

1 Prevalence and architecture of developmental disorders caused by *de*  
2 *novo* mutation

3  
4 The Deciphering Developmental Disorders Study

5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24

25 **Abbreviations**

26 PTV: Protein-Truncating Variant  
27 DNM: *De Novo* Mutation  
28 DD: Developmental Disorder  
29 DDD: Deciphering Developmental Disorders study  
30

31 **Key Words**

32 *De novo* mutation; Developmental Disease; Seizures; Intellectual Disability; PhenIcons; Average  
33 Faces; ANKRD11; ARID1B; KMT2A; DDX3X; ADNP; MED13L; DYRK1A; EP300; SCN2A; SETD5;  
34 KCNQ2; MECP2; SYNGAP1; ASXL3; SATB2; TCF4; CDK13; CREBBP; DYNC1H1; FOXP1; PPP2R5D;  
35 PURA; CTNBN1; KAT6A; SMARCA2; STXBP1; EHMT1; ITPR1; KAT6B; NSD1; SMC1A; TBL1XR1;  
36 CASK; CHD2; CHD4; HDAC8; USP9X; WDR45; AHDC1; CSNK2A1; GNAI1; GNAO1; HNRNPU;  
37 KANSL1; KIF1A; MEF2C; PACS1; SLC6A1; CNOT3; CTCF; EEF1A2; FOXG1; GATAD2B; GRIN2B;  
38 IQSEC2; POGZ; PUF60; SCN8A; TCF20; BCL11A; BRAF; CDKL5; NFIX; PTPN11; AUTS2; CHAMP1;  
39 CNKSR2; DNM1; KCNH1; NAA10; PPM1D; ZBTB18; ZMYND11; ASXL1; COL4A3BP; KCNQ3; MSL3;

40 MYT1L; PDHA1; PPP2R1A; SMAD4; TRIO; WAC; CHD8; GABRB3; KDM5B; PTEN; QRI1; SET;  
41 ZC4H2; ALG13; SCN1A; SUV420H1; SLC35A2  
42  
43

## 44 Summary

45 Individuals with severe, undiagnosed developmental disorders (DDs) are enriched for damaging  
46 *de novo* mutations (DNMs) in developmentally important genes. We exome sequenced 4,293  
47 families with individuals with DDs, and meta-analysed these data with another 3,287 individuals  
48 with similar disorders. We show that the most significant factors influencing the diagnostic yield  
49 of DNMs are the sex of the affected individual, the relatedness of their parents, whether close  
50 relatives are affected and parental ages. We identified 94 genes enriched for damaging DNMs,  
51 including 14 without previous compelling evidence. We have characterised the phenotypic  
52 diversity among these disorders. We estimate that 42% of our cohort carry pathogenic DNMs in  
53 coding sequences, and approximately half disrupt gene function, with the remainder resulting  
54 in altered-function. We estimate that developmental disorders caused by DNMs have an  
55 average birth prevalence of 1 in 213 to 1 in 448, depending on parental age. Given current  
56 global demographics, this equates to almost 400,000 children born per year.

## 57 Main text

58 Approximately 2-5% of children are born with major congenital malformations and/or manifest  
59 severe neurodevelopmental disorders during childhood<sup>1,2</sup>. While diverse mechanisms can cause  
60 such developmental disorders, including gestational infection and maternal alcohol  
61 consumption, damaging genetic variation in developmentally important genes has a major  
62 contribution. Several recent studies have identified a substantial causal role for DNMs not  
63 present in either parent<sup>3-16</sup>. Despite the identification of many developmental disorders caused  
64 by DNMs, it is generally accepted that many more such disorders await discovery<sup>15</sup>, and the  
65 overall contribution of DNMs to developmental disorders is not known. Moreover, some  
66 pathogenic DNMs completely ablate the function of the encoded protein, whereas others alter  
67 the function of the encoded protein<sup>17</sup>; the relative contributions of these two mechanistic  
68 classes is also not known.

69  
70 We recruited 4,293 individuals to the Deciphering Developmental Disorders (DDD) study<sup>15</sup> via  
71 genetics services of the UK National Health Service and Republic of Ireland. Each of these  
72 individuals was referred with a severe undiagnosed developmental disorder and most were the  
73 only affected family member. Most (81%) individuals had been screened for large pathogenic  
74 deletions and duplications. We systematically phenotyped these individuals and sequenced the  
75 exomes of these individuals and their parents. Growth measurements, family history, and  
76 developmental milestones were collected, and detailed clinical phenotypes were captured  
77 using Human Phenotype Ontology (HPO) terms. Analyses of 1,133 of these trios were described  
78 previously<sup>15,18</sup>. We generated a high sensitivity set of 8,361 candidate DNMs in coding or  
79 splicing sequence (mean of 1.95 DNMs per proband), while removing systematic erroneous  
80 calls (Supplementary Table 1). This rate of candidate DNMs per proband is higher than other  
81 studies<sup>3-15</sup>, because we wish to maintain high sensitivity, and can address lower specificity via  
82 subsequent validation. 1,624 genes contained two or more DNMs in unrelated individuals.

83  
84 Twenty-three percent of individuals had likely pathogenic protein-truncating or missense DNMs  
85 within the clinically curated set of genes robustly associated with dominant developmental

86 disorders<sup>18</sup>. We investigated factors associated with whether an individual had a likely  
87 pathogenic DNM in these curated genes (Figure 1a, b, Supplementary Table 1). We observed  
88 that males had a lower chance of carrying a likely pathogenic DNM ( $P = 1.6 \times 10^{-4}$ ; OR 0.75, 0.65  
89 - 0.87 95% CI), as has also been observed in autism<sup>19</sup>. We also observed increased likelihood of  
90 having a pathogenic DNM with the extent of speech delay ( $P = 0.00115$ ), but not other  
91 indicators of severity relative to the rest of the cohort. Individuals with other affected family  
92 members were less likely to have pathogenic DNMs (affected siblings:  $P = 7.3 \times 10^{-18}$ , affected  
93 parents:  $P = 5.7 \times 10^{-9}$ ), and individuals who were from self-declared consanguineous unions  
94 were less likely to have a pathogenic DNM ( $P = 8.0 \times 10^{-11}$ ). Furthermore, the total genomic  
95 extent of autozygosity (due to parental relatedness) was negatively correlated with the  
96 likelihood of having a pathogenic DNM ( $P = 1.7 \times 10^{-7}$ ), for every  $\log_{10}$  increase in autozygous  
97 length, the probability of having a pathogenic DNM dropped by 7.5%, likely due to increasing  
98 burden of recessive causation (Figure 1c). Nonetheless, 6% of individuals with autozygosity  
99 equivalent to a first cousin union or greater had a plausibly pathogenic DNM, underscoring the  
100 importance of considering *de novo* causation in all families.

101  
102 Paternal age has been shown to be the primary factor influencing the number of DNMs in a  
103 child<sup>20,21</sup>, and thus is expected to be a risk factor for pathogenic DNMs. Paternal age was only  
104 weakly associated with likelihood of having a pathogenic DNM ( $P = 0.016$ ). However, focusing  
105 on the minority of DNMs that were truncating and missense variants in known DD-associated  
106 genes limits our power to detect such an effect. Analysing all 8,409 high confidence exonic and  
107 intronic autosomal DNMs confirmed a strong paternal age effect ( $P = 1.4 \times 10^{-10}$ , 1.53  
108 DNMs/year, 1.07-2.01 95% CI), as well as highlighting a weaker, independent, maternal age  
109 effect ( $P = 0.0019$ , 0.86 DNMs/year, 0.32-1.40 95% CI, Figure 1d,e), as has recently been  
110 described in whole genome analyses<sup>22</sup>. These genome-wide estimates were scaled from exome-  
111 based estimates, of 0.0306 DNMs/year paternal effect and 0.0172 DNMs/year maternal effect.

112  
113 We identified genes significantly enriched for damaging DNMs by comparing the observed  
114 gene-wise DNM count to that expected under a null mutation model<sup>23</sup>, as described  
115 previously<sup>15</sup>. We combined this analysis with 4,224 published DNMs in 3,287 affected  
116 individuals from thirteen exome or genome sequencing studies (Supplementary Table 2)<sup>3-14</sup> that  
117 exhibited a similar excess of DNMs in our curated set of DD-associated genes (Extended Data  
118 Figure 1). We found 93 genes with genome-wide significance ( $P < 5 \times 10^{-7}$ , Figure 2), 80 of which  
119 had prior evidence of DD-association (Supplementary Table 3). We have developed visual  
120 summaries of the phenotypes associated with each gene to facilitate clinical use. In addition,  
121 we created anonymised average face images from individuals with DNMs in genome-wide  
122 significant genes (Figure 2) from ordinary (2D) clinical photos using previously validated  
123 software<sup>24</sup>. These images highlight facial dysmorphologies specific to certain genes. After  
124 careful review by two experienced clinical geneticists, average face images for twelve genes  
125 were determined to be truly anonymised and of sufficient quality. To assess any increase in  
126 power to detect novel DD-associated genes, we excluded individuals with likely pathogenic  
127 variants in known DD-associated genes<sup>15</sup>, leaving 3,158 probands from our cohort, along with  
128 2,955 probands from the meta-analysis studies. In this subset, fourteen genes for which no  
129 statistically-compelling prior evidence for DD causation was available achieved genome-wide

130 significance: *CDK13*, *CHD4*, *CNOT3*, *CSNK2A1*, *GNAI1*, *KCNQ3*, *MSL3*, *PPM1D*, *PUF60*, *QRICH1*,  
131 *SET*, *SUV420H1*, *TCF20*, and *ZBTB18* ( $P < 5 \times 10^{-7}$ , Table 1, Extended Data Figure 4). The clinical  
132 features associated with these newly confirmed disorders are summarised in Extended Data  
133 Figure 2, Extended Data Figure 3 and Supplementary Information. *QRICH1* would not achieve  
134 genome-wide significance without excluding individuals with likely pathogenic variants in DD-  
135 associated genes. In addition to discovering novel DD-associated genes, we identified several  
136 new disorders linked to known DD-associated genes, but with different modes of inheritance or  
137 molecular mechanisms. We found *USP9X* and *ZC4H2* had a genome-wide significant excess of  
138 DNMs in female probands, indicating these genes have X-linked dominant modes of inheritance  
139 in addition to previously reported X-linked recessive mode of inheritance in males<sup>25,26</sup>. In  
140 addition, we found truncating mutations in *SMC1A* were strongly associated with a novel  
141 seizure disorder ( $P = 6.5 \times 10^{-19}$ ), while in-frame/missense mutations in *SMC1A* with dominant  
142 negative effects<sup>27</sup> are a known cause of Cornelia de Lange Syndrome (CdLS). Individuals with  
143 truncating mutations in *SMC1A* lacked the characteristic facial dysmorphology of CdLS.

144  
145 We then explored two approaches for integrating phenotypic data into disease gene  
146 association: statistical assessment of Human Phenotype Ontology (HPO) term similarity  
147 between individuals sharing candidate DNMs in the same gene (as we described previously<sup>28</sup>)  
148 and phenotypic stratification based on specific clinical characteristics. Combining genetic  
149 evidence and HPO term similarity increased the significance of some known DD-associated  
150 genes. However, significance decreased for a larger number of genes causing severe DD but  
151 associated with non discriminative HPO terms (Extended Data Figure 5a). Although we did not  
152 incorporate categorical phenotypic similarity in the gene discovery analyses described above,  
153 the systematic acquisition of phenotypic data on affected individuals within DDD enabled  
154 aggregate representations to be created for each gene achieving genome-wide significance. We  
155 present these in the form of icon-based summaries of growth and developmental milestones  
156 (PhenIcons), heatmaps of the recurrently coded HPO terms and, where photos for at least ten  
157 children with mutations in the same gene were available, an anonymised average facial  
158 representation (Supplementary Information).

159  
160 Twenty percent of individuals had HPO terms which indicated seizures and/or epilepsy. We  
161 compared analysis within this phenotypically stratified group with gene-wise analyses of the  
162 entire cohort, to see if it increased power to detect known seizure-associated genes (Extended  
163 Data Figure 5b). Fifteen seizure-associated genes were genome-wide significant in both the  
164 seizure-only and the entire-cohort analyses. Nine seizure-associated genes were genome-wide  
165 significant in the entire cohort but not in the seizure subset. Of the 285 individuals with  
166 truncating or missense DNMs in known seizure-associated genes, 56% of individuals had no  
167 coded terms related to seizures/epilepsy. These findings suggest that the power of increased  
168 sample size far outweighs specific phenotypic expressivity due to the shared genetic etiology  
169 between individuals with and without epilepsy in our cohort. Despite this, nearly three times as  
170 many individuals with seizures had a DNM in a seizure-associated gene compared to individuals  
171 without seizures (Extended Data Figure 5c). At matched sample sizes, more genes exceeded  
172 genomewide significance in seizure samples than in unstratified samples (Extended Data Figure  
173 5d). This highlights the cost-benefit of recruiting a phenotypically more homogenous cohort.

174  
175 The large number of genome-wide significant genes identified in the analyses above allows us  
176 to compare empirically different experimental strategies for novel gene discovery in a  
177 genetically heterogeneous cohort. We compared the power of exome and genome sequencing  
178 to detect genome-wide significant genes, assuming that budget and not samples are limiting,  
179 under different scenarios of cost ratios and sensitivity ratios (Extended Data Figure 6a). At  
180 current cost ratios (exome costs 30-40% of a genome) and with a plausible sensitivity  
181 differential (genome detects 5% more exonic variants than exome<sup>29</sup>) exome sequencing detects  
182 more than twice as many genome-wide significant genes. These empirical estimates were  
183 consistent with power simulations for identifying dominant loss-of-function genes (Extended  
184 Data Figure 6b). In summary, while genome sequencing gives greatest sensitivity to detect  
185 pathogenic variation in a single individual (or outside of the coding region), exome sequencing  
186 is more powerful for novel disease gene discovery (and, analogously, likely delivers lower cost  
187 per diagnosis currently).

