1

**Copy number variation profile-based genomic subtyping of premenstrual dysphoric disorder in Chinese**

Hong Xue[1,2,3*], Zhenggang Wu[2,3], Xi Long[2], Ata Ullah[2], Si Chen[2], Wai-Kin Mat[2], Peng Sun[1], Ming-Zhou Gao[1], Jie-Qiong Wang[1], Hai-Jun Wang[1], Xia Li[1], Wen-Jun Sun[1], and Ming-Qi Qiao[1*]

[1]Shandong University of Traditional Chinese Medicine, Jinan, Shandong, People's Republic of China

[2]Division of Life Science, Hong Kong University of Science and Technology, Hong Kong, People's Republic of China

[3]School of Basic Medicine and Clinical Pharmacy, China Pharmaceutical University, Nanjing, Jiangsu, People's Republic of China

* Correspondence:

Hong Xue

hxue@ust.hk

Mingqi Qiao

qmingqi@163.com

## Abstract

23  Premenstrual dysphoric disorder (PMDD) affects nearly 5% women of reproductive age. The
24
25  symptomatic heterogeneity, along with largely unknown genetics, of PMDD have greatly

26  hindered its effective treatment. In the present study, 127 Chinese PMDD patients of the

27  'invasion' and 'depression' subtypes clinically differentiated by us earlier were analyzed

28  together with 108 non-PMDD controls for genome-wide copy number variations (CNVs).

29  Germline genomic DNA samples from white blood cells were subjected to AluScan

30  sequencing-based CNV profiling, which enabled clustering of patient samples readily into the

31  V and D groups, dominated by the "invasion" and "depression" clinical subtypes,

32  respectively; the CNVs obtained with 100-kb windows yielded two clusters that were

33  correlated with these subtypes with a consistency of up to 89.8%. Diagnostic correlation- and

34  frequency-based CNV features of either CNV-gain (CNVG) or CNV-loss (CNVL) that could

35  differentiate between V and D subtypes were selected and analyzed. CNVG features located

36  preferentially in S2-phase replicating regions and enriched with steroid hormone biosynthesis

37  pathway of genes were found protective against PMDD. Moreover, machine learning

38  employing the correlation-based CNV features could predict with >80% accuracy whether a

39  genomic sample was D-type, V-type or control. In terms of their CNV profiles, the D- and V-

40  types differed more from one another than from the controls, thereby providing a genomic

41  basis for the clinical D-V subtyping of PMDD. Genome-wide profiling of CNVs, as a new

42  approach to complex disease genetics, has revealed recurrent CNVs and genomic features

43  beyond individual genes and mutations underlying PMDD clinical diversity.

44

## Introduction

46  Premenstrual dysphoric disorder (PMDD) is a syndrome that afflicts 5-10% of women in

47  their reproductive years (1). The severity of the syndrome is typically highest just before the

48    menstruation period, suggesting that the symptoms were linked to hormonal changes. This

49    has been confirmed by the findings of premenstrual neurosteroid fluctuations, and alterations

50    in the sensitivity of GABA$_A$ receptors to neurosteroids giving rise to mood instability (2, 3).

51    Cortical gamma-aminobutyric acid (GABA) levels also declined during the menstrual cycle

52    in healthy women but increased in women with PMDD from the follicular phase to the mid-

53    luteal and late luteal stages (4). Furthermore, PMDD has been associated with the estrogen

54    receptor alpha gene *ESR1* (5), and the ESC/E(Z) genes affecting the interactions of sex

55    hormones with other genes (6). Five major contributors to the etiology of PMDD include: (1)

56    genetic susceptibility; (2) progesterone and its metabolite ALLO; (3) estrogen, serotonin and

57    brain-derived neurotrophic factor (BDNF); (4) brain structure and function; and (5) the

58    hypothalamic-pituitary-adrenal axis and hypothalamic-pituitary-gonadal axis (7). The

59    schizophrenia-associated SNPs in *GABRB2*, located in introns 8 and 9 near an AluYi6

60    insertion, have been associated with both schizophrenia and bipolar disorder (8, 9), heroin

61    addiction (10), altruism (11), autism and mental retardation (12). Deletion of gabrb2 genes

62    from knockout mice also brought about schizophrenic symptoms that were alleviated by the

63    antipsychotic Risperidone (13). Recently, analysis of germline copy-number-variations

64    (CNVs) at the nsv1177513 site in Exon 11, and the esv2730987 site in Intron 6, of *GABRB2*

65    in PMDD and schizophrenia patients showed that CNV alterations at both esv2730987 and

66    nsv1177513 were significantly associated with schizophrenia in Chinese and Germans as

67    well as PMDD in Chinese (14). Moreover, subjects with different levels of susceptibility to

68    cancer could be distinguished by means of diagnostic CNV marker features selected from the

69    germline genomes with the application of machine learning (15).

70

71    It is recognized that the symptoms of PMDD are consistent with multiple clinical subtypes. A

72    Delphi survey led to the proposal of three symptoms-based types of PMDD, *viz.* a

73    predominantly physical type, a predominantly emotional type, and a mixed type (16); and

74    DSM-V proposed that PMDD is defined by one or more of the symptoms of marked affective

75    lability, marked irritability or anger, marked depressed mood and hopelessness, and marked

76    anxiety and tension, plus at least one of seven other symptoms. At the School of Basic

77    Medicine, Shandong University of Traditional Chinese Medicine, the medical records also

78    pointed to at least two major types of PMDD, viz. an irritability-marked 'invasion' type

79    (58.9%) and a depressive mood-marked 'depression' type (27.5%) (17). In view of the

80    spectrum of PMDD symptoms, the objective of the present study was to enquire whether the

81    two major clinical subtypes of PMDD could be corelated with genomic profiles. Through

82    genome-wide CNV profiling by AluScan next-generation sequencing (18, 19), the results

83    revealed two large clusters of CNV profiles that were highly correlated with the clinical

84    "depression" and "invasion" subtypes. Furthermore, CNV-gain (CNVG) and CNV-loss

85    (CNVL) features diagnostic of PMDD or each of the two clinical subtypes were uncovered

86    among CNVs called from sequence windows of different sizes, which were variously

87    distributed in genomic regions of different replication timing and overlapped with genes in

88    various genetic pathways of potential clinical relevance. These results provided genomic

89    verification for the invasion- and depression-subtypes employed by us previously (17), which

90    corresponded to part of the complex symptoms stipulated by DSM-V (20) as diagnostic

91    criteria for PMDD.

92

93    **Methods**

94    **Clinical assessments**

95    Clinical diagnosis of PMDD patients (P-type subjects) from asymptomatic controls (C-type

96    subjects) was performed in accordance to the protocol in Diagnostic and Statistical Manual of

97    Mental Disorders (DSM-IV) by two psychiatrists independently. The identifications of

4

98    'depression-type' and 'invasion-type' subjects were carried out as previously described (17).

99

**Genomic DNA samples**

101    Peripheral white blood cell DNA samples were collected from PMDD patients and non-

102    PMDD control subjects with approval by the institutional ethic committee of Shandong

103    University of Chinese Medicine. The patients and healthy volunteers who participated in this

104    study all signed the informed consent form. The samples consist of a control cohort of 108

105    subjects and a PMDD cohort of 127 cases. The latter cohort was further divided into the

106    depression-subtype (71 cases) and invasion-subtype (56 cases). The subtypings of the 127

107    PMDD cases were given in Table S1.

108

**AluScan sequencing and CNV calling**

110    Samples of ~0.1μg DNA were subjected to inter-Alu PCR amplification using the four Alu-

111    consensual primers AluY278T18, AluY66H21, R12A/267 and L12A/8 (18). The 200 bp to

112    ~6 kb amplicons in each sample were employed to build a library for sequencing on the

113    Illumina platform with 100 bp paired-end reads. According to the standard framework, all the

114    reads were mapped to reference human genome hg19 downloaded from UCSC by BWA,

115    followed by base recalibration and local realignment by GATK (21). CNVs were called from

116    the AluScan sequences with the method of AluScanCNV2 (19, 22) based on sequence

117    windows of 50-500 kb in 50-kb increments on the 22 autosomes and the X chromosome. The

118    CNV profiles of all 108 control and 127 PMDD subjects were available in Table S2.

119

**Clustering and grouping of patient samples based on CNV profiles**

121    The profiles of CNVG and CNVL called from the 127 P-group samples using different CNV-

122    calling window sizes were separately subjected to correlation analysis and hierarchical

123    clustering with 1,000 bootstraps using the 'pvclust' R package (23). The derived correlation

124    heatmaps as well as the CNVG-based and CNVL-based dendrograms obtained for each

125    window size were employed to determine the two subgroups of CNV profiles using two

126    different grouping methods for cross validation.

127

128    In the first method, *viz*. the straightforward '*cutree*' method, the sub-clustering was carried

129    out using the '*cutree*' function from the 'dendextend' R package (24) to cut each dendrogram

130    into 2-8 sub-clusters (Figure S1). The DNA samples located in the sub-cluster populated with

131    the highest number of clinical depression-type samples among all the sub-clusters was

132    referred as D-type genomic samples; and the DNA samples located in the remaining sub-

133    clusters were combined and referred as V-type genomic samples. In the second, or '*semi-*

134    *supervised*' method, some branches on the dendrograms were first rotated around their

135    respective nodes to bring the closely co-localized samples into tightly knit sub-clusters

136    enclosed by black square boxes on the diagonal of each heatmap. Thereupon, all the samples

137    within the same block box were all designated as D-type or V-type genomic samples

138    depending on whether the majority clinical subtype of the samples were depression-type or

139    invasion-type. The designated D- and V-type genomic samples derived using the two

140    grouping methods for ten different window sizes are shown in Table S3 and exemplified by

141    the blue and red branches in Figure 1 and Figure S2 respectively for the 100-kb CNV profiles.

142

143    For either the '*cutree*' method or the '*semi-supervised*' method, let the number of CNV-based

144    D-type samples that also belonged to the clinical depression-subtype be represented by $True_D$,

145    and the number of CNV-based V-type samples that also belonged to the clinical invasion-

146    subtypes be represented by $True_V$. Accordingly, the consistency ($Y$) between CNV-based

147    classification and the clinical classification of PMDD patient samples could be estimated by:

$$Y = (True_D + True_V)/127$$

148     On this basis, the levels of consistency between CNV-based and clinical subtypings for the

149     different CNVs called using different window sizes for both the '*cutree*' and '*semi-supervised*'

150     methods are shown in Table S4.

151

152     **Selection of diagnostic CNV features**

153     The selection of diagnostic CNV features was performed using either (a) correlation-based

154     method or (b) frequency-based method as described (15). CfsSubsetEval from the Weka

155     package was employed together with BestFirst search method to select the correlation-based

156     diagnostic CNV features. Fisher's exact tests were employed to select the frequency-based

157     CNV features that showed significantly different occurrence frequencies between a pair of

158     sample groups (e.g. P-vs-C or D-vs-V) with a false discovery rate (FDR) less than 0.01.

159

160     **Predictive subtyping of genomic samples by machine learning**

161     Earlier, diagnostic germline CN-gains and CN-losses from leucocyte DNA samples of

162     subjects with or without past episodes of cancers in tissues other than leucocytes were found

163     to provide a useful basis to predict the propensity of the subject to cancer (15). Since the 127

164     P-group and 108 C-group DNA samples from PMDD and control subjects consisted of a

165     mixture of D-type, V-type and C-type DNAs, the question arose whether it was possible to

166     predict the typing of DNA samples between the D-vs-V, D-vs-C and V-vs-C choices

167     employing the diagnostic CNVG and CNVL features obtained with the correlation-based

168     method.

