# A combined RNA-seq and whole genome sequencing approach for identification of non-coding pathogenic variants in single families.

Revital Bronstein,[1] Elizabeth E. Capowski,[2] Sudeep Mehrotra,[1] Alex D. Jansen,[2] Daniel Navarro-Gomez,[1] Mathew Maher,[1] Emily Place,[1] Riccardo Sangermano,[1] Kinga M.Bujakowska,[1] David M. Gamm,[3] Eric A. Pierce,[1]*.

1. Ocular Genomics Institute, Massachusetts Eye and Ear Infirmary, Harvard Medical School, Boston, MA, 02114, USA.

2. University of Wisconsin-Madison, Waisman Center Stem Cell Research Program, Madison, WI, 53705, USA.

3. Department of Ophthalmology and Visual Sciences, University of Wisconsin-Madison, McPherson Eye Research Institute, Waisman Center Stem Cell Research Program, Madison, WI, 53705, USA.

*Correspondence: eric_pierce@meei.harvard.edu

# Abstract

Inherited retinal degenerations (IRDs) are at the focus of current genetic

therapeutic advancements. For a genetic treatment such as gene therapy to be

successful an accurate genetic diagnostic is required. Genetic diagnostics relies

on the assessment of the probability that a given DNA variant is pathogenic.

Non-coding variants present a unique challenge for such assessments as

compared to coding variants. For one, non-coding variants are present at much

higher number in the genome than coding variants. In addition, our

understanding of the rules that govern the non-coding regions of the genome is

less complete than our understanding of the coding regions. Methods that allow

for both the identification of candidate non-coding pathogenic variants and their

functional validation may help overcome these caveats allowing for a greater

number of patients to benefit from advancements in genetic therapeutics. We

present here an unbiased approach combining whole genome sequencing

(WGS) with patient induced pluripotent stem cell (iPSC) derived retinal organoids

(ROs) transcriptome analysis. With this approach we identified and functionally

validated a novel pathogenic non-coding variant in a small family with a

previously unresolved genetic diagnosis.

## Introduction

38

39    Inherited retinal degenerations (IRDs) are a leading cause of blindness,

40    altogether affecting >2million people worldwide. IRDs are characterized by

41    progressive degeneration of photoreceptor and/or retinal pigment epithelial

42    (RPE) cells of the retina with variable age of onset and rates of degeneration[1].

43    Despite tremendous ongoing efforts and research into various therapeutics,

44    treatments for IRDs remain limited. Currently there are two regulatory agency-

45    approved treatment approaches, retinal prosthesis implants and  gene

46    augmentation therapy for IRD caused by mutations in the *RPE65* gene [1–7]. The

47    eye is a prime candidate for gene therapy approaches due to its relative immune-

48    privilege, surgical accessibility and ease of non-invasive monitoring. In addition,

49    IRDs are Mendelian disorders caused by mutations in single genes in the vast

50    majority of cases. Owing to these favorable circumstances several gene/genetic

51    therapy clinical trials have been initiated for IRDs, including those caused by

52    mutations in the *ABCA4, CEP290, CHM, CNGA3, CNGB3*, *MYO7A, RPGR,*

53    *RS1,* and *USH2A* [6].

54    As each genetic therapy targets a specific gene, for a patient to be considered for

55    treatment they must obtain a reliable genetic diagnosis. Difficulties inherent to

56    genetic diagnostics are rooted in the fact that every individual carries millions of

57    DNA variants in their genome [8,9]. The large majority of the DNA variants are

58    found in non-coding regions of the genome such as intergenic and intronic

59    regions. Since non-coding sequences can better tolerate sequence variation

3

60  compared to coding sequences, most of these variants are benign and do not

61  lead to disease. Still, some non-coding variants are found to be pathogenic by

62  altering gene expression and/or splicing patterns [10–12]. As a result, non-coding

63  variants are among the hardest to classify and thus under-diagnosed [13,14].

64  Algorithms exist that predict the effect of a non-coding variant on gene

65  expression or splicing based on analysis of the DNA sequence alone, but their

66  accuracy for diagnostic purposes remains undetermined [15–18].

67  In order to functionally test the effects of non-coding variants one needs to

68  quantify the level of gene expression and analyze the splicing patterns of the

69  presumably affected genes. When multiple variants in multiple genes need to be

70  evaluated, advanced methods for whole transcriptome analysis are

71  advantageous. Indeed, previous studies successfully utilized large RNA-seq

72  datasets from tissue biopsies to identify novel non-coding pathogenic variants

73  [19,20]. For example, Evrony *et al* narrowed down a linkage analysis in a very large

74  pedigree to a single non-coding variant using RNA-seq. This non-coding variant

75  was shown to cause intron retention in the *DONSON* gene and is most likely the

76  genetic cause of microcephaly-micromelia syndrome (MMS) in this population [19].