188  
189 Our previous simulations suggested that analysis of a cohort of 4,293 DDD families ought to be  
190 able to detect approximately half of all haploinsufficient DD-associated genes at genome-wide  
191 significance<sup>15</sup>. Empirically, we have identified 47% (50/107) of haploinsufficient genes  
192 previously robustly associated with neurodevelopmental disorders<sup>18</sup>. We hypothesised that  
193 genetic testing prior to recruitment into our study may have depleted the cohort of the most  
194 clinically recognisable disorders. Indeed, we observed that the genes associated with the most  
195 clinically recognisable disorders were associated with a significant, three-fold lower enrichment  
196 of truncating DNMs than other DD-associated genes (~40-fold enrichment vs ~120-fold  
197 enrichment, Figure 3a). Removing these most recognisable disorders from the analysis, we  
198 identified 55% (42/76) of the remaining haploinsufficient DD-associated genes. The known DD-  
199 associated haploinsufficient genes that did not reach genome-wide significance were clearly  
200 enriched for those with lower mutability, which we would expect to lower power to detect in  
201 our analyses. We identified DD-associated genes (e.g. *NRXN2*) with high mutability, low clinical  
202 recognisability and yet no signal of enrichment for DNMs in our cohort, as assessed by  $\Delta_{AIC}$   
203 (Extended Data Figure 7, Supplementary Table 4). Our analyses call into question whether these  
204 genes really are associated with haploinsufficient neurodevelopmental disorders and highlights  
205 the potential for well-powered gene discovery analyses to refute prior credence regarding  
206 disease gene associations or prior inferences regarding an underlying haploinsufficient  
207 mechanism.

208  
209 We estimated the likely prevalence of pathogenic missense and truncating DNMs within our  
210 cohort by increasing the stringency of called DNMs until the observed synonymous DNMs  
211 equated that expected under the null mutation model (Extended Data Figure 8a), then  
212 quantifying the excess of observed missense and truncating DNMs across all genes (Figure 3b).  
213 We observed an excess of 576 truncating and 1,220 missense mutations, suggesting 41.8%  
214 (1,796/4,293) of the cohort has a pathogenic DNM. This estimate of the number of excess  
215 missense and truncating DNMs in our cohort is robust to varying the stringency of DNM calling  
216 (Extended Data Figure 8b). The vast majority of synonymous DNMs are likely to be benign, as  
217 evidenced by them being distributed uniformly (Figure 3d) among genes irrespective of their

218 tolerance of truncating variation in the general population (as quantified by the probability of  
219 being LoF-intolerant (pLI) metric<sup>30</sup>). By contrast, missense and truncating DNMs are significantly  
220 enriched in genes with the highest probabilities of being intolerant of truncating variation  
221 (Figure 3d). The pLI-based distributions were similar to distributions which used functional  
222 constraint (Extended Data Figure 9)<sup>31</sup>. Only 51% (923/1,796) of these excess missense and  
223 truncating DNMs are located in DD-associated dominant genes, with the remainder likely to  
224 affect genes not yet associated with DDs. A much higher proportion of the excess truncating  
225 DNMs (71%) than missense DNMs (42%) affected known DD-associated genes. This suggests  
226 that whereas most haploinsufficient DD-associated genes have already been identified, many  
227 DD-associated genes characterised by pathogenic missense DNMs remain to be discovered.

228  
229 Understanding the mechanism of action of a monogenic disorder is an important prerequisite  
230 for designing therapeutic strategies<sup>32</sup>. We sought to estimate the relative proportion of altered-  
231 function and loss-of-function mechanisms among the excess DNMs in our cohort, by assuming  
232 that the vast majority of truncating mutations operate by a loss-of-function mechanism and  
233 using two independent approaches to estimate the relative contribution of the two  
234 mechanisms among the excess missense DNMs (Methods). First, we used the observed ratio of  
235 truncating and missense DNMs within haploinsufficient DD-associated genes to estimate the  
236 proportion of the excess missense DNMs that likely act by loss-of-function (Figure 3c). This  
237 approach estimated that 59% (55 - 64% 95% CI) of excess missense and truncating DNMs  
238 operate by loss-of-function, and 41% by altered-function. Second, we took advantage of the  
239 different population genetic characteristics of known altered-function and loss-of-function DD-  
240 associated genes. Specifically, we observed that these two classes of DD-associated genes are  
241 differentially depleted of truncating variation in individuals without overt developmental  
242 disorders (pLI metric<sup>30</sup>). We modelled the observed pLI distribution of excess missense DNMs as  
243 a mixture of the pLI distributions of known altered-function and loss-of-function DD-associated  
244 genes (Figure 3e, f), and estimated that 63% (50 - 76% 95% CI) of excess missense DNMs likely  
245 act by altered-function mechanisms. Incorporating the truncating DNMs operating by a loss-of-  
246 function mechanism, this approach estimated that 57% (48 - 66% 95% CI) of excess missense  
247 and truncating DNMs operate by loss-of-function and 43% by altered-function.

248  
249 We estimated the birth prevalence of monoallelic developmental disorders by using the  
250 germline mutation model to calculate the expected cumulative germline mutation rate of  
251 truncating DNMs in haploinsufficient DD-associated genes and scaling this upwards based on  
252 the composition of excess DNMs in the DDD cohort described above (see Methods), correcting  
253 for disorders that are under-represented in our cohort as a result of prior genetic testing (e.g.  
254 clinically-recognisable disorders and large pathogenic CNVs identified by prior chromosomal  
255 microarray analysis). This gives a mean prevalence estimate of 0.34% (0.31-0.37 95% CI), or 1 in  
256 295 births. By factoring in the paternal and maternal age effects on the mutation rate (Figure 1)  
257 we modelled age-specific estimates of birth prevalence (Figure 4) that range from 1 in 448  
258 (both mother and father aged 20) to 1 in 213 (both mother and father aged 45). Assuming a  
259 yearly global birth rate of 18.6 live births/1000 individuals, and a mean age when giving birth of  
260 26.6 years, nearly 400,000 of the 140 million annual births will have a developmental disorder  
261 caused by a DNM.

262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282

In summary, we have shown that *de novo* mutations account for approximately half of the genetic architecture of severe developmental disorders, and are split roughly equally between loss-of-function and altered-function. Whereas most haploinsufficient DD-associated genes have already been identified, currently many activating and dominant negative DD-associated genes have eluded discovery. This elusiveness likely results from these disorders being individually rarer, being caused by a relatively small number of missense mutations within each gene. It would be valuable to estimate the penetrance of *de novo* mutations in the genes we identified exceeding genome-wide significance, but we cannot formally assess penetrance with our data. Future evaluations could integrate depletion of damaging variation in large healthy populations with patterns of segregation in affected families. Discovery of the remaining dominant developmental disorders requires larger studies and novel, more powerful, analytical strategies for disease-gene association that leverage gene-specific patterns of population variation, specifically the observed depletion of damaging variation. The integration of accurate and complete quantitative and categorical phenotypic data into the analysis will improve the power to identify ultrarare DD with distinctive clinical presentations. We have estimated the mean birth prevalence of dominant monogenic developmental disorders to be around 1 in 295, which is greater than the combined impact of trisomies 13, 18 and 21<sup>33</sup> and highlights the cumulative population morbidity and mortality imposed by these individually rare disorders.

283    **References**

284    1.     Sheridan, E. *et al.* Risk factors for congenital anomaly in a multiethnic birth cohort: an  
285     analysis of the Born in Bradford study. *Lancet* **382**, 1350-9 (2013).

286    2.     Ropers, H.H. Genetics of early onset cognitive impairment. *Annu Rev Genomics Hum*  
287     *Genet* **11**, 161-87 (2010).

288    3.     De Ligt, J. *et al.* Diagnostic exome sequencing in persons with severe intellectual  
289     disability. *The New England Journal of Medicine* **367**, 1921-9 (2012).

290    4.     De Rubeis, S. *et al.* Synaptic, transcriptional and chromatin genes disrupted in autism.  
291     *Nature* **515**, 209-215 (2014).

292    5.     Epi4K Consortium & Epilepsy Phenome/Genome Project. De novo mutations in epileptic  
293     encephalopathies. *Nature* **501**, 217-21 (2013).

294    6.     EuroEPINOMICS-RES Consortium, Epilepsy Phenome/Genome Project & Epi4K  
295     Consortium. De novo mutations in synaptic transmission genes including DNMT1 cause  
296     epileptic encephalopathies. *Am J Hum Genet* **95**, 360-70 (2014).

297    7.     Fromer, M. *et al.* De novo mutations in schizophrenia implicate synaptic networks.  
298     *Nature* **506**, 179-184 (2014).

299    8.     Gilissen, C. *et al.* Genome sequencing identifies major causes of severe intellectual  
300     disability. *Nature* **511**, 344-7 (2014).

301    9.     Iossifov, I. *et al.* The contribution of de novo coding mutations to autism spectrum  
302     disorder. *Nature* **515**, 216-221 (2014).

303    10.    Iossifov, I. *et al.* De Novo Gene Disruptions in Children on the Autistic Spectrum. *Neuron*  
304     **74**, 285-299 (2012).

305    11.    O’Roak, B.J. *et al.* Sporadic autism exomes reveal a highly interconnected protein  
306     network of de novo mutations. *Nature* **485**, 1-7 (2012).

307    12.    Rauch, A. *et al.* Range of genetic mutations associated with severe non-syndromic  
308     sporadic intellectual disability: an exome sequencing study. *Lancet* **380**, 1674-82 (2012).

309    13.    Sanders, S.J. *et al.* De novo mutations revealed by whole-exome sequencing are strongly  
310     associated with autism. *Nature* **485**, 237-41 (2012).

311    14.    Zaidi, S. *et al.* De novo mutations in histone-modifying genes in congenital heart disease.  
312     *Nature* **498**, 220-3 (2013).

313    15.    The Deciphering Developmental Disorders Study. Large-scale discovery of novel genetic  
314     causes of developmental disorders. *Nature* **519**, 223-228 (2015).

315    16.    de Ligt, J., Veltman, J.A. & Vissers, L.E.L.M. Point mutations as a source of de novo  
316     genetic disease. *Current Opinion in Genetics & Development* **23**, 257-263 (2013).

317    17.    Wilkie, A.O. The molecular basis of genetic dominance. *J Med Genet* **31**, 89-98 (1994).

318    18.    Wright, C.F. *et al.* Genetic diagnosis of developmental disorders in the DDD study: a  
319     scalable analysis of genome-wide research data. *The Lancet* (2014).

320    19.    Jacquemont, S. *et al.* A higher mutational burden in females supports a "female  
321     protective model" in neurodevelopmental disorders. *Am J Hum Genet* **94**, 415-25  
322     (2014).

323    20.    Kong, A. *et al.* Rate of de novo mutations and the importance of father's age to disease  
324     risk. *Nature* **488**, 471-5 (2012).

- 325 21. Rahbari, R. *et al.* Timing, rates and spectra of human germline mutation. *Nat Genet* **48**,  
326 126-33 (2016).
- 327 22. Wong, W.S. *et al.* New observations on maternal age effect on germline de novo  
328 mutations. *Nat Commun* **7**, 10486 (2016).
- 329 23. Samocha, K.E. *et al.* A framework for the interpretation of de novo variation in human  
330 disease. *Nature Genetics* **46**, 944-950 (2014).
- 331 24. Ferry, Q. *et al.* Diagnostically relevant facial gestalt information from ordinary photos.  
332 *eLife* **3**, e02020-e02020 (2014).
- 333 25. Hirata, H. *et al.* ZC4H2 mutations are associated with arthrogryposis multiplex congenita  
334 and intellectual disability through impairment of central and peripheral synaptic  
335 plasticity. *Am J Hum Genet* **92**, 681-95 (2013).
- 336 26. Homan, C.C. *et al.* Mutations in USP9X are associated with X-linked intellectual disability  
337 and disrupt neuronal cell migration and growth. *Am J Hum Genet* **94**, 470-8 (2014).
- 338 27. Liu, J. *et al.* SMC1A expression and mechanism of pathogenicity in probands with X-  
339 Linked Cornelia de Lange syndrome. *Hum Mutat* **30**, 1535-42 (2009).
- 340 28. Akawi, N. *et al.* Discovery of four recessive developmental disorders using probabilistic  
341 genotype and phenotype matching among 4,125 families. *Nature Genetics* **47**, 1363-  
342 1369 (2015).
- 343 29. Meynert, A.M., Ansari, M., FitzPatrick, D.R. & Taylor, M.S. Variant detection sensitivity  
344 and biases in whole genome and exome sequencing. *BMC Bioinformatics* **15**, 247 (2014).
- 345 30. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**,  
346 285-291 (2016).
- 347 31. Petrovski, S., Wang, Q., Heinzen, E.L., Allen, A.S. & Goldstein, D.B. Genic intolerance to  
348 functional variation and the interpretation of personal genomes. *PLoS Genet* **9**,  
349 e1003709 (2013).
- 350 32. Boycott, K.M., Vanstone, M.R., Bulman, D.E. & Mackenzie, A.E. Rare-disease genetics in  
351 the era of next-generation sequencing: discovery to translation. *Nature Reviews*  
352 *Genetics* **14**, 681-91 (2013).
- 353 33. Springett, A. *et al.* Congenital Anomaly Statistics 2011: England and Wales. (2013).
- 354 34. Bragin, E. *et al.* DECIPHER: database for the interpretation of phenotype-linked plausibly  
355 pathogenic sequence and copy-number variation. *Nucleic Acids Res* **42**, D993-D1000  
356 (2014).
- 357 35. Köhler, S. *et al.* Clinical diagnostics in human genetics with semantic similarity searches  
358 in ontologies. *American Journal of Human Genetics* **85**, 457-464 (2009).
- 359 36. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler  
360 transform. *Bioinformatics* **25**, 1754-1760 (2009).
- 361 37. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing  
362 next-generation DNA sequencing data. *Genome Res* **20**, 1297-303 (2010).
- 363 38. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,  
364 2078-2079 (2009).
- 365 39. Ramu, A. *et al.* DeNovoGear: de novo indel and point mutation discovery and phasing.  
366 *Nature Methods* **10**, 985-7 (2013).
- 367 40. Abecasis, G.R. *et al.* An integrated map of genetic variation from 1,092 human genomes.  
368 *Nature* **491**, 56-65 (2012).