169

170     For example, in a choice between the P-vs-C types, a mixture of P- and C-type samples were

171     randomly separated into a labeled Learning Band and an unlabeled Test Band, with equal or

172   near equal number of samples in the two bands. Diagnostic CNVG and CNVL features were

173   selected from the labeled Combined Learning Band with machine learning using the

174   correlation-based method and employed to estimate the risk factor $R$ for each DNA sample in

175   the Test Band according to Eqn. 1.

176
$$R = \log\left(\frac{Pr(PMDD|Features)}{Pr(Control|Features)}\right) \qquad \text{Eqn. 1}$$

$$Pr(PMDD|Features) = Pr\left(Features|PMDD\right) \times Pr\left(PMDD\right)/Pr\left(Features\right)$$

$$Pr(Control|Features) = Pr\left(Features|Control\right) \times Pr\left(Control\right)/Pr\left(Features\right)$$

177   where Pr(PMDD|Features) was the posterior probability of membership in the PMDD group

178   given the CNV data of a particular Test Band sample; Pr(Control|Features) was its posterior

179   probability of membership in the Control group given the same test CNV data;

180   Pr(Features|PMDD) was the likelihood function of the test CNV data given membership in

181   the PMDD group; Pr(Features|Control) was the likelihood function of the test CNV data

182   given membership in the Control group; Pr(PMDD) and Pr(Control) were the prior

183   distributions of PMDD and Control samples respectively within the Learning Band; and

184   Pr(Features) was the prior distribution of CNV-features among all the CNVs within the

185   Learning Band.

186

187   For every sample in the Test Band, its value of R estimated using Eqn. 1 would predict

188   whether the sample belonged to the Control group or PMDD group: it would predictively

189   belong to Control group (*viz.* 'non-PMDD') if $R < 0$; belong to PMDD group if $R > 0$; or no

190   prediction could be made if $R = 0$. For every PMDD sample in the Test Band, $R > 0$

191   represented a 'true' prediction whereas $R < 0$ represented a 'not true' prediction. On the other

192   hand, for any Control sample in the Test Band, $R > 0$ represented a 'not true' prediction

193   whereas $R < 0$ represented a 'true' prediction. Accuracy of prediction was therefore given by:

194
$$\text{Accuracy} = \frac{[\text{True predictions of Control}]+[\text{True predictions of PMDD}]}{[\text{Total predictions of Control}]+[\text{Total predictions of PMDD}]} \times 100\% \qquad \text{Eqn. 2}$$

195

196  Repetition of this procedure 1,000 times would yield 1,000 Accuracy estimates, and in turn

197  the Average Accuracy regarding the P-vs-C typing.

198

199  **Functional annotation of genes overlapping with diagnostic CNV features**

200  By comparing the genomic coordinates of all the frequency-based diagnostic CNV features to

201  those of the known genes retrieved from the R package

202  'TxDb.Hsapiens.UCSC.hg19.knownGene' version 3.2.2 (25), and considering any gene to be

203  'overlapping' with a CNV feature if any proportion of its sequence (from > 0% to 100% in 10%

204  increments) coincided with part or all of the CNV feature, the list of CNV-overlapping genes

205  obtained was uploaded to DAVID Bioinformatics Resources as a *test-list* employing the

206  'RDAVIDWebService' R package (26). All the known genes on chromosomes 1-22 and X

207  were also uploaded as the *background-list*. Comparison of the two lists using the

208  'getFunctionalAnnotationChart' of the 'RDAVIDWebService' R package revealed gene

209  pathways or categories, as defined in the GO, KEGG and INTERPRO databases, that were

210  enriched with the *test-list* of genes among the *background-list* of genes. The pathways or

211  categories yielding <0.05 Benjamini-corrected *p*-values were regarded to be significantly

212  enriched in the genes on the *test-list* (Table S5).

213

214  **Genomic-feature content of diagnostic CNV features in different replication phases**

215  DNA sequences on 22 autosomes and chromosome X were subject to replication-time

216  segmentation according to Long & Xue (27). Briefly speaking, experiment-assessed

217  replication timing of all 1-kb sequence windows in the genomes of fifteen human cell lines

218  were retrieved from the 'UW Repli-seq track' in the UCSC Table Browser (28), and the

219  representative replication phase of each sequence window was identified as one of the six

220 types of sequencing segments (viz. G1b, S1, S2, S3, S4 and G2) based on their experiment-

221 assessed replication timing in all fifteen human cell lines.

222

223 The density or intensity of genomic features were quantified as described in Ng et al (29) in

224 the diagnostic CNV features and the non-diagnostic-CNV regions in each type of replication

225 phase. The genomic feature content of diagnostic CNV features is indicated by the fold

226 change of the density or intensity of the genomic feature in diagnostic CNV features relative

227 to the non-diagnostic-CNV regions.

228

229 **Statistical analysis**

230 All comparisons of CNV frequencies were conducted using Fisher's exact tests, and the *p*-

231 values were adjusted by false discovery rate for multiple comparisons. In functional

232 annotation of genes, *p*-values from DAVID web service were subject to Benjamini-correction

233 for multiple comparisons. When annotating the genes that overlapped with any diagnostic

234 CNV feature, empirical *p*-values were estimated using Monte Carlo methods with 1,000

235 simulations to validate the significant gene pathways/categories based on the 50-, 100- or

236 450-kb size groups of CNV features. In each round of simulation, sequence windows of the

237 same size as the targeted group of CNV features were randomly selected from chromosomes

238 1-22 and X, with the number of selected windows being equal to the average number of CNV

239 features in the different type-comparisons to be analysed (see Table S6). For each simulation,

240 the genes that overlapped with any of the selected sequence windows were functionally

241 annotated. The empirical *p*-value of a targeted pathway was given by $(r+1)/(n+1)$, where $n =$

242 1,000 and $r$ = number of simulations that displayed significant enrichment (<0.05 Benjamini-

243 corrected *p*-values) in the targeted pathway.

244

245 **Software for data processing and visualization**

246 Data processing tasks were carried out using custom R codes, except that tasks requiring

247 machine learning were processed using Weka package. All figures were drawn under R

248 environment using the 'ggplot2' (30), 'pheatmap' (31) and 'quantsmooth' (32) packages,

249 except for Figure 3 which was drawn using http://bioinformatics.psb.ugent.be/webtools/Venn/,

250 and Figure 5 using Integrative Genomics Viewer 2.3.69 (33).

251

## Results

253 **Correlations between clinical diagnosis and CNV profiles**

254 In order to examine whether there might be significant correlation between the clinical

255 symptoms of PMDD patients and their germline CNV profiles, the CNVGs and CNVLs

256 called from different sizes of sequence windows on the 127 P-type DNA samples, were

257 subjected to hierarchical clustering in each instance. The CNVGs and CNVLs called from

258 100-kb sequence windows of the 71 depression-subtype and 56 invasion-subtype patient

259 samples were segregated using the cutree and semi-supervised methods into distinct D-type

260 and V-type clusters in the dendrograms as shown in Figure S2 and Figure 1 respectively. The

261 clusters obtained from the CNVG dendrograms were designated as $D_G$ and $V_G$ clusters, and

262 the clusters obtained from the CNVL dendrograms were designated as $D_L$ and $V_L$ clusters.

263 Notably, the cutree method yielded 72 $V_G$-type and 55 $D_G$-type CNVG profiles with 81.10%

264 consistency between the invasion-vs-depression clinical classification and the V-vs-D

265 CNVG-based classification (Figure S2A); whereas the semi-supervised method yielded 61

266 $V_G$-type and 66 $D_G$-type CNVG profiles with 89.76% consistency between the invasion-vs-

267 depression clinical classification and the V-vs-D CNVG-based classification (Figure 1A).

268 Therefore, using either the cutree method or the semi-supervised method, the CNVG-based

269 classification was highly correlated with the clinical symptom-based classification of the

11

270    PMDD genomes; this was likewise the case with the $D_L$-type and $V_L$-type CNVLs.

271    Altogether, for the CNVGs and CNVLs in, the 50-500 kb window sizes, the cutree method

272    yielded consistencies of 68-91%, and the semi-supervised method yielded consistencies of

273    88-98%, between the CNV-based and symptom-based classifications. The semi-supervised

274    classifications of P-type samples based on CNVs called from 50-500 kb window sizes were

275    available in Figure S3. These results demonstrated that both the CNVGs and CNVLs

276    contributed to the etiology of the depression-type and the invasion-type symptoms. Moreover,

277    the comparable results obtained using the cutree and semi-supervised methods confirmed the

278    robustness of the CNV-symptom correlations. When the CNVGs or CNVLs called from 100-

279    kb sequence windows of the 108 C-type control samples were subject to hierarchical

280    clustering along with the P-type samples, ~40% of the C-type CNV profiles formed a tight

281    sub-cluster and ~60% were dispersely distributed in the dendrogram, forming sub-clusters

282    with the depression-subtype or invasion-subtype PMDD samples (Figure S4).

283

284    **Use of diagnostic CNV-features for predictive subtyping**

285    The correlation between germline CNV profiles and clinical subtypes of PMDD suggests that

286    it would be practicable to predict from the germline CNVs of women their propensity to

287    develop PMDD, as well as the likely subtype of the PMDD clinical condition. Toward this

288    objective, the method developed earlier by us through the use of diagnostic CNV-features

289    selected with machine learning to assess a subject's propensity for cancer (15) could be

290    employed as described in 'Selection of diagnostic CNV features' under Method. Figure 2

291    shows the diagnostic CNV features selected by either the correlation method or the frequency

292    method for prediction the propensity of a test subject's germline CNVs for which of the P, C,

293    $V_G$, $D_G$, $V_L$ and $D_L$ genomic groups: the P-type and C-type outcomes would be assessed

294    based on PMDD symptoms; $V_G$ and $D_G$ would be based on the distinction between the V and

295    D clusters in the CNVG dendrogram in Figure 1A; and $V_L$ and $D_L$ would be based on the

296    distinction between the V and D clusters in the CNVL dendrogram in Figure 1B. The

297    diagnostic CNVG and CNVL features selected using the correlation and frequency methods

298    are given in Table S7.

299

300    Figure 2A shows the sets of diagnostic CNV features selected using the correlation-based

301    (red triangles) or frequency-based (black circles) method to enable a choice between a pair of

302    genomic groups. For example, the $D_G$-vs-C panel of Figure 2A contained a mixture of 66 $D_G$-

303    type samples and 108 C-type samples. The diagnostic CNV features selected from the total of

304    144 samples by means of either the correlation method (red triangles) or the frequency

305    method (grey circles) were distributed in a crescent near the y-axis and another crescent near

306    the x-axis. Accordingly, any DNA sample in the mixture that was enriched with near-y

307    diagnostic CNV features would be predicted to be endowed with a greater propensity for C-

308    type over $D_G$-type, whereas any DNA sample that was enriched with near-x diagnostic CNV

309    features would be predicted to be endowed with a propensity for $D_G$-type over C-type. In the

310    $D_G$-vs-C panel of Figure 2B, diagnostic CNV features selected using the correlated method

311    was employed to predict the $D_G$-vs-C nature in the 174-sample mixture as described under

312    the 'Predictive subtyping of genomic samples by machine learning' section in Methods. After

313    1,000 trial runs, each with a random partition of the samples into an 87-sample Learning

314    Band and an 87-sample Test Band, the average prediction accuracy obtained was 83.0%.

315    Altogether, the seven panels in Figure 2B yielded average prediction accuracies ranging from

316    81.0% to 88.4%. Interestingly, the list of correlation-based CNV features useful for

317    differentiating between the propensities toward the D and V subtypes (Table S8) showed that

318    the CNV features biased in favor of V-type samples were mostly CNVL features (27/42 for

319    $V_G$ and 14/17 for $V_L$). The accuracies of sample-classification predictions derived from the

320    cutree method are available in Figure S5.