77  Cummings *et al* used 184 skeletal muscle RNA-seq samples available through

78  Genotype-Tissue Expression resource (GTEx) [21] as a reference panel for 50

79  patients with undiagnosed muscle disorders. This comparison led to a genetic

80  diagnosis for 17 previously unsolved families and identification of several splice

81  altering variants [20].

82    Such studies depend on the availability of large, publicly available RNA-seq

83    datasets and/or on a large cohort of patients. They also require biopsy samples

84    from a clinically relevant tissue or cell type. Both gene expression and splicing

85    are tissue-specific [22,23] owing to restricted availability of transcription and splicing

86    factors with variable usage of regulatory DNA sequences [24–27]. Consequently,

87    DNA variants in such regulatory sequences can have tissue-specific outcomes

88    on gene expression and splicing [18,28]. Thus, analyzing RNA from a clinically

89    relevant tissue or cell type is crucial to obtain a more focused and reliable

90    diagnostic result. When a clinically relevant tissue is not accessible, *ex vivo*

91    surrogate models can sometimes suffice. Indeed, a study aimed at discerning the

92    genetic causality of patients with monogenetic neuromuscular disorders found

93    that t-myotubes, skeletal myotubes derived by myoD overexpression in

94    fibroblasts, accurately reflected the muscle transcriptome and faithfully revealed

95    pathogenic variants [29].

96    In this study, we aimed to develop a pipeline that would detect putative non-

97    coding pathogenic mutations in a small family. We present a pilot study

98    performed in a five member family with two siblings affected by cone dysfunction

99    syndrome in which we successfully identify and functionally validate a novel deep

100   intronic variant without the use of large reference datasets. Obtaining a clinically

101   relevant tissue from IRD patients is not possible as the retina cannot be safely

102   biopsied. To overcome this limitation, we have made use of patient-derived

103   induced pluripotent stem cells (iPSCs) that were differentiated *in vitro* to form

104   retinal organoids (ROs)[30]. ROs have been shown to recapitulate many aspects of

105    human retinal structure and function [30,31,40,32–39]. We show that the RO

106    transcriptome is much closer to the transcriptome of normal human retina than

107    other more readily available diagnostic tissues. More importantly, we show for

108    the first time that analysis of a patient-derived RO transcriptome can successfully

109    detect pathogenic deep intronic variants that activate cryptic splice sites, leading

110    to a new genetic diagnosis. Our approach can lead to a larger number of patients

111    to be eligible for genetic therapies.

## Materials and Methods

112

### Human Subjects

113

114    The study was approved by the institutional review board at the Massachusetts

115    Eye and Ear (Human Studies Committee MEE in USA) and adhered to the

116    Declaration of Helsinki. Informed consent was obtained from all individuals on

117    whom genetic testing and further molecular evaluations were performed.

### Pluripotent Stem Cell Induction

118

119    Tissue samples were obtained with written informed consent in adherence with

120    the Declaration of Helsinki and with approval from institutional review boards at

121    the University of Wisconsin-Madison and Massachusetts Eye and Ear Infirmary.

122    Blood samples from 4 individuals from family OGI-081 (197, 198, 200 and 340)

123    were collected and reprogrammed by Cellular Dynamics, Inc (now FUJIFILM

124    Cellular Dynamics, Inc) as custom MyCell products. Three independent clones

125    from each individual were karyotypically normal, expressed pluripotency markers

126    and successfully differentiated to retinal organoids ([41]: lines 1579, 1580, 1581

127  and 1582). Stem cells were maintained on Matrigel (ThermoFisher) in either

128  mTeSR1 (WiCell) or Stemflex (ThermoFisher) and passaged with either Versene

129  or ReLeSR (STEMCELL Technologies).

