
Goldmann JM, Seplyarskiy VB, Wong WSW, Vilboux T, Neerincx PD, Bodian DL, Solomon BD, Veltman JA, Deeken JF, Gilissen C, Niederhuber JE.

[Germline de novo mutation clusters arise during oocyte aging in genomic regions with high double-strand-break incidence.](#)

Nature Genetics 2018,

<https://doi.org/10.1038/s41588-018-0071-6>

Copyright:

This is the authors accepted manuscript of an article that has been published in its final form by Nature Publishing Group, 2018.

DOI link to article:

<https://doi.org/10.1038/s41588-018-0071-6>

Date deposited:

06/03/2018

Embargo release date:

05 September 2018

1 Germline *de novo* mutation clusters arise during 2 oocyte aging in genomic regions with increased 3 double-strand break incidence

4 Jakob M. Goldmann^{1*}, Vladimir B. Seplyarskiy^{2,3*}, Wendy S.W. Wong^{4*}, Thierry Vilboux⁴,
5 Pieter B. Neerincx^{5,6}, Dale L. Bodian⁴, Benjamin D. Solomon^{7,8}, Joris A. Veltman^{9,10}, John F.
6 Deeken⁴, Christian Gilissen^{9#}, John E. Niederhuber^{4,11#}

7
8 ¹Department of Human Genetics, Radboud Institute for Molecular Life Sciences, Radboud
9 University Medical Center, Geert Grooteplein 10, 6525 GA Nijmegen, the Netherlands

10 ²Division of Genetics, Brigham & Women's Hospital, Harvard Medical School, Boston,
11 Massachusetts, USA.

12 ³Institute for Information Transmission Problems of the Russian Academy of Sciences
13 (Kharkevich Institute), Bolshoi Karetny pereulok 19, Moscow 127994, Russia

14 ⁴Inova Translational Medicine Institute (ITMI), Inova Health Systems, Falls Church, VA, USA

15 ⁵Department of Genetics, University of Groningen, University Medical Center Groningen,
16 Groningen, The Netherlands

17 ⁶Genomics Coordination Center, University of Groningen, University Medical Center
18 Groningen, Groningen, The Netherlands

19 ⁷Department of Pediatrics, Inova Children's Hospital, Inova Health System, Falls Church, VA,
20 USA

21 ⁸Department of Pediatrics, Virginia Commonwealth University School of Medicine, 1201 E
22 Marshall St, Richmond, VA, USA

23 ⁹Department of Human Genetics, Donders Centre for Neuroscience, Radboud University
24 Medical Center , Geert Grooteplein 10, 6525 GA Nijmegen, the Netherlands

25 ¹⁰Institute of Genetic Medicine, International Centre for Life, Newcastle University,
26 Newcastle upon Tyne, United Kingdom

27 ¹¹Johns Hopkins University School of Medicine, 733 North Broadway Street, Baltimore, MD,
28 USA

29
30 *These authors contributed equally.

31 # These authors jointly supervised this work.

32

33 To whom correspondence should be addressed: Christian.gilissen@radboudumc.nl and
34 John.Niederhuber@inova.org

35

36 Clustering of mutations has been observed in cancer genomes as well as for germline *de*
37 *novo* mutations (DNMs). We identified 1,796 clustered DNMs (cDNMs) within whole-
38 genome sequencing data from 1,291 parent-offspring trios to investigate their patterns and
39 inferred a mutational mechanism. We found that the number of clusters on the maternal
40 allele was positively correlated with maternal age and that these consist of more individual
41 mutations with larger inter-mutational distances compared to paternal clusters. More than
42 50% of maternal clusters were located on chromosomes 8, 9 and 16, in regions with an
43 overall increased maternal mutation rate. Maternal clusters in these regions showed a
44 distinct mutation signature characterized by C>G transversions. Finally, we found that
45 maternal clusters associate with processes involving double-stranded-breaks (DSBs) such as
46 meiotic gene conversions and *de novo* deletions events. This suggest accumulation of DSB-
47 induced mutations throughout oocyte aging as an underlying mechanism for maternal
48 mutation clusters.

49

50

51 *De novo* mutations (DNMs) arise spontaneously in parental gametes and result in
52 approximately 50-100 germline mutations in their offspring¹⁻⁴. As such, DNMs are both
53 drivers of evolution as well as a common cause of sporadic disorders. The total number of
54 DNMs is highly correlated with paternal age and, to a lesser extent, with maternal age²⁻⁴. The
55 paternal age effect, giving rise to about one additional DNM in the offspring per year of life
56 of the father before conception, is thought to be due to the higher number of cell divisions
57 that spermatogonial cells of older men have undergone prior to this period^{5,6}. The
58 mechanisms underlying the maternal age effect, giving rise to about one additional DNM per
59 4 years of life of the mother, are still unknown. Approximately 2-3% of all DNMs in the
60 offspring occur in close spatial proximities (below 20kb) as clustered mutations^{4,7-11}. These
61 clustered DNMs (cDNMs) have a distinct nucleotide substitution spectrum with an
62 enrichment of C>G mutations, suggesting mutational mechanisms different from
63 unclustered DNMs^{4,9,10,12,13}. The precise composition of the mutation spectrum also varies
64 with the inter-mutational distances of the clusters^{10,14}. Contrary to unclustered DNMs, no
65 paternal bias has been observed for the number of cDNMs^{4,9,12}. Here, we investigated
66 cDNMs, their potential contribution to the paternal and maternal age effect on the total
67 number of DNMs, and the possible mechanisms underlying their occurrence.