321

322    Favorable diagnostic CNV-features were often shared by more than one PMDD types, as

323    indicated by the overlaps between the colored circles for the P-vs-C (blue), D-vs-C (red) and

324    V-vs-C (green) comparisons in the Venn diagrams (Figure 3A and B). A range of CNV

325    features were shared by all three kinds of circles, suggesting that they represented key CNV

326    features differentiating between the control and PMDD patient samples (Table S9). Notably

327    also, in all the panels in Figure 3, there was no CNV feature was shared only by the red

328    circles for D-vs-C and the green circles for V-vs-C, which suggests that the CNV-features

329    favoring the D-type genomes differed diametrically from the CNV-features favoring the V-

330    type genomes. As well, there were more D-favoring CNVG features than CNVL features, but

331    more V-favoring CNVL features than CNVG features.

332

333    **Genome-wide distribution of diagnostic CNV features**

334    In order to have a global view of CNV profiles, the locations and replication timing of all

335    frequency-based diagnostic CNV features, whether overlapping with any known genes or not,

336    were plotted on Figure S6. The results showed that the CNV features were widely spread on

337    all the somatic chromosomes and chromosome X. Chromosomes 4, 13, 18 21 and X were

338    particularly abundant in CNV features that replicated in the G2 phase. Given the correlation

339    between the clinical symptom-based typing of PMDD cases and the clustering of germline

340    diagnostic CNV features, these CNV features could be useful guides in a search for some of

341    genomic sites underlying PMDD.

342

14

343     In Figure 4, the distributions of the CNV features among DNA regions replicating at different

344     cell cycle phases exhibited a number of characteristics: (a) In terms of the number of CNV

345     features that differed between a pair of CNV-types, the P-vs-C panel (viz. P>C or P<C) gave

346     rise to the smallest difference, whereas the D-vs-V pair ($D_G>V_G$ or $D_G<V_G$) gave rise to the

347     largest difference; (b) the ratio of CNVL features relative to CNVG features (viz. L/G on

348     chart) that favored the C-type over P-type were 1.32 for 50-kb CNV features, 2.13 for 100-kb

349     ones and 2.05 for 450-kb ones, all greater than unity (Figure 4A); (c) the P-vs-C comparisons

350     were suggestive of protective effects of smaller size CNVLs in the early replication phases

351     and larger CNVLs in the later phases (Figure 4A); (d) the CNVLs captured by 50-kb

352     windows included significantly more V-favoring than either C-favoring or D-favoring ones

353     (L/G = 2.41 in Figure 4C and 1.80 in Figure 4D); (e) the CNVGs were significantly enriched

354     in D-favoring features compared to C-favoring or V-favoring ones, whereas CNVLs were

355     significantly enriched in V-favoring features compared to C-favoring or D-favoring ones; (f)

356     D-vs-V comparisons suggest that V-type PMDD was correlated with smaller CNVG features

357     belonging to the early replication phases and large CNVGs belonging to the later phases; (g)

358     Large CNVG features were enriched in the G2-phase replicating sequences, especially

359     among the features selected for the D-vs-C and D-vs-V comparisons (see G2-phase columns

360     marked with red asterisks in Figure 4B and D); (h) More than half of the large G2-phase

361     CNVG features in the $C>D_G$ group are identical to those of the $V_G>D_G$ group, suggesting the

362     shared genetic variations in G2 phase underlying V and C types; (i) Large CNVL feature

363     were enriched in S3-phase replicating sequences in the $C>V_G$ and $D_G>V_G$ groups but not in

364     the $V_G>C$ or $V_G>D_G$ groups. The replication-phase distributions of CNV features obtained

365     based on the $D_L$- or $V_L$-type samples derived from the CNVL dendrogram in Figure 1B were

366     available in Figure S7.

367

**Pathways and genes enriched in diagnostic CNV features**

A wide range of genes showed sequence overlaps with the frequency-based diagnostic CNVG and CNVL features of a range of KEGG pathways in PMDD and its subtypes (Table 1) which pointed to their possible contributions to the PMDD disorder, and some major genes were contained in more than one pathway (Table 2). It was striking that, as indicated in lines 1-5 of Table 2, the control C-type was favored by high frequencies of CNVG features relative to the diseased P-, D- or V-type, suggesting that a major causal factor of the PMDD disorder could be decreased levels of the CNVG features overlapping with the steroid hormone biosynthesis pathway, with the involvement of *CYP*- and *UGT*-genes replicating in phases S2 and S1. As shown in lines 9-17 of Table 2, the C-type and V-type profiles were favored over the D-type by high frequencies of CNVL features in the *GRI*-genes, which were involved in pathways of nicotine addiction, circadian entrainment, serotonergic synapse, dopaminergic synapse and cAMP signaling. The chromosomal sites of these genes and their overlaps with the 100-kb CNV features are shown in Figure 5.

The 50-kb CNV features overlapped with the genes in the glutamatergic-synapse, alcoholism, and systemic lupus erythematosus pathway genes, as well as steroid hormone biosynthesis pathway genes replicating in S2 and S3 (Table S5). On the other hand, the 450-kb CNV features overlapped with chemokine signaling pathway genes (Table S5). Because high regional density of genes could impact on gene annotations by yielding false-positive co-localizations when a CNV feature incidentally captured a gene cluster belonging to a pathway, empirical *p*-values based on Monte Carlo simulations were also estimated for the 100-kb CNV features (Table S10), which provided additional support for some of the pathway in Table 2 through the elimination of such false positives (see 'Statistical analysis' in Methods).

16

**Genomic features enriched in diagnostic CNV features**

Co-localization analysis revealed various associations between 100-kb frequency-based CNV features and a wide spectrum of genomic features in different replication phases (Table 3 and Table S11). D versus V differences in genomic feature contents can be identified from the thermal scale plots of co-localization scores illustrated in Figure 6 and Figure S8. The genomics features apparently differed between D and V types included: (1) In terms of retrotransposons, D-favoring CNVG features enriched with more of the subfamily of evolutionarily very young short transposons SVAef, while V-favoring CNVG and D-favoring CNVL features enriched with the very young long transposon subfamily, L1vy. (2) With respect to genetic markers, P-favoring, especially D-favoring CNVL features were enriched with recombination events as well as genetic variation hotspots and clusters (27). GWAS reported markers were co-localized with D-favoring CNVGs in S1, V-favoring CNVLs in S4 and C-favoring CNVLs G1b. As well, ClinVar markers were enriched in V-favoring CNVG of S2 phase and D-favoring CNVL of S1. (3) In respect to the group of CpG-related genomic features, the main difference between the two types was that D-favoring CNVL features were more enriched with CpG features such as MeBS in S4 replicating sequences. Compared with D- and V-favoring, the C-favoring CNV features were more prominently enriched with CpG features, especially for C-favoring CNVG in S3 and CNVL in G2 and S4 phases. (4) In regard to non-coding RNA, LINC was enriched in V-favoring CNVL as well as C-favoring CNV features, but not in D-favoring features. (5) To a lesser extent, the enrichment of histone binding sites in D-favoring CNVG features of G2 and S2 phases. In contrast, histone sites were enriched in V-favoring CNVL features of S4 phase. This trend was clearly visible from Figure 6, where twelve kind histone binding sites were analyzed separately and displayed side-by-side.

418     Some of the strongly enriched genomic features with great than one-fold enrichment was

419     listed in Table 3. For example, enrichment of DNase I hypersensitive sites (DNase) was

420     found in C>P CNVL and C>$V_L$ CNVG features in G2-phase replicating sequences and P>C

421     (as well as D>C and V>C) CNVL features in S4-phase replicating sequences. Regulatory

422     elements isolated by formaldehyde (FAIRE) were found to enrich in C>V, D>C and D>V

423     CNVL features that located in the late-replicating S4 and G2 phases. Disease- or trait-

424     associated SNPs identified by genome-wide association studies (GWAS) were enriched in

425     C>V CNVLs in G1b phase, and P>C CNVGs in S1 phase reaching a fold-change of 5.5

426     relative to non-diagnostic-CNV regions in S1 phase. The C-favoring (C>P and C>$D_L$) CNVG

427     features and C-favoring CNVL (C>P and C>$V_G$) features tend to co-localize with CpG

428     islands (CpGi) in median to late-replicating S3-G2 phases. A range of methylation-related

429     features (Me450, MeBS, and MeMRE) were found to enrich in C-favoring CNVG features

430     mainly in early to median G1b-S3 phases, and C-favoring CNVL features in late-replicating

431     S4-G2 phases. Long intergenic non-coding RNAs (LINC) were found to be enriched in C-

432     favoring CNVGs mainly in S2 phase or CNVLs mainly in G2 phase, D-favoring CNVGs in

433     G1b phase, and V-favoring CNVLs in S3-G2 phases.

434

435     **Discussion**

436     Application of either the cutree method or the semi-supervised method to the hierarchically

437     clustered CNVGs or CNVLs in the germline genomes of PMDD subjects enabled the

438     distinction between the D-type and V-type CNV profiles. The high degree of consistency

439     between the clinical depression-subtype and D-type CNV profiles, and between the clinical

440     invasion-subtype and V-type CNV profiles, indicated that the two clinical PMDD subtypes

441     were intrinsically correlated with the two dissimilar types of CNV profiles. This was further

442     conformed when diagnostic CNVG and CNVL features were selected by means of machine

443     learning using the correlation method, and employed as abundance markers to predict

444     whether a given germline genomic sample belonged to the control group, the PMDD group,

445     the V-type CNV group or the D-type CNV group, yielding average accuracies of prediction

446     of 81.0-88.4% (Figure 3), which in turn validated the use of diagnostic CNVG and CNVL

447     features to identify the genes and pathways that overlapped with such diagnostic features as

448     potential contributors to the PMDD disorder.

449

450     In this regard, there exists overall accord between the cutree and the semi-supervised

451     methods in terms of diagnostic CNV features identified, replication-phase distribution and

452     pathway enrichments (Table S5, S12 and S13). As indicated in DSM-V, PMDD is defined by

453     a complex system of symptoms. In the present study, limited data allowed the analysis of

454     only the depression-type and invasion-type symptoms. Nevertheless, the Venn diagrams in

455     Figure 3 clearly showed that the CNVs underlying the D-type and V-type CNV profiles were

456     strikingly more divergent from one another than their separate divergences from the CNVs

457     underlying the C-type. This finding was also consistent with the results in Figure 2A, which

458     showed that there were more correlation-based or frequency-based CNV features that could

459     be employed to distinguish between $D_G$-vs-$V_G$ or $D_L$-vs-$V_L$ compared to CNV-features that

460     could distinguish between $D_G$-vs-C, $V_G$-vs-C, $D_L$-vs-C or $V_L$-vs-C. As well, the mixed

461     distribution of control CNV profiles among depression- or invasion-type CNV profiles

462     (Figure S3) indicated that the difference between the CNVs in the two subtypes of PMDD

463     was larger than their individual differences from the control. This surprising genome

464     condition, as illustrated in Figure 7, raises the question of whether the depression-type and

465     invasion-type conditions of PMDD might represent two distinct clinical disorders.