130  <u>Retinal Organoid (RO) Differentiation</u>

131  Differentiation of iPSCs was performed as previously described [41]. Briefly,

132  embryoid bodies (EB) were lifted with either 2 mg/ml dispase or ReLeSR and

133  weaned into Neural Induction Media (NIM: DMEM:F12 1:1, 1% N2 supplement,

134  1x MEM nonessential amino acids (MEM NEAA), 1x GlutaMAX and 2 mg/ml

135  heparin (Sigma)) over the course of 4 days.  On day 6, 1.5 nM BMP4 (R&D

136  Systems) was added to fresh NIM and on day 7, EBs were plated on Matrigel at

137  a density of 200 EBs per well of a 6-well plate.  Half the media was replaced with

138  fresh NIM on days 9, 12 and 15 to gradually dilute the BMP4 and on day16, the

139  media was changed to Retinal Differentiation Media (RDM: DMEM:F12 3:1, 2%

140  B27 supplement, MEM NEAA, 1X antibiotic, anti-mycotic and 1x GlutaMAX). On

141  days 25-30, optic vesicle-like structures were manually dissected and maintained

142  as free floating organoids in poly HEMA (Sigma)- coated flasks with twice weekly

143  feeding of 3D-RDM (RDM + 5% FBS (WiCell), 100 $\mu$M taurine (Sigma) and

144  1:1000 chemically defined lipid supplement) to which 1 $\mu$M all-trans retinoic acid

145  (Sigma) was added until d100. Live cultures were imaged on a Nikon Ts2-FL

146  equipped with a DS-fi3 camera.

147  <u>Immunocytochemistry and Microscopy</u>

148  Organoids were fixed in 4% paraformaldehyde at room temperature for 40 min,

149  cryopreserved in 15% sucrose followed by equilibration in 30% sucrose, and

150 sectioned on a cryostat. Sections were blocked for 1 hr at room temperature (RT)

151 in 10% normal donkey serum, 5% BSA, 1% fish gelatin and 0.5% Triton then

152 incubated overnight at 4$^{o}$C with primary antibodies diluted in block. Table S1 lists

153 primary antibodies, dilutions and sources. Slides were incubated with species-

154 specific fluorophore-conjugated secondary antibodies diluted 1:500 in block, for

155 30 minutes in the dark at RT (Alexa Fluor 488, AF546 and AF647). Sections

156 were imaged on a Nikon A1R-HD laser scanning confocal microscope (Nikon

157 Corporation, Tokyo, Japan).

158 <u>DNA Sequencing</u>

159 DNA was extracted from venous blood using the DNeasy Blood and Tissue Kit

160 (Qiagen, Hilden, Germany). OGI-081-197 underwent GEDi sequencing as

161 described previously [42]. All five family members underwent whole exome and

162 PCR-free whole genome sequencing. Sequencing was done at the Genomics

163 Core at Massachussets Eye and Ear as described previously [43].

164 <u>RNA Sequencing</u>

165 For transcriptome analysis, ROs from at least 2 different clones per individual

166 were harvested at approximately day160 (early stage 3), lysed in 350 µl buffer

167 RLT+ME from the RNeasy mini kit (Qiagen), snap frozen on dry ice and stored at

168 -80C. At a later time samples were defrosted on ice and passed through

169 QIAshredder columns (Qiagen). Subsequently, Total RNA was extracted per the

170 manufacturer's instructions. RNA quality and quantity was assessed on an

171 Agilent 2100 Bioanalyzer, RIN number ranged between 9.6-9.9. For each sample

172 1µg of total RNA spiked with 1.2ng Sequins (v2) controls [44] was used to generate

173   RNA-seq paired-end libraries with the Illumina TruSeq Stranded Total RNA kit.

174   Ribosomal RNA was removed with the Ribo-Zero Human/Mouse/Rat kit.

175   Libraries were multiplexed and sequenced on an Illumina HiSeq 2500 instrument

176   for 101 cycles.

177   <u>Bioinformatics</u>

178   Whole exome sequence data was analyzed in house [43] and whole genome data

179   was analyzed in collaboration with the Broad Institute of MIT and Harvard using

180   methodology described previously [20]. Briefly, BWA was used for alignment.

181   GATK was used for single nucleotide polymorphism and insertion/deletion calls.

182   Additional variant annotation was performed using the Variant Effect Predictor

183   (VEP) [45]. Variants of interest were limited to polymorphisms with less than 0.005

184   allelic frequency in the gnomAD and ExAC databases [8,9]. Whole genome copy

185   number analysis, with consideration of structural changes, was done using

186   Genome STRiP 2.0 [46].

187   For analyses of RNA-seq data, read quality was assessed with FastQC v0.11.3

188   (Babraham  Bioinformatics, Cambridge, UK)  and MultiQC v1.2 [47]. Reads were

189   aligned to the human genome version GRCh37 by the STAR v2.5.3a [48] aligner in

190   two-pass mode within the sample and across replicates for each sample sets.