68 Whole genomes of 1,291 parent-offspring trios from the Inova Translational Medicine
69 Institute longitudinal childhood study cohort were sequenced using Illumina HiSeq2000 with
70 average 40x coverage by Illumina services (La Jolla, USA; **Table 1, Supplementary Table 1**).
71 This cohort represents a sample of the general population of average health giving birth at a
72 single hospital¹⁵. After quality control, we identified 73,755 high-confidence DNMs using a
73 random forest classifier (**Online Methods, Supplementary Table 2**). We defined cDNMs as
74 DNMs within the same individual with all pair-wise inter-mutational distances smaller than
75 20kb. In total we identified 1,796 cDNMs (2.4% of all DNMs) distributed across 799 clusters,
76 with 2-10 mutations per cluster, of which 678 clusters (85%) consisted of exactly two
77 mutations (**Supplementary Tables 3-6**). 144 cDNMs in 72 clusters were directly adjacent. By
78 performing read-phasing, we successfully identified the parent-of-origin for 700 cDNMs
79 (39.0% of all cDNMs) across 400 clusters (**Table 1, Supplementary Table 7-8**). In 98.0%
80 (204/208) of the fully phased clusters, all cDNMs arose on the same allele, which is in line
81 with our simulations of the false detection rate of cluster detection (**Supplementary Table**
82 **9**). In contrast to unclustered DNMs, we did not observe an excess of cDNMs on the paternal
83 allele (202 maternal clusters and 198 paternal clusters, chi-square goodness-of-fit $p=0.84$). In
84 addition, we created a validation dataset based on four independent studies with phased
85 DNMs from whole-genome sequencing (WGS)^{4,9,10,12}, resulting in a total of 1,643 cDNMs
86 across 745 clusters, with limited information on parental ages (**Table 1, Supplementary**
87 **Table 10**).

88 To investigate the contribution of cDNMs to the parental age effects, we used a linear
89 regression model to correlate the age of each parent with the number of phased cDNMs in
90 the offspring. Although the number of paternal cDNMs did not show a significant correlation
91 with the paternal age ($p=0.087$), we found a highly significant correlation of maternal cDNMs
92 with maternal age ($p<10^{-10}$). This effect was similar in our replication cohort (maternal
93 $p=0.00155$ and paternal $p=0.319$, **Supplementary Figures 1,2**). In the primary cohort, the
94 cDNMs accounted for 23% (95% c.i. 7-38%) of the maternal age effect (p -value for maternal
95 age effect of unclustered DNMs $p=1.5\times 10^{-19}$). For the clusters where only a subset of cDNMs
96 could be phased, we extrapolated the parent-of-origin. Based on this extrapolation, we also

97 observe a significant paternal age effect of a smaller amplitude than the maternal age effect
98 (paternal effect size $p=0.026/\text{year}$, $p=8\times 10^{-7}$, maternal effect size $0.041/\text{year}$, $p=3\times 10^{-11}$).
99 While in the primary cohort, only 5% of the probands with the youngest mothers had one or
100 more maternal cDNMs per genome, this was more than 5 times higher (risk ratio test,
101 $p=1.4\times 10^{-11}$; c.i. 3.0-9.4) in probands from the oldest mothers (27% having a maternal cDNM,
102 **Figure 1a**). This difference was not significant for the paternal cDNMs (13% vs 19%; risk ratio
103 test $p=0.08$; risk ratio 1.42; 95% c.i. 0.95-2.12). In the replication cohort, the risk ratio was
104 3.02 for maternal cDNMs (c.i. 1.22-7.45; $p=0.011$, **Supplementary Figure 3**) and 0.60 (c.i.
105 0.30-1.22; $p=0.15$) for paternal cDNMs.

106 We found that this maternal age effect of clusters stems mostly from clusters with inter-
107 mutational distances greater than 1kb (**Figures 1b,c, Supplementary Tables 11 and 12,**
108 **Supplementary Figure 3**). Strikingly, the maximum number of DNMs in the phased clusters
109 of an individual correlates positively with maternal age ($p<10^{-10}$, replication cohort $p<10^{-4}$),
110 but is correlated only marginally significant with paternal age ($p=0.050$, replication cohort
111 $p=0.408$, **Figure 1d,e, Supplementary Figure 3**). Clusters with more than two mutations were
112 4.2 times more likely to contain maternal cDNMs than paternal cDNMs (95% c.i. 2.5 – 7.6;
113 $p=1.7\times 10^{-7}$). These results show that maternal clusters contain more cDNMs with larger
114 inter-mutational distances.

115 We previously observed that maternal DNMs are enriched within specific genomic regions
116 on chromosomes 8 and 16⁴. In this study, we found that 58.4% of maternal cDNMs localize
117 to chromosomes 8, 9 and 16 ($p<10^{-16}$, replication cohort $p<10^{-16}$, Chi-square test; **Figure 2a,**
118 **Supplementary Figures 4 and 5**). This in contrast to paternal cDNMs for which the number
119 correlates with chromosome length ($R^2=0.72$, $p=6*10^{-7}$, replication cohort $R^2=0.43$, $p=0.001$).
120 The maternal cDNMs on these three chromosomes occur specifically in regions that are also
121 enriched for maternal unclustered DNMs (**Figure 2b, Supplementary Figures 6 and 7,**
122 **Supplementary Note 1**) and their mutation spectrum is strongly enriched for C>G
123 substitutions compared to other maternal cDNMs (**Figure 2c,d**, bootstrapping $p=0.022$).
124 These observations are further supported by the patterns of clusters with more than two
125 cDNMs, which are more likely to be on the maternal allele. These clusters are also enriched
126 on the chromosomes 8, 9 and 16 (Chi-square test $p=3\times 10^{-09}$), and show an excess of C>G
127 substitutions (Chi-square test $p=4.5\times 10^{-11}$). Taken together, this suggests a different
128 mutational mechanism for maternal cDNMs in these regions compared to the rest of the
129 genome.