466

467     A faithful temporal order of DNA replication is fundamental to normal cellular function, and

19

468  aberrant replication timings were observed in complex diseases including cancers (34, 35).

469  Accordingly, the relative abundances of diagnostic CNVG and CNVL features among

470  genomic DNA sequence regions preferentially replicating in each one of the six phases of cell

471  cycle, namely G1b, S1, S2, S3, S4 and G2 were examined in Figure 4. The peaks of C-

472  favoring CNVL features in P-vs-C comparisons (downward hollow green bars in Figure 4A)

473  shifted clearly from the early S1 phase among the 50-kb features to the late G2 phase among

474  the 450-kb features, pointing to the enrichment of some small CNVL features in the early-

475  replicating regions and larger CNVL features in the late-replicating regions among the

476  determinants of the C-type, *viz.* in the prevention of PMDD occurrence. As well, more D-

477  favoring CNVG features were located in G2-replicating sequences compared to genomic

478  DNA sequences replicating in other cell-cycle phases within the D>C and D>V groups,

479  which was particularly notable in view of the enrichment of G2 phase-replicating sequences

480  in non-coding sequences (27, 29). In addition, the abundance of V-favoring 50-kb CNVL

481  features in the V>C and V>D comparisons (Figure 4B and D) suggests that small-size

482  CNVLs also played important roles in the development of V-type PMDD.

483

484  When diagnostic CNVs were analyzed for their genomic feature enrichment with reference to

485  replication phases, interesting observations were obtained (Figure 6; Table 3). It has been

486  revealed that the late-replicating S4-G2 phases in the gene-distal zones are found to be

487  depleted of functional genomic features (27). However, the present study observed

488  associations of open chromatin signals, regulatory elements and epigenetic regulation sites

489  with the diagnostic CNV features in these late-replicating sequences (Table 3), indicating that

490  the diagnostic CNV features might represent pivotal genomic sites in the late-replicating

491  sequences that sequence alterations may give raise to functional perturbations underlying

492  PMDD and its two subtypes. As illustrated herein, genomic feature content analysis,

493 implemented with replication phase information, has pointed to the likelihood of genomic

494 events underlying the subtyping of PMDD and hence a genomic nature of the disorder and its

495 clinical diversity. Since genomic features included in the analysis were broad in spectrum and

496 well beyond the boundary of known genes, the feature enrichment analysis performed herein

497 may complement with and surpass genetic pathway analysis as a powerful tool for genomic

498 studies on complex traits and disorders.

499

500 Previously, a number of genes was proposed to be PMDD suspectable genes, including those

501 of steroid hormone biosynthesis (2, 3), and estrogen signaling (36, 37), and these proposals

502 were supported by the presence of these genes in Table 1. The overlaps of genes of nicotine

503 addiction, glutamatergic synapses, olfactory transduction, alcoholism, systemic lupus

504 erythematosus, hypogonadism, premature ovarian failure, and breast cancer with PMDD

505 might be suggestive of hitherto hidden aspects of central nervous system or endocrine system

506 involvements with PMDD. The *GRIA4* gene, overlapping with the 100-kb CNV features for

507 the C>D and V>D comparisons, groups, has also been found to be associated with

508 schizophrenia (38), in accordance with the shared CNVs between schizophrenia and PMDD

509 (14).

510

511 In conclusion, through CNV profiling, the present study provided evidence for strong

512 correlation of the clinical depression-subtype or invasion-subtype with the D-type and V-type

513 germline genomes, marked by the overlaps between their CNVs and the machine-selected

514 diagnostic CNV features that favored one or another type of genomes. On account of this

515 correlation, the diagnostic CNV features could be employed as frequency markers to predict

516 the propensity to PMDD and one of its clinical subtypes, as well as position markers to

517 identify candidate PMDD genes and pathways. Moreover, the genetic difference between the

518  depression-favoring and invasion-favoring CNV profiles was found to exceed their individual

519  divergences from the normal controls (Figure 7), raising the question of how this outcome

520  might have been evolved. Future studies will be required to determine how many of the array

521  of PMDD symptoms besides the depression-subtype and invasion-subtype ones could be

522  significantly correlated with CNVs, and what complex diseases other than PMDD would

523  embody CNV-symptom correlations as strong as those encountered with PMDD.

524

## Acknowledgements

535

## Conflict of Interest

537  The authors declare that the research was conducted in the absence of any commercial or

538  financial relationships that could be construed as a potential conflict of interest.

539

## Author Contributions

541  HX and MQ conceived and designed the experiments, ZW, XLo, AU, SC and WM performed

542  the AluScan sequencing related experiments and analysis of the sequencing data. PS, MG, JW,

543     HW, XLi, WS and MQ coordinated the collection of PMDD and control cohorts, and HX,

544     ZW, XLo, SC and MQ wrote the paper.

545

## Supplementary Information

547     Supplementary materials are available online, including Figures S1-S8 and Table S1-S15.

548

## References

550     1.    Cohen LS, Soares CN, Otto MW, Sweeney BH, Liberman RF, Harlow BL. (2002):

551           Prevalence and predictors of premenstrual dysphoric disorder (PMDD) in older

552           premenopausal women - The Harvard Study of Moods and Cycles. *J Affect Disorders.*

553           70(2):125-132.

554     2.    Smith SS, Ruderman Y, Frye C, Homanics G, Yuan M. (2006): Steroid withdrawal in

555           the mouse results in anxiogenic effects of 3alpha,5beta-THP: a possible model of

556           premenstrual dysphoric disorder. *Psychopharmacology (Berl).* 186(3):323-333.

557     3.    Turkmen S, Backstrom T, Wahlstrom G, Andreen L, Johansson IM. (2011):

558           Tolerance to allopregnanolone with focus on the GABA-A receptor. *Br J Pharmacol.*

559           162(2):311-327.

560     4.    Epperson CN, Haga K, Mason GF, Sellers E, Gueorguieva R, Zhang W, et al. (2002):

561           Cortical gamma-aminobutyric acid levels across the menstrual cycle in healthy

562           women and those with premenstrual dysphoric disorder: a proton magnetic resonance

563           spectroscopy study. *Arch Gen Psychiatry.* 59(9):851-858.

564     5.    Huo L, Straub RE, Roca C, Schmidt PJ, Shi K, Vakkalanka R, et al. (2007): Risk for

565           premenstrual dysphoric disorder is associated with genetic variation in ESR1, the

566           estrogen receptor alpha gene. *Biol Psychiatry.* 62(8):925-933.

23

567   6.   Dubey N, Hoffman JF, Schuebel K, Yuan Q, Martinez PE, Nieman LK, et al. (2017):

568        The ESC/E(Z) complex, an effector of response to ovarian steroids, manifests an

569        intrinsic difference in cells from women with premenstrual dysphoric disorder. *Mol*

570        *Psychiatry.* 22(8):1172-1184.

571   7.   Raffi ER, Freeman MP. (2017): The etiology of premenstrual dysphoric disorder: 5

572        interwoven pieces. Current Psychiatry. 16(9):20-28.

573   8.   Lo WS, Lau CF, Xuan Z, Chan CF, Feng GY, He L, et al. (2004): Association of

574        SNPs and haplotypes in GABAA receptor beta2 gene with schizophrenia. *Mol*

575        *Psychiatry.* 9(6):603-608.

576   9.   Chen J, Tsang SY, Zhao CY, Pun FW, Yu Z, Mei L, et al. (2009): GABRB2 in

577        schizophrenia and bipolar disorder: disease association, gene expression and clinical

578        correlations. *Biochem Soc Trans.* 37(Pt 6):1415-1418.

579   10.  Kim YS, Yang M, Mat WK, Tsang SY, Su ZH, Jiang XF, et al. (2015): GABRB2

580        Haplotype Association with Heroin Dependence in Chinese Population. *PLoS One.*

581        10(11).

582   11.  Tsang SY, Zhong SF, Mei LL, Chen JH, Ng SK, Pun FW, et al. (2013): Social

583        Cognitive Role of Schizophrenia Candidate Gene GABRB2. *PLoS One.* 8(4).

584   12.  Lew AR, Kellermayer TR, Sule BP, Szigeti K. (2018): Copy Number Variations in

585        Adult-onset Neuropsychiatric Diseases. *Curr Genomics.* 19(6):420-430.

586   13.  Yeung RK, Xiang ZH, Tsang SY, Li R, Ho TYC, Li Q, et al. (2018): GABRB2

587        knockout mice displayed schizophrenia-like and comorbid phenotypes with

588        interneuron-astrocyte-microglia dysregulation. *Transl Psychiatry.* 8:128.

589   14.   Ullah A, Long X, Mat WK, Hu T, Khan MI, Hui L, et al. (2020): Highly recurrent

590         copy number variations in GABRB2 associated with schizophrenia and premenstrual

591         dysphoric disorder. *Front Psychiatry.* 11:572.

592   15.   Ding X, Tsang SY, Ng SK, Xue H. (2014): Application of machine learning to

593         development of copy number variation-based prediction of cancer risk. *Genomics*

594         *Insights.* 7:1-11.

595   16.   Ismail KM, Nevatte T, O'Brien S, Paschetta E, Backstrom T, Dennerstein L, et al.

596         (2013): Clinical subtypes of core premenstrual disorders: a Delphi survey. *Arch*

597         *Womens Ment Health.* 16(3):197-201.

598   17.   Qiao M, Sun P, Wang H, Wang Y, Zhan X, Liu H, et al. (2017): Epidemiological

599         distribution and subtype analysis of premenstrual dysphoric disorder syndromes and

600         symptoms based on TCM theories. *Biomed Res Int.* 2017:4595016.

601   18.   Mei L, Ding X, Tsang SY, Pun FW, Ng SK, Yang J, et al. (2011): AluScan: a method

602         for genome-wide scanning of sequence and structure variations in the human genome.

603         *BMC Genomics.* 12:564.

604   19.   Yang JF, Ding XF, Chen L, Mat WK, Xu MZ, Chen JF, et al. (2014): Copy number

605         variation analysis based on AluScan sequences. *J Clin Bioinforma.* 4(1):15.

606   20.   American Psychiatric Association. (2013): Diagnostic and statistical manual of

607         mental disorders (Fifth Edition, DSM-5TM). American Psychiatric Publishing,

608         Washington, DC and London England. p. 171-174.

609   21.   McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al.

610         (2010): The Genome Analysis Toolkit: a MapReduce framework for analyzing next-

611         generation DNA sequencing data. *Genome Res.* 20(9):1297-1303.

612  22.  Hu T, Chen S, Ullah A, Xue H. (2019): AluScanCNV2: An R package for copy
613        number variation calling and cancer risk prediction with next-generation sequencing
614        data. *Genes Dis*. 6(1):43-46.

615  23.  Suzuki R, Shimodaira H. pvclust: Hierarchical Clustering with P-Values via
616        Multiscale Bootstrap Resampling. R package version 2.0-0. ed2015.

617  24.  Galili T. (2015): dendextend: an R package for visualizing, adjusting and comparing
618        trees of hierarchical clustering. *Bioinformatics*. 31(22):3718-3720.

619  25.  Carlson M, Maintainer BP. (2015): TxDb.Hsapiens.UCSC.hg19.knownGene:
620        Annotation package for TxDb object(s).