191   Annotations were derived from the Human GENCODE v19 (Ensemble74).

192   FeatureCount v1.5.2 [49] from the Subread package, was used to generate gene

193   expression matrix  with the following non-default settings, reads must be paired,

194   both the pairs must be mapped, use only uniquely  mapped  reads, multi-mapped

195   reads are not counted, chimeric reads are not counted and strand specificity

196    turned on. Anaquin [50] was further used to evaluate alignment sensitivity and gene

197    expression. Here sensitivity indicates the fraction of annotated regions covered

198    by alignments of the reads by STAR (Table S2). No limit of quantification or limit

199    of detection was reported.

200    For discovery of novel or known alternative splicing events we used a

201    combination of CASH v2.2.1 [51] and MAJIQ v1.1.3a [52]. CASH was operated with

202    default settings. MAJIQ was run with a minimum of 5 reads for junction detection

203    and 10 reads for the calculation of delta percent spliced-in (dPSI). EdgeR

204    (v3.2.2) [53] was used to perform differential gene expression with default settings.

205    Data normalization was performed using trimmed mean of M-values (TMM).

206    Next, the raw read counts were converted to transcript per million (TPM)

207    expression values. The Picard tools v1.87 and RSeQC v2.6.4 [54] were used to

208    calculate  mean fragment length. The approach implemented in Kallisto [55] was

209    used to covert raw reads to TPM values. An average TPM of the third lowest

210    Sequins between test and control samples was calculated and used as cutoff.

211    TPM values for the GTEx samples used in figure3 were downloaded from the

212    GTEx portal. The human normal retina (HNR) samples [56] and the ENCODE skin

213    samples were reanalyzed as described above. The ENCODE skin samples were

214    used for the analysis performed with MAJIQ as they were generated with a

215    stranded total RNA-seq library same as our RO samples. In contrast the GTEx

216    samples were generated with an mRNA non-stranded library.

217

218    <u>RT-PCR and cloning</u>

219    RT-PCR was conducted using SuperScript IV first-Strand synthesis system

220    (ThermoFisher Scientific, Waltham, MA). Exon 14b was amplified with primers F:

221    GACATGTTGCTAAGATTGAAATCCGT from exon 14 and R:

222    GACCCAGCTTTCAGAGTAACCAGAAC from exon 15 using Phusion

223    polymerase (NEB, Ipswich, MA). The longer band containing exon14b was then

224    excised from the gel and purified using Zymoclean gel DNA recovery kit

225    (Zymoresearch, Irvine, CA) and cloned using pGEM-T Easy Vector System

226    (Promega, Madison, WI). The plasmid was used in a transformation into

227    Subcloning Efficiency DH5α Competent Cells (Invitrogen, Carlsbad, CA). The

228    plasmid was isolated with Zyppy Plasmid Miniprep Kit (Zymoresearch, Irvine,

229    CA). All procedures described in this section were conducted according to the

230    manufacturer instructions. The DNA sequence of Exon 14b was found to be:

231    GCCAGGTGCAGTGGCTCACGACTGTAATTCCAACACTTTGGGAGGCCAAGG

232    TGGCAGGATCACATAAGTCCAGGAGTTCAAGACAAGCCTGGACAACATG.


233

# Results

## Unresolved genetic analysis of family OGI-081

The pilot study reported here was conducted on a five member family with two siblings shown by clinical testing to be affected by a cone dysfunction syndrome (Figure1a). Both affected patients had nystagmus and decreased vision from infancy and at age 8, OGI-081-197 was also noted to have photophobia. Visual acuity for both affected patients measured 20/150-200 and remained stable for 3 years. Full field electroretinogram (ffERG) testing of retinal function for OGI-081-197 was significant for reduced and delayed cone photoreceptor responses, with normal rod photoreceptor response amplitudes. Optical coherence tomography (OCT) imaging of the retina showed retinal degenerative changes in the fovea (Figure1b). In addition to their retinal disease, both affected patients were found to have Chiari malformations. Interestingly, vision phenotypes such as photophobia, vision loss and nystagmus have been reported as accompanying symptoms in some forms of Chiari malformations[57–60]. Unfortunately, despite evidence that Chiari malformations have a hereditable component [61–63], the genes involved are not yet well defined [59,64,65]. For that reason we could not rule out the possibility that the vision phenotypes and Chiari malformations share a common genetic causality.