130 To confirm these findings, we created a dataset of (unphased) clustered SNP variants based
131 on publically available population-based genetic data¹⁶ (**Online Methods**). This resulted in
132 1,146,891 clustered SNPs (cSNPs) across 522,487 clusters (**Supplementary Table 13**). We
133 found that cSNPs on chromosomes that are enriched for maternal cDNMs are enriched for
134 C>G substitutions (bootstrapping test, see **Online Methods**, $p<0.001$, **Figure 2e**). To further
135 investigate this association, we calculated a genome-wide score for C>G cSNP enrichment
136 (**Supplementary Methods**) and found that the number of maternal cDNMs in a region is
137 significantly correlated with high C>G scores (Poisson regression $p<10^{-16}$ for maternal
138 cDNMs, $p=0.33$ for paternal cDNMs, **Supplementary Figure 8**). Using this method we also
139 identified an additional region on chromosome 2 that is enriched for maternal cDNMs
140 (**Figure 2f**). This strong association between C>G scores of cSNPs with maternal cDNMs

141 highlights the maternal clusters' profound contribution to population polymorphisms in
142 these regions.

143 The observed age-effect of maternal cDNMs suggests underlying mechanisms that are active
144 during oocyte aging, a process that has been associated with the decreasing efficiency of
145 double-strand break (DSB) repair¹⁷⁻¹⁹. We therefore hypothesized that the maternal-aging
146 associated clusters arise via a DSB-associated mechanism and investigated the occurrence of
147 cDNMs at regions that are associated with DSBs. As proxies for DSB sites we used (1) sites of
148 *de novo* meiotic gene conversion (MGC), (2) the flanking regions of *de novo* CNV breakpoints
149 in our cohort, and (3) known recombination hotspots²⁰.

150 We used MGC sites from Halldorsson et al.²¹ and found that these events co-localize with
151 maternal cDNMs significantly more often than expected by chance ($p=0.002$, permutation
152 testing, **Figure 3a, Supplementary Table 14**). This association is not significant for paternal
153 MGCs with paternal cDNMs ($p=0.609$).

154 In our primary cohort, we identified 45 high-quality *de novo* CNVs, of which 5 have a total of
155 17 DNMs within 100kb flanking the breakpoints (**Figure 3b, Supplementary Methods**).
156 Exactly 12 of these 17 DNMs are cDNMs, which constitutes a high enrichment ($p = 2.58 \times 10^{-16}$,
157 Fisher's exact test). For 6 of these DNMs the parent-of-origin was resolved and in all cases
158 the DNMs arose from the maternal allele ($p=0.03$, Fisher's exact test). In concordance with
159 this, all 5 CNVs are deletions of the maternal allele (**Supplementary Table 15**). An
160 arrangement of several DNMs and a *de novo* deletion on the same allele within the same
161 generation is very unlikely to occur by chance and suggests a single event as a common
162 cause. In our replication cohort, we also discovered 5 *de novo* deletion events. Two of these
163 CNVs have a total of 4 DNMs from the same individual within 100kb of the CNV breakpoints,
164 and two of these are within 20kb of each other (**Supplementary Figure 9**), again showing an
165 enrichment of cDNMs ($p=0.002$, Fisher's exact test). Interestingly, cSNPs were significantly
166 closer to CNV breakpoints than expected by chance (**Figure 3c**, Mann-Whitney test on 1% of
167 data $p < 10^{-9}$), corroborating the co-segregation of CNV events and clustered mutations.

168 Finally, we used gender specific recombination scores²⁰ to assess whether cDNMs occur
169 more often at regions of high recombination. We did not find a significant overlap of
170 maternal cDNMs with regions of high maternal recombination ($p=0.204$ permutation testing,
171 **Figure 3d, Supplementary Table 14**). Nevertheless, genomic regions with maternal cDNMs
172 had higher sex-matched recombination scores than regions with only unclustered maternal
173 DNMs (primary cohort $p=0.019$, replication cohort $p=0.13$) and higher than regions of
174 paternal cDNMs (primary cohort $p=0.004$, replication cohort $p=0.29$; **Supplementary Figure**
175 **10**). In addition, genomic regions with cSNPs have significantly higher recombination rates
176 than genomic regions without cSNPs ($p=3.91 \times 10^{-49}$). Interestingly, our observed association
177 of cDNMs with recombination rates is much smaller than the observed association with
178 MGCs. This is in line with the maternal age effect of MGCs being larger compared to the
179 maternal age effect of the crossover rate^{21,22}. Campbell et al. found that, with increasing
180 maternal age, recombination occurs more frequently outside of recombination hotspots²³. In
181 addition, these events were increasingly deregulated, appearing in closer proximity of each
182 other than expected based on models of crossover interference. The fact that recombination
183 events have shown to be mutagenic²⁴⁻²⁶ suggests that this increase in deregulated
184 recombination events may be the underlying cause of cDNM formation. In this study, we
185 found that that chromosomes 8, 9 and 16 are heavily enriched for maternal clusters and

186 strikingly these chromosomes also have the highest degree of cross-over events escaping
187 interference²³.

188 Additionally, cDNM mutational spectra, and in particular those of maternal cDNMs, are very
189 similar to the previously identified signature of somatic mutations caused by deficiency in
190 homologous recombination repair of DSBs^{27,28} (Signature 3, **Supplementary Figure 11**). The
191 proband's parents are very unlikely to suffer from DNA repair deficiencies such as those
192 underlying cancer mutation profiles, therefore this finding is in agreement with a key role for
193 imperfect DSB repair after unregulated recombination in the formation of maternal
194 mutation clusters. However, we found no statistical association between variants in genes
195 involved in homologous recombination repair or in establishing recombination sites²⁹⁻³¹
196 (**Supplementary Tables 16 and 17**).

197 Although the formation of clustered mutations has the potential to be highly deleterious,
198 there seems to be selection in favor of high recombination rates in ageing oocytes^{32,33}. It has
199 been argued that these high recombination rates provide protection against aneuploidies³³⁻
200 ³⁴, the risk of which increases with maternal age. Taken together, our results show that
201 deregulated recombination is a likely cause for DNMs clusters, whereas replicative errors are
202 not a likely cause. A recent paper that studied genome-wide *de novo* mutations in a cohort
203 of 1,548 Icelanders also found that clustered mutations increase faster with maternal than
204 paternal age³⁵. In addition, the authors observed a non-uniform distribution of these events
205 across the genome³⁵ corresponding with the regions that we reported here.

