exported unique elements from each network that are

196  significantly associated to that GO term (FDR < 0.05) for comparison purposes. Networks were

197  constructed for direct interactions between selected objects and filtered for brain tissue and

198  human species.

199      It is noteworthy that "neuron differentiation (FDR<0.001)" was a common GO process

200  associated with genes from both PRSs and ePGS genes. However, this finding was due to

201  different element networks in each score (**Figure 3**). In ePGS, "neuron differentiation" was

202  mapped to elements such as "Nestin", which is present in neural stem and progenitor cells and

203  directly involved in differentiation process[33]. In PRS, "neuron differentiation" was mapped to

204  elements such as "olfactory receptor" and less connections are seen between elements. Taken

205  together, the findings depicted in **Figure 3** suggest that while both ePGS and PRSs are linked to

206  processes related to neuron projection development, these relations occur via unique and specific

207  mechanisms. The unique elements related to the ePGS score, in these examples, are richer and

8

208    more connected, suggesting that variations in the ePGS score possibly represent variation on

209    these specific biological processes.

210

211    **ePGS genes represent co-expression networks that are preserved across species.**

212    Since our example ePGSs were originally informed by co-expression networks identified in mice

213    (**Supplemental Methods**), we examined whether ePGS genes would also represent co-

214    expression networks in humans and compare brain co-expression patterns between ePGS genes

215    and traditional PRS genes. We used PFC gene expression data in human post-mortem brain

216    tissue from the BrainSpan database (from embryonic to adulthood, N= 42)[34] and analyzed the

217    correlation between the expression levels in the PFC for the ePGS and PRS gene lists. It is

218    important to note that in this comparison the gene list used for the ePGS originates from mouse,

219    whereas that for the PRS is from human data sets. Our results show that ePGS gene networks, in

220    the examples given here, have greater PFC gene co-expression percentage in humans in

221    comparison to PRS gene lists (**Figure 4**).  For the *SLC6A3* PFC ePGS, 40% of the gene pairs had

222    an absolute expression correlation r>=0.5 and 80% of the correlations were significant at P<0.05.

223    However, when using the genes of a traditional PRS for broad depression, a lower percentage of

224    co-expression was observed with 17% of the gene pairs had an absolute expression correlation

225    r>= 0.5 and only 62% of the correlations were significant at P<0.05. The same comparisons were

226    done for the *DRD2* PFC ePGS and its respective comparable size broad depression PRS, and

227    more robust co-expression patterns were consistently observed in ePGS in comparison to PRSs

228    for broad depression (see **Figure 4**). The results from these examples indicate that ePGSs

229    informed by mice RNAseq data represent brain gene co-expression networks also in humans, and

230    these gene networks are more tightly connected than those represented by genes that constitute

231     the traditional PRS in the examples seen here. This finding demonstrates a successful cross

232     species translation of genome functional annotation into the ePGS scores.

233

234     **ePGS reflects tissue specific co-expression networks.**

235     The ePGS calculation is informed by RNAseq data, which quantifies genome-wide tissue-

236     specific transcription (**Supplemental figure 2**). Therefore, the ePGS is based on tissue-specific

237     gene co-expression data to identify the gene network.  The tissue-specific genotype-gene

238     expression association from GTEx is then to weight the ePGS SNPs. Thus, both the selection of

239     the genes and their weighting are derived from tissue specific data sets. In contrast, a PRS is

240     based on the genotype, which is the same across different cells and tissue types.

241         To exemplify the importance of tissue specificity, we compared two gene networks built

242     on the same gene as the initial anchor, *SLC6A3*, in the PFC and the striatum. Please note the

243     differences in visualization of the *SLC6A3* PFC (total number of genes = 262) and *SLC6A3*

244     Striatum (total number of genes = 346) networks (**Supplemental Figure 3a**). We identified 53

245     genes in common between the networks (**Supplemental figure 3b**), which represents a small

246     percentage of the total number of genes from both regions (21% for *SLC6A3* PFC ePGS and 15%

247     for *SLC6A3* Striatum ePGS). This finding highlights the considerable tissue specificity of the

248     networks, even when based on the same initial gene as the anchor, which demonstrates the

249     ability of the ePGS to represent tissue specific information[35].

250     **ePGS interacts with environmental variation**

251         Despite a broadly-held conviction that genotype – phenotype relations can be context

252     specific, the demonstration of gene x environment interactions has been controversial. The

253     controversy was focused largely on candidate gene approaches that commonly failed to replicate

10

254    and generally flew in the face of the polygenic nature of the target phenotypes. Unfortunately,

255    despite its polygenic nature, investigations using polygenic scores derived from GWASs show

256    only modest success in revealing gene-environment interactions[36][37][38]. This is actually

257    unsurprising. A PRS is based on a GWAS using the most significant SNPs representing genetic

258    variants strongly associated with a condition or trait. The considerable strength of the PRS

259    method is the ability to capture polygenetically-determined predispositions for phenotypic

260    outcomes as simple main effects using a continuous measure. A PRS is thus an ideal tool for the

261    study of main effects of genetic variation. However, the reliance on SNPs that pass a designated

262    level of statistical association with the phenotype of interest biases in favor of those variants that

263    exhibit minimal environmental dependency. The implication is that SNPs in genes that that are

264    highly dependent upon environmental context are less likely to emerge as significant as main

265    effects in a GWAS, considering the rigorous GWAS-level of statistical significance for main

266    effects. **Figure 5** shows a Manhattan plot for the broad depression GWAS[39]. SNPs in green are

267    those included in the *SLC6A3* PFC ePGS, demonstrating that the variants included in the ePGS

268    lie well below the GWAS significance level. This difference would be expected if SNP's

269    comprising an ePGS are context dependent. This may explain why the ePGS may be more suited

270    to identify GxE interaction effects[40] as documented below.