The phenotypes segregation in OGI-081 is indicative of a recessive mode of inheritance. Thus, for genetic testing we searched for genes that have putatively pathogenic variants in both alleles (Figure1a). Selective exon capture based

12

256   genetic diagnostic testing was performed using the Genetic Eye Disease (GEDi)

257   test [42]. Since this did not identify a clear cause of disease, whole exome

258   sequencing (WES) for the five members of the family was performed. Both the

259   GEDi testing and WES identified a single rare variant in the *CNGB3* gene

260   (c.1148delC, p.Thr383IlefsTer13) which has been reported to be pathogenic [66–68],

261   but a second rare variant in *CNGB3* was not identified, nor were other potential

262   causative genetic variants forthcoming for the two affected members of the

263   family. *CNGB3* mutations are among the most common causes of cone

264   dysfunction syndrome, but to the best of our knowledge, Chiari malformations

265   have not been reported as an accompanying symptom in *CNGB3* patients [69,70].

266   Moreover, the $1.75e^{-3}$ gnomAD allele frequency of the p.Thr383IlefsTer13 variant

267   is higher than expected for recessive pathogenic variants and two homozygous

268   individuals are reported in the gnomAD database. Since it has been proposed

269   that up to 1 in 4-5 individuals in the general population may be a carrier of null

270   mutations in IRD genes [71] it was possible that the presence of variant

271   p.Thr383IlefsTer13 was an incidental finding.

272   We therefore decided to test for two possible disease scenarios. One is that

273   *CNGB3* accounts for the cone dysfunction syndrome and the second allele is a

274   non-coding variant. In this case the Chiari malformation has a separate,

275   unrelated causality. The second possibility is that a novel disease-causing gene

276   is responsible for both the cone dysfunction syndrome and the Chiari

277   malformation. To test these hypotheses, we performed both whole genome

278   sequencing (WGS) and RNA-seq analysis of a surrogate retinal tissue to

13

279    determine whether the combination of these orthogonal investigations could yield

280    a clear genetic solution.

281    Analysis of DNA variants detected by WGS identified 3268 segregating rare

282    variants that could be sorted into 8191 allelic pairs in 642 genes (TableS3). The

283    variant ranked at the top of the list of potential causes of disease remained the

284    known pathogenic variant c.1148delC; p.Thr383IlefsTer13 in the *CNGB3* gene.

285    However, a second coding variant once again was not found in this gene. Next,

286    we set to establish a clinically relevant surrogate transcriptome for the human

287    retina.

288    <u>The iPSC-derived retinal organoid (RO) transcriptome can be used as a</u>

289    <u>surrogate for a human retinal biopsy</u>

290    Patient-derived iPSCs were generated from peripheral blood monocytes of all

291    members of family OGI-081 excluding OGI-081-199. The iPSCs were subjected

292    to an *in vitro* differentiation process to generate ROs with attached RPE (Figure

293    2a) [30]. RPE was specifically retained in the ROs so as to concurrently identify

294    potential mutations in RPE genes as well as photoreceptor genes (Figure2b).

295    The ROs were kept in culture for 160 days (early stage 3 [30]), a time point at

296    which outer segments are visible by light microscopy and cone and rod cells are

297    clearly distinguished by immunocytochemistry (Figure2b, c-h). At da160, ROs

298    were harvested for total RNA isolation, library preparation, and RNA-seq analysis

299    (Table1).

14

300     In order to evaluate if day 160 ROs can be used as informative surrogates for the

301     adult human retinal transcriptome, we examined the expression levels of 270

302     known IRD genes reported in the RetNet database. For this analysis, we

303     compared an in-house dataset composed of 3 post-mortem human normal

304     retinas (HNR, N=3)[56] to ROs derived from the unaffected sibling (N=5) or skin or

305     whole blood samples taken from the GTEx database (Figure 3 and Table S4)[21].

306     Skin and blood represent tissues that are more readily available – and thus

307     commonly used – for surrogate diagnostic testing. We found 224 IRD genes to

308     be expressed in HNR (TPM>1). Interestingly we were able to detect expression

309     of 254 IRD disease genes in the RO samples (Figure 3 and TableS4). The higher

310     detection rate of disease-causing genes in ROs compared to HNR is most likely

311     because of overall higher TPM values in RO samples (Figure 3a&b), possibly

312     due to higher RNA quality compared to post-mortem HNR samples. In addition,

313     the presence of RPE and photoreceptors with varied maturation statuses in the

314     RO samples could be contributing factors to this finding. As expected, skin and

315     blood expressed lower numbers of IRD genes (188 & 130 respectively) at much

316     lower TPM values (Figure 3a&b and Table S4). Since IRD genes were very

317     poorly represented in the blood samples, we excluded blood from further

318     analysis.