- 369 41. McLaren, W. *et al.* Deriving the consequences of genomic variants with the Ensembl API  
370 and SNP Effect Predictor. *Bioinformatics* **26**, 2069-70 (2010).
- 371 42. Felzenszwalb, P.F., Girshick, R.B., McAllester, D. & Ramanan, D. Object detection with  
372 discriminatively trained part-based models. *IEEE transactions on pattern analysis and*  
373 *machine intelligence* **32**, 1627-45 (2010).
- 374 43. Xiong, X. & De la Torre, F. Supervised Descent Method and Its Applications to Face  
375 Alignment. in *2013 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*  
376 532-539 (IEEE, Portland, OR, 2013).
- 377 44. Cooper, G.M. *et al.* A copy number variation morbidity map of developmental delay. *Nat*  
378 *Genet* **43**, 838-46 (2011).
- 379 45. Sagoo, G.S. *et al.* Array CGH in patients with learning disability (mental retardation) and  
380 congenital anomalies: updated systematic review and meta-analysis of 19 studies and  
381 13,926 subjects. *Genet Med* **11**, 139-46 (2009).
- 382 46. Central Intelligence Agency. The World Factbook. Vol. 2016 (2016).
- 383 47. The World Bank. Fertility rate, total (births per woman). in *World Development*  
384 *Indicators* Vol. 2016 (2016).
- 385 48. Copen, C.E., Thoma, M.E. & Kirmeyer, S. Interpregnancy Intervals in the United States:  
386 Data From the Birth Certificate and the National Survey of Family Growth. in *National*  
387 *Vital Statistics Reports* Vol. 64 (National Center for Health Statistics, Hyattsville, MD,  
388 2015).
- 389

## 390 Acknowledgments

391 We thank the families for their participation and patience. We are grateful to the Exome  
392 Aggregation Consortium for making their data available. The DDD study presents independent  
393 research commissioned by the Health Innovation Challenge Fund (grant HICF-1009-003), a  
394 parallel funding partnership between the Wellcome Trust and the UK Department of Health,  
395 and the Wellcome Trust Sanger Institute (grant WT098051). The views expressed in this  
396 publication are those of the author(s) and not necessarily those of the Wellcome Trust or the  
397 UK Department of Health. The study has UK Research Ethics Committee approval  
398 (10/H0305/83, granted by the Cambridge South Research Ethics Committee and GEN/284/12,  
399 granted by the Republic of Ireland Research Ethics Committee). The research team  
400 acknowledges the support of the National Institutes for Health Research, through the  
401 Comprehensive Clinical Research Network. The authors wish to thank the Sanger Human  
402 Genome Informatics team, the Sample Management team, the Illumina High-Throughput team,  
403 the New Pipeline Group team, the DNA pipelines team and the Core Sequencing team for their  
404 support in generating and processing the data. D.R.F. is funded through an MRC Human  
405 Genetics Unit program grant to the University of Edinburgh. Finally we gratefully acknowledge  
406 the contribution of two esteemed DDD clinical collaborators, John Tolmie and Louise Brueton,  
407 who died in the course of the study.

## 408 Author Contributions

409 Jeremy F McRae<sup>1</sup>, Stephen Clayton<sup>1</sup>, Tomas W Fitzgerald<sup>1</sup>, Joanna Kaplanis<sup>1</sup>, Elena Prigmore<sup>1</sup>,  
410 Diana Rajan<sup>1</sup>, Alejandro Sifrim<sup>1</sup>, Stuart Aitken<sup>2</sup>, Nadia Akawi<sup>1</sup>, Mohsan Alvi<sup>3</sup>, Kirsty Ambridge<sup>1</sup>,  
411 Daniel M Barrett<sup>1</sup>, Tanya Bayzetinova<sup>1</sup>, Philip Jones<sup>1</sup>, Wendy D Jones<sup>1</sup>, Daniel King<sup>1</sup>, Netravathi  
412 Krishnappa<sup>1</sup>, Laura E Mason<sup>1</sup>, Tarjinder Singh<sup>1</sup>, Adrian R Tivey<sup>1</sup>, Munaza Ahmed<sup>4</sup>, Uruj Anjum<sup>5</sup>,  
413 Hayley Archer<sup>6</sup>, Ruth Armstrong<sup>7</sup>, Jana Awada<sup>1</sup>, Meena Balasubramanian<sup>8</sup>, Siddharth Banka<sup>9</sup>,  
414 Diana Baralle<sup>4</sup>, Angela Barnicoat<sup>10</sup>, Paul Batstone<sup>11</sup>, David Baty<sup>12</sup>, Chris Bennett<sup>13</sup>, Jonathan  
415 Berg<sup>12</sup>, Birgitta Bernhard<sup>14</sup>, A Paul Bevan<sup>1</sup>, Maria Bitner-Glindzicz<sup>10</sup>, Edward Blair<sup>15</sup>, Moira  
416 Blyth<sup>13</sup>, David Bohanna<sup>16</sup>, Louise Bourdon<sup>14</sup>, David Bourn<sup>17</sup>, Lisa Bradley<sup>18</sup>, Angela Brady<sup>14</sup>,  
417 Simon Brent<sup>1</sup>, Carole Brewer<sup>19</sup>, Kate Brunstrom<sup>10</sup>, David J Bunyan<sup>4</sup>, John Burn<sup>17</sup>, Natalie  
418 Canham<sup>14</sup>, Bruce Castle<sup>19</sup>, Kate Chandler<sup>9</sup>, Elena Chatzimichali<sup>1</sup>, Deirdre Cilliers<sup>15</sup>, Angus Clarke<sup>6</sup>,  
419 Susan Clasper<sup>15</sup>, Jill Clayton-Smith<sup>9</sup>, Virginia Clowes<sup>14</sup>, Andrea Coates<sup>13</sup>, Trevor Cole<sup>16</sup>, Irina  
420 Colgiu<sup>1</sup>, Amanda Collins<sup>4</sup>, Morag N Collinson<sup>4</sup>, Fiona Connell<sup>20</sup>, Nicola Cooper<sup>16</sup>, Helen Cox<sup>16</sup>,  
421 Lara Cresswell<sup>21</sup>, Gareth Cross<sup>22</sup>, Yanick Crow<sup>9</sup>, Mariella D'Alessandro<sup>11</sup>, Tabib Dabir<sup>18</sup>,  
422 Rosemarie Davidson<sup>23</sup>, Sally Davies<sup>6</sup>, Dylan de Vries<sup>1</sup>, John Dean<sup>11</sup>, Charu Deshpande<sup>20</sup>, Gemma  
423 Devlin<sup>19</sup>, Abhijit Dixit<sup>22</sup>, Angus Dobbie<sup>13</sup>, Alan Donaldson<sup>24</sup>, Dian Donnai<sup>9</sup>, Deirdre Donnelly<sup>18</sup>,  
424 Carina Donnelly<sup>9</sup>, Angela Douglas<sup>25</sup>, Sofia Douzgou<sup>9</sup>, Alexis Duncan<sup>23</sup>, Jacqueline Eason<sup>22</sup>, Sian  
425 Ellard<sup>19</sup>, Ian Ellis<sup>25</sup>, Frances Elmslie<sup>5</sup>, Karenza Evans<sup>6</sup>, Sarah Everest<sup>19</sup>, Tina Fendick<sup>20</sup>, Richard  
426 Fisher<sup>17</sup>, Frances Flinter<sup>20</sup>, Nicola Foulds<sup>4</sup>, Andrew Fry<sup>6</sup>, Alan Fryer<sup>25</sup>, Carol Gardiner<sup>23</sup>, Lorraine  
427 Gaunt<sup>9</sup>, Neeti Ghali<sup>14</sup>, Richard Gibbons<sup>15</sup>, Harinder Gill<sup>26</sup>, Judith Goodship<sup>17</sup>, David Goudie<sup>12</sup>,  
428 Emma Gray<sup>1</sup>, Andrew Green<sup>26</sup>, Philip Greene<sup>2</sup>, Lynn Greenhalgh<sup>25</sup>, Susan Gribble<sup>1</sup>, Rachel  
429 Harrison<sup>22</sup>, Lucy Harrison<sup>4</sup>, Victoria Harrison<sup>4</sup>, Rose Hawkins<sup>24</sup>, Liu He<sup>1</sup>, Stephen Hellens<sup>17</sup>, Alex  
430 Henderson<sup>17</sup>, Sarah Hewitt<sup>13</sup>, Lucy Hildyard<sup>1</sup>, Emma Hobson<sup>13</sup>, Simon Holden<sup>7</sup>, Muriel Holder<sup>14</sup>,  
431 Susan Holder<sup>14</sup>, Georgina Hollingsworth<sup>10</sup>, Tessa Homfray<sup>5</sup>, Mervyn Humphreys<sup>18</sup>, Jane Hurst<sup>10</sup>,

432 Ben Hutton<sup>1</sup>, Stuart Ingram<sup>8</sup>, Melita Irving<sup>20</sup>, Lily Islam<sup>16</sup>, Andrew Jackson<sup>2</sup>, Joanna Jarvis<sup>16</sup>, Lucy  
433 Jenkins<sup>10</sup>, Diana Johnson<sup>8</sup>, Elizabeth Jones<sup>9</sup>, Dragana Josifova<sup>20</sup>, Shelagh Joss<sup>23</sup>, Beckie  
434 Kaemba<sup>21</sup>, Sandra Kazembe<sup>21</sup>, Rosemary Kelsell<sup>1</sup>, Bronwyn Kerr<sup>9</sup>, Helen Kingston<sup>9</sup>, Usha Kini<sup>15</sup>,  
435 Esther Kinning<sup>23</sup>, Gail Kirby<sup>16</sup>, Claire Kirk<sup>18</sup>, Emma Kivuva<sup>19</sup>, Alison Kraus<sup>13</sup>, Dhavendra Kumar<sup>6</sup>,  
436 V.K Ajith Kumar<sup>10</sup>, Katherine Lachlan<sup>4</sup>, Wayne Lam<sup>2</sup>, Anne Lampe<sup>2</sup>, Caroline Langman<sup>20</sup>, Melissa  
437 Lees<sup>10</sup>, Derek Lim<sup>16</sup>, Cheryl Longman<sup>23</sup>, Gordon Lowther<sup>23</sup>, Sally A Lynch<sup>26</sup>, Alex Magee<sup>18</sup>, Eddy  
438 Maher<sup>2</sup>, Alison Male<sup>10</sup>, Sahar Mansour<sup>5</sup>, Karen Marks<sup>5</sup>, Katherine Martin<sup>22</sup>, Una Maye<sup>25</sup>, Emma  
439 McCann<sup>27</sup>, Vivienne McConnell<sup>18</sup>, Meriel McEntagart<sup>5</sup>, Ruth McGowan<sup>11</sup>, Kirsten McKay<sup>16</sup>,  
440 Shane McKee<sup>18</sup>, Dominic J McMullan<sup>16</sup>, Susan McNerlan<sup>18</sup>, Catherine McWilliam<sup>11</sup>, Sarju  
441 Mehta<sup>7</sup>, Kay Metcalfe<sup>9</sup>, Anna Middleton<sup>1</sup>, Zosia Miedzybrodzka<sup>11</sup>, Emma Miles<sup>9</sup>, Shehla  
442 Mohammed<sup>20</sup>, Tara Montgomery<sup>17</sup>, David Moore<sup>2</sup>, Sian Morgan<sup>6</sup>, Jenny Morton<sup>16</sup>, Hood  
443 Mugalaasi<sup>6</sup>, Victoria Murday<sup>23</sup>, Helen Murphy<sup>9</sup>, Swati Naik<sup>16</sup>, Andrea Nemeth<sup>15</sup>, Louise Nevitt<sup>8</sup>,  
444 Ruth Newbury-Ecob<sup>24</sup>, Andrew Norman<sup>16</sup>, Rosie O'Shea<sup>26</sup>, Caroline Ogilvie<sup>20</sup>, Kai-Ren Ong<sup>16</sup>,  
445 Soo-Mi Park<sup>7</sup>, Michael J Parker<sup>8</sup>, Chirag Patel<sup>16</sup>, Joan Paterson<sup>7</sup>, Stewart Payne<sup>14</sup>, Daniel  
446 Perrett<sup>1</sup>, Julie Phipps<sup>15</sup>, Daniela T Pilz<sup>23</sup>, Martin Pollard<sup>1</sup>, Caroline Pottinger<sup>27</sup>, Joanna Poulton<sup>15</sup>,  
447 Norman Pratt<sup>12</sup>, Katrina Prescott<sup>13</sup>, Sue Price<sup>15</sup>, Abigail Pridham<sup>15</sup>, Annie Procter<sup>6</sup>, Hellen  
448 Purnell<sup>15</sup>, Oliver Quarrell<sup>8</sup>, Nicola Ragge<sup>16</sup>, Raheleh Rahbari<sup>1</sup>, Josh Randall<sup>1</sup>, Julia Rankin<sup>19</sup>, Lucy  
449 Raymond<sup>7</sup>, Debbie Rice<sup>12</sup>, Leema Robert<sup>20</sup>, Eileen Roberts<sup>24</sup>, Jonathan Roberts<sup>7</sup>, Paul Roberts<sup>13</sup>,  
450 Gillian Roberts<sup>25</sup>, Alison Ross<sup>11</sup>, Elisabeth Rosser<sup>10</sup>, Anand Saggar<sup>5</sup>, Shalaka Samant<sup>11</sup>, Julian  
451 Sampson<sup>6</sup>, Richard Sandford<sup>7</sup>, Ajoy Sarkar<sup>22</sup>, Susann Schweiger<sup>12</sup>, Richard Scott<sup>10</sup>, Ingrid Scurr<sup>24</sup>,  
452 Ann Selby<sup>22</sup>, Anneke Seller<sup>15</sup>, Cheryl Sequeira<sup>14</sup>, Nora Shannon<sup>22</sup>, Saba Sharif<sup>16</sup>, Charles Shaw-  
453 Smith<sup>19</sup>, Emma Shearing<sup>8</sup>, Debbie Shears<sup>15</sup>, Eamonn Sheridan<sup>13</sup>, Ingrid Simonic<sup>7</sup>, Roldan  
454 Singzon<sup>14</sup>, Zara Skitt<sup>9</sup>, Audrey Smith<sup>13</sup>, Kath Smith<sup>8</sup>, Sarah Smithson<sup>24</sup>, Linda Sneddon<sup>17</sup>, Miranda  
455 Splitt<sup>17</sup>, Miranda Squires<sup>13</sup>, Fiona Stewart<sup>18</sup>, Helen Stewart<sup>15</sup>, Volker Straub<sup>17</sup>, Mohnish Suri<sup>22</sup>,  
456 Vivienne Sutton<sup>25</sup>, Ganesh Jawahar Swaminathan<sup>1</sup>, Elizabeth Sweeney<sup>25</sup>, Kate Tatton-Brown<sup>5</sup>,  
457 Cat Taylor<sup>8</sup>, Rohan Taylor<sup>5</sup>, Mark Tein<sup>16</sup>, I Karen Temple<sup>4</sup>, Jenny Thomson<sup>13</sup>, Marc Tischkowitz<sup>7</sup>,  
458 Susan Tomkins<sup>24</sup>, Audrey Torokwa<sup>4</sup>, Becky Treacy<sup>7</sup>, Claire Turner<sup>19</sup>, Peter Turnpenny<sup>19</sup>, Carolyn  
459 Tysoe<sup>19</sup>, Anthony Vandersteen<sup>14</sup>, Vinod Varghese<sup>6</sup>, Pradeep Vasudevan<sup>21</sup>, Parthiban  
460 Vijayarangakannan<sup>1</sup>, Julie Vogt<sup>16</sup>, Emma Wakeling<sup>14</sup>, Sarah Wallwark<sup>7</sup>, Jonathon Waters<sup>10</sup>, Astrid  
461 Weber<sup>25</sup>, Diana Wellesley<sup>4</sup>, Margo Whiteford<sup>23</sup>, Sara Widaa<sup>1</sup>, Sarah Wilcox<sup>7</sup>, Emily Wilkinson<sup>1</sup>,  
462 Denise Williams<sup>16</sup>, Nicola Williams<sup>23</sup>, Louise Wilson<sup>10</sup>, Geoff Woods<sup>7</sup>, Christopher Wragg<sup>24</sup>,  
463 Michael Wright<sup>17</sup>, Laura Yates<sup>17</sup>, Michael Yau<sup>20</sup>, Chris Nellåker<sup>28,29,30</sup>, Michael J Parker<sup>31</sup>, Helen V  
464 Firth<sup>1,7,32</sup>, Caroline F Wright<sup>1,32</sup>, David R FitzPatrick<sup>1,2,32</sup>, Jeffrey C Barrett<sup>1,32</sup>, Matthew E  
465 Hurles<sup>1,32</sup>