621  26.  Fresno C, Fernandez EA. (2013): RDAVIDWebService: a versatile R interface to
622        DAVID. *Bioinformatics*. 29(21):2810-2811.

623  27.  Long X, Xue H. (2020): Genetic-variant hotspots and hotspot clusters in the human
624        genome facilitating adaptation while increasing instability. *bioRxiv* doi:
625        https://doi.org/10.1101/2020.10.16.342188.

626  28.  Hansen RS, Thomas S, Sandstrom R, Canfield TK, Thurman RE, Weaver M, et al.
627        (2010): Sequencing newly replicated DNA reveals widespread plasticity in human
628        replication timing. *Proc Natl Acad Sci U S A*. 107(1):139-144.

629  29.  Ng SK, Hu T, Long X, Chan CH, Tsang SY, Xue H. (2016): Feature co-localization
630        landscape of the human genome. *Sci Rep*. 6:20650.

631  30.  Wickham H. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New
632        York; 2009.

633  31.  Kolde R. pheatmap: Pretty Heatmaps. 2015.

634   32.   Oosting J, Eilers P, Menezes R. quantsmooth: Quantile smoothing and genomic
635         visualization of array data. 2014.

636   33.   Robinson JT, Thorvaldsdottir H, Winckler W, Guttman M, Lander ES, Getz G, et al.
637         (2011): Integrative genomics viewer. *Nat Biotechnol.* 29(1):24-26.

638   34.   Macheret M, Halazonetis TD. (2018): Intragenic origins due to short G1 phases
639         underlie oncogene-induced DNA replication stress. *Nature.* 555(7694):112-116.

640   35.   Watanabe Y, Maekawa M. (2010): Spatiotemporal regulation of DNA replication in
641         the human genome and its association with genomic instability and disease. *Curr Med*
642         *Chem.* 17(3):222-233.

643   36.   Huo L, Straub RE, Roca C, Schmidt PJ, Shi K, Vakkalanka R, et al. (2007): Risk for
644         premenstrual dysphoric disorder is associated with genetic variation in ESR1, the
645         estrogen receptor alpha gene. *Biol Psychiat.* 62(8):925-933.

646   37.   Dubey N, Hoffman JF, Schuebel K, Yuan Q, Martinez PE, Nieman LK, et al. (2017):
647         The ESC/E(Z) complex, an effector of response to ovarian steroids, manifests an
648         intrinsic difference in cells from women with premenstrual dysphoric disorder. *Mol*
649         *Psychiatr.* 22(8):1172-1184.

650   38.   Makino C, Fujii Y, Kikuta R, Hirata N, Tani A, Shibata A, et al. (2003): Positive
651         association of the AMPA receptor subunit GluR4 gene (GRIA4) haplotype with
652         schizophrenia: linkage disequilibrium mapping using SNPs evenly distributed across
653         the gene region. *Am J Med Genet B Neuropsychiatr Genet.* 116B(1):17-22.

654

## Figure Legends

656   **Figure 1**. **Hierarchical clustering of PMDD samples based on their pairwise similarities**
657   **in genome wide CNV profiles.** For all 127 P-group samples, all CNVs were identified from

658     AluScan sequencing data with 100-kb non-overlapping scanning windows across the genome

659     and used in the plots of similarity scores for CNVGs (A) and CNVL (B), respectively. The

660     dendrograms on top of the heat maps were bootstrapped 1,000 times. The color of each

661     square in the heat map indicates the correlation coefficient (r) of a pair of samples according

662     to the blue-red thermal scale. The semi-supervised classification of samples based on (A)

663     CNVGs and (B) CNVLs was indicated by the red dendrogram branches for V-type and blue

664     ones for D-type genomes. The bands below the dendrograms and on the left-hand side of the

665     heat maps portrayed the subtyping of PMDD samples based on clinical symptoms, with

666     purple bands representing the clinically determined invasion subtype (n = 56) and orange

667     bands the depression subtype (n = 71). Each of the square diagonal boxes in panels (A) and

668     (B) enclosed a group of genomes with close correlations between each other in the group,

669     such that they could be identified as a coherent block of genomes belonging to either the V-

670     type or D-type CNV profiles depending on their enrichment in the invasion- or depression-

671     subtype samples (see 'Clustering of patient samples based on CNV profiles' in Methods).

672     Comparable heat maps obtained using sequence window sizes of 50 to 500 kb for CNV-

673     calling are shown in Figure S3.

674

675     **Figure 2. Occurrence frequencies of diagnostic CNV features and their prediction**

676     **accuracies for seven pairs of sample groups.** Panel (A) shows the frequency distribution of

677     diagnostic CNV features for different pairs of sample groups. The x-axis represents the

678     frequency of CNVs in the first-named group (Group 1 as shown on x-axis), and y-axis the

679     frequency of CNVs in the second-named group (Group 2 as shown on y-axis) in a given pair

680     of sample groups. Diagnostic CNV features with higher frequencies in Group 1 relative to

681     Group 2 (located in lower right crescent) are referred to as 'Group 1-favoring' features,

682     whereas diagnostic CNV features with higher frequencies in Group 2 relative to Group 1

683     (located in upper left crescent) are 'Group 2-favoring' features. Black circles are CNV

684     features selected using the frequency-based method with FDR < 0.01 (Fisher's exact tests),

685     and red triangles are CNV features selected using the correlation-based method. Panel (B)

686     shows the prediction accuracies (estimated using Eqn.2 in Methods) of sample classification

687     in seven sample-pairs based on CNV features selected using the correlation method. For each

688     of the seven pairs, prediction accuracy was estimated 1,000 times and the average accuracy

689     (Av.) was given in the pertinent panel. Subscript G denotes that the D- or V-type samples

690     were derived from the dendrogram of CNVGs (Figure 1A), while subscript L denotes that the

691     D- or V-type samples were derived from the dendrogram of CNVLs (Figure 1B).

692

693     **Figure 3. Overlaps between the diagnostic CNV features differentiating the two**

694     **subtypes of PMDD collectively and individually from the control.** CNV features identified

695     using (A) correlation-based method, and (B) frequency-based method. Circled 'G' indicates

696     CNVG features and circled 'L' indicates CNVL features. The '>' and '<' signs portray the

697     relative frequencies of the CNV features for a pair of sample groups, e.g. P>C represents

698     diagnostic CNV features that occurred in higher frequencies in P-group compared to C-

699     group. Subscript G denotes that the D- or V-type samples were derived from the CNVG

700     dendrogram in Figure 1A, whereas subscript L denotes that the D- or V-type samples were

701     derived from the CNVL dendrogram in Figure 1B.

702

703     **Figure 4. Distribution of frequency-based diagnostic CNV features among genomic**

704     **sequences of different DNA replication phases.** Number of base pairs of the CNV features

705     called using 50, 100 and 450-kb windows for (A) P-vs-C, (B) $D_G$-vs-C, (C) $V_G$-vs-C, and (D)

706     $D_G$-vs-$V_G$ groups. The solid bars represent CNVG features and hollow bars represent CNVL

707     features in each panel. The replication phases G1b to G2 are color coded as shown. The '>' or

708    '<' sign portrays larger or smaller frequencies of the CNV features in favor of the first-named

709    group over the second-named one. L/G represents the ratio of the number of CNVLs over the

710    number of CNVGs. Significant enrichment of CNV features in a particular replication phase

711    in the genome is indicated by asterisks that are color coded according to the replication

712    phase, or in black asterisks for comparison between an L/G value in the upper half of a panel

713    and an L/G value in the lower half (Bonferroni-corrected, *** $p < 0.005$, ** $p < 0.01$, * $p <$

714    0.05). Numerical $p$-values are shown in Table S14. Subscript G denotes that the D- or V-type

715    samples were derived from the CNVG dendrogram in Figure 1A. See Figure S7 for the

716    results obtained based on the D- or V-type samples derived from the CNVL dendrogram in

717    Figure 1B.

718

719    **Figure 5. Selected genes overlapping with frequency-based diagnostic CNV features.**

720    Expanded views of chromosomal segments on (A) chromosomes 2 and 7 for steroid

721    biosynthesis pathway genes, (B) chromosomes 5, 11, 12, 16 and 17 for *GRI*-genes of the

722    glutamatergic synapse and nicotine addiction pathways, and (C) chromosomes 6 and X for

723    the non-pathway *TRERF1* and *POF1B* genes with color-coded representation of the DNA

724    replication phase in the 'Phase' track, and aligned gene sequence(s) in blue (e.g. *UGT1A8* or

725    *TRERF1*) as described in RefSeq Genes in UCSC Genome Browser. Green rectangular boxes

726    either below the genes indicate the presence of diagnostic CNVG or CNVL feature(s). Inside

727    each box, colored stripes are indicative of CNVL features(s), and solid coloring is indicative

728    of CNVG features(s): purple for predominantly V-favoring features, orange for D-favoring

729    features, and green for C-favoring features.

730

731    **Figure 6. Enrichment analysis of genomic-feature contents in different replication**

732    **phases for control and PMDD subtypes.** Frequency-based CNV features diagnostic for C

733    group, i.e., control, as well as that for D and C groups of PMDD samples clustered by CNVG

734    dendrogram, identified with 100-kb scanning windows, were used in the analysis.

735    Enrichment analysis results were plotted for CNVG features in the upper two panels and that

736    for CNVL features in the bottom two panels. A similar analysis performed in parallel for

737    clustered by CNVL dendrogram can be found in Figure S8. Fold-change of each genomic

738    feature in the diagnostic CNV features relative to the non-diagnostic-CNV regions was

739    estimated according to 'Genomic-feature content of diagnostic CNV features in different

740    replication phases' in Methods, and was color-coded based on the thermal scale. Fold-change

741    greater than 2-fold was capped at 2 in the heat map. 'Group 1' indicated the first-named

742    group and 'Group 2' the second-named group in a given pair of samples. Genomic features

743    were grouped into Retrotransposon (SVAef, SVAcd, SVAab, AluYvy, AluYy,s AluS, AluJ,

744    FLAM, L1vy, L1y, L1m, L1o, MIR, L2), Genetic markers (RecD, RecH, RecK, GWAS,

745    ClinVar, GV hotspot, Cluster, CNVG), Regulatory sites (H3k27me3, H4k20me1, H3k9me1,

746    H2az, H3k79me2, H3k36me3, H3k4me3, H3k9ac, H3k4me2, H3k27ac, H3k4me1, H3k9me3,

747    MeMRE, MeDIP, MeBS, CpGi, CpGe, Me450, TFBS, REG, FAIRE, DNase) and

748    Gene/Transcription (Gene, EXPS, LRNA+, LRNA-, LINC) groups on the x-axis based on

749    their sequence and functional properties. The descriptions of genomic features and numeric

750    data were available in Table S11.

751

752    **Figure 7. Genetic distances between the two subtypes of PMDD and the control.**

753    Pairwise distances were estimated based on the abundance of diagnostic CNV features

754    between C-, D- and V-type genomes. The numbers of frequency-based CNV features were

755    employed as an approximate index of the genetic distance between the D-vs-C, V-vs-C or D-

756    vs-V sample pairs in Table S15, which comprised the 50-500 kb frequency-based CNV

757     features. Notably, the D-vs-V distance was larger than the D-vs-C distance or the V-vs-C

758     distance.