206

207 **URLs**

208 goleft indexcov: <https://github.com/brentp/goleft/tree/master/indexcov>

209 agg gvcf aggregation tool: <https://github.com/Illumina/agg>

210 **Acknowledgements**

211 This study was funded by the Inova Health System with support from Fairfax County and the
212 philanthropic support from the Odeen family. We thank the Inova translational medicine
213 institute staff for supporting the study. We also thank the families who participated in the
214 genomic studies that made this research possible. This work was partly financially supported
215 by grants from the Netherlands Organization for Scientific Research (916-14-043 to
216 C.Gilissen and 918-15-667 to J. Veltman), and the European Research Council (ERC Starting
217 grant DENOVO 281964 to J. Veltman).

218 This study makes use of data generated by the Genome of the Netherlands Project. A full list
219 of the investigators is available from www.nlgenome.nl. Funding for the project was
220 provided by the Netherlands Organization for Scientific Research under award number
221 184021007, dated July 9, 2009 and made available as a Rainbow Project of the Biobanking
222 and Biomolecular Research Infrastructure Netherlands (BBMRI-NL). The sequencing was
223 carried out in collaboration with the Beijing Institute for Genomics (BGI).

224 **Author contributions**

225 C.G. and J.E.N. designed the study. J.M.G., V.B.S., and W.S.W.W. performed the data
226 analyses. W.S.W.W carried out QC, and *de novo* mutations calling. T.V. performed the Sanger
227 validation. B.D.S., J.F.D, and J.E.N. supervised the data collection, sequencing and writing of
228 the manuscript. D.B. assisted in data analyses and interpretation. J.M.G., V.B.S., W.S.W.W.,
229 J.A.V. and C.G. drafted the manuscript. P.B.N. acquired part of the replication data. All
230 authors contributed to the final version of the paper.

231 **Competing Financial Interests**

232 The authors do not declare any competing financial interests.

233

234 References

235

- 236 1. Veltman, J.A. & Brunner, H.G. De novo mutations in human genetic disease. *Nature*
237 *reviews. Genetics* **13**, 565-75 (2012).
- 238 2. Kong, A. *et al.* Rate of de novo mutations and the importance of father's age to
239 disease risk. *Nature* **488**, 471-5 (2012).
- 240 3. Wong, W.S.W. *et al.* New observations on maternal age effect on germline de novo
241 mutations. *Nature communications* **7**, 10486 (2016).
- 242 4. Goldmann, J.M. *et al.* Parent-of-origin-specific signatures of de novo mutations.
243 *Nature Genetics* (2016).
- 244 5. Crow, J.F. The origins, patterns and implications of human spontaneous mutation.
245 *Nature Reviews Genetics* **1**, 40-47 (2000).
- 246 6. Segurel, L., Wyman, M.J. & Przeworski, M. Determinants of mutation rate variation in
247 the human germline. *Annu Rev Genomics Hum Genet* **15**, 47-70 (2014).
- 248 7. Michaelson, Jacob J. *et al.* Whole-Genome Sequencing in Autism Identifies Hot Spots
249 for De Novo Germline Mutation. *Cell* **151**, 1431-1442 (2012).
- 250 8. Schrider, D.R., Hourmozdi, J.N. & Hahn, M.W. Pervasive multinucleotide mutational
251 events in eukaryotes. *Current biology : CB* **21**, 1051-4 (2011).
- 252 9. Yuen, R.K. *et al.* Genome-wide characteristics of de novo mutations in autism. *npj*
253 *Genomic Medicine* **1**, 16027 (2016).
- 254 10. Besenbacher, S. *et al.* Multi-nucleotide de novo Mutations in Humans. *PLOS Genetics*
255 **12**, e1006315 (2016).
- 256 11. Terekhanova, N.V., Bazykin, G.A., Neverov, A., Kondrashov, A.S. & Seplyarskiy, V.B.
257 Prevalence of Multinucleotide Replacements in Evolution of Primates and Drosophila.
258 *Molecular Biology and Evolution* **30**, 1315-1325 (2013).
- 259 12. Francioli, L.C. *et al.* Genome-wide patterns and properties of de novo mutations in
260 humans. *Nature Genetics advance on*(2015).
- 261 13. Rahbari, R. *et al.* Timing, rates and spectra of human germline mutation. *Nat Genet*
262 **48**, 126-133 (2016).
- 263 14. Harris, K. & Nielsen, R. Error-prone polymerase activity causes multinucleotide
264 mutations in humans. *Genome research* **24**, 1445-54 (2014).
- 265 15. Bodian, D.L. *et al.* Utility of whole-genome sequencing for detection of newborn
266 screening disorders in a population cohort of 1,696 neonates. *Genetics in medicine : official journal of the American College of Medical Genetics* (2015).
- 267 16. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74
268 (2015).
- 270 17. Titus, S. *et al.* Impairment of BRCA1-related DNA double-strand break repair leads to
271 ovarian aging in mice and humans. *Science translational medicine* **5**, 172ra21 (2013).
- 272 18. White, R.R. & Vijg, J. Do DNA Double-Strand Breaks Drive Aging? *Molecular Cell* **63**,
273 729-738 (2016).
- 274 19. Oktay, K., Turan, V., Titus, S., Stobezki, R. & Liu, L. BRCA Mutations, DNA Repair
275 Deficiency, and Ovarian Aging. *Biology of reproduction* **93**, 67 (2015).
- 276 20. Kong, A. *et al.* Fine-scale recombination rate differences between sexes, populations
277 and individuals. *Nature* **467**, 1099-103 (2010).
- 278 21. Halldorsson, B.V. *et al.* The rate of meiotic gene conversion varies by sex and age.
279 *Nature Genetics* (2016).