466  
467 <sup>1</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, CB10  
468 1SA, UK

469 <sup>2</sup>MRC Human Genetics Unit, MRC IGMM, University of Edinburgh, Western General Hospital,  
470 Edinburgh, EH4 2XU, UK

471 <sup>3</sup>Department of Engineering Science, University of Oxford, Parks Road, Oxford, OX1 3PJ, UK

472 <sup>4</sup>Wessex Clinical Genetics Service, University Hospital Southampton, Princess Anne Hospital,  
473 Coxford Road, Southampton, SO16 5YA, UK and Wessex Regional Genetics Laboratory,  
474 Salisbury NHS Foundation Trust, Salisbury District Hospital, Odstock Road, Salisbury, Wiltshire,

475 SP2 8BJ, UK and Faculty of Medicine, University of Southampton, Building 85, Life Sciences  
476 Building, Highfield Campus, Southampton, SO17 1BJ, UK  
477 <sup>5</sup>South West Thames Regional Genetics Centre, St George's Healthcare NHS Trust, St George's,  
478 University of London, Cranmer Terrace, London, SW17 0RE, UK  
479 <sup>6</sup>Institute Of Medical Genetics, University Hospital Of Wales, Heath Park, Cardiff, CF14 4XW, UK  
480 and Department of Clinical Genetics, Block 12, Glan Clwyd Hospital, Rhyl, Denbighshire, LL18  
481 5UJ, UK  
482 <sup>7</sup>East Anglian Medical Genetics Service, Box 134, Cambridge University Hospitals NHS  
483 Foundation Trust, Cambridge Biomedical Campus, Cambridge, CB2 0QQ, UK  
484 <sup>8</sup>Sheffield Regional Genetics Services, Sheffield Children's NHS Trust, Western Bank, Sheffield,  
485 S10 2TH, UK  
486 <sup>9</sup>Manchester Centre for Genomic Medicine, St Mary's Hospital, Central Manchester University  
487 Hospitals NHS Foundation Trust, Manchester Academic Health Science Centre, Manchester  
488 M13 9WL, UK  
489 <sup>10</sup>North East Thames Regional Genetics Service, Great Ormond Street Hospital for Children NHS  
490 Foundation Trust, Great Ormond Street Hospital, Great Ormond Street, London, WC1N 3JH,  
491 UK  
492 <sup>11</sup>North of Scotland Regional Genetics Service, NHS Grampian, Department of Medical Genetics  
493 Medical School, Foresterhill, Aberdeen, AB25 2ZD, UK  
494 <sup>12</sup>East of Scotland Regional Genetics Service, Human Genetics Unit, Pathology Department, NHS  
495 Tayside, Ninewells Hospital, Dundee, DD1 9SY, UK  
496 <sup>13</sup>Yorkshire Regional Genetics Service, Leeds Teaching Hospitals NHS Trust, Department of  
497 Clinical Genetics, Chapel Allerton Hospital, Chapeltown Road, Leeds, LS7 4SA, UK  
498 <sup>14</sup>North West Thames Regional Genetics Centre, North West London Hospitals NHS Trust, The  
499 Kennedy Galton Centre, Northwick Park And St Mark's NHS Trust Watford Road, Harrow, HA1  
500 3UJ, UK  
501 <sup>15</sup>Oxford Regional Genetics Service, Oxford Radcliffe Hospitals NHS Trust, The Churchill Old  
502 Road, Oxford, OX3 7LJ, UK  
503 <sup>16</sup>West Midlands Regional Genetics Service, Birmingham Women's NHS Foundation Trust,  
504 Birmingham Women's Hospital, Edgbaston, Birmingham, B15 2TG, UK  
505 <sup>17</sup>Northern Genetics Service, Newcastle upon Tyne Hospitals NHS Foundation Trust, Institute of  
506 Human Genetics, International Centre for Life, Central Parkway, Newcastle upon Tyne, NE1  
507 3BZ, UK  
508 <sup>18</sup>Northern Ireland Regional Genetics Centre, Belfast Health and Social Care Trust, Belfast City  
509 Hospital, Lisburn Road, Belfast, BT9 7AB, UK  
510 <sup>19</sup>Peninsula Clinical Genetics Service, Royal Devon and Exeter NHS Foundation Trust, Clinical  
511 Genetics Department, Royal Devon & Exeter Hospital (Heavitree), Gladstone Road, Exeter, EX1  
512 2ED, UK  
513 <sup>20</sup>South East Thames Regional Genetics Centre, Guy's and St Thomas' NHS Foundation Trust,  
514 Guy's Hospital, Great Maze Pond, London, SE1 9RT, UK  
515 <sup>21</sup>Leicestershire Genetics Centre, University Hospitals of Leicester NHS Trust, Leicester Royal  
516 Infirmary (NHS Trust), Leicester, LE1 5WW, UK  
517 <sup>22</sup>Nottingham Regional Genetics Service, City Hospital Campus, Nottingham University Hospitals  
518 NHS Trust, The Gables, Hucknall Road, Nottingham NG5 1PB, UK

519 <sup>23</sup>West of Scotland Regional Genetics Service, NHS Greater Glasgow and Clyde, Institute Of  
520 Medical Genetics, Yorkhill Hospital, Glasgow, G3 8SJ, UK  
521 <sup>24</sup>Bristol Genetics Service (Avon, Somerset, Gloucs and West Wilts), University Hospitals Bristol  
522 NHS Foundation Trust, St Michael's Hospital, St Michael's Hill, Bristol, BS2 8DT, UK  
523 <sup>25</sup>Merseyside and Cheshire Genetics Service, Liverpool Women's NHS Foundation Trust,  
524 Department of Clinical Genetics, Royal Liverpool Children's Hospital Alder Hey, Eaton Road,  
525 Liverpool, L12 2AP, UK  
526 <sup>26</sup>National Centre for Medical Genetics, Our Lady's Children's Hospital, Crumlin, Dublin 12,  
527 Ireland  
528 <sup>27</sup>Department of Clinical Genetics, Block 12, Glan Clwyd Hospital, Rhyl, Denbighshire, Wales,  
529 LL18 5UJ, UK  
530 <sup>28</sup>Nuffield Department of Obstetrics & Gynaecology, University of Oxford, Level 3, Women's  
531 Centre, John Radcliffe Hospital, Oxford, OX3 9DU, UK  
532 <sup>29</sup>Institute of Biomedical Engineering, Department of Engineering Science, University of Oxford,  
533 Old Road Campus Research Building, Oxford, OX3 7DQ, UK  
534 <sup>30</sup>Big Data Institute, University of Oxford, Roosevelt drive, Oxford, OX3 7LF, UK  
535 <sup>31</sup>The Ethox Centre, Nuffield Department of Population Health, University of Oxford, Old Road  
536 Campus, Oxford, OX3 7LF, UK  
537 <sup>32</sup>These authors jointly supervised this work.  
538

539 Patient recruitment and phenotyping: M. Ahmed, U.A., H.A., R.A., M. Balasubramanian, S.  
540 Banka, D. Baralle, A. Barnicoat, P.B., D. Baty, C. Bennett, J. Berg, B.B., M.B-G., E.B., M. Blyth,  
541 D. Bohanna, L. Bourdon, D. Bourn, L. Bradley, A. Brady, C. Brewer, K.B., D.J.B., J. Burn, N.  
542 Canham, B.C., K.C., D.C., A. Clarke, S. Clasper, J.C-S., V.C., A. Coates, T.C., A. Collins, M.N.C.,  
543 F.C., N. Cooper, H.C., L.C., G.C., Y.C., M.D., T.D., R.D., S. Davies, J.D., C. Deshpande, G.D., A.  
544 Dixit, A. Dobbie, A. Donaldson, D. Donnai, D. Donnelly, C. Donnelly, A. Douglas, S. Douzgou,  
545 A. Duncan, J.E., S. Ellard, I.E., F.E., K.E., S. Everest, T.F., R.F., F.F., N.F., A. Fry, A. Fryer, C.G., L.  
546 Gaunt, N.G., R.G., H.G., J.G., D.G., A.G., P.G., L. Greenhalgh, R. Harrison, L. Harrison, V.H., R.  
547 Hawkins, S. Hellens, A.H., S. Hewitt, E.H., S. Holden, M. Holder, S. Holder, G.H., T.H., M.  
548 Humphreys, J.H., S.I., M.I., L.I., A.J., J.J., L.J., D. Johnson, E.J., D. Josifova, S.J., B. Kaemba, S.K.,  
549 B. Kerr, H.K., U.K., E. Kinning, G.K., C.K., E. Kivuva, A.K., D. Kumar, V.A.K., K.L., W.L., A.L., C.  
550 Langman, M.L., D.L., C. Longman, G.L., S.A.L., A. Magee, E. Maher, A. Male, S. Mansour, K.  
551 Marks, K. Martin, U.M., E. McCann, V. McConnell, M.M., R.M., K. McKay, S. McKee, D.J.M., S.  
552 McNerlan, C.M., S. Mehta, K. Metcalfe, Z.M., E. Miles, S. Mohammed, T.M., D.M., S. Morgan,  
553 J.M., H. Mugalaasi, V. Murday, H. Murphy, S.N., A. Nemeth, L.N., R.N-E., A. Norman, R.O.,  
554 C.O., K-R.O., S-M.P., M.J. Parker, C. Patel, J. Paterson, S. Payne, J. Phipps, D.T.P., C. Pottinger,  
555 J. Poulton, N.P., K.P., S. Price, A. Pridham, A. Procter, H.P., O.Q., N.R., J. Rankin, L. Raymond,  
556 D. Rice, L. Robert, E. Roberts, J. Roberts, P.R., G.R., A.R., E. Rosser, A. Sagggar, S. Samant, J.S.,  
557 R. Sandford, A. Sarkar, S. Schweiger, R. Scott, I. Scurr, A. Selby, A. Seller, C.S., N.S., S. Sharif,  
558 C.S-S., E. Shearing, D.S., E. Sheridan, I. Simonic, R. Singzon, Z.S., A. Smith, K.S., S. Smithson,  
559 L.S., M. Splitt, M. Squires, F.S., H.S., V. Straub, M. Suri, V. Sutton, E. Sweeney, K.T-B., C.  
560 Taylor, R.T., M. Tein, I.K.T., J.T., M. Tischkowitz, S.T., A.T., B.T., C. Turner, P.T., C. Tysoe, A.V.,  
561 V.V., P. Vasudevan, J.V., E. Wakeling, S. Wallwark, J.W., A.W., D. Wellesley, M. Whiteford, S.  
562 Wilcox, D. Williams, N.W., L.W., G.W., C.W., M. Wright, L.Y., M.Y., H.V.F., D.R.F.

563  
564 Sample and data processing: S. Clayton, T.W.F., E.P., D. Rajan, K.A., D.M.B., T.B., P.J., N.K.,  
565 L.E.M., A.R.T., A.P.B., S. Brent, E.C., I.C., E.G., S.G., L. Hildyard, B.H., R.K., D.P., M.P., J. Randall,  
566 G.J.S., S. Widaa, E. Wilkinson  
567  
568 Validation experiments: J.F.M., E.P., D. Rajan, A. Sifrim, N.K., C.F.W.  
569  
570 Study design: M.J. Parker, H.V.F., C.F.W., D.R.F., J.C.B., M.E.H.  
571  
572 Method development and data analysis: J.F.M., S. Clayton, T.W.F., J.K., E.P., D. Rajan, A. Sifrim,  
573 S.A., N.A., M. Alvi, P.J., W.D.J., D. King, T.S., J.A., D.d.V., L. He, R.R., G.J.S., P.  
574 Vijayarangakannan, C.N., H.V.F., C.F.W., D.R.F., J.C.B., M.E.H.  
575  
576 Data interpretation: J.F.M., H.V.F., C.F.W., D.R.F., J.C.B., M.E.H.  
577 Writing: J.F.M., C.F.W., D.R.F., M.E.H.  
578  
579 Experimental and analytical supervision: M.J. Parker, H.V.F., C.F.W., D.R.F., J.C.B., M.E.H.  
580  
581 Project Supervision: M.E.H.