271        The results of analyses using the ePGS method have consistently revealed significant and

272    gene x environment interactions. What is essential to appreciate is the high degree of replication

273    of these findings across highly diverse populations, including those of different ancestry. There

274    are now a number of published studies that demonstrate the capacity of the ePGS to identify

275    gene-environment interactions. Importantly, these analyses use a variety of measures of

276    environmental quality and phenotypic outcomes. For example, De Lima et al (2022) described

277    that PFC ePGS based on the leptin receptor gene moderated the effect of postnatal adversity on

11

278  child eating behaviour[41]. This was an example of a hypothesis-driven analysis based on prior

279  knowledge of leptin receptor activity in appetite regulation. Dalmaz et al (2021) showed that a

280  network of genes co-expressed with the synaptic protein VAMP1 gene in the PFC moderates the

281  influence of the early environment on cognitive function in children[42]. Miguel et al (2019) found

282  a significant association between history of exposure to perinatal hypoxic ischemic conditions

283  and children's cognitive flexibility, but this was moderated by the PFC *SLC6A3* ePGS [43].

284  In a study that used a WGCNA approach to define the ePGS, Arcego et al[44] provided

285  evidence for a hippocampal glucocorticoid-sensitive gene network as a moderated of the effect

286  of early life adversity on later mental health in two distinct populations. The ePGS was based on

287  a gene network derived from RNAseq with hippocampus in non-human primates using WGCNA

288  to identify the glucocorticoid-sensitive module. Interestingly, the authors also used parallel

289  independent component analysis to identify brain regions significantly associated with the

290  glucocorticoid-sensitive gene network. In sum, an increasing evidence suggests that the ePGS is

291  an appropriate method to identify GxE interaction effects.

292

293  **ePGS has high trans-ancestry portability of genetic data**

294  Allele frequency varies across ancestries[45] and the lack of proper diverse populations

295  representation in current genetic association studies hampers the translation of findings into

296  clinical applications[46]. Efforts are being made to identify genetic variations common and unique

297  to different populations, such as the 1000 Genomes Project that identified novel SNPs[47] and the

298  HapMap consortium[48]. Nevertheless the level of precision currently available for European

299  ancestry is still not uniformly available for other ancestries[49]. In PRS, the SNP list is derived

300  from the GWAS and the same variants are included in the calculation of the polygenic score in

301    diverse populations, which challenges PRS trans-ancestry portability[46,50,51]. The calculation of a

302    PRS relies on SNPs, a level of analysis at which ancestral differences are greatest. In contrast, as

303    the ePGS calculation emerges from a gene list, the SNPs included in the same ePGS may differ

304    across ancestries but will still represent the same gene list and the same co-expression network.

305            The use of genetic scores that perform functional annotation or that consider genes as the

306    first level of information, instead of SNPs, may have advantages for trans-ancestry application of

307    genetic data[52,53], as is the case of the ePGS method. Indeed, we see high trans-ancestry

308    portability and replicability of findings using ePGS[9,15-17,42,43,54]. To illustrate the differences

309    between the traditional PRS and the ePGS in terms of score composition and trans-ancestry

310    portability, we calculated PRSs of comparable size to ePGS (*SLC6A3* or *DRD2*) in the 1000

311    Genomes Project dataset. The scores were calculated separately for each ancestry to account for

312    ancestry-specific allele frequencies and linkage disequilibrium. Ancestries include African,

313    American, East Asian, European and South Asian (**Supplemental Methods**). The same number

314    of SNPs present in each ePGS for each ancestry was selected from the most significant variants

315    described in the reference GWAS (broad depression[39]), and subjected to linkage disequilibrium

316    clumping ($r^2<0.2$) for calculation of PRS separately in each ancestry. Next, the SNPs derived

317    from the calculated PRSs for each ancestry were assigned to genes and compared with ePGSs

318    gene list. **Figure 6** shows the gene overlap between the five different ancestries for each ePGS

319    and their respective comparable size PRS. The ePGS has a higher percentage of gene overlap

320    between different ancestries in comparison to PRS scores in the examples seen here. These

321    results could explain the performance of the ePGS in terms of replication seen in studies across

322    ancestries using the ePGS method[9,15-17,42,43] since ePGS preserves more information (number of

323    genes) across ancestries in comparison to PRS. We also compared the score distribution density

324   across ancestries (**Supplemental Figure 4**). Overall, the ePGS has a greater density overlap

325   between ancestries than the PRS.

326

327   **Future steps and perspectives in ePGS research**

328   The ePGS calculation is initiated by the definition of a biologically relevant gene

329   network, and this can be done in multiple ways. The examples provided here utilized co-

330   expression data from mice anchored in specific genes for the identification of co-expression

331   networks (*SLC6A3* or *DRD2*). However, other types of data and levels of information can also be

332   used to inform the calculation of ePGS, such as protein-protein interactions, DNA methylation

333   data, or differently expressed gene lists[17]. A promising venue currently being used in our lab

334   consist of utilizing weighted gene correlation network analyses (WGCNA)[18] in RNAseq data to

335   identify co-expression gene networks significantly associated with an exposure or condition in

336   controlled animal model experiments or in postmortem human tissue, in a data driven manner,

337   thus completely abandoning the hypothesis-driven approach. This perspective is well aligned

338   with the complex system in biology paradigm, and it is an anticipated improvement of the

339   method. Arcego et al (2023) is a demonstration of this improvement as the authors used

340   WGCNA to identify a hippocampal network of genes responsive to glucocorticoid treatment in

341   macaques and then calculated an ePGS in humans based on this identified gene network[44].

342   After the selection of the gene network, the list of genes can be filtered by diverse

343   parameters. Adding filters allow the integration of additional information such as the

344   developmental period, by filtering the gene selection for genes upregulated during a certain stage

345   using Brainspan[9,34,55]. Chromosome conformation information can also be added[56], by using data

346   from high-throughput sequencing (Hi-C) and assigning noncoding SNPs to their cognate genes

14

347     based on long-range interactions using H-MAGMA[57] input files that describe gene–SNP pairs

348     based on brain Hi-C data[58]. FIMO[59] can also be used to include variants affecting transcription

349     factor binding motifs from the genes of the network. Finally, candidate regulatory variants can be

350     added by mapping available SNPs on promoter regions (up to 4kb upstream of the transcription

351     start site) of the genes that compose the network. Lastly, the weight attributed to each SNP in the

352     ePGS calculation can be derived from different GWASs. In the current examples, a GWAS for

353     gene expression (GTEx[19]) was used, thus reflecting individual variations in gene expression of

354     the network in the specific brain region. All these parameters can be accommodated to

355     contemplate different research questions. Finally, adaptation of the ePGS technique for the use of

356     single-cell and spatial transcriptomics will add still increased resolution and specificity to the

357     polygenic scores.