319     We next examined the complexity of the HNR, RO and skin transcriptomes,

320     which is reflected by the multitude of isoforms that are produced from each gene

321     locus [72,73]. Isoform diversity that occurs via alternative splicing of the pre-mRNA

322     can be represented by the splice junctions detected in each gene locus. We used

15

323   the MAJIQ [52] algorithm to detect splice junctions in HNR and RO samples of the

324   OGI-081 unaffected sibling, and the ENCODE database [74,75] to detect splice

325   junctions in corresponding skin samples. The skin samples from the ENCODE

326   database were more suited for this analysis then the GTEx samples due to RNA-

327   seq library type (see Materials and Methods). We found a comparable number of

328   splice junctions in the HNR and RO samples (41,121 and 31,535 (76%)

329   respectively) but a much lower number in the skin samples (16,713 (40%))

330   (Figure 3c and TableS5). This result is probably due to the fact that HNR and

331   ROs are composed of a more diverse cell population than skin. More importantly,

332   in IRD genes, the gap in complexity between the HNR (946 junctions) and ROs

333   (739 (78%)) as compared to skin (227(24%)) is even greater (Figure 3d & Table

334   S5). Thus, the RO transcriptome provided a close facsimile of the human retinal

335   and IRD transcriptomes at the gene expression and splicing pattern levels,

336   whereas the skin transcriptome did not.

337   <u>Detection of a novel non-coding pathogenic variant in *CNGB3*</u>

338   In order to find the underlying genetic cause of the retinal degeneration in OGI-

339   081 affected patients, we conducted differential splicing and gene expression

340   analyses of the RNA-seq data obtained from day 160 ROs of affected versus

341   unaffected siblings. The differential splicing analysis was conducted with CASH

342   [51] and MAJIQ [52] algorithms, and we used edgeR algorithm [53] for differential gene

343   expression. CASH detected 106 differential splicing events in 101 genes (Table

344   S6), while MAJIQ detected 522 differential splicing junctions in 260 genes (Table

345   S7). A comparison to the 642 genes with DNA variant pairs indicated that only

16

346  two genes, *CNGB3* and *NCALD,* had altered splicing patterns and a DNA variant

347  in each of their alleles that segregated according to disease status in OGI-081

348  (Figure 4a).

349  For the *NCALD* gene, the alternative splicing events identified by CASH and

350  MAJIQ were isoform switching events between minor isoforms (data not shown)

351  whose biological significances are unclear. In addition, the DNA variants in this

352  gene are located 15-50kb away from the nearest alternative junction (Figure 4b

353  and Table S8). Given such large distances, it is not likely that these variants

354  could cause the alternative splicing events. Moreover, the differential gene

355  expression analysis did not find the *NCALD* gene to be differentially expressed

356  between the affected and unaffected siblings (Table S9). Taken together, these

357  results indicate that the *NCALD* gene variants are most likely not the genetic

358  cause for the disease phenotypes observed in OGI-081.

359  We next examined the *CNGB3* gene locus. Based on segregation in the family,

360  we found one allele carrying the intronic variant chr8:g.87618576G>A while the

361  second allele carried the known pathogenic variant c.1148delC;

362  p.Thr383IlefsTer13, and a second intronic variant chr8:g.87676221T>C (Table

363  S3). Both MAJIQ and CASH detected an alternative splicing event spanning

364  variant chr8:g.87618576G>A. In contrast, no alternative splicing events were

365  found to span variant chr8:g.87676221T>C (Figure4b and TableS6-8). Close

366  examination of the alternative splicing event spanning variant

367  chr8:g.87618576G>A revealed that it incorporated a cryptic exon into *CNGB3* in

368  RNA samples taken from both affected siblings but not from the unaffected

369  sibling (Figure 5a&b). The cryptic exon is spliced between canonical exon 14 and

370  exon 15 and will therefore be termed exon14b from hereon. The inclusion of

371  exon 14b was validated by RT-PCR and subsequent cloning and Sanger

372  sequencing of the novel longer isoform from the two affected siblings (Figure 5b).

373  Both the addition of exon 14b to the *CNGB3* transcript as a result of variant

374  chr8:g.87618576G>A, and the single base pair deletion in the second allele

375  carrying variant c.1148delC; p.Thr383IlefsTer13, lead to a frame shift and

376  subsequent premature termination. We therefore expected both alleles to

377  undergo nonsense mediated decay (NMD) with down regulation of *CNGB3*

378  mRNA levels in the affected siblings as compared to the unaffected sibling.