## 582 Author Information

583 Exome sequencing data are accessible via the European Genome-phenome Archive (EGA) under  
584 accession EGAS00001000775. Details of DD-associated genes are available at  
585 [www.ebi.ac.uk/gene2phenotype](http://www.ebi.ac.uk/gene2phenotype). M.E.H. is a co-founder of, and holds shares in, Congenica Ltd,  
586 a genetics diagnostic company. Correspondence and requests for materials should be  
587 addressed to M.E.H (meh@sanger.ac.uk).  
588

589 **Tables**

590 Table 1: Genes achieving genome-wide significant statistical evidence without previous compelling  
 591 evidence for being developmental disorder genes. The numbers of unrelated individuals with  
 592 independent *de novo* mutations (DNMs) are given for protein truncating variants (PTV) and missense  
 593 variants. Counts of individuals in other cohorts are given in brackets if present. The *P*-value reported is  
 594 the minimum *P*-value from the testing of the DDD dataset or the meta-analysis dataset. The subset  
 595 providing the *P*-value is also listed. Mutations are considered clustered if the *P*-value from proximity  
 596 clustering of DNMs is less than 0.01.

Gene	Missense	PTV	P-value	Test	Clustering
<i>CDK13</i>	10	1	$3.2 \times 10^{-19}$	DDD	Yes
<i>GNAI1</i>	7 (1)	1	$2.1 \times 10^{-13}$	DDD	No
<i>CSNK2A1</i>	7	0	$1.4 \times 10^{-12}$	DDD	Yes
<i>PPM1D</i>	0	5 (1)	$6.3 \times 10^{-12}$	Meta	No
<i>CNOT3</i>	5	2 (1)	$5.2 \times 10^{-11}$	DDD	Yes
<i>MSL3</i>	0	4	$2.2 \times 10^{-10}$	DDD	No
<i>KCNQ3</i>	4 (3)	0	$3.4 \times 10^{-10}$	Meta	Yes
<i>ZBTB18</i>	1 (1)	4	$1.4 \times 10^{-9}$	DDD	No
<i>PUF60</i>	4 (1)	3	$2.6 \times 10^{-9}$	DDD	No
<i>TCF20</i>	1	5	$2.7 \times 10^{-9}$	DDD	No
<i>SUV420H1</i>	0 (2)	2 (3)	$2.9 \times 10^{-9}$	Meta	No
<i>CHD4</i>	8 (1)	1	$7.6 \times 10^{-9}$	DDD	No
<i>SET</i>	0	3	$1.2 \times 10^{-7}$	DDD	No
<i>QRICH1</i>	0	3 (1)	$3.6 \times 10^{-7}$	Meta	No

597

598

599 Extended Data Tables

600

601 Extended Data Table 1: Phenotypes tested for association with having a pathogenic *de novo* mutation.

602

Category	Phenotype	Type	Value	95% CI	P-value
Post-natal	abnormal cranial MRI	Odds ratio	1.365	1.125 – 1.656	0.002
	feeding problems	Odds ratio	1.176	1.01 – 1.369	0.039
	neonatal intensive care	Odds ratio	0.896	0.762 – 1.054	0.190
	anticonvulsant drugs	Odds ratio	0.582	0.246 – 1.377	0.270
Pre-natal	bleeding	Odds ratio	0.892	0.714 – 1.114	0.346
	maternal illness	Odds ratio	0.908	0.764 – 1.079	0.278
	maternal diabetes	Odds ratio	0.787	0.504 – 1.229	0.341
	abnormal scan	Odds ratio	0.839	0.692 – 1.017	0.078
	assisted reproduction	Odds ratio	0.868	0.554 – 1.36	0.584
	increased nuchal translucency	Odds ratio	1.432	0.903 – 2.271	0.126
Family history	consanguinity	Odds ratio	0.234	0.138 – 0.397	$8.0 \times 10^{-11}$
	similar phenotype parents	Odds ratio	0.295	0.184 – 0.474	$5.7 \times 10^{-9}$
	similar phenotype relatives	Odds ratio	0.553	0.402 – 0.761	$1.5 \times 10^{-4}$
	similar phenotype siblings	Odds ratio	0.311	0.23 – 0.421	$7.3 \times 10^{-18}$
	only patient affected	Odds ratio	2.478	2.001 – 3.068	$3.9 \times 10^{-19}$
	X-linked inheritance	Odds ratio	0.839	0.436 – 1.613	0.752
	Multiple births	Beta	0.043	-0.058 – 0.144	0.403
	History of pregnancy loss	Beta	-0.039	-0.155 – 0.078	0.516
Developmental milestones	first words	Beta	0.205	0.081 – 0.328	0.001
	walked independently	Beta	0.125	0.016 – 0.235	0.025
	sat independently	Beta	0.050	-0.069 – 0.17	0.408
	social smile	Beta	0.072	-0.066 – 0.211	0.305
Growth	height	Beta	0.008	-0.111 – 0.126	0.897
	birthweight	Beta	-0.018	-0.135 – 0.098	0.756
	OFC	Beta	-0.094	-0.215 – 0.026	0.125
	weight	Beta	-0.331	-1.278 – 0.615	0.493
Age	age at assessment	Beta	0.116	0.015 – 0.217	0.025
	gestation	Beta	0.079	-0.033 – 0.19	0.167
	father's age	Beta	0.137	0.027 – 0.247	0.015
	mother's age	Beta	0.108	-0.003 – 0.219	0.056
Other	phenotypic terms (n)	Beta	0.104	0.004 – 0.203	0.041
	autozygosity length	Beta	-0.185	-0.254 – -0.115	$1.6 \times 10^{-7}$
	sex (male)	Odds ratio	0.750	0.646 – 0.87	$1.6 \times 10^{-4}$

## 603 Supplementary Table Legends

604 Note: These are included in the supplementary info, but are required here for the auto-  
605 numbering.

606

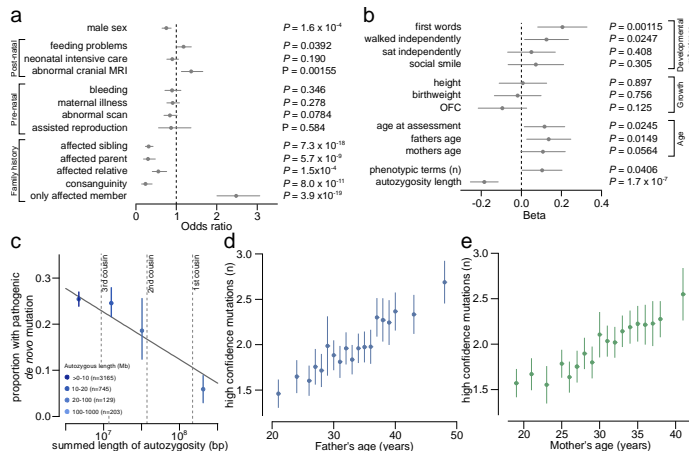
607 Supplementary Table 1: Table of *de novo* mutations (DNM) in the 4,293 DDD individuals. The table  
608 includes sex, chromosome, position, reference and alternate alleles, HGNC symbol, VEP consequence,  
609 posterior probability of DNM and validation status where available. Individual IDs are available on  
610 request. This list excludes the sites that failed validations, but includes sites that passed validation  
611 (confirmed), sites that were uncertain (uncertain), and sites that were not tested by secondary  
612 validation (NA). Genome positions are given as GRCh37 coordinates.

613 Supplementary Table 2: Details of cohorts used in meta-analyses. This includes numbers of individuals  
614 by sex and publication details.

615 Supplementary Table 3: Genes with genome-wide significant statistical evidence to be developmental  
616 disorder genes. The numbers of unrelated individuals with independent *de novo* mutations (DNMs) are  
617 given for protein truncating variants (PTV) and missense variants. If any additional individuals were in  
618 other cohorts, that number is given in brackets. The *P*-value reported is the minimum *P*-value from the  
619 testing of the DDD dataset or the meta-analysis dataset. The subset providing the *P*-value is also listed.  
620 Mutations are considered clustered if the *P*-value proximity clustering of DNMs is less than 0.01.

621 Supplementary Table 4: Comparison of known haploinsufficient (HI) neurodevelopment genes to HI and  
622 non-HI enrichment models. Genes are ranked by difference in the Akaike's Information Criterion  
623 computed for models where the genes match either expected non-HI PTV enrichment (model 1), or  
624 expected HI protein-truncating variant (PTV) enrichment (model 2).

625



628

629

630

631

632

633

634

635

636

637

638

639

640

641

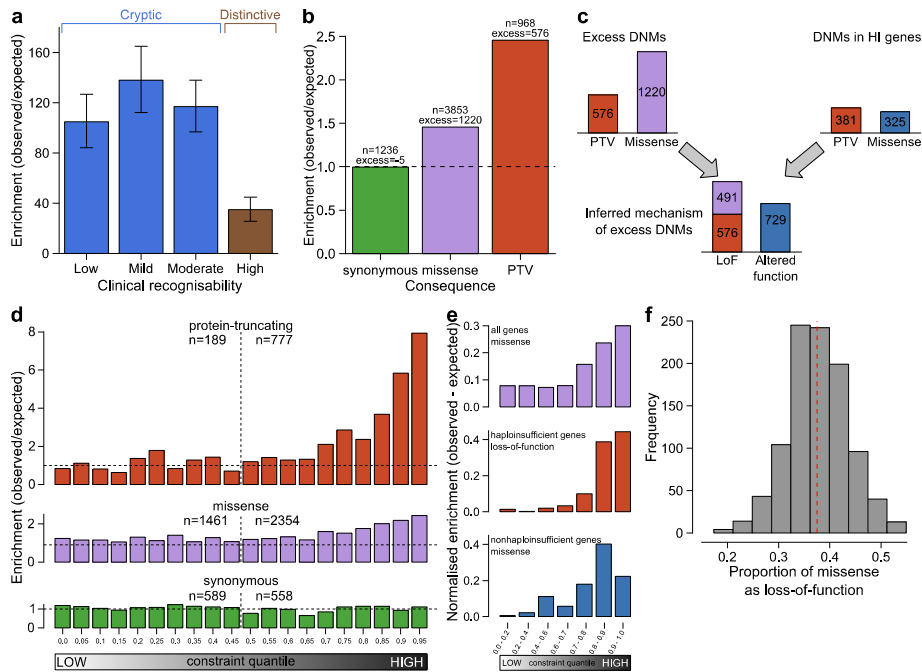
642

643

644

Figure 1: Association of phenotypes with presence of likely pathogenic *de novo* mutations (DNMs). **a**, Odds ratios for binary phenotypes. Positive odds ratios are associated with increased risk of pathogenic DNMs when the phenotype is present. P-values are given for a Fisher's Exact test. **b**, Beta coefficients from logistic regression of quantitative phenotypes versus presence of a pathogenic DNM. All phenotypes aside from length of autozygous regions were corrected for gender as a covariate. The developmental milestones (age to achieve first words, walk independently, sit independently and social smile) were log-scaled before regression. The growth parameters (height, birthweight and occipitofrontal circumference (OFC)) were evaluated as absolute distance from the median. **c**, Relationship between length of autozygous regions chance of having a pathogenic DNM. The regression line is plotted as the dark gray line. The 95% confidence interval for the regression is shaded gray. The autozygosity lengths expected under different degrees of consanguineous unions are shown as vertical dashed lines. n, number of individuals in each autozygosity group. **d**, Relationship between age of fathers at birth of child and number of high confidence DNMs. n, number of high confidence DNMs. **e**, Relationship between age of mothers at birth of child and number of high confidence DNMs. Error bars indicate 95% c.i. n, number of high confidence DNMs.





654

655

656

657

658

659

660

661

662

663

664

665

666

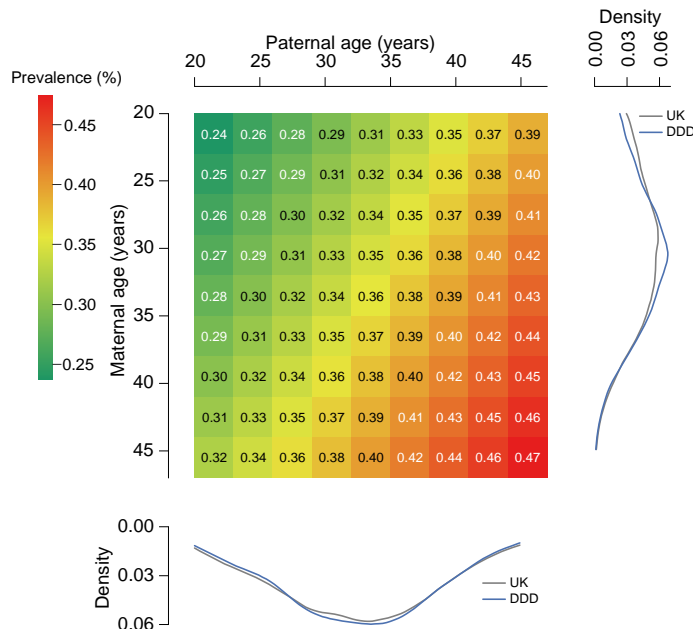
667

668

669

Figure 3: Excess of *de novo* mutations (DNMs). **a**, Enrichment ratios of observed to expected loss-of-function DNMs by clinical recognisability for dominant haploinsufficient neurodevelopmental genes as judged by two consultant clinical geneticists. Error bars indicate 95% CI. **b**, Enrichment of DNMs by consequence normalised relative to the number of synonymous DNMs. **c**, Proportion of excess DNMs with loss-of-function or altered-function mechanisms. Proportions are derived from numbers of excess DNMs by consequence, and numbers of excess truncating and missense DNMs in dominant haploinsufficient genes. **d**, Enrichment ratios of observed to expected DNMs by pLI constraint quantile for loss-of-function, missense and synonymous DNMs. Counts of DNMs in each lower and upper half of the quantiles are provided. **e**, Normalised excess of observed to expected DNMs by pLI constraint quantile. This includes missense DNMs within all genes, loss-of-function including missense DNMs in dominant haploinsufficient genes and missense DNMs in dominant nonhaploinsufficient genes (genes with dominant negative or activating mechanisms). **f**, Proportion of excess missense DNMs with a loss-of-function mechanism. The red dashed line indicates the proportion in observed excess DNMs at the optimal goodness-of-fit. The histogram shows the frequencies of estimated proportions from 1000 permutations, assuming the observed proportion is correct.