358     **Discussion**

359            Aligned with the idea of incorporating functional genomics information to PRS

360     technology, we have developed the expression based polygenic score (ePGS). While both PRS

361     and ePGS summarize the small effects of multiple SNPs using the genotype information, the use

362     of tissue specific gene expression data in the ePGS technique transforms the polygenic score into

363     a functional genomic tissue-specific measure. The ePGS also reflects the combined biological

364     function of gene networks.

365            Here we demonstrated the consequences of rethinking SNP selection and incorporating

366     other levels of information to polygenic scores, such as gene expression and tissue specific data.

367     We compare ePGS and PRS features and score content. The ePGS reflects cohesive gene

368     networks, demonstrating a high level of co-expression between the genes. This could be

369     explained by ePGS considering only exon DNA sequences and being built from gene co-

370     expression information. It is important to highlight that since genes do not work in isolation, but

15

371    rather in networks[5], the use of a gene network perspective has the potential to better reflect

372    biological functions associated with these genes. We demonstrated that the ePGS and PRS reflect

373    different biological processes, when comparing unique elements that are related to a common

374    gene ontology term. The ePGS unique elements, in the examples demonstrated here, appear to be

375    richer and more connected, suggesting that variations in the ePGS score may represent variation

376    on a specific biological process. We also demonstrated that ePGS based gene networks represent

377    tissue specific co-expression networks in humans. The possibility of reflecting functional

378    genomics information in a tissue specific manner is one of the strengths of the ePGS,

379    demonstrated here by the uniqueness of the *SLC6A3* PFC gene network in comparison to the

380    *SLC6A3* Striatum gene network. As a consequence of these above-mentioned features, the ePGS

381    is suited to test gene by environment effects, evidenced by previous published studies[9,16,42-44].

382    The content of ePGS on different ancestries seem consistent when comparing the ePGS and PRS

383    score gene overlap. This is expected since the use of genome functional annotation has the power

384    to improve prediction of complex traits within and between ancestries[60] and the incorporation of

385    functional markers, such as gene expression, improves trans-ancestry portability of genomic

386    data[61]. The ePGS uses genome functional annotation in two steps of its calculation; in the co-

387    expression basis and by weighing the SNPs using GTEx genotype-gene expression association.

388          An advantage of using a gene network approach like the ePGS is the possibility of

389    integrating other data modalities also represented by networks or with high dimensionality. For

390    example, the integration of genetic and neuroimage information by parallel independent

391    component analysis, which estimates the maximum independent components within each data

392    modality separately while also maximizing the association between modalities using an entropy

393    term based on information theory [62]. Studies using pICA and the ePGS have found interesting

16

394    results linking both data modalities and informing on the neuroanatomical basis of the effects of

395    variations in the gene network expression[9,42,43,63].

396         In conclusion, the ePGS method is purely based on biological, co-expression data and no

397    information on association with outcomes of interest (e.g. GWAS for diseases) is used. The

398    differences between conventional PRSs and ePGSs presented here, may explain the successful

399    ePGS performance in gene by environment interaction models and across ancestries, suggesting

400    that the ePGS is an interesting method to capture individual biological variation in response to

401    environmental changes[7,17], and may profoundly influence the way we study human disease

402    biology.

17

## References

1    Uffelmann, E. *et al.* Genome-wide association studies. *Nature Reviews Methods Primers* **1** (2021). https://doi.org:10.1038/s43586-021-00056-9

2    Wray, N. R., Lee, S. H., Mehta, D., Vinkhuyzen, A. A., Dudbridge, F. & Middeldorp, C. M. Research review: Polygenic methods and their application to psychiatric traits. *J Child Psychol Psychiatry* **55**, 1068-1087 (2014). https://doi.org:10.1111/jcpp.12295

3    Dudbridge, F. Power and predictive accuracy of polygenic risk scores. *PLoS genetics* **9**, e1003348 (2013). https://doi.org:10.1371/journal.pgen.1003348

4    Schadt, E. E. Molecular networks as sensors and drivers of common human diseases. *Nature* **461**, 218-223 (2009). https://doi.org:10.1038/nature08454

5    Gaiteri, C., Ding, Y., French, B., Tseng, G. C. & Sibille, E. Beyond modules and hubs: the potential of gene coexpression networks for investigating molecular mechanisms of complex brain disorders. *Genes, Brain, Behavior* **13**, 13-24 (2014).

6    Silveira, P. P. *et al.* Cumulative prenatal exposure to adversity reveals associations with a broad range of neurodevelopmental outcomes that are moderated by a novel, biologically informed polygenetic score based on the serotonin transporter solute carrier family C6, member 4 (SLC6A4) gene expression. *Dev Psychopathol* **29**, 1601-1617 (2017). https://doi.org:10.1017/S0954579417001262

7    Hari Dass, S. A. *et al.* A biologically-informed polygenic score identifies endophenotypes and clinical conditions associated with the insulin receptor function on specific brain regions. *EBioMedicine* **42**, 188-202 (2019). https://doi.org:10.1016/j.ebiom.2019.03.051

8    Miguel, P. M. *et al.* Prefrontal Cortex Dopamine Transporter Gene Network Moderates the Effect of Perinatal Hypoxic-Ischemic Conditions on Cognitive Flexibility and Brain Gray Matter Density in Children. *Biol Psychiatry* **86**, 621-630 (2019). https://doi.org:10.1016/j.biopsych.2019.03.983

9    de Lima, R. M. S. *et al.* Amygdala 5-HTT Gene Network Moderates the Effects of Postnatal Adversity on Attention Problems: Anatomo-Functional Correlation and Epigenetic Changes. *Front Neurosci* **14**, 198 (2020). https://doi.org:10.3389/fnins.2020.00198