**Table 1.** Selected genes overlapping with 100-kb frequency-based CNVG and CNVL features[1] with adjusted *p*-values less than 0.01.

| Gene | Ratio | H%[2] | L%[3] | CNV | CNV location (x 100 kb)[4] | *p*-value[5] | CNV in replication phase (%)[6] | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | G1b | S1 | S2 | S3 | S4 | G2 |
| **GABA_A Receptor Family** (KEGG: Nicotine addiction) | | | | | | | | | | | | |
| *GABRR1, 2* | $C > D_G$ | 26 | 2 | L | 6:899-900 | 2.05E-04 | 0 | 0 | 99 | 0 | 0 | 0 |
| *GABRG3* | $C > D_G$ | 61 | 27 | L | 15:275-276 | 5.55E-04 | 0 | 0 | 0 | 0 | 0 | 100 |
| *GABRG3* | $V_G > D_G$ | 59 | 27 | L | 15:275-276 | 3.46E-03 | 0 | 0 | 0 | 0 | 0 | 100 |
| **Glutamate Metabotropic Receptor** (KEGG: Glutamatergic synapse) | | | | | | | | | | | | |
| *GRM4* | $D_L > V_L$ | 80 | 51 | L | 6:339-340 | 6.99E-03 | 0 | 87 | 0 | 0 | 0 | 0 |
| *GRM8* | $D_L > V_L$ | 91 | 45 | L | 7:1,267-1,268 | 1.08E-06 | 0 | 0 | 0 | 37 | 45 | 0 |
| *GRM5* | $V_G > D_G$ | 16 | 0 | L | 11:882-883 | 3.98E-03 | 0 | 90 | 0 | 0 | 0 | 0 |
| **Glutamate Ionotropic Receptor** (KEGG: Glutamatergic synapse, Nicotine addiction) | | | | | | | | | | | | |
| *GRIA1* | $V_G > D_G$ | 31 | 0 | L | 5:1,530-1,531 | 6.10E-06 | 0 | 0 | 0 | 12 | 78 | 0 |
| *GRIK2* | $V_G > D_G$ | 16 | 0 | L | 6:1,025-1,026 | 3.29E-03 | 0 | 0 | 0 | 0 | 0 | 100 |
| *GRIA4* | $V_G > D_G$ | 25 | 3 | L | 11:1,057-1,058 | 3.98E-03 | 0 | 0 | 0 | 0 | 100 | 0 |
| *GRIN2B* | $V_G > D_G$ | 23 | 0 | L | 12:138-139 | 2.56E-04 | 0 | 0 | 0 | 0 | 100 | 0 |
| *GRIN2A* | $V_G > D_G$ | 39 | 6 | L | 16:98-99 | 1.61E-04 | 0 | 0 | 0 | 77 | 16 | 0 |
| *GRIN2A* | $V_G > D_G$ | 26 | 5 | L | 16:99-100 | 6.84E-03 | 0 | 0 | 0 | 94 | 0 | 0 |
| *GRIN2C* | $V_G > D_G$ | 39 | 12 | L | 17:728-729 | 4.36E-03 | 3 | 92 | 0 | 0 | 0 | 0 |
| **UDP Glucuronosyltransferase 1 Family** (KEGG: Steroid hormone biosynthesis) | | | | | | | | | | | | |
| *UGT1A1, 3-10* | $V_L > D_L$ | 36 | 11 | L | 2:2,346-2,347 | 6.62E-03 | 0 | 0 | 98 | 0 | 0 | 0 |
| **Cytochrome P450** (KEGG: Steroid hormone biosynthesis) | | | | | | | | | | | | |
| *CYP3A4, 5, 7* | $C > D_G$ | 30 | 5 | G | 7:993-994 | 3.05E-04 | 0 | 100 | 0 | 0 | 0 | 0 |
| *CYP3A4, 5, 7* | $C > D_L$ | 30 | 9 | G | 7:993-994 | 6.15E-03 | 0 | 100 | 0 | 0 | 0 | 0 |
| *CYP3A4, 5, 7* | $V_L > D_L$ | 55 | 9 | G | 7:993-994 | 6.37E-07 | 0 | 100 | 0 | 0 | 0 | 0 |
| *CYP11B1, 2* | $V_L > D_L$ | 19 | 0 | G | 8:1,439-1,440 | 6.55E-04 | 0 | 1 | 68 | 0 | 0 | 0 |
| *CYP11A1* | $V_L > D_L$ | 51 | 22 | L | 15:746-747 | 7.80E-03 | 100 | 0 | 0 | 0 | 0 | 0 |
| **Premature Ovarian Failure Protein 1B** | | | | | | | | | | | | |
| *POF1B* | $C > P$ | 29 | 9 | G | X:845-846 | 2.13E-03 | 0 | 0 | 0 | 0 | 6 | 92 |
| *POF1B* | $C > D_G$ | 29 | 0 | G | X:845-846 | 2.19E-06 | 0 | 0 | 0 | 0 | 6 | 92 |
| *POF1B* | $V_G > D_G$ | 18 | 0 | G | X:845-846 | 1.20E-03 | 0 | 0 | 0 | 0 | 6 | 92 |
| *POF1B* | $D_G > C$ | 23 | 5 | G | X:846-847 | 2.94E-03 | 0 | 0 | 0 | 0 | 0 | 100 |
| *POF1B* | $D_G > V_G$ | 23 | 2 | G | X:846-847 | 1.67E-03 | 0 | 0 | 0 | 0 | 0 | 100 |
| **Transcriptional Regulating Factor 1** (Breast cancer anti-estrogen resistance 2) | | | | | | | | | | | | |
| *TRERF1* | $C > D_G$ | 20 | 0 | L | 6:423-424 | 5.55E-04 | 8 | 92 | 0 | 0 | 0 | 0 |
| *TRERF1* | $V_G > D_G$ | 16 | 0 | L | 6:423-424 | 3.98E-03 | 8 | 92 | 0 | 0 | 0 | 0 |
| *TRERF1* | $D_L > V_L$ | 42 | 15 | G | 6:424-425 | 7.63E-03 | 97 | 0 | 0 | 0 | 0 | 0 |

| | Ratio | [2] | [3] | | [4] | [5] | [6] | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Opioid Binding Protein/Cell Adhesion Molecule Like** (Hypogonadotropic Hypogonadism 14) | | | | | | | | | | | | |
| *OPCML* | $D_G > C$ | 100 | 79 | L | 11:1,331-1,332 | 3.17E-04 | 0 | 0 | 0 | 0 | 0 | 100 |
| *OPCML* | $D_G > V_G$ | 100 | 85 | L | 11:1,331-1,332 | 7.57E-03 | 0 | 0 | 0 | 0 | 0 | 100 |
| *OPCML* | $D_G > C$ | 61 | 24 | L | 11:1,332-1,333 | 1.07E-04 | 0 | 0 | 0 | 0 | 0 | 100 |
| *OPCML* | $D_G > V_G$ | 61 | 15 | L | 11:1,332-1,333 | 6.24E-06 | 0 | 0 | 0 | 0 | 0 | 100 |
| **MACRO Domain Containing 2** (Mono-ADP Ribosylhydrolase 2, Hypogonadotropic Hypogonadism 21) | | | | | | | | | | | | |
| *MACROD2* | $C > D_G$ | 21 | 2 | G | 20:158-159 | 7.43E-04 | 0 | 0 | 0 | 31 | 34 | 0 |
| *MACROD2* | $D_G > C$ | 97 | 77 | G | 20:151-152 | 1.45E-03 | 0 | 0 | 0 | 100 | 0 | 0 |
| *MACROD2* | $C > D_G$ | 64 | 35 | G | 20:144-145 | 2.07E-03 | 0 | 0 | 0 | 74 | 15 | 0 |
| *MACROD2* | $C > D_G$ | 64 | 24 | L | 20:156-157 | 2.82E-05 | 0 | 0 | 0 | 75 | 23 | 0 |
| *MACROD2* | $D_G > V_G$ | 97 | 62 | G | 20:151-152 | 8.43E-06 | 0 | 0 | 0 | 100 | 0 | 0 |
| *MACROD2* | $V_G > D_G$ | 31 | 2 | G | 20:158-159 | 2.79E-05 | 0 | 0 | 0 | 31 | 34 | 0 |
| *MACROD2* | $V_G > D_G$ | 21 | 2 | G | 20:145-146 | 2.03E-03 | 0 | 0 | 0 | 100 | 0 | 0 |
| *MACROD2* | $V_G > D_G$ | 64 | 35 | G | 20:144-145 | 6.56E-03 | 0 | 0 | 0 | 74 | 15 | 0 |
| *MACROD2* | $V_G > D_G$ | 74 | 24 | L | 20:156-157 | 1.17E-06 | 0 | 0 | 0 | 75 | 23 | 0 |
| *MACROD2* | $V_G > D_G$ | 15 | 0 | L | 20:146-147 | 7.25E-03 | 0 | 0 | 0 | 100 | 0 | 0 |

[1] See Table S7 for data on 50-500 kb CNV features;

[2] CNV frequency in first-named, higher frequency group (H), e.g. C-group of C>P pair;

[3] CNV frequency in second-named, lower frequency group (L), e.g. P-group of C>P pair;

[4] Chromosome number with start and end coordinates (to be multiplied by 100 kb);

[5] FDR-corrected *p*-value obtained using Fisher's exact test on counts of CNV features in the two groups compared, as specified in the 'Ratio' column;

[6] % base pairs in 6 replication phases; the % in genomic regions with unknown replication timing are not shown.

**Table 2.** Representative pathways enriched in 100-kb frequency-based CNV features.

| No. | Group[1] H | Group[1] L | CNV[2] | KEGG pathway[3] | Chromosome[4] | G1b | S1 | S2 | S3 | S4 | G2 | Proportion (%)[6] | p-value[7] |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | C | P | G | Steroid hormone biosynthesis | 2 | - | - | ++++++ | - | - | - | > 0 - 100 | [5.7E-04, 1.5E-09] |
| 2 | C | $V_G$ | G | Steroid hormone biosynthesis | 2 | - | - | ++++++ | - | - | - | > 0 - 100 | [2.7E-04, 9.3E-10] |
| 3 | C | $V_L$ | G | Steroid hormone biosynthesis | 2 | - | - | ++++++ | - | - | - | > 0 - 100 | [2.7E-04, 1.3E-09] |
| 4 | C | $D_G$ | G | Steroid hormone biosynthesis | 2,7 | - | ++ | ++++ | - | - | - | > 0 - 100 | [5.6E-05, 1.4E-11] |
| 5 | C | $D_L$ | G | Steroid hormone biosynthesis | 2,7 | - | ++ | ++++ | - | - | - | > 0 - 100 | [3.1E-05, 7.6E-12] |
| 6 | $V_L$ | $D_L$ | L | Steroid hormone biosynthesis | 2,15 | + | - | +++++ | - | - | - | > 0 - 100 | [4.8E-02, 6.7E-05] |
| 7 | $V_L$ | $D_L$ | G | Steroid hormone biosynthesis | 7,8,1 | - | +++ | ++ | - | - | - | 60 | [4.9E-02, 4.9E-02] |
| 8 | $D_G$ | $V_G$ | L | Ovarian steroidogenesis | 15,10,14,16,7 | +++ | + | + | + | - | - | 40 | [4.8E-02, 4.8E-02] |
| 9 | C | $D_G$ | L | Nicotine addiction | 6,11,12,15,16,17,5 | - | + | ++ | + | ++ | + | > 0 - 20 | [3.2E-02, 1.9E-02] |
| 10 | C | $D_L$ | L | Nicotine addiction | 6,11,12,15,16,17,5 | - | + | ++ | + | ++ | + | > 0 - 20 | [4.9E-02, 1.2E-02] |
| 11 | C | $D_G$ | L | Circadian entrainment | 11,17,12,16,20,5,8 | + | + | - | + | ++++ | + | 20 | [2.9E-02, 2.9E-02] |
| 12 | $V_G$ | $D_G$ | L | Serotonergic synapse | 1,11,15,12,17,20,21,5,6,7 | + | + | + | +++ | - | + | > 0 - 30 | [3.8E-02, 1.2E-02] |
| 13 | $V_G$ | $D_G$ | L | Glutamatergic synapse | 11,15,17,1,12,16,20,5,6 | + | ++ | + | + | ++ | + | > 0, 10 | [4.6E-02, 3.6E-02] |
| 14 | $V_G$ | $D_G$ | L | Nicotine addiction | 15,11,12,16,17,5 | - | + | - | ++ | +++ | + | 10 | [4.6E-02, 4.6E-02] |
| 15 | $V_G$ | $D_G$ | L | Dopaminergic synapse | 1,11,12,16,17,2,20,21,4,5,7,8 | + | - | ++ | ++ | ++ | + | > 0 - 20 | [4.4E-02, 3.0E-02] |
| 16 | $V_G$ | $D_G$ | L | Circadian entrainment | 1,11,17,12,16,20,21,4,5 | + | + | ++ | + | ++ | + | > 0 - 20 | [3.0E-02, 1.7E-02] |
| 17 | $V_G$ | $D_G$ | L | cAMP signaling pathway | 1,11,5,10,12,16,17,3,4,6,7 | + | + | + | ++ | ++ | - | 20 | [3.3E-02, 3.3E-02] |