670



671

672 Figure 4: Prevalence of live births with developmental disorders caused by dominant *de novo* mutations  
 673 (DNMs). The prevalence within the general population is provided as percentage for combinations of  
 674 parental ages, extrapolated from the maternal and paternal rates of DNMs. Distributions of parental  
 675 ages within the DDD cohort and the UK population are shown at the matching parental axis.

676

677

## 678 Methods

### 679 Family recruitment

680 At 24 clinical genetics centers within the United Kingdom (UK) National Health Service and the  
681 Republic of Ireland, 4,293 patients with severe, undiagnosed developmental disorders and their  
682 parents (4,125 families) were recruited and systematically phenotyped. The study has UK  
683 Research Ethics Committee approval (10/H0305/83, granted by the Cambridge South Research  
684 Ethics Committee and GEN/284/12, granted by the Republic of Ireland Research Ethics  
685 Committee). Families gave informed consent for participation.

686  
687 Clinical data (growth measurements, family history, developmental milestones, etc.) were  
688 collected using a standard restricted-term questionnaire within DECIPHER<sup>34</sup>, and detailed  
689 developmental phenotypes for the individuals were entered using Human Phenotype Ontology  
690 (HPO) terms<sup>35</sup>. Saliva samples for the whole family and blood-extracted DNA samples for the  
691 probands were collected, processed and quality controlled as previously described<sup>15</sup>.

### 693 Exome sequencing

694 Genomic DNA (approximately 1 µg) was fragmented to an average size of 150 base-pairs (bp)  
695 and subjected to DNA library creation using established Illumina paired-end protocols. Adaptor-  
696 ligated libraries were amplified and indexed via polymerase chain reaction (PCR). A portion of  
697 each library was used to create an equimolar pool comprising eight indexed libraries. Each pool  
698 was hybridized to SureSelect ribonucleic acid (RNA) baits (Agilent Human All-Exon V3 Plus with  
699 custom ELID C0338371 and Agilent Human All-Exon V5 Plus with custom ELID C0338371) and  
700 sequence targets were captured and amplified in accordance with the manufacturer's  
701 recommendations. Enriched libraries were subjected to 75-base paired-end sequencing  
702 (Illumina HiSeq) following the manufacturer's instructions.

### 704 Alignment and calling single nucleotide variants, insertions and deletions

705 Mapping of short-read sequences for each sequencing lanelet was carried out using the  
706 Burrows-Wheeler Aligner (BWA; version 0.59)<sup>36</sup> backtrack algorithm with the GRCh37 1000  
707 Genomes Project phase 2 reference (also known as hs37d5). Sample-level BAM improvement  
708 was carried out using the Genome Analysis Toolkit (GATK; version 3.1.1)<sup>37</sup> and SAMtools  
709 (version 0.1.19)<sup>38</sup>. This consisted of a realignment of reads around known and discovered indels  
710 followed by base quality score recalibration (BQSR), with both steps performed using GATK.  
711 Lastly, SAMtools calmd was applied and indexes were created.

712  
713 Known indels for realignment were taken from the Mills Devine and 1000 Genomes Project  
714 Gold set and the 1000 Genomes Project phase low-coverage set, both part of the GATK  
715 resource bundle (version 2.2). Known variants for BQSR were taken from dbSNP 137, also part  
716 of the GATK resource bundle. Finally, single nucleotide variants (SNVs) and indels were called  
717 using the GATK HaplotypeCaller (version 3.2.2); this was run in multisample calling mode using  
718 the complete data set. GATK Variant Quality Score Recalibration (VQSR) was then computed on  
719 the whole data set and applied to the individual-sample variant calling format (VCF) files.

720 DeNovoGear (version 0.54)<sup>39</sup> was used to detect SNV, insertion and deletion *de novo* mutations  
721 (DNMs) from child and parental exome data (BAM files).

722

### 723 Variant annotation

724 Variants in the VCF were annotated with minor allele frequency (MAF) data from a variety of  
725 different sources. The MAF annotations used included data from four different populations of  
726 the 1000 Genomes Project<sup>40</sup> (AMR, ASN, AFR and EUR), the UK10K cohort, the NHLBI GO Exome  
727 Sequencing Project (ESP), the Non-Finnish European (NFE) subset of the Exome Aggregation  
728 Consortium (ExAC) and an internal allele frequency generated using unaffected parents from  
729 the cohort.

730

731 Variants in the VCF were annotated with Ensembl Variant Effect Predictor (VEP)<sup>41</sup> based on  
732 Ensembl gene build 76. The transcript with the most severe consequence was selected and all  
733 associated VEP annotations were based on the predicted effect of the variant on that particular  
734 transcript; where multiple transcripts shared the same most severe consequence, the canonical  
735 or longest was selected. We included an additional consequence for variants at the last base of  
736 an exon before an intron, where the final base is a guanine, since these variants appear to be as  
737 damaging as a splice donor variant<sup>28</sup>.

738

739 We categorized variants into three classes by VEP consequence:

- 740 1. protein-truncating variants (PTV): splice donor, splice acceptor, stop gained, frameshift,  
741 initiator codon, and conserved exon terminus variant.
- 742 2. missense variants: missense, stop lost, inframe deletion, inframe insertion, coding  
743 sequence, and protein altering variant.
- 744 3. silent variants: synonymous.

745

### 746 *De novo* mutation filtering

747 We filtered candidate DNM calls to reduce the false positive rate but maximize sensitivity,  
748 based on prior results from experimental validation by capillary sequencing of candidate  
749 DNMs<sup>15</sup>. Candidate DNMs were excluded if not called by GATK in the child, or called in either  
750 parent, or if they had a maximum MAF greater than 0.01. Candidate DNMs were excluded  
751 when the forward and reverse coverage differed between reference and alternative alleles,  
752 defined as  $P < 10^{-3}$  from a Fisher's exact test of coverage from orientation by allele summed  
753 across the child and parents.

754

755 Candidate DNMs were also excluded if they met two of the three following three criteria: 1) an  
756 excess of parental alternative alleles within the cohort at the DNMs position, defined as  $P < 10^{-3}$   
757 under a one-sided binomial test given an expected error rate of 0.002 and the cumulative  
758 parental depth; 2) an excess of alternative alleles within the cohort in DNMs in a gene, defined  
759 as  $P < 10^{-3}$  under a one-sided binomial test given an expected error rate of 0.002 and the  
760 cumulative depth, or 3) both parents had one or more reads supporting the alternative allele.

761

762 If, after filtering, more than one variant was observed in a given gene for a particular trio, only  
763 the variant with the highest predicted functional impact was kept (protein truncating >  
764 missense > silent).

765

#### 766 *De novo* mutation validation

767 For candidate DNMs of interest, primers were designed to amplify 150-250 bp products  
768 centered around the site of interest. Default primer3 design settings were used with the  
769 following adjustments: GC clamp = 1, human mispriming library used. Site-specific primers were  
770 tailed with Illumina adapter sequences. PCR products were generated with JumpStart AccuTaq  
771 LA DNA polymerase (Sigma Aldrich), using 40 ng genomic DNA as template. Amplicons were  
772 tagged with Illumina PCR primers along with unique barcodes enabling multiplexing of 96  
773 samples. Barcodes were incorporated using Kapa HiFi mastermix (Kapa Biosystems). Samples  
774 were pooled and sequenced down one lane of the Illumina MiSeq, using 250 bp paired end  
775 reads. An in-house analysis pipeline extracted the read count per site and classified inheritance  
776 status per variant using a maximum likelihood approach (see Supplementary Note).

777

#### 778 Individuals with likely pathogenic variants

779 We previously screened 1,133 individuals for variants that contribute to their disorder<sup>15,18</sup>. All  
780 candidate variants in the 1,133 individuals were reviewed by consultant clinical geneticists for  
781 relevance to the individuals' phenotypes. Most diagnosable pathogenic variants occurred *de*  
782 *nov*o in dominant genes, but a small proportion also occurred in recessive genes or under other  
783 inheritance modes. DNMs within dominant DD-associated genes were very likely to be  
784 classified as the pathogenic variant for the individuals' disorder. Due to the time required to  
785 review individuals and their candidate variants, we did not conduct a similar review in the  
786 remainder of the 4,293 individuals. Instead we defined likely pathogenic variants as candidate  
787 DNMs found in autosomal and X-linked dominant DD-associated genes, or candidate DNMs  
788 found in hemizygous DD-associated genes in males. 1,136 individuals in the 4,293 cohort had  
789 variants either previously classified as pathogenic<sup>15,18</sup>, or had a likely pathogenic DNM.

790

#### 791 Gene-wise assessment of DNM significance

792 Gene-specific germline mutation rates for different functional classes were computed<sup>15,23</sup> for  
793 the longest transcript in the union of transcripts overlapping the observed DNMs in that gene.  
794 We evaluated the gene-specific enrichment of PTV and missense DNMs by computing its  
795 statistical significance under a null hypothesis of the expected number of DNMs given the gene-  
796 specific mutation rate and the number of considered chromosomes<sup>23</sup>.

797

798 We also assessed clustering of missense DNMs within genes<sup>15</sup>, as expected for DNMs operating  
799 by activating or dominant negative mechanisms. We did this by calculating simulated  
800 dispersions of the observed number of DNMs within the gene. The probability of simulating a  
801 DNM at a specific codon was weighted by the trinucleotide sequence-context<sup>15,23</sup>. This allowed  
802 us to estimate the probability of the observed degree of clustering given the null model of  
803 random mutations.

804

805 Fisher's method was used to combine the significance testing of missense + PTV DNM  
806 enrichment and missense DNM clustering. We defined a gene as significantly enriched for  
807 DNMs if the PTV enrichment  $P$ -value or the combined missense  $P$ -value less than  $7 \times 10^{-7}$ , which  
808 represents a Bonferroni corrected  $P$ -value of 0.05 adjusted for  $4 \times 18500$  tests ( $2 \times$  consequence  
809 classes tested  $\times$  protein coding genes).

810

### 811 [Composite face generation](#)

812 Families were given the option to have photographs of the affected individual(s) uploaded  
813 within DECIPHER<sup>34</sup>. Using images of individuals with DNMs in the same gene we generated de-  
814 identified realistic average faces (composite faces). Faces were detected using a discriminately  
815 trained deformable part model detector<sup>42</sup>. The annotation algorithm identified a set of 36  
816 landmarks per detected face<sup>43</sup> and was trained on a manually annotated dataset of 3100  
817 images<sup>24</sup>. The average face mesh was created by the Delaunay triangulation of the average  
818 constellation of facial landmarks for all patients with a shared genetic disorder.

819

820 The averaging algorithm is sensitive to left-right facial asymmetries across multiple patients. For  
821 this purpose, we use a template constellation of landmarks based on the average constellations  
822 of 2000 healthy individuals<sup>24</sup>. For each patient, we align the constellation of landmarks to the  
823 template with respect to the points along the middle of the face and compute the Euclidean  
824 distances between each landmark and its corresponding pair on the template. The faces are  
825 mirrored such that the half of the face with the greater difference is always on the same side.

826

827 The dataset used for this work may contain multiple photos for one patient. To avoid biasing  
828 the average face mesh towards these individuals, we computed an average face for each  
829 patient and use these personal averages to compute the final average face. Finally, to avoid any  
830 image in the composite dominating from variance in illumination between images, we  
831 normalised the intensities of pixel values within the face to an average value across all faces in  
832 each average. The composite faces were assessed visually to confirm successful ablation of any  
833 individually identifiable features. Visual assessment of the composite photograph by two  
834 experienced clinical geneticists, alongside the individual patient photos, was performed for all  
835 93 genome-wide significant DD-associated genes for which clinical photos were available for  
836 more than one patient, to remove potentially identifiable composite faces as well as quality  
837 control on the automated composite face generation process. Eighty-one composite faces were  
838 excluded leaving the twelve de-identified composite faces that are shown in Figure 2 and  
839 Extended Data Figure 3. Each of the twelve composite faces that passed de-identification and  
840 quality control was generated from photos of ten or more patients.

841

### 842 [Assessing power of incorporating phenotypic information](#)

843 We previously described a method to assess phenotypic similarity by HPO terms among groups  
844 of individuals sharing genetic defects in the same gene<sup>28</sup>. We examined whether incorporating  
845 this statistical test improved our ability to identify dominant genes at genome-wide  
846 significance. Per gene, we tested the phenotypic similarity of individuals with DNMs in the  
847 gene. We combined the phenotypic similarity  $P$ -value with the genotypic  $P$ -value per gene (the

848 minimum P-value from the DDD-only and meta-analysis) using Fisher's method. We examined  
849 the distribution of differences in P-value between tests without the phenotypic similarity P-  
850 value and tests that incorporated the phenotypic similarity P-value.

851

852 Many (854, 20%) of the DDD cohort experience seizures. We investigated whether testing  
853 within the subset of individuals with seizures improved our ability to find associations for  
854 seizure specific genes. A list of 102 seizure-associated genes was curated from three sources, a  
855 gene panel for Ohtahara syndrome, a currently used clinical gene panel for epilepsy and a panel  
856 derived from DD-associated genes<sup>18</sup>. The P-values from the seizure subset were compared to P-  
857 values from the complete cohort.