- 280 22. Martin, H.C. *et al.* Multicohort analysis of the maternal age effect on recombination.
281 *Nature Communications* **6**, 7846 (2015).
- 282 23. Campbell, C.L. *et al.* Escape from crossover interference increases with maternal age.
283 *Nature Communications* **6**, 6260 (2015).
- 284 24. Arbeithuber, B., Betancourt, A.J., Ebner, T. & Tiemann-Boege, I. Crossovers are
285 associated with mutation and biased gene conversion at recombination hotspots.
286 *Proceedings of the National Academy of Sciences* **112**, 2109-2114 (2015).
- 287 25. Lercher, M.J. & Hurst, L.D. Human SNP variability and mutation rate are higher in
288 regions of high recombination. *Trends Genet* **18**, 337-40 (2002).
- 289 26. Webster, M.T. & Hurst, L.D. Direct and indirect consequences of meiotic
290 recombination: implications for genome evolution. *Trends Genet* **28**, 101-9 (2012).
- 291 27. Alexandrov, L.B. *et al.* Signatures of mutational processes in human cancer. *Nature*
292 **500**, 415-21 (2013).
- 293 28. Záborszky, J. *et al.* Loss of BRCA1 or BRCA2 markedly increases the rate of base
294 substitution mutagenesis and has distinct effects on genomic deletions. *Oncogene*
295 **36**, 746-755 (2017).
- 296 29. Moynahan, M.E., Chiu, J.W., Koller, B.H. & Jasin, M. Brca1 controls homology-
297 directed DNA repair. *Molecular cell* **4**, 511-8 (1999).
- 298 30. Patel, K.J. *et al.* Involvement of Brca2 in DNA repair. *Molecular cell* **1**, 347-57 (1998).
- 299 31. Baudat, F. *et al.* PRDM9 Is a Major Determinant of Meiotic Recombination Hotspots
300 in Humans and Mice. *Science* **327**, 836-840 (2010).
- 301 32. Kong, A. *et al.* Recombination rate and reproductive success in humans. *Nature*
302 *Genetics* **36**, 1203-1206 (2004).
- 303 33. Ottolini, C.S. *et al.* Genome-wide maps of recombination and chromosome
304 segregation in human oocytes and embryos show selection for maternal
305 recombination rates. *Nature Genetics* **47**, 727 (2015).
- 306 34. Middlebrooks, C.D. *et al.* Evidence for dysregulation of genome-wide recombination
307 in oocytes with nondisjoined chromosomes 21. *Human Molecular Genetics* **23**, 408-
308 417 (2014).
- 309 35. Jonsson, H. *et al.* Parental influence on human germline de novo mutations in 1,548
310 trios from Iceland. *Nature* **549**, 519-522 (2017).

311
312

313 **Figure 1:** Differences between maternal and paternal cDNMs (a) The fraction of probands
314 with maternal and paternal clustered mutations (y-axis), grouped by parental age quantiles.
315 Error bars indicate the binomial 95% confidence intervals. Labels on the lower axis indicate
316 age ranges of the respective groups and group sizes. See **Supplementary Figure 1** for graphs
317 and regression lines. (b) The number of paternal and maternal cDNMs (y-axis) stratified by
318 the distance to the nearest other cDNM (x-axis). (c) The size of paternal and maternal age
319 effect of clusters with at least one phased cDNM (y-axis) by inter-mutational distance (x-
320 axis). Whiskers indicate the 95% confidence interval. (d) Age of fathers at conception and (e)
321 age of the mothers at conception (y-axis) by the number of mutations in the offspring's
322 largest mutation cluster (x-axis). We considered only clusters where at least one cDNM is on
323 the allele from the respective parent (paternal allele for **d** and maternal allele for **e**).
324 Numbers indicate the size of each group. Boxplot compartments: box: interquartile range;
325 line: median; whiskers: extreme values $<1.5 \times$ interquartile ranges from box borders).

326 **Figure 2:** Patterns of cDNMs across the chromosomes. (a) The proportion of phased cDNMs
327 per chromosome. Error bars indicate the binomial 95% confidence intervals. (b) Overview of
328 chromosome 16 region enriched for maternal cluster mutations. X-axis and ideogram
329 indicate chromosomal position. The red and blue histograms indicate the number of
330 maternal cDNMs and paternal cDNMs identified in this study, respectively. The pale red and
331 pale blue histograms indicate the number of maternal and paternal unclustered DNMs
332 (ucDNMs). The lowest track indicates normalized cSNP C>G score, which is predictive for
333 maternal DNMs. (c) The nucleotide substitution spectrum of maternal and paternal clusters
334 and unclustered DNMs. The star indicates a significant difference assessed by bootstrapping
335 (**Online Methods**). Error bars indicate the binomial 95% confidence intervals. (d) The
336 nucleotide substitution spectrum of cDNMs by location. Error bars indicate the binomial 95%
337 confidence intervals. (e) The nucleotide substitution spectrum of polymorphism-derived
338 clustered mutation by location. The star indicates a significant difference assessed by
339 bootstrapping. Error bars indicate the binomial 95% confidence intervals. (f) Region with
340 increased maternal mutation rate on chromosome 2 (region displayed chr2:1-
341 100,000,000bp; region with maternal cDNMs chr2:40,000,000-60,000,000).

342 **Figure 3:** cDNMs and sites likely affected by DSBs. (a) Z-scores of expected and observed
343 overlaps of cDNM clusters in our cohort and sex-matched meiotic gene conversion in
344 another cohort²¹. Diamonds: observed values, boxplot compartments: box: interquartile
345 range; line: median; whiskers: extreme values $<1.5 \times$ interquartile ranges from box borders.
346 (b) DNMs detected close to sites of *de novo* CNVs. Data of DNMs is listed in **Supplementary**
347 **Table 15**. (c) cSNP density close to CNV breakpoints (**Online Methods**). (d) Z-scores of
348 expected and observed overlap cDNM clusters and sex-matched recombination hotspots.
349 Symbols and boxplots as in (a).