10   Morgunova, A. *et al.* DCC gene network in the prefrontal cortex is associated with total brain volume in childhood. *Journal of psychiatry & neuroscience : JPN* **46**, E154-E163 (2020). https://doi.org:10.1503/jpn.200081

11   Potter-Dickey, A. *et al.* Associations Among Parental Caregiving Quality, Cannabinoid Receptor 1 Expression-Based Polygenic Scores, and Infant-Parent Attachment: Evidence for Differential Genetic Susceptibility? *Frontiers in neuroscience* **15**, 704392 (2021). https://doi.org:10.3389/fnins.2021.704392

12   Dalmaz, C. *et al.* Prefrontal cortex VAMP1 gene network moderates the effect of the early environment on cognitive flexibility in children. *Neurobiol Learn Mem* **185**, 107509 (2021). https://doi.org:10.1016/j.nlm.2021.107509

13   Selenius, J. S. *et al.* The relationship between health-related quality of life and melancholic depressive symptoms is modified by brain insulin receptor gene network. *Sci Rep* **11**, 21588 (2021). https://doi.org:10.1038/s41598-021-00631-w

14   de Mendonca Filho, E. J. *et al.* Cognitive Development and Brain Gray Matter Susceptibility to Prenatal Adversities: Moderation by the Prefrontal Cortex Brain-Derived Neurotrophic Factor Gene Co-expression Network. *Frontiers in neuroscience* **15**, 744743 (2021). https://doi.org:10.3389/fnins.2021.744743

15   Restrepo-Lozano, J. M. *et al.* Corticolimbic DCC gene co-expression networks as predictors of impulsivity in children. *Molecular psychiatry* **27**, 2742-2750 (2022). https://doi.org:10.1038/s41380-022-01533-7

| 451 | 16 | de Lima, R. M. S., Barth, B., Arcego, D. M., de Mendonça Filho, E. J., Patel, S., Wang, Z., Pokhvisneva, I., Parent, C. Levitan, R.D., Kobor, M.S., Bittencourt, A.P.S.V, Meaney, M.J. Dalmaz, C., Silveira, P.P. Leptin receptor co-expression gene network moderates the effect of early life adversity on eating behavior in children. *Communications Biology* **Accepted** (2022). |

451  16   de Lima, R. M. S., Barth, B., Arcego, D. M., de Mendonça Filho, E. J., Patel, S., Wang, Z.,
452       Pokhvisneva, I., Parent, C. Levitan, R.D., Kobor, M.S., Bittencourt, A.P.S.V, Meaney, M.J. Dalmaz,
453       C., Silveira, P.P. Leptin receptor co-expression gene network moderates the effect of early life
454       adversity on eating behavior in children. *Communications Biology* **Accepted** (2022).
455  17   Latsko, M. S., Wang, Z., Zhang, T.Y., Parent, C., O'Toole, N., Pokhvisneva I., Kee, M. Z. L., Wen, X.,
456       Craig, K., Boyce, W.T., Meaney, M. J., Silveira, P. P. A translational polygenic score of biological
457       sensitivity to context. *Am J Psychiat* **Accepted** (2022).
458  18   Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network analysis.
459       *BMC bioinformatics* **9**, 559 (2008).
460  19   Consortium, G. T. The Genotype-Tissue Expression (GTEx) project. *Nat Genet* **45**, 580-585 (2013).
461       https://doi.org:10.1038/ng.2653
462  20   Mulligan, M. K., Mozhui, K., Prins, P. & Williams, R. W. GeneNetwork: A Toolbox for Systems
463       Genetics. *Methods Mol Biol* **1488**, 75-120 (2017). https://doi.org:10.1007/978-1-4939-6427-7_4
464  21   Ramasamy, A. *et al.* Genetic variability in the regulation of gene expression in ten regions of the
465       human brain. *Nat Neurosci* **17**, 1418-1428 (2014). https://doi.org:10.1038/nn.3801
466  22   Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping identifiers for the integration of
467       genomic datasets with the R/Bioconductor package biomaRt. *Nature protocols* **4**, 1184 (2009).
468  23   Durinck, S. *et al.* BioMart and Bioconductor: a powerful link between biological databases and
469       microarray data analysis. *Bioinformatics* **21**, 3439-3440 (2005).
470       https://doi.org:10.1093/bioinformatics/bti525
471  24   Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74 (2015).
472       https://doi.org:10.1038/nature15393
473  25   Gay, N. R. *et al.* Impact of admixture and ancestry on eQTL analysis and GWAS colocalization in
474       GTEx. *Genome Biology* **21**, 233 (2020). https://doi.org:10.1186/s13059-020-02113-0
475  26   Warde-Farley, D. *et al.* The GeneMANIA prediction server: biological network integration for
476       gene prioritization and predicting gene function. *Nucleic acids research* **38**, W214-W220 (2010).
477  27   Montojo, J. *et al.* GeneMANIA Cytoscape plugin: fast gene function predictions on the desktop.
478       *Bioinformatics* **26**, 2927-2928 (2010). https://doi.org:10.1093/bioinformatics/btq562
479  28   Scardoni, G., Petterlini, M. & Laudanna, C. Analyzing biological network parameters with
480       CentiScaPe. *Bioinformatics* **25**, 2857-2859 (2009).
481       https://doi.org:10.1093/bioinformatics/btp517
482  29   Shannon, P. *et al.* Cytoscape: a software environment for integrated models of biomolecular
483       interaction networks. *Genome Res* **13**, 2498-2504 (2003). https://doi.org:10.1101/gr.1239303
484  30   Sakharkar, M. K., Chow, V. T. & Kangueane, P. Distributions of exons and introns in the human
485       genome. *In silico biology* **4**, 387-393 (2004).
486  31   Bartonicek, N. *et al.* Intergenic disease-associated regions are abundant in novel transcripts.
487       *Genome Biology* **18**, 241 (2017). https://doi.org:10.1186/s13059-017-1363-3
488  32   Zhang, F. & Lupski, J. R. Non-coding genetic variants in human disease. *Human molecular*
489       *genetics* **24**, R102-R110 (2015).
490  33   Bernal, A. & Arranz, L. Nestin-expressing progenitor cells: function, identity and therapeutic
491       implications. *Cell Mol Life Sci* **75**, 2177-2195 (2018). https://doi.org:10.1007/s00018-018-2794-z
492  34   Miller, J. A. *et al.* Transcriptional landscape of the prenatal human brain. *Nature* **508**, 199 (2014).
493  35   Hari Dass, S. A. *et al.* A biologically-informed polygenic score identifies endophenotypes and
494       clinical conditions associated with the insulin receptor function on specific brain regions.
495       *EBioMedicine* **42**, 13 (2019).
496  36   Mullins, N. *et al.* Polygenic interactions with environmental adversity in the aetiology of major
497       depressive disorder. *Psychol Med* **46**, 759-770 (2016).
498       https://doi.org:10.1017/S0033291715002172