858

### 859 [Assessing power of exome vs genome sequencing](#)

860 We compared the expected power of exome sequencing versus genome sequencing to identify  
861 disease genes. Within the DDD cohort, 55 dominant DD-associated genes achieve genome-wide  
862 significance when testing for enrichment of DNMs within genes. We did not incorporate  
863 missense DNM clustering due to the large computational requirements for assessing clustering  
864 in many replicates.

865

866 We assumed a cost of 1,000 USD per individual for genome sequencing. We allowed the cost of  
867 exome sequencing to vary relative to genome sequencing, from 10-100%. We calculated the  
868 number of trios that could be sequenced under these scenarios. Estimates of the improved  
869 power of genome sequencing to detect DNMs in the coding sequence are around 1.05-fold<sup>29</sup>  
870 and we increased the number of trios by 1.0–1.2-fold to allow this.

871

872 We sampled as many individuals from our cohort as the number of trios and counted which of  
873 the 55 DD-associated genes still achieved genome-wide significance for DNM enrichment. We  
874 ran 1000 simulations of each condition and obtained the mean number of genome-wide  
875 significant genes for each condition.

876

### 877 [Associations with presence of likely pathogenic \*de novo\* mutations](#)

878 We tested whether phenotypes were associated with the likelihood of having a likely  
879 pathogenic DNM. We analysed all collected phenotypes which could be coded in either a binary  
880 or quantitative format. Categorical phenotypes (e.g. sex coded as male or female) were tested  
881 by Fisher's exact test while quantitative phenotypes (e.g. duration of gestation coded in weeks)  
882 were tested with logistic regression, using sex as a covariate.

883

884 We investigated whether having autozygous regions affected the likelihood of having a  
885 diagnostic DNM. Autozygous regions were determined from genotypes in every individual, to  
886 obtain the total length per individual. We fitted a logistic regression for the total length of  
887 autozygous regions on whether individuals had a likely pathogenic DNM. To illustrate the  
888 relationship between length of autozygosity and the occurrence of a likely pathogenic DNM, we  
889 grouped the individuals by length and plotted the proportion of individuals in each group with a  
890 DNM against the median length of the group.

891  
892 The effects of parental age on the number of DNMs were assessed using 8,409 high confidence  
893 (posterior probability of DNM > 0.5) unphased coding and noncoding DNMs in 4,293  
894 individuals. A Poisson multiple regression was fit on the number of DNMs in each individual  
895 with both maternal and paternal age at the child's birth as covariates. The model was fit with  
896 the identity link and allowed for overdispersion. This model used exome-based DNMs, and the  
897 analysis was scaled to the whole genome by multiplying the coefficients by a factor of 50, based  
898 on ~2% of the genome being well covered in our data (exons + introns).

#### 900 [Excess of \*de novo\* mutations by consequence](#)

901 We identified the threshold for posterior probability of DNM at which the number of observed  
902 candidate synonymous DNMs equalled the number of expected synonymous DNMs. Candidate  
903 DNMs with scores below this threshold were excluded. We also examined the likely sensitivity  
904 and specificity of this threshold based on validation results for DNMs within a previous  
905 publication<sup>15</sup> in which comprehensive experimental validation was performed on 1,133 trios  
906 that comprise a subset of the families analysed here.

907  
908 The numbers of expected DNMs per gene were calculated per consequence from expected  
909 mutation rates per gene and the 2,407 male and 1,886 females in the cohort. We calculated the  
910 excess of DNMs for missense and PTVs as the ratio of numbers of observed DNMs versus  
911 expected DNMs, as well as the difference of observed DNMs minus expected DNMs.

#### 913 [Ascertainment bias within dominant neurodevelopmental genes](#)

914 We identified 150 autosomal dominant haploinsufficient genes that affect neurodevelopment  
915 within our curated developmental disorder gene set. Genes affecting neurodevelopment were  
916 identified where the affected organs included the brain, or where HPO phenotypes linked to  
917 defects in the gene included either an abnormality of brain morphology (HP:0012443) or  
918 cognitive impairment (HP:0100543) term.

919  
920 The 150 genes were classified for ease of clinical recognition of the syndrome from gene  
921 defects by two consultant clinical geneticists. Genes were rated from 1 (least recognisable) to 5  
922 (most recognisable). Categories 1 and 2 contained 5 and 22 genes respectively, and so were  
923 combined in later analyses. The remaining categories had more than 33 genes per category.  
924 The ratio of observed loss-of-function DNMs to expected loss-of-function DNMs was calculated  
925 for each recognisability category, along with 95% confidence intervals from a Poisson  
926 distribution given observed counts.

927  
928 We estimated the likelihood of obtaining the observed number of PTV DNMs under two  
929 models. Our first model assumed no haploinsufficiency, and mutation counts were expected to  
930 follow baseline mutation rates. Our second model assumed fully penetrant haploinsufficiency,  
931 and scaled the baseline PTV mutation expectations by the observed PTV enrichment in our  
932 known haploinsufficient neurodevelopmental genes, stratified by clinical recognisability into  
933 low (containing genes with our "low", "mild" and "moderate" labels) and high categories. We

934 calculated the likelihoods of both models per gene as the Poisson probability of obtaining the  
935 observed number of PTVs, given the expected mutation rates. We computed the Akaike's  
936 Information Criterion for each model and ranked them by the difference between model 1 and  
937 model 2 ( $\Delta_{AIC}$ ).  
938

### 939 Proportion of *de novo* mutations with loss-of-function mechanism

940 The observed excess of missense/inframe indel DNMs is composed of a mixture of DNMs with  
941 loss-of-function mechanisms and DNMs with altered-function mechanisms. We found that the  
942 excess of PTV DNMs within dominant haploinsufficient DD-associated genes had a greater skew  
943 towards genes with high intolerance for loss-of-function variants than the excess of missense  
944 DNMs in dominant non-haploinsufficient genes. We binned genes by the probability of being  
945 loss-of-function intolerant<sup>30</sup> constraint decile and calculated the observed excess of missense  
946 DNMs in each bin. We modelled this binned distribution as a two-component mixture with the  
947 components representing DNMs with a loss-of-function or function-altering mechanism. We  
948 identified the optimal mixing proportion for the loss-of-function and altered-function DNMs  
949 from the lowest goodness-of-fit (from a spline fitted to the sum-of-squares of the differences  
950 per decile) to missense/inframe indels in all genes across a range of mixtures.  
951

952 The excess of DNMs with a loss-of-function mechanism was calculated as the excess of DNMs  
953 with a VEP loss-of-function consequence, plus the proportion of the excess of missense DNMs  
954 at the optimal mixing proportion.  
955

956 We independently estimated the proportions of loss-of-function and altered-function. We  
957 counted PTV and missense/inframe indel DNMs within dominant haploinsufficient genes to  
958 estimate the proportion of excess DNMs with a loss-of-function mechanism, but which were  
959 classified as missense/inframe indel. We estimated the proportion of excess DNMs with a loss-  
960 of-function mechanism as the PTV excess plus the PTV excess multiplied by the proportion of  
961 loss-of-function classified as missense.  
962

### 963 Prevalence of developmental disorders from dominant *de novo* mutations

964 We estimated the birth prevalence of monoallelic developmental disorders by using the  
965 germline mutation model. We calculated the expected cumulative germline mutation rate of  
966 truncating DNMs in 238 haploinsufficient DD-associated genes. We scaled this upwards based  
967 on the composition of excess DNMs in the DDD cohort using the ratio of excess DNMs (n=1816)  
968 to DNMs within dominant haploinsufficient DD-associated genes (n=412). Around 10% of DDs  
969 are caused by *de novo* CNVs<sup>44,45</sup>, which are underrepresented in our cohort as a result of prior  
970 genetic testing. If included, the excess DNM in our cohort would increase by 21%, therefore we  
971 scaled the prevalence estimate upwards by this factor.  
972

973 Mothers aged 29.9 and fathers aged 29.5 have children with 77 DNMs per genome on  
974 average<sup>21</sup>. We calculated the mean number of DNMs expected under different combinations of  
975 parental ages, given our estimates of the extra DNMs per year from older mothers and fathers.  
976 We scaled the prevalence to different combinations of parental ages using the ratio of expected

977 mutations at a given age combination to the number expected at the mean cohort parental  
978 ages.

979

980 To estimate the annual number of live births with developmental disorders caused by DNMs,  
981 we obtained country population sizes, birth rates, age at first birth<sup>46</sup>, and calculated global birth  
982 rate (18.58 live births/1000 individuals) and age at first birth (22.62 years), weighted by  
983 population size. We calculated the mean age when giving birth (26.57 years) given a total  
984 fertility rate of 2.45 children per mother<sup>47</sup>, and a mean interpregnancy interval of 29 months<sup>48</sup>.  
985 We calculated the number of live births given our estimate of DD prevalence caused by DNMs  
986 at this age (0.00288), the global population size (7.4 billion individuals) and the global birth rate.

987

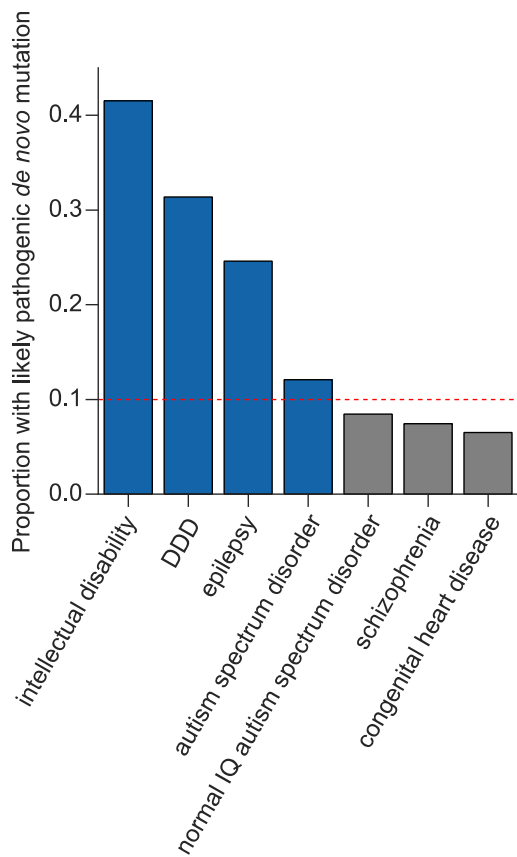
#### 988 [Code availability](#)

989 Source code for filtering candidate DNMs, testing DNM enrichment, DNM clustering and  
990 phenotypic similarity can be found here: <https://github.com/jeremymcrae/denovoFilter>,  
991 <https://github.com/jeremymcrae/mupit>, <https://github.com/jeremymcrae/denovonear>,  
992 [https://github.com/jeremymcrae/hpo\\_similarity](https://github.com/jeremymcrae/hpo_similarity)

993

994 Extended Data Figures

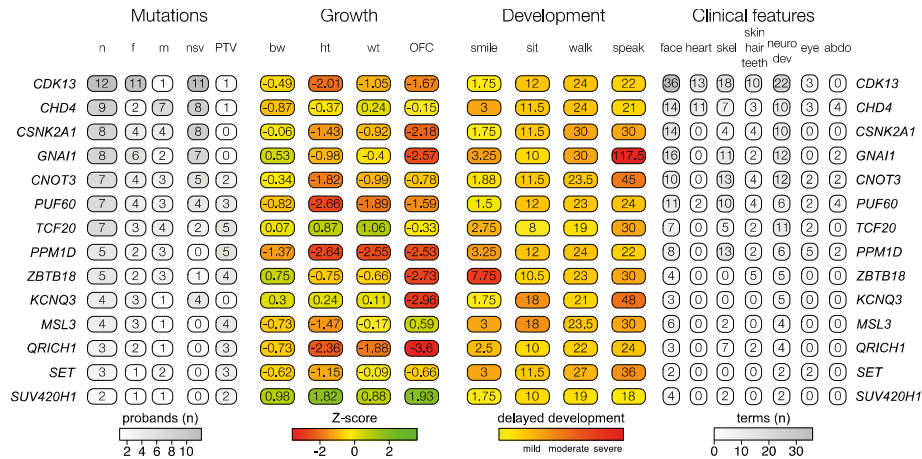
995



996

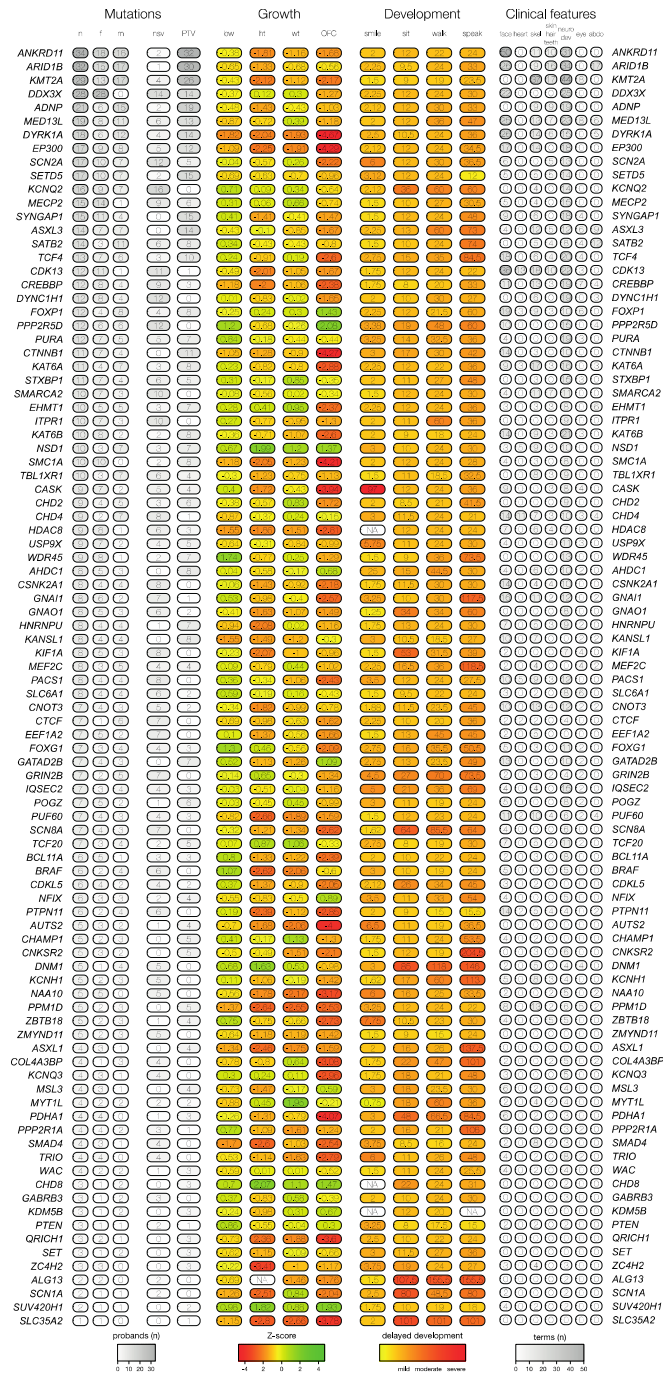
997 Extended Data Figure 1: Proportion of individuals with a *de novo* mutation (DNM) likely to be  
998 pathogenic. These only included individuals with protein altering or protein truncating DNMs in  
999 dominant or X-linked dominant developmental disorder (DD) associated genes, or males with DNMs in  
1000 hemizygous DD-associated genes. The proportions given are for those individuals with any DNMs rather  
1001 than the total number of individuals in each subset. Cohorts included in the DNM meta-analyses are  
1002 shaded blue.