[1] Significant difference in CNV frequencies between compared groups, with 'H' and 'L' indicating higher- and lower-frequency group respectively. The subscripts 'G' and 'L' indicate sample groups clustered based on CNVG and CNVL respectively;

[2] 'G' indicates copy-number-gains and 'L' indicates copy-number-losses;

[3] KEGG pathway IDs are, in order of appearance in table, hsa00140, hsa04913, hsa05033, hsa04713, hsa04726, hsa04724, hsa04728 and hsa04024;

[4] Chromosomes where pathway genes overlapped with CNV feature(s);

[5] Approximate distribution of pathway genes in different replication phases, with '-' indicating 0%, and one '+' indicating 0-20%, up to six '+' indicating 100%;

[6] Proportion of gene sequence overlapping with the CNV feature(s) ranging from > 0% to 100%;

[7] Range of Benjamini-adjusted p-values of pathway enrichment pertaining to bottom and top figures referred to in footnote 6.

**Table 3.** Selected genomic features in 100-kb frequency-based diagnostic CNV features with fold-change greater than 1 in at least one replication phase(s).

| Group[1] H | L | CNV[2] | Fold-change in different replication phases[3] G1b | S1 | S2 | S3 | S4 | G2 |
|---|---|---|---|---|---|---|---|---|
| *CpGi (CpG island)* | | | | | | | | |
| C | P | G | -0.32 | 0.06 | -0.06 | **1.62** | -0.18 | **1.15** |
| C | P | L | -0.18 | 0.05 | 0.39 | 0.46 | **1.10** | **2.20** |
| C | $D_L$ | G | -0.45 | 0.02 | -0.08 | 0.59 | **1.18** | 0.57 |
| C | $V_G$ | L | 0.14 | 0.06 | 0.69 | 0.71 | 0.47 | **1.42** |
| P | C | L | -0.34 | 0.12 | -0.02 | 0.38 | -0.23 | **1.08** |
| *Me450 (Methylation status using HumanMethylation450)* | | | | | | | | |
| C | P | G | -0.51 | 0.54 | **1.04** | 0.88 | -0.01 | 0.43 |
| C | P | L | -0.29 | 0.23 | -0.02 | 0.33 | **1.16** | **1.05** |
| C | $V_L$ | G | **1.00** | 0.04 | -0.02 | -0.35 | -0.48 | -0.25 |
| *MeBS (cytosine methylation using bisulfite sequencing)* | | | | | | | | |
| C | P | L | -0.12 | 0.45 | 0.12 | 0.49 | **1.53** | **1.55** |
| C | $D_L$ | G | -0.53 | 0.10 | 0.06 | 0.64 | **1.65** | 0.59 |
| C | $V_G$ | L | 0.27 | 0.26 | 0.69 | 0.94 | **1.51** | 0.29 |
| $D_G$ | C | L | **1.03** | 0.54 | 0.35 | 0.24 | **1.20** | -0.39 |
| $D_L$ | C | L | **1.38** | 0.20 | 0.25 | 0.37 | 0.13 | -0.34 |
| $D_G$ | $V_G$ | L | **1.15** | 0.41 | 0.37 | 0.46 | **1.69** | -0.24 |
| $D_L$ | $V_L$ | L | 0.82 | 0.45 | 0.53 | 0.59 | **1.47** | -0.15 |
| $V_L$ | $D_L$ | G | 0.15 | -0.17 | -0.14 | 0.45 | **1.35** | **1.21** |
| *MeMRE (Methylation using MRE-* | | | | | | | | |
| C | P | G | -0.26 | 0.23 | -0.30 | **2.00** | -0.12 | -0.19 |
| C | P | L | -0.23 | -0.06 | 0.31 | 0.47 | **1.44** | 0.90 |
| P | C | L | -0.13 | 0.10 | 0.10 | 0.33 | -0.22 | **1.70** |
| *DNase (DNase I hypersensitive sites)* | | | | | | | | |
| C | P | L | -0.12 | 0.02 | 0.11 | 0.55 | 0.46 | **1.02** |
| C | $V_L$ | G | 0.22 | -0.30 | 0.25 | -0.55 | 0.04 | **1.14** |
| P | C | L | -0.10 | 0.08 | 0.51 | 0.38 | **1.80** | -0.13 |
| $D_G$ | C | L | 0.28 | 0.27 | 0.24 | -0.05 | **1.37** | -0.52 |
| $D_L$ | C | L | 0.47 | 0.22 | 0.39 | -0.12 | **1.77** | -0.14 |
| $V_G$ | C | L | -0.16 | 0.22 | 0.05 | 0.00 | **1.38** | -0.57 |
| $V_L$ | C | L | -0.15 | 0.14 | 0.25 | 0.09 | **1.08** | -0.66 |
| *FAIRE (Formaldehyde-assisted isolation of regulatory elements)* | | | | | | | | |
| C | P | L | -0.18 | 0.17 | -0.10 | 0.15 | 0.71 | **2.69** |
| C | $V_G$ | L | -0.23 | 0.12 | 0.03 | 0.14 | 0.46 | **2.97** |
| C | $V_L$ | L | 0.02 | 0.08 | -0.04 | 0.04 | 0.33 | **2.54** |
| P | C | L | -0.09 | 0.16 | 0.53 | -0.03 | 0.86 | **1.56** |
| $D_G$ | C | L | 0.32 | 0.28 | 0.16 | -0.17 | **1.55** | 0.14 |
| $D_L$ | C | L | 0.34 | 0.29 | 0.33 | -0.24 | **2.11** | **1.25** |
| $D_G$ | $V_G$ | L | 0.22 | 0.17 | 0.11 | 0.11 | **1.15** | 0.92 |
| $D_L$ | $V_L$ | L | 0.20 | 0.21 | 0.02 | 0.03 | 0.77 | **2.08** |

| Group[1] H | L | CNV[2] | Fold-change in different replication phases[3] G1b | S1 | S2 | S3 | S4 | G2 |
|---|---|---|---|---|---|---|---|---|
| *LINC (Large intergenic non-coding RNA)* | | | | | | | | |
| *Continued* | | | | | | | | |
| C | P | G | -0.97 | -0.52 | **18.30** | **4.23** | 0.32 | -0.11 |
| C | P | L | -0.99 | **1.05** | 0.17 | **1.02** | -0.73 | 0.92 |
| C | $V_L$ | G | 0.11 | -0.96 | 0.03 | **1.75** | -0.55 | **2.42** |
| C | $V_G$ | L | -1.00 | -0.77 | -0.79 | -0.93 | -0.47 | **1.59** |
| C | $V_L$ | L | -0.98 | -0.74 | -0.72 | -0.87 | -0.50 | **1.06** |
| $V_G$ | C | L | 0.48 | -0.70 | 0.32 | **2.14** | **1.61** | **1.36** |
| $V_L$ | C | L | 0.31 | -0.75 | -0.02 | **1.80** | **1.01** | 0.74 |
| $V_G$ | $D_G$ | L | 0.37 | -0.67 | **1.19** | 0.52 | 0.33 | **1.34** |
| *RecD (Sex-averaged rates of recombination)* | | | | | | | | |
| P | C | L | **1.31** | 0.39 | 0.26 | 0.44 | 0.36 | 0.97 |
| $D_G$ | C | L | 0.23 | 0.16 | 0.43 | 0.42 | 0.35 | **1.24** |
| $D_L$ | C | L | 0.40 | 0.30 | 0.52 | 0.50 | 0.51 | **1.56** |
| $D_L$ | $V_L$ | L | 0.25 | 0.13 | 0.24 | 0.22 | 0.45 | **1.06** |
| *GWAS (GWAS-identified SNPs)* | | | | | | | | |
| C | $V_G$ | L | **2.08** | -0.01 | -0.02 | -0.35 | 0.18 | 0.60 |
| C | $V_L$ | L | **1.63** | -0.13 | -0.01 | -0.23 | 0.10 | 0.38 |
| P | C | G | -0.26 | **5.48** | -0.38 | 0.27 | -0.26 | -0.34 |
| $D_G$ | C | G | -0.16 | **1.61** | -0.13 | -0.28 | -0.05 | -0.01 |
| $D_L$ | C | G | -0.08 | **1.78** | -0.17 | -0.22 | -0.03 | -0.16 |
| $V_G$ | C | L | -0.34 | -0.20 | -0.22 | -0.23 | **2.78** | -0.14 |
| $V_L$ | C | L | -0.28 | -0.37 | -0.29 | -0.17 | **2.00** | 0.01 |
| $D_G$ | $V_G$ | G | -0.16 | **1.25** | -0.02 | 0.33 | -0.28 | -0.25 |
| $D_L$ | $V_L$ | G | -0.08 | **1.27** | -0.03 | 0.24 | -0.24 | -0.11 |
| *GV hotspot (Density-based genetic-variant hotspot)* | | | | | | | | |
| C | $V_G$ | G | **1.97** | -0.34 | -0.21 | -0.27 | -0.15 | **1.22** |
| C | $V_L$ | G | **1.63** | 0.01 | -0.27 | -0.09 | -0.28 | 0.69 |
| P | C | L | 0.12 | -0.01 | -0.11 | 0.70 | **1.58** | 0.29 |
| $D_G$ | $V_G$ | G | -0.03 | -0.23 | -0.44 | 0.14 | 0.20 | **1.03** |
| $D_L$ | $V_L$ | G | 0.01 | -0.31 | -0.43 | 0.28 | 0.00 | **1.06** |
| $D_G$ | $V_G$ | L | 0.30 | 0.27 | 0.09 | 0.42 | 0.61 | **1.22** |
| *Cluster (Cluster of genetic-variant hotspots)* | | | | | | | | |
| C | $D_L$ | L | 0.03 | -0.84 | 0.38 | -0.87 | -0.30 | **2.42** |
| C | $V_G$ | G | **3.68** | -1.00 | -1.00 | 0.42 | -0.15 | -0.29 |
| C | $V_L$ | G | **2.90** | -1.00 | -1.00 | 0.37 | -0.21 | -0.52 |
| C | $V_G$ | L | -0.23 | **1.42** | 0.03 | 0.24 | 0.26 | 0.83 |
| P | C | L | -1.00 | -1.00 | -0.27 | **4.33** | **4.26** | -1.00 |
| $D_G$ | C | L | **2.15** | 0.73 | -0.05 | 0.57 | **1.34** | -1.00 |
| $D_L$ | C | L | **2.46** | 0.90 | 0.07 | 0.73 | **1.11** | -1.00 |