| 499 | 37 | Peyrot, W. J. *et al.* Does Childhood Trauma Moderate Polygenic Risk for Depression? A Meta-analysis of 5765 Subjects From the Psychiatric Genomics Consortium. *Biol Psychiatry* **84**, 138-147 (2018). https://doi.org:10.1016/j.biopsych.2017.09.009 |
|---|---|---|
| 502 | 38 | Trotta, A. *et al.* Interplay between Schizophrenia Polygenic Risk Score and Childhood Adversity in First-Presentation Psychotic Disorder: A Pilot Study. *PloS one* **11**, e0163319 (2016). https://doi.org:10.1371/journal.pone.0163319 |
| 505 | 39 | Howard, D. M. *et al.* Genome-wide meta-analysis of depression identifies 102 independent variants and highlights the importance of the prefrontal brain regions. *Nature Neuroscience* **22**, 343-352 (2019). https://doi.org:10.1038/s41593-018-0326-7 |
| 508 | 40 | Silveira, P. P. & Meaney, M. J. Examining the biological mechanisms of human mental disorders resulting from gene-environment interdependence using novel functional genomic approaches. *Neurobiol Dis* **178**, 106008 (2023). https://doi.org:10.1016/j.nbd.2023.106008 |
| 511 | 41 | de Lima, R. M. S. *et al.* Leptin receptor co-expression gene network moderates the effect of early life adversity on eating behavior in children. *Communications Biology* **5**, 1092 (2022). |
| 513 | 42 | Dalmaz, C. *et al.* Prefrontal cortex VAMP1 gene network moderates the effect of the early environment on cognitive flexibility in children. *Neurobiology of Learning and Memory*, 107509 (2021). |
| 516 | 43 | Miguel, P. M. *et al.* Prefrontal cortex dopamine transporter gene network moderates the effect of perinatal hypoxic-ischemic conditions on cognitive flexibility and brain gray matter density in children. *Biological psychiatry* **86**, 621-630 (2019). |
| 519 | 44 | Arcego, D. M. *et al.* A Glucocorticoid-Sensitive Hippocampal Gene Network Moderates the Impact of Early-Life Adversity on Mental Health Outcomes. *Biological Psychiatry* (2023). https://doi.org:https://doi.org/10.1016/j.biopsych.2023.06.028 |
| 522 | 45 | Choudhury, A. *et al.* Population-specific common SNPs reflect demographic histories and highlight regions of genomic plasticity with functional relevance. *BMC genomics* **15**, 1-20 (2014). |
| 524 | 46 | Sirugo, G., Williams, S. M. & Tishkoff, S. A. The missing diversity in human genetic studies. *Cell* **177**, 26-31 (2019). |
| 526 | 47 | McVean, G. A. *et al.* An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56-65 (2012). https://doi.org:10.1038/nature11632 |
| 528 | 48 | Consortium, I. H. A haplotype map of the human genome. *Nature* **437**, 1299 (2005). |
| 529 | 49 | Fitipaldi, H. & Franks, P. W. Ethnic, gender and other sociodemographic biases in genome-wide association studies for the most burdensome non-communicable diseases: 2005-2022. *Hum Mol Genet* **32**, 520-532 (2023). https://doi.org/10.1093/hmg/ddac245 |
| 532 | 50 | Martin, A. R., Kanai, M., Kamatani, Y., Okada, Y., Neale, B. M. & Daly, M. J. Clinical use of current polygenic risk scores may exacerbate health disparities. *Nature genetics* **51**, 584-591 (2019). https://doi.org:10.1038/s41588-019-0379-x |
| 535 | 51 | Lam, M. *et al.* Comparative genetic architectures of schizophrenia in East Asian and European populations. *Nature genetics* **51**, 1670-1678 (2019). https://doi.org:10.1038/s41588-019-0512-x |
| 537 | 52 | Liang, Y. *et al.* Polygenic transcriptome risk scores (PTRS) can improve portability of polygenic risk scores across ancestries. *Genome biology* **23**, 23 (2022). https://doi.org:10.1186/s13059-021-02591-w |
| 540 | 53 | Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nature genetics* **47**, 1228-1235 (2015). https://doi.org:10.1038/ng.3404 |
| 543 | 54 | Restrepo-Lozano, J. M. *et al.* Corticolimbic DCC gene co-expression networks as predictors of impulsivity in children. *Molecular Psychiatry* **27**, 2742-2750 (2022). https://doi.org:10.1038/s41380-022-01533-7 |
| 546 | 55 | de Mendonça Filho, E. J., Barth, B., Bandeir, D. R., de Lima, R. M. S., Arcego, D. M. Dalmaz, C., Pokhvisneva, I., Sassi, R. B., Hall, G. B. C., Meaney, M. J., Silveira, P. P. . Cognitive Development |