1003



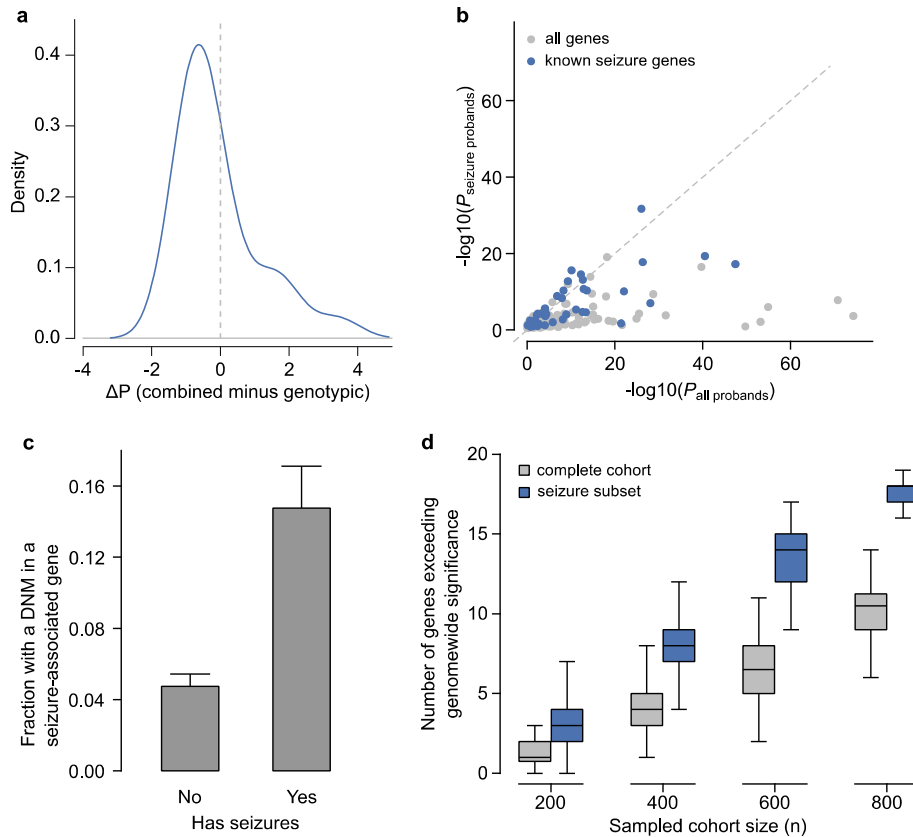
1004  
1005 Extended Data Figure 2: Phenotypic summary of genes without previous compelling evidence.  
1006 Phenotypes are grouped by type. The first group indicates counts of individuals with DNMs per gene by  
1007 sex (m: male, f: female), and by functional consequence (nsv: nonsynonymous variant, PTV: protein-  
1008 truncating variant). The second group indicates mean values for growth parameters: birthweight (bw),  
1009 height (ht), weight (wt), occipitofrontal circumference (OFC). Values are given as standard deviations  
1010 from the healthy population mean derived from ALSPAC data. The third group indicates the mean age  
1011 for achieving developmental milestones: age of first social smile, age of first sitting unassisted, age of  
1012 first walking unassisted and age of first speaking. Values are given in months. The final group  
1013 summarises Human Phenotype Ontology (HPO)-coded phenotypes per gene, as counts of HPO-terms  
1014 within different clinical categories.

1015



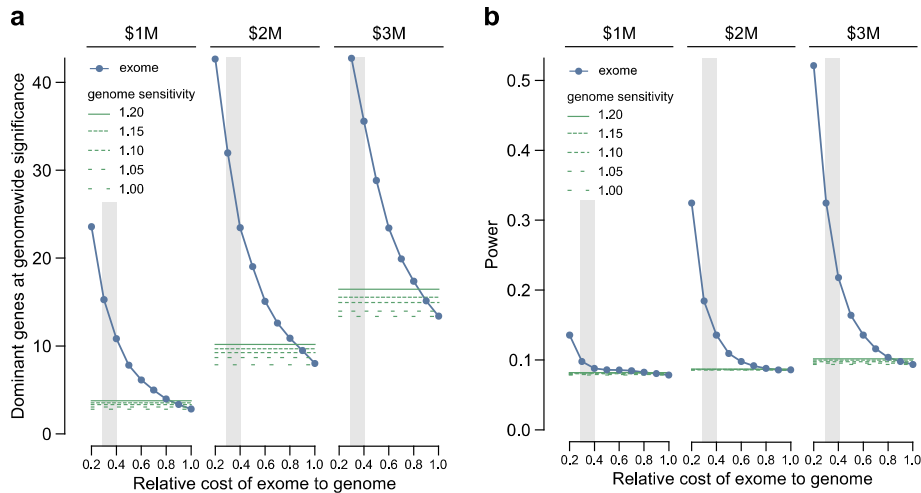
1016  
 1017 Extended Data Figure 3: Phenotypic summary of individuals with *de novo* mutations in genes achieving  
 1018 genome-wide significance. Phenotypes are grouped by type. The first group indicates counts of  
 1019 individuals with DNMs per gene by sex (m: male, f: female), and by functional consequence (nsv:  
 1020 nonsynonymous variant, PTV: protein-truncating variant). The second group indicates mean values for  
 1021 growth parameters: birthweight (bw), height (ht), weight (wt), occipitofrontal circumference (OFC).  
 1022 Values are given as standard deviations from the healthy population mean derived from ALSPAC data.  
 1023 The third group indicates the mean age for achieving developmental milestones: age of first social smile,  
 1024 age of first sitting unassisted, age of first walking unassisted and age of first speaking. Values are given in  
 1025 months. The final group summarises Human Phenotype Ontology (HPO)-coded phenotypes per gene, as  
 1026 counts of HPO-terms within different clinical categories.





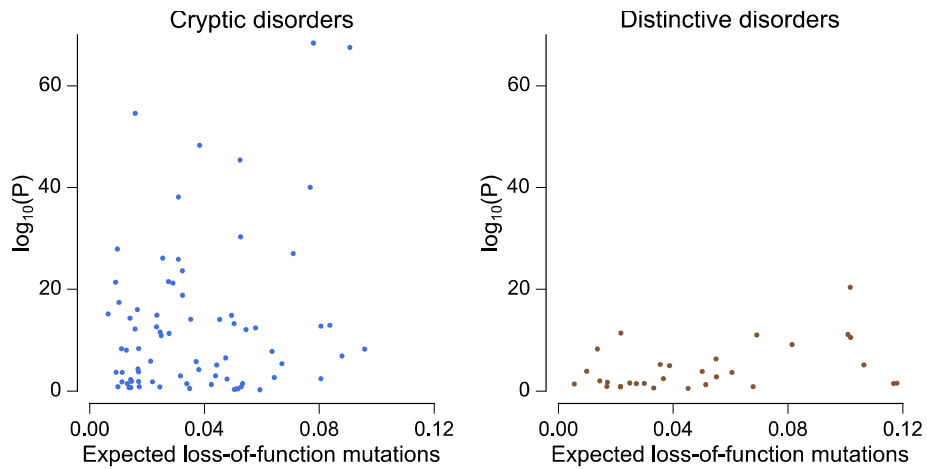
1032  
 1033 Extended Data Figure 5: Effect of clustering by phenotype on the ability to identify genomewide  
 1034 significant genes. **a**, Comparison of P-values derived from genotypic information alone versus P-values  
 1035 that incorporate genotypic information and phenotypic similarity. **b**, Comparison of P-values from tests  
 1036 in the complete DDD cohort versus tests in the subset with seizures. Genes that were previously linked  
 1037 to seizures are shaded blue. **c**, Proportion of cohort with a *de novo* mutation (DNM) in a seizure-  
 1038 associated gene, stratified by whether seizure-affected status. Bars indicate 95% CI. **d**, Comparison of  
 1039 power to identify genomewide significant genes in probands with seizures, versus the unstratified  
 1040 cohort, at matched sample sizes.

1041



1042  
 1043 Extended Data Figure 6: Power of genome versus exome sequencing to discover dominant genes  
 1044 associated with developmental disorders. **a**, the number of genes exceeding genome-wide significance  
 1045 was estimated at three different fixed budgets (1 million (M) USD, 2M and 3M) and a range of relative  
 1046 sensitivities for genomes versus exomes to detect *de novo* mutations. The number of genes identifiable  
 1047 by exome sequencing are shaded blue, whereas the number of genes identifiable by genome  
 1048 sequencing are shaded green. The regions where exome sequencing costs 30-40% of genome  
 1049 sequencing are shaded with a grey background, which corresponds to the price differential in 2016. **b**,  
 1050 simulated estimates of power to detect loss-of-function genes in the genome at different cohort sizes,  
 1051 given fixed budgets.

1052



1053

1054

Extended Data Figure 7: Gene-wise significance of neurodevelopmental genes versus the expected number of mutations per gene. Points are shaded by clinical recognisability classification. Genes have

1055

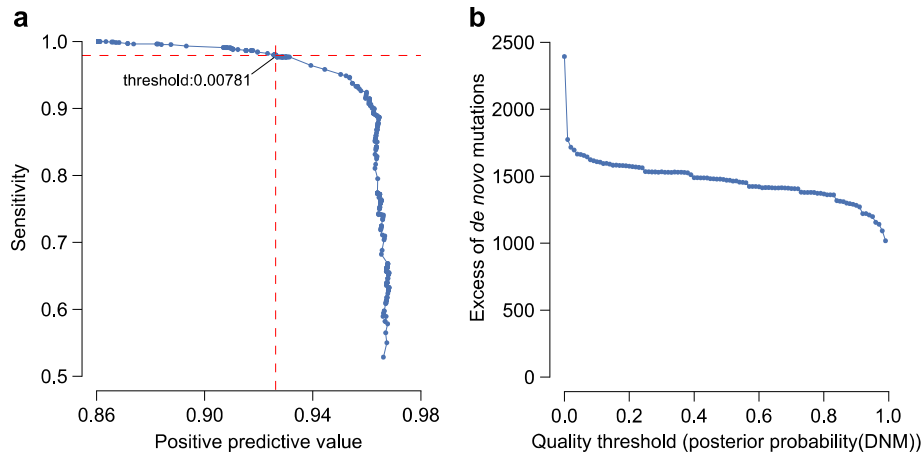
been separated into two plots, one plot with genes for cryptic disorders with low, mild or moderate

1056

clinical recognisability, and one plot with genes for distinctive disorders with high clinical recognisability.

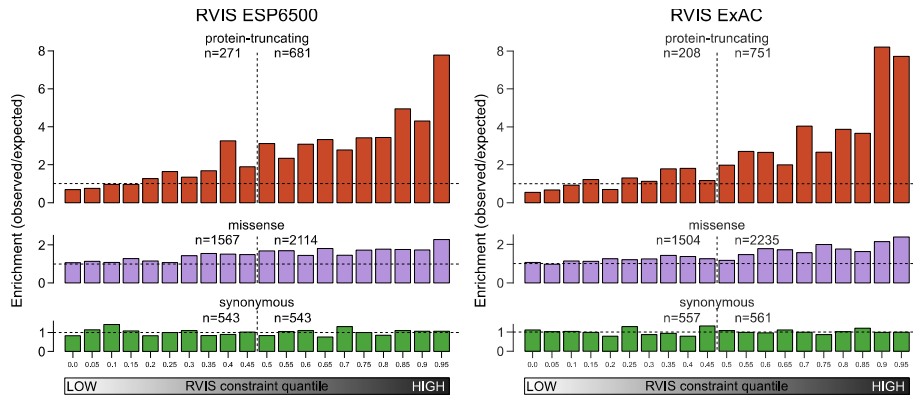
1058

1059



1060  
 1061 Extended Data Figure 8: Stringency of *de novo* mutation (DNM) filtering. **a**, Sensitivity and specificity of  
 1062 DNM validations within sets filtered on varying thresholds of DNM quality (posterior probability of  
 1063 DNM). The analysed DNMs were restricted to sites identified within the earlier 1133 trios<sup>15</sup>, where all  
 1064 candidate DNMs underwent validation experiments. The labelled value is the quality threshold at which  
 1065 the number of candidate synonymous DNMs equals the number of expected synonymous mutations  
 1066 under a null germline mutation rate. **b**, Excess of missense and loss-of-function DNMs at varying DNM  
 1067 quality thresholds. The DNM excess is adjusted for the sensitivity and specificity at each threshold.

1068



1069

1070 Extended Data Figure 9: Enrichment of de novo mutations by consequence type, across RVIS functional  
 1071 constraint quantiles. A comparison of enrichment for RVIS values generated from ESP650 data versus  
 1072 ExAC data is provided.

1073