350

351

352 **Main Tables**

353 **Table 1: Overview of cohorts**

Cohort		Total number	Paternal number	Maternal number
Primary cohort	Probands	1,291		
	DNMs	73,755	20,196	5,547
	cDNMs	1,796	323	377
	Clusters	799	110 (+88)	94 (+108)
Replication cohort	Probands	1,557		
	DNMs	74,395	9,466	2,796
	cDNMs	1,643	133	195
	Clusters	745	40 (+49)	67 (+46)

354 Numbers of probands, DNMs, cDNMs and clusters of the cohorts used in this study. The
 355 numbers in brackets indicate clusters where not all cDNMs could be phased for the
 356 respective parent.

357

358 Online Methods

359 Cohort

360 The cohort used in this study is from Inova Translational Medicine Institute's
361 Longitudinal Childhood Genome Study (previously referred to as the First 1,000 Days of
362 Life and Beyond study), which represents a population cohort in good general health^{4,15}.
363 The study was conducted by the Inova Translational Medicine Institute and approved by
364 both the Inova and Western Institutional Review Boards (study 20120204). Parents and
365 the newborns were recruited at Inova Fairfax Hospital between 2012 and 2014. A
366 summary of participants' ages is given in **Supplementary Table 1**.

367 Whole genome sequencing

368 Sample preparation, processing and whole-genome sequencing (WGS) have been
369 previously described^{4,15}. Briefly, DNA was extracted from peripheral blood obtained
370 from each family member. Whole genome sequencing using paired-end 100bp reads
371 (median fragment length is 375) at an average 40X coverage was performed by Illumina
372 Services (San Diego, CA). The sequenced reads were aligned to the hg19 human
373 reference genome by the ISAAC aligner³⁶ with the Illumina Whole Human Genome
374 Sequencing Service Informatics Pipeline version 2.01 - 2.03.

375 To systematically analyze the data quality of all sequencing reactions, a principal
376 component analysis on scaled summary statistics was performed (**Supplementary**
377 **Figure 12, Supplementary Table 18**). The first principal component is highly
378 correlated to average sequencing coverage; a group of outlying points refers to a group
379 of sequencing reactions with average genome coverage above 70x. The second principal
380 component is associated with the date of sequencing and the version of the software
381 used for analysis, respectively. The third principal component is related to the estimated
382 ancestries of the sequenced individuals.

383 DNM calling and quality control

384 Callable regions of each sample were determined by CallableLoci in GATK version 3.1.
385 The number of callable bases by batch is shown in **Supplementary Figure 13**. The
386 batch number does not significantly influence the number of DNMs called
387 (**Supplementary Table 19**). Joint calling using HaplotypeCaller, PhaseByTransmission
388 and ReadBackPhasing in GATK version 3.1 were performed on each of the 1,315 trios in
389 the canonical autosomes³⁷. The putative *de novo* mutations were generated from taking
390 PASS filter calls with heterozygous in the proband and homozygous reference in both
391 parents in the PhaseByTransmission results in each trio. We have previously analyzed
392 816 trios⁴, of which, 65 trios were also sequenced by the Illumina services with pipeline
393 version 2.0.0-2.0.1, and are not part of this cohort. These 65 trios sequenced by Illumina
394 have gone through the same pipeline to generate a set of putative DNMs. We defined the
395 positive set as those putative DNMs that overlap with previous identified DNMs
396 identified using Complete Genomics (CG) technology (2,670), as well as those that were
397 validated by Sanger sequencing (34), the total number in the true positive set is 2,704.
398 The negative set consists of 50 random putative DNMs in each of the 65 trios that are not
399 in the previously identified set by CG (50*65=3,250), as well as 4 false positive sites
400 identified by Sanger, the total number of negative sites is 3,254. We note that some of
401 the sites in the negative set are true positives but the number is likely to be low. The test

402 set which consists of the positive and negative sets was split by 90:10 ratio into training
403 set and test set. The R libraries randomForest version 4.6.10 and caret version 6.0.52
404 were used to train the random forest classifier. The OOB estimate of error rate on
405 training set is 1.77% and the error rate in the test set is 2.18%. The features used in the
406 classifier and their relative importances are shown in **Supplementary Table 20**. The
407 confusion matrix for the test set is shown in **Supplementary Table 21**.

408 In order to minimize the bias due to mapping errors and coverage differences, we
409 further filtered the predicted DNMs by (1) callable regions in the cohort: A site is in the
410 callable region if at least 90% of the samples has the PASS status by GATK CallableLoci³⁷,
411 (2) good mappability regions, where mappable is defined according to the CRG 100mer
412 (file wgEncodeCrgMapabilityAlign100mer.bw from UCSC Table Browser) being equal to
413 1³⁸, sites also called by the Illumina Isaac Small Variant Caller, and sites with FS
414 (FisherStrand test score) ≥ 20 , and sites with exceptionally high or low PL values
415 (**Supplementary Table 22**). An overview of the filtering procedure is given in
416 **Supplementary Table 2**.

417 In the initial sequencing cohort, there were 12 monozygotic twin pairs, 29 dizygotic twin
418 pairs and a family of three trizygotic siblings. In order to assess the consistency in *de*
419 *novo* calling, we investigated the concordance percentages of monozygotic and dizygotic
420 families (**Supplementary Table 23** and **Supplementary Table 24**). DNM calls in
421 monozygotic twins are on average 95% concordant, the dizygotic average concordance
422 is 0.1%. This is similar to concordance ratios observed previously⁴.

423 We removed 1 trio with an exceptional high number of DNM calls, 8 trios with a large
424 chromosomal anomaly in either the proband or one of the parents and removed
425 (arbitrarily) one of the monozygotic twins in each set. After performing simple multiple
426 linear regression, 3 samples have a significant Bonferroni p-value for studentized
427 residuals (Bonferroni corrected $p < 0.05$) and are removed from the cohort, resulting in
428 1,291 trios (**Supplementary Table 2**). We investigated the effect of average genome
429 coverage on the filtered data. The results are shown in **Supplementary Figure 14**.