20

548         and Brain Gray Matter Susceptibility to Prenatal Adversities: Moderation by the Prefrontal
549         Cortex Brain-Derived Neurotrophic Factor Gene Co-expression. *Frontiers in neuroscience* (2021).
550         https://doi.org:10.3389/fnins.2021.744743
551   56   Dalmaz, C., Pokhvisneva, I., Wang, Z., Barth, B., Patel, S., de Lima, R.M.S, de Mendonça Filho,
552         E.J., Arcego, D. M., Kobor, M. S., O'Donnell, K. J., Meaney, M. J., Silveira, P. P. . Syntaxin-1A
553         gene network moderates the vulnerability/resilience to early life trauma-induced depressive
554         symptoms in women. (Submitted).
555   57   Sey, N. Y. A. *et al.* A computational tool (H-MAGMA) for improved prediction of brain-disorder
556         risk genes by incorporating brain chromatin interaction profiles. *Nat Neurosci* **23**, 583-593
557         (2020). https://doi.org:10.1038/s41593-020-0603-0
558   58   Wang, D. *et al.* Comprehensive functional genomic resource and integrative model for the
559         human brain. *Science* **362** (2018). https://doi.org:10.1126/science.aat8464
560   59   Bailey, T. L., Johnson, J., Grant, C. E. & Noble, W. S. The MEME Suite. *Nucleic acids research* **43**,
561         W39-W49 (2015). https://doi.org:10.1093/nar/gkv416
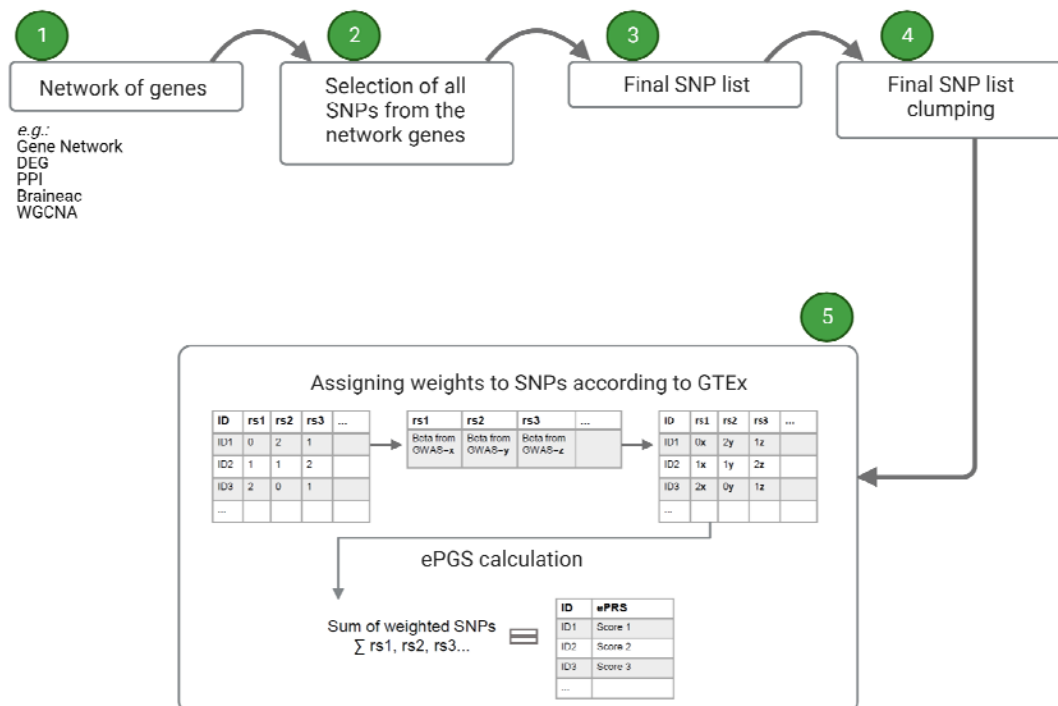562   60   Zheng, Z. *et al.* Leveraging functional genomic annotations and genome coverage to improve
563         polygenic prediction of complex traits within and between ancestries. *bioRxiv*,
564         2022.2010.2012.510418 (2022). https://doi.org:10.1101/2022.10.12.510418
565   61   Amariuta, T. *et al.* Improving the trans-ancestry portability of polygenic risk scores by prioritizing
566         variants in predicted cell-type-specific regulatory elements. *Nature Genetics* **52**, 1346-1354
567         (2020). https://doi.org:10.1038/s41588-020-00740-8
568   62   Pearlson, G. D., Liu, J. & Calhoun, V. D. An introductory review of parallel independent
569         component analysis (p-ICA) and a guide to applying p-ICA to genetic data and imaging
570         phenotypes to identify disease-associated biological pathways and systems in common complex
571         disorders. *Front Genet* **6**, 276 (2015). https://doi.org:10.3389/fgene.2015.00276
572   63   de Mendonça Filho, E. J. *et al.* Cognitive Development and Brain Gray Matter Susceptibility to
573         Prenatal Adversities: Moderation by the Prefrontal Cortex Brain-Derived Neurotrophic Factor
574         Gene Co-expression Network. *Frontiers in Neuroscience* **15** (2021).
575         https://doi.org:10.3389/fnins.2021.744743

576

**Figures and legends**



**Figure 1. Schematic figure representing the key steps to calculate the ePGS. 1)** Construction of a network of genes that is defined by a set of genes that interact in a biologically meaningful way. Some examples are co-expression of transcripts from animal models (GeneNetwork), as used in the current study, and different expression analysis (DEG). Additionally, it can be defined by protein-protein interaction (PPI), co-expression of transcripts from human samples (Braineac) and by weighted gene co-expression network analysis (WGCNA). At this step, tissue specificity can be defined by selecting transcript data from specific tissues of interest. The list of genes can also be filtered by a specific developmental time point, for example, by using publicly available databases such as the BrainSpan[34]. Furthermore, the list of genes can be filtered by other conditions and interests. **2)** Selection of all existing SNPs from the gene network was done using biomaRt package. From this list we retained common SNPs with a) SNPs from the study sample genotyping data and b) SNPs present in GTEx (which is a genome-wide analysis that has gene expression as the outcome; GTEx was chosen to weight the selected SNPs in the examples provided here). The common SNPs represent the final SNP list that is subjected to linkage disequilibrium clumping ($r^2 > 0.2$). **5)** Weight the SNPs: the number of effect alleles (genotype information from the study sample) at a given SNP is multiplied by the effect size of the association between SNPs and the gene expression (GTEx). The sum of all weighted SNPs for each individual corresponds to the individual ePGS.