379  However, contrary to our expectations, *CNGB3* expression was not significantly

380  down regulated in our RNA-seq dataset, as analyzed using the edgeR program.

381  A comparison of each of the affected siblings with the unaffected sibling yielded

382  $\log_2$FC values of -1.05 and -1.13, indicating slightly lower expression levels in the

383  affected siblings that did not reach statistical significance (p-values of 0.13 and

384  0.16 respectively (Table S9)). These two frame shift alleles are predicted to

385  encode truncated proteins. The protein encoded by the exon 14b including

386  isoform is predicted to maintaining a full transmembrane domain but lack the

387  ligand binding domain of CNGB3 (Figure 5c). Similarly, the protein encoded by

388  the isoform carrying the known pathogenic variant c.1148delC;

389  p.Thr383IlefsTer13 is predicted to have a truncated transmembrane domain in

390  addition to lacking the ligand binding domain (Figure 5c). In order to determine

391  whether the truncated CNGB3 proteins are being translated, we performed

18

392    immunohistochemistry on ROs from the exon 14b allele carrier parent (OGI-081-

393    200) and one affected sibling (OGI-081-197, Figure 6). CNGB3 is a subunit of the

394    cone cyclic nucleotide-gated (CNG) channel, which localizes to cone

395    photoreceptor outer segments in chicken and mice [76,77]. We have also validated

396    the localization of human CNGB3 to cone photoreceptor outer segments in the

397    human retina (Figure S2). We therefore immunostained ROs for CNGB3 and ML

398    opsin, the latter serving as a marker for photoreceptor outer segments. For these

399    studies, stage 3 ROs were kept in culture for a total of 262 days, allowing cones

400    full opportunity to mature and localize ML opsin and CNGB3 to the photoreceptor

401    outer segments. As expected, CNGB3 co-localized with ML opsin in cone

402    photoreceptor outer segments in the parent (Figure 6c), with weaker staining

403    observed in inner segments as well (Figure 6b&c), presumably due to

404    mislocalization of truncated CNGB3 protein produced by the exon 14b including

405    allele. In ROs from the affected sibling, where both alleles are predicted to result

406    in truncated proteins, CNGB3 was only observed diffusely in the cell body and in

407    inner segments; i.e., no co-localization with ML opsin was observed in cone

408    photoreceptor outer segments (Figure 6e&f). Taken together, results from the

409    differential splicing analysis indicate that the likely cause for the inherited retinal

410    degeneration in OGI-081 is two pathogenic alleles in *CNGB3* - the known

411    pathogenic allele p.Thr383IlefsTer13 and the novel deep intronic allele

412    chr8:g.87618576G>A; NM_019098.3:c.1663 – 2137C>T; pLeu524 IlefsTer50.


413

414 <u>Splicing prediction algorithms</u>

415 With the identification of the non-coding pathogenic variant in *CNGB3,* we set out

416 to examine the mechanism by which it promotes the inclusion of exon 14b. We

417 analyzed the splice junctions surrounding exon14b with the variant analysis tool

418 Alamut Visual. Alamut Visual incorporates three splicing predictors capable of

419 analyzing deep intronic variants, SpliceSiteFinder-like (SSF)[78], MaxEntScan [79]

420 and NNSPLICE[80]. All three algorithms predicted chr8:g.87618576G>A to

421 strengthen a cryptic donor splice site (DSS) (Table 2). All three algorithms also

422 detected a potential acceptor splice site (ASS) at position c.1663-2238, exactly

423 where our Sanger sequencing indicated the acceptor site of exon 14b resides.

424 Interestingly, exon 14b ASS is a stronger than the one located at exon 15 (Table

425 2). It is plausible that the availability of this acceptor site and its ability to compete

426 with the acceptor site of exon 15 contributed to the effect of variant

427 chr8:g.87618576G>A on the splicing pattern of *CNGB3* in the affected siblings. In

428 addition we noticed the presence of a second even stronger alternative ASS

429 54bp upstream of exon 15 ASS. It is possible that this secondary competitor

430 further weakens the exon 15 ASS thus enhancing the effects of variant

431 chr8:g.87618576G>A.

432 Next, using OGI-081 as a true positive case, we tested whether splicing

433 prediction algorithms could be used to prioritize candidate non-coding splicing

434 altering variants, circumventing the need for RNA-seq analysis. We annotated

435 the 3,268 rare variants with allelic pairs identified in OGI-081 with two splicing

436 prediction programs. (i). Alamut Batch that makes its prediction by the combined

437    calculations of the same splicing predictions algorithms as Alamut visual but is

438    capable of calculating the effects of multiple variants. (ii). SpliceAI, a deep neural

439    network tool, to predict splice junctions from pre-mRNA transcript sequence [18].