430 The method for determining the parent-of-origin of DNMs with Illumina WGS trio data
431 was previously described^{3,4}. Briefly, GATK PhaseByTransmission was used to assign
432 parent-of-origin to informative heterozygous SNPs in the proband, GATK
433 ReadBackPhasing was used to link DNMs to these informative SNPs. If contradictory
434 markers were linked to the same DNM, it would not be assigned a parent-of-origin.
435 Overall, 227 of the 25,970 filtered DNMs are linked to contradictory markers (0.87%).

436 **Clustered DNMs**

437 We defined cDNMs as DNMs on the same chromosome of the same individual within
438 20kb of each other. In order to estimate the chance of two DNMs being closer than 20kb
439 on the same chromosome, we simulated 70,000 mutations at random positions within
440 the callable and mappable genome. The randomized positions were given sample IDs as
441 in the set of observed DNMs and the distances were calculated. We found that the false
442 discovery rate of cluster detection is 0.0375 at a threshold of 20kb (**Supplementary**
443 **Table 9**). Statistics on the number of cDNMs per cluster are given in **Supplementary**
444 **Table 3**.

445 For analyses on clusters we extrapolated the parent-of-origin by considering all cDNM to
446 originate from the same allele.

447 Sanger validation

448 We performed Sanger validation on 163 clustered DNM sites on the proband and his or
449 her parents, of which 62 are on chromosomes 8, 9 and 16 (**Supplementary Table 25A**).
450 Overall, 91.3% of the DNMs are validated, 92.7% on chromosomes 8,9 and 16 vs. 90.4%
451 on other chromosomes. The number of sites validated in each pipeline version is
452 proportional to the number of trios sequenced in each pipeline (**Supplementary Table**
453 **25B**). There is no significant difference in the proportion of sites validated in each
454 pipeline ($P = 0.92$, Fisher's exact test). No evidence of the mutations was found at any
455 site in the parents. All of the invalidated sites were due to lack of evidence in the
456 proband.

457 Clustered polymorphism variants

458 We use polymorphism data from the 1000 Genomes Project Consortium¹⁶. We only
459 considered non-singleton variants with below 1% derived allele frequency, using the
460 ancestral variant determined by The 1000 Genomes Project Consortium. Clusters were
461 defined as two or more SNPs at distances between 10-1000 nucleotides from each other,
462 such that all the genotypes carrying the derived allele for one of the SNPs also carry the
463 derived allele for any other SNP within the cluster. We show that cSNP spectra are
464 similar to cDNM spectra: enriched by C>G mutations and depleted by CpG>TpG
465 mutations, compared to unclustered DNMs. We restricted ourselves to distances
466 between cSNPs shorter than 1000 nucleotides, because of two reasons. First, the
467 probability of recombination scales with the distance between SNP positions and thus
468 longer clusters are more frequently disrupted. Second, the probability to observe two
469 independent mutations on the same haplotype would be ~20 fold higher in 1-20 kb
470 range than in 0-1 kb range. In contrast we observe 806 cDNMs in 0-1 kb range and 990
471 cDNMs in 1-20 kb range. Therefore, we expect a higher noise to signal ratio for larger
472 distances. In line with this, the spectra of larger clusters are progressively less similar to
473 cDNMs (**Supplementary Figure 15**). For analyzing the density of cSNPs around CNV
474 breakpoints, we calculated the distances between cSNPs and CNVs on the chromosomes
475 of each individual. We only considered cSNPs flanking CNVs, but not within its body.
476 These distances were compared to the distances between cSNPs and the CNVs on the
477 same chromosome of a random other individual.

478 Statistical assessment of the maternal age effect

479 For analyzing the parental age effects on both the number of clusters as well as the
480 number of cDNMs, linear models were fitted using the R statistical environment version
481 3.3.3 with standard settings. The reported p-values reflect the difference from zero of
482 the respective age effect.

483 Extrapolations of DNM phasing were done by assigning a cluster's unphased DNMs the
484 same allele as the phased ones. In order to correct for the false detection rate of 3.75%
485 (**Supplementary Table 9**), we sampled 1,000 subsets of $100\% - 3.75\% = 96.25\%$ of
486 cDNMs and calculated the age effects on all of them. We report the median effect size
487 and the median p-value.

488 For comparing proband groups' risks for having DNM clusters we used risk ratio
489 statistics as implemented in the R package "epitools". For assessing the enrichment of
490 C>G substitutions on chromosomes 8, 9 and 16, we re-sampled the chromosome
491 annotation 1,000 times and compared the difference of the fractions of C>G mutations
492 on the special chromosomes and the remaining autosomes to the observed value.

493 **Statistical assessment of nucleotide substitution profiles**

494 The significance of differences between nucleotide substitution profiles was assessed by
495 bootstrapping: We resampled the grouping variable 1,000 times and compared the
496 resulting random groups to the observed groups. For assessing C>G enrichment we
497 calculated p-values by counting the number of random groups where the difference in
498 C>G fractions between the groups is equal to or larger than in the observed set and
499 dividing by the number of samplings.

500 **Statistical assessment of DSB proxy regions overlap**

501 For calculating distributions on the expected number of overlaps between DNM clusters
502 and DSB proxy regions we used permutation testing as implemented in the R library
503 RegioneR³⁹. DNM cluster regions were defined as the positions of cDNMs and the space
504 between them. Recombination hotspots were defined as genomic sites with a
505 recombination-score above 10^{20} . Meiotic gene conversions were filtered for non-
506 crossover gene conversions only detected in the chip dataset²¹. In absence of knowledge
507 about the exact boundaries of the conversion streak and confronted with the majority of
508 meiotic gene conversions being observed only in one SNP, we defined the positions of
509 meiotic gene conversions as the distance between the two SNPs adjacent to the SNPs
510 affected by conversion. The cluster regions were randomized 500 times to genomic
511 positions where at least 1000 base pairs were within the callable and mergable subset of
512 the genome. For every randomization round the number of cluster positions overlapping
513 DSB proxy regions was compared to the observed number of overlaps. For the
514 calculation of z-scores of an overlap count the mean number of overlaps was subtracted
515 before division by the standard deviation of the number of overlaps.