22

**Figure 2. Network visualization comparison of *SLC6A3* derived ePGS and comparable size PRSes. a)** *SLC6A3* PFC ePGS gene network; **b)** Broad depression PRS gene network comparable size with *SLC6A3* PFC ePGS; **c)** One-way ANOVA of total connectivity (total degree values) for ePGS and PRS comparable size. Gene co-expression interactions were obtained from GeneMANIA (http://genemania.org) and used to generate the networks with Cytoscape® application, which specifies amount of interactions between pairs of genes based on their co-expression, represented by the number of edges (gray lines) in the networks. The Centiscape plug-in in Cytoscape® was used to calculate the centrality of the genes in each network, defining the degree (number of connections with other nodes, represented by node size, in which bigger nodes indicates more connections with other nodes) and betweenness (number of times a node lies on the shortest path between other nodes, represented by node's color in which darker colors indicate higher betweenness in the networks) for the components of the networks.
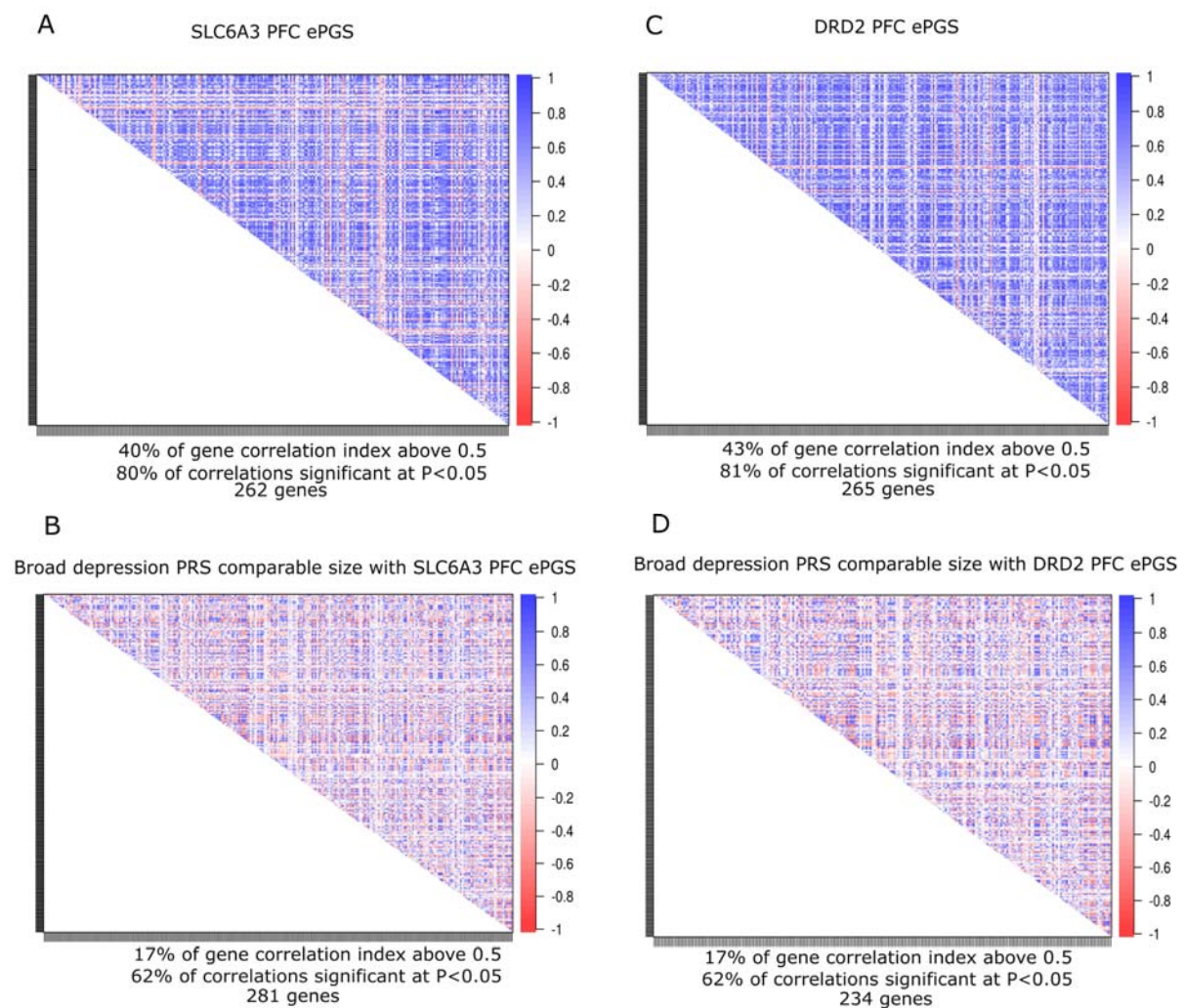
23

**Figure 3. Unique elements for 'neuron differentiation', a common gene ontology enrichment analysis term for both ePGS and PRS.** Gene ontology (GO) enrichment analysis was performed using Metacore®. The function "compare experiments" was used to obtain common significant (FDR <0.05) GO terms between the gene networks while also identifying the unique elements from each network that are significantly associated to the GO term. Networks were plotted in MetaCore® using the unique elements of each network for the GO enrichment term selected. Figures **a, b, c,** and **d** show visual comparisons of the different contributions of ePGS and PRS to the GO term. The details of the legends of the network's figures can be found in https://portal.genego.com/legends/MetaCoreQuickReferenceGuide.pdf.