36

***LINC (Large intergenic non-coding RNA)***

| | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| C | $D_G$ | G | 0.50 | -0.64 | **7.56** | **1.38** | **1.46** | 0.11 | $V_G$ | C | G | -1.00 | -0.26 | **1.48** | 0.46 | 0.12 | -0.83 |
| C | $D_L$ | G | -0.98 | -0.76 | **8.90** | **1.53** | **1.51** | 0.16 | $V_L$ | C | G | -1.00 | **1.45** | 0.72 | 0.19 | -0.08 | -0.38 |
| C | $D_G$ | L | -0.86 | 0.44 | 0.46 | -0.79 | -0.88 | **3.10** | $V_G$ | C | L | -0.41 | -0.72 | -0.70 | **1.68** | 0.46 | -0.97 |
| $D_G$ | C | G | **1.84** | -0.92 | -0.68 | 0.21 | -0.99 | -0.80 | $V_L$ | C | L | -0.21 | -0.78 | -0.41 | **1.20** | 0.27 | -0.89 |
| $D_G$ | C | L | -0.85 | -0.84 | -0.95 | **1.89** | -1.00 | -1.00 | $D_G$ | $V_G$ | G | 0.02 | -0.65 | -0.55 | -0.50 | **1.52** | 0.89 |
| $D_L$ | C | L | -0.81 | -0.82 | -0.95 | **1.85** | -1.00 | -1.00 | $D_L$ | $V_L$ | G | -0.05 | -0.64 | -0.55 | -0.53 | **1.41** | 0.95 |
| $V_G$ | $D_G$ | G | 0.87 | -0.74 | **4.34** | -0.66 | 0.81 | 0.16 | $D_G$ | $V_G$ | L | 0.91 | 0.42 | 0.43 | 0.64 | **1.71** | **1.25** |
| $V_L$ | $D_L$ | G | 0.76 | -0.75 | **3.45** | -0.69 | 0.65 | -0.08 | $V_G$ | $D_G$ | L | -0.64 | -0.36 | **3.00** | 0.50 | 0.12 | 0.20 |
| | | | | | | | | | $V_L$ | $D_L$ | L | -0.44 | -0.52 | -0.15 | 0.41 | -0.38 | **1.07** |

[1] Significant difference in CNV frequencies between compared groups, with 'H' and 'L' indicating higher- and lower-frequency group respectively. The subscripts 'G' and 'L' indicate sample groups clustered based on CNVG and CNVL respectively;

[2] 'G' indicates copy-number-gains and 'L' indicates copy-number-losses;

[3] Fold-change (> 1-fold in bold) of genomic feature density or intensity in diagnostic CNV features relative to non-diagnostic-CNV regions in replication phase.

**A** Clustered by 100-kb CNVG

**B** Clustered by 100-kb CNVL

$V_G = 61$  $D_G = 66$

$V_L = 53$  $D_L = 74$

**A** **Correlation-based CNV features**

Clustered by CNVG ($V_G = 61$, $D_G = 66$)   Clustered by CNVL ($V_L = 53$, $D_L = 74$)

*Patient > Control*

Ⓖ  P>C  $D_G$>C  $V_G$>C: 10, 10, 26, 0, 5, 0, 8

Ⓛ  P>C  $D_G$>C  $V_G$>C: 7, 2, 10, 0, 8, 0, 10

Ⓖ  P>C  $D_L$>C  $V_L$>C: 10, 13, 20, 0, 2, 0, 7

Ⓛ  P>C  $D_L$>C  $V_L$>C: 8, 6, 10, 0, 3, 0, 9

*Control > Patient*

Ⓖ  C>P  C>$D_G$  C>$V_G$: 10, 4, 9, 0, 3, 0, 11

Ⓛ  C>P  C>$D_G$  C>$V_G$: 27, 6, 26, 1, 4, 0, 14

Ⓖ  C>P  C>$D_L$  C>$V_L$: 10, 4, 7, 2, 1, 0, 13

Ⓛ  C>P  C>$D_L$  C>$V_L$: 31, 5, 14, 1, 1, 0, 15

**B** **Frequency-based CNV features**

Clustered by CNVG ($V_G = 61$, $D_G = 66$)   Clustered by CNVL ($V_L = 53$, $D_L = 74$)

*Patient > Control*

Ⓖ  P>C  $D_G$>C  $V_G$>C: 4, 155, 479, 11, 79, 0, 204

Ⓛ  P>C  $D_G$>C  $V_G$>C: 2, 58, 229, 11, 56, 0, 300

Ⓖ  P>C  $D_L$>C  $V_L$>C: 3, 154, 445, 16, 76, 0, 298

Ⓛ  P>C  $D_L$>C  $V_L$>C: 2, 60, 204, 11, 54, 0, 383

*Control > Patient*

Ⓖ  C>P  C>$D_G$  C>$V_G$: 17, 117, 322, 37, 24, 0, 126

Ⓛ  C>P  C>$D_G$  C>$V_G$: 74, 190, 411, 60, 91, 0, 236

Ⓖ  C>P  C>$D_L$  C>$V_L$: 20, 108, 262, 40, 27, 0, 168

Ⓛ  C>P  C>$D_L$  C>$V_L$: 66, 183, 374, 65, 101, 0, 371

**A** P > C features

50-kb   L/G = 0.75
100-kb   L/G = 0.51
450-kb   L/G = 0.64

Base pair (x 10⁷)

L/G = 1.52
L/G = 2.13
L/G = 2.05

C > P features

**B** D_G > C features

50-kb   L/G = 0.36
100-kb   L/G = 0.46
450-kb   L/G = 0.79

Base pair (x 10⁷)

L/G = 1.13
L/G = 1.39
L/G = 1.21

C > D_G features

CNV type:
CNVG: solid column
CNVL: open column

**C** V_G > C features

50-kb   L/G = 2.41
100-kb   L/G = 1.25
450-kb   L/G = 1.03

Base pair (x 10⁷)

L/G = 1.12
L/G = 2.07
L/G = 1.93

C > V_G features

**D** D_G > V_G features

50-kb   L/G = 0.56
100-kb   L/G = 0.75
450-kb

Base pair (x 10⁷)

L/G = 1.80
L/G = 1.24

V_G > D_G features

Replication Phase
G1b
S1
S2
S3
S4
G2

**A**

**Chromosome 2: Steroid hormone biosynthesis**

234,540 kb    234,580 kb    234,620 kb    234,660 kb

Phase

*UGT1A8*
*UGT1A10*
*UGT1A9*
*UGT1A7*
*UGT1A6*
*UGT1A5*
*UGT1A4*
*UGT1A3*
*UGT1A1*

**CNVL for $V_L > D_L$** ➜

**CNVG for $C > D_G$, $C > V_G$, $C > P$, $C > D_L$, $C > V_L$** ➜

**Chromosome 7: Steroid hormone biosynthesis**

99,200 kb    99,300 kb    99,400 kb    99,500 kb

Phase

*CYP3A5*                                    *CYP3A4*

*CYP3A7-CYP3AP1*

*CYP3A7*

➜ **CNVG for $V_L > D_L$, $C > D_G$, $C > D_L$**

Replication Phase:

G1b
S1
S2
S3
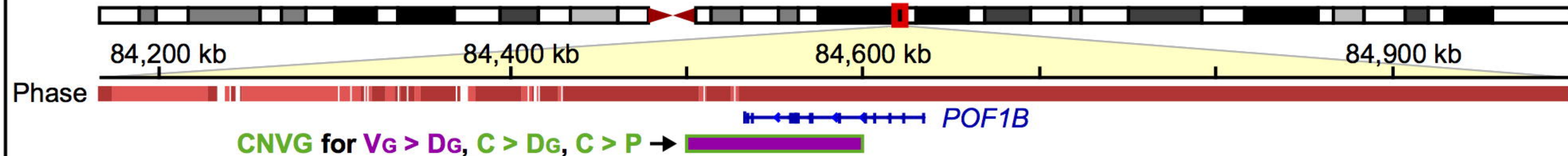S4
G2

**B**

**Chromosome 5: Glutamatergic synapse, Nicotine addiction**

152,900 kb   153,000 kb   153,100 kb   153,200 kb

Phase

*GRIA1*

CNVG for **D$_G$ > V$_G$**, **D$_G$ > C**, **P > C**, **D$_L$ > V$_L$**, **D$_L$ > C**    CNVL for **V$_G$ > D$_G$**, **C > D$_G$**, **V$_L$ > D$_L$**, **C > D$_L$**    CNVG for **V$_G$ > C**

**Chromosome 11: Glutamatergic synapse, Nicotine addiction**

105,400 kb   105,500 kb   105,600 kb   105,700 kb   105,800 kb   105,900 kb

Phase

*GRIA4*

CNVL for **C > D$_G$**, **C > V$_G$**, **C > P**, **C > D$_L$**, **C > V$_L$** →    ← CNVG for **V$_G$ > D$_G$**

**Chromosome 12: Glutamatergic synapse, Nicotine addiction**

13,600 kb   13,700 kb   13,800 kb   13,900 kb   14,000 kb   14,100 kb   14,200 kb

Phase

*GRIN2B*

← CNVL for **V$_G$ > D$_G$**, **C > D$_G$**, **V$_L$ > D$_L$**, **C > D$_L$**

Replication Phase:
G1b
S1
S2
S3
S4
G2

**Chromosome 16: Glutamatergic synapse, Nicotine addiction**

9,900 kb   10,000 kb   10,100 kb   10,200 kb

Phase

*GRIN2A*

CNVG for **D$_G$ > V$_G$** →

CNVL for **V$_G$ > D$_G$**, **C > D$_G$**, **V$_L$ > D$_L$**    CNVL for **C > D$_G$**, **C > V$_G$**, **C > P**, **C > D$_L$**, **C > V$_L$**

CNVL for **V$_G$ > D$_G$**, **C > D$_G$**, **C > D$_L$**, **C > P**

**Chromosome 17: Glutamatergic synapse, Nicotine addiction**

72,780 kb   72,820 kb   72,860 kb   72,900 kb   72,940 kb

Phase

*GRIN2C*

CNVL for **V$_G$ > D$_G$**, **C > D$_G$**, **C > D$_L$**, **C > P**

**C**

**Chromosome 6:**

41,900 kb · 42,100 kb · 42,300 kb · 42,500 kb · 42,700 kb

Phase

*TRERF1*

**CNVL for VG > DG, VL > DL, C > DG, C > DL** → ← **CNVG for DL > VL**

**Chromosome X:**

84,200 kb · 84,400 kb · 84,600 kb · 84,900 kb

Phase

*POF1B*

**CNVG for VG > DG, C > DG, C > P** →

Replication Phase:
- G1b
- S1
- S2
- S3
- S4
- G2