516 **De novo CNVs**

517 In the primary cohort, we called *de novo* CNVs using both coverage-based method
518 FREEC⁴⁰ and read-pair based method Manta⁴¹. We also calculated window based
519 normalized coverage with "goleft indexcov". For each proband, we called CNVs using the
520 default options in FREEC with the proband as the case and one of the parents as control.
521 We then required the CNVs subtracted from each parent to have 90% reciprocal overlap,
522 with copy number equals 1 or 3, both parents have the mean normalized coverage
523 between 0.85 and 1.15 in the region, the proband have mean normalized coverage
524 smaller than 0.85 or greater than 1.15 in the region, with length greater or equal to
525 10kb. We performed joint calling for each trio with Manta using default options. We then
526 filter for SV type being DEL or DUP, proband with GT equals to 0/1 and both parents
527 with GT equal to 0/0, proband's PR and SR for ALT allele ≥ 3 and the proportion of PR
528 and SR for ALT ≥ 0.2 , parents' proportion of PR and SR for ALT ≤ 0.05 .

529 In the complete genomics data in the replication cohort, we required the *de novo* CNV to
530 be called by both coverage-based and read-based methods. For the coverage-based
531 method, we first subtracted CNVs in the proband from one of the parents using the

532 cnvSegmentsDiploidBeta files, and then we intersect the two putative *de novo* CNV files
533 subtracted from each parent, with 90% overlap, and size >9999. For the read-based
534 method, we subtracted highConfidenceSvEventsBeta file from the proband from
535 allSvEventsBeta file from each of the parents, and intersected the two subtracted files
536 requiring 90% overlap. The final list of *de novo* CNVs is generated by intersecting the
537 coverage-based and read-based files from the same proband, requiring 90% overlap.
538 Bedtools 2.22.0 was used to carry out region subtractions and intersections⁴².

539 **Mutation signatures**

540 A large set of mutational signatures is known from cancer studies²⁷, some of which are
541 well annotated with mutational influences. To fit the patterns of our DNMs to these
542 signatures we used an algorithm similar to the one described in⁴³: a non-negative least-
543 squares algorithm finds the mixture of known signatures that describes best the
544 observed pattern. In order to get an indication of the robustness of the fitted mixture of
545 signatures, a bootstrapping analysis was done. The mutations of a group were
546 resampled 1,000 times with replacement and the standard deviation as well as the 95%
547 confidence intervals of each fitted signature were calculated.

548 **Single variant association study of parents genotype in *BRCA1*, *BRCA2* and 549 *PRDM9* with number of phased cDNMs in the proband**

550 The small variants in the autosomes were merged using “agg” with Illumina genome VCF
551 files using default parameters. No sample had a call rate <90%. In this analysis, only
552 those variants in *BRCA1*, *BRCA2*, *PRDM9*, and marker rs2914276, with call rate >90%, no
553 significant deviation from Hardy-Weinberg equilibrium ($P>0.001$), and with minor allele
554 frequency >0.005 were included. No LD pruning was performed. If a parent has more
555 than 1 offspring in the cohort (twins or siblings), only one of the sibling’s number of
556 phased cDNMs is kept as the phenotype for the respective parent. The association
557 analysis was performed with Plink v.1.90b⁴⁴ with additive model on paternal genotypes
558 with paternal number of cDNMs, using paternal age at conception, and father’s first 3
559 PCs as covariates; and on maternal genotypes with maternal number of cDNMs, using
560 maternal age at conception, and mother’s first 3 PCs as covariates, respectively. The
561 association study included 1,247 fathers and 1,247 mothers. No variant reached
562 significance ($P<0.05$) after Bonferroni correction.

563

564 **Data availability**

565 *De novo* mutation calls used in this manuscript will be available in dbGaP, by the
566 accession code phs001522.v1.p1.

567 **Code availability**

568 Code available upon request.

569

570

571

572 **Methods-only references**

573

- 574 36. Raczky, C. *et al.* Isaac: ultra-fast whole-genome secondary analysis on Illumina
575 sequencing platforms. *Bioinformatics* **29**, 2041-2043 (2013).
- 576 37. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for
577 analyzing next-generation DNA sequencing data. *Genome research* **20**, 1297-303
578 (2010).
- 579 38. Derrien, T. *et al.* Fast Computation and Applications of Genome Mappability. *PLoS*
580 *ONE* **7**, e30377 (2012).
- 581 39. Gel, B. *et al.* regioneR: an R/Bioconductor package for the association analysis of
582 genomic regions based on permutation tests. *Bioinformatics* **11**, btv562 (2015).
- 583 40. Boeva, V. *et al.* Control-FREEC: a tool for assessing copy number and allelic content
584 using next-generation sequencing data. *Bioinformatics (Oxford, England)* **28**, 423-5
585 (2012).
- 586 41. Chen, X. *et al.* Manta: rapid detection of structural variants and indels for germline
587 and cancer sequencing applications. *Bioinformatics (Oxford, England)* **32**, 1220-2
588 (2016).
- 589 42. Quinlan, A.R. & Hall, I.M. BEDTools: a flexible suite of utilities for comparing genomic
590 features. *Bioinformatics (Oxford, England)* **26**, 841-2 (2010).
- 591 43. Blokzijl, F., Janssen, R., Van Boxtel, R. & Cuppen, E. MutationalPatterns: an
592 integrative R package for studying patterns in base substitution catalogues. *bioRxiv*
593 (2016).
- 594 44. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based
595 linkage analyses. *American journal of human genetics* **81**, 559-75 (2007).
- 596 45. Glusman, G., Caballero, J., Mauldin, D.E., Hood, L. & Roach, J.C. Kaviar: an accessible
597 system for testing SNV novelty. *Bioinformatics (Oxford, England)* **27**, 3216-7 (2011).

598

599