24

**Figure 4. Correlation matrix of gene expression for ePGS gene networks and PRS gene networks based on BrainSpan human post-mortem brain tissue (from embryonic to adulthood, N=42). a)** *SLC6A3* PFC ePGS gene network: 40% of the gene correlations was above 0.5 and 80% of the correlations are significant at P<0.05; **b)** Broad depression PRS gene network comparable size with the *SLC6A3* PFC ePGS: 17% of the ge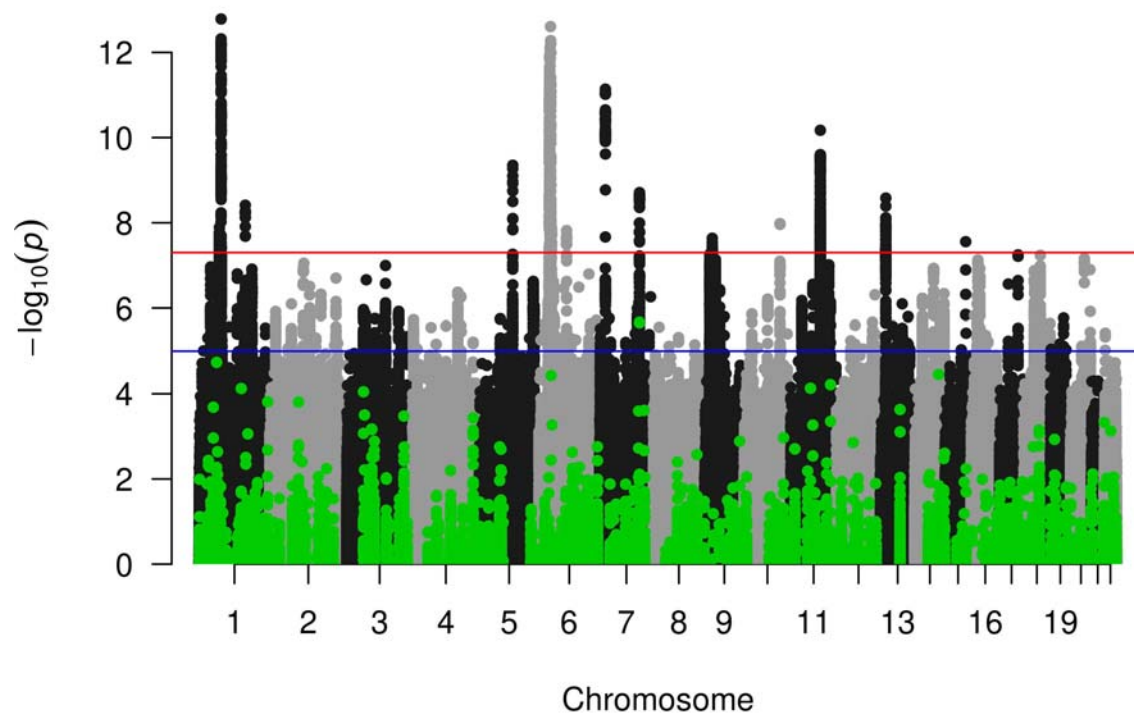ne correlations above 0.5; 62% correlations significant at P<0.05; **c)** *DRD2* PFC ePGS gene network: 43% of gene correlations above 0.5 and 81% of correlations significant at P<0.05; **d)** Broad depression PRS gene network comparable size with the *DRD2* PFC ePGS: 17% of the gene correlations above 0.5; 62% correlations significant at P<0.05.

631

**Figure 5. Manhattan plot for Howard (2019) broad depression GWAS results and *SLC6A3*
PFC ePGS SNPs.** Gray and black dots represent -log10(p) from the broad depression GWAS.
Green dots represent -log10(p) from GTEx for the SNPs included in *SLC6A3* PFC ePGS. It
demonstrates that all SNPs from the ePGS are not statistically significant at the genome wide
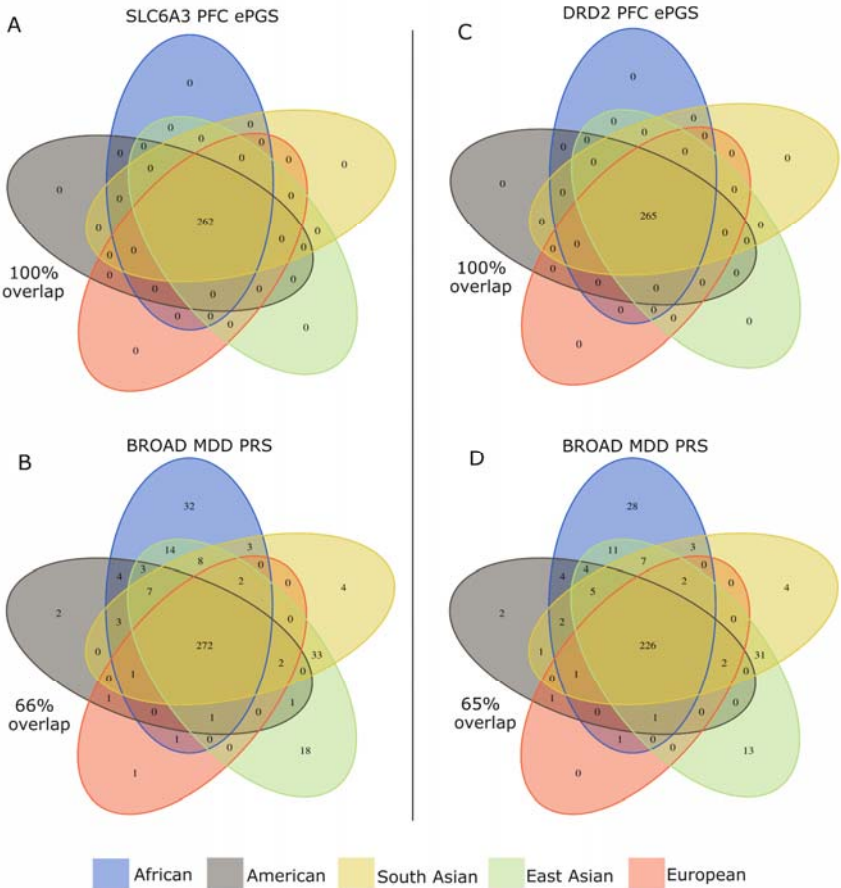level.

637

**Figure 6. Venn diagrams of gene overlap for ePGSes and PRSes calculated based on the ePGS and PRS in the 1000 Genomes Project dataset.** Gene overlap between the five different ancestries for *SLC6A3* and *DRD2* ePGS and their respective comparable size PRS. It demonstrates that the ePGS have more common genes between different ancestries in comparison to PRS scores.

**Acknowledgments**

**Author contributions**

BB and PPS designed the study, BB, EJMF, DMA and IP conducted data analysis, BB, EJMF, DMA and IP generated the figures and BB and PPS wrote the manuscript with editing of MJM and input from all authors. PPS and MJM supervised the research. All authors read and approved the final manuscript.

**Competing interest declaration**

The authors declare no competing interests.