440    Alamut batch calculated a high probability for altering splicing for 532 variants in

441    315 genes (Table S10). Although variant chr8:g.87618576G>A, the novel

442    pathogenic variant identified in this study, was predicted by Alamut Batch to

443    strongly activate a cryptic donor site the large number of additional candidate

444    variants make this tool too cumbersome for identification of candidate non-coding

445    pathogenic variants. For SpliceAI, to identify synonymous exonic, near intronic,

446    and deep intronic variants predicted to affect splicing at a validation rate of 40%

447    the authors used Δ Score greater than or equal to 0.2, 0.2, and 0.5 respectively.

448    Out of the variants segregating in OGI-081only eight had scores 0.2<0.5 and only

449    one, variant Chr9:g.86536129C>T, received a Δ Score >0.5 (Table S11). Variant

450    chr8:g.87618576G>A, the novel pathogenic variant identified in this study as

451    activating a cryptic donor splice site, was calculated by SpliceAI to have a donor

452    gain Δ Score of 0.3 well below the 0.5 cutoff for deep intronic variants. Thus, had

453    we used SpliceAI splicing predictions as a filter to identify potential causal

454    variants for functional validations, variant chr8:g.87618576G>A would have been

455    overlooked. The ASS of exon 14b was not identified by either algorithm and

456    therefore could not have been used to highlight variant chr8:g.87618576G>A as

457    a more plausible pathogenic variant.

458

## Discussion

460  We present here an unbiased approach based on the combination of WGS and

461  RNA-seq data to identify and functionally validate pathogenic non-coding variants

462  without the use of large datasets. We show that *ex vivo* models, such as iPSC

463  derived ROs, can serve as a surrogate source of a patient's own retinal tissue for

464  RNA and protein analyses. IRDs are currently at the focus of gene therapy

465  advances and several clinical trials are underway, including a trial for *CNGB3*

466  gene augmentation therapy [6]. This work was aimed at expanding the number of

467  patients eligible for clinical trials and forthcoming therapies. Indeed, our findings

468  here make the two affected siblings of OGI-081 eligible to participate in ongoing

469  clinical trials for *CNGB3* gene therapy. Our approach is applicable to any

470  inherited disease, both WGS and RNA-seq techniques are commercially

471  available, gold standards are being established and the analysis tools are readily

472  accessible [47,48,82,49–55,81]. *Ex vivo* organoid models are being developed for a

473  multitude of tissues including brain [83,84], kidney [85], liver [86,87] and lung [88].

474  Non-coding variants present a challenge for a correct genetic diagnosis that is

475  imperative for a successful genetic therapy. The combination of WGS and RNA-

476  seq methodologies allows us to both detect non-coding variants and evaluate

477  their functionality throughout the genome. Indeed, a similar approach has already

478  been successfully employed to diagnose diseases where RNA could be

479  harvested from biopsies of disease-relevant tissue [19,20]. These studies relied on

480  the availability of large control datasets of RNA-seq samples from unaffected

481 individuals and/or a large cohort of patients [19,20]. Our work shows that the correct

482 diagnosis of non-coding variants is possible without reliance on such resources.

483 WGS analysis of all five members of OGI-081 and segregation analysis of the

484 variants within the family narrowed down the search from tens of thousands of

485 variants to a few hundred with allelic pairs. We then used RNA-seq analysis

486 comparing two affected siblings to an unaffected one as an orthogonal approach

487 to identify genes with altered splicing or expression in disease. Thus identifying

488 the deep intronic variant chr8:g.87618576G>A as a novel pathogenic variant in

489 the *CNGB3* gene. The iPSC derived ROs served both as a source of disease

490 relevant transcriptome and as a system for functional validation of the truncated

491 proteins. In future studies, for families with a single patient the parents may serve

492 as control samples, so that each parent controls for the effect of the allele

493 inherited from the other parent making our approach applicable even for ultra-

494 rare diseases.

495 Once the deep intronic variant was detected and validated we were able to use

496 that prior knowledge to identify additional factors that may have contributed to the

497 inclusion of exon 14b such as the availability of the cryptic acceptor site of exon

498 14b and the comparative weakness of the exon 15 acceptor site. Such complex

499 dependencies are a prime example as to why sequence based predictions of

500 splicing patterns are hard to compute. Still, several splicing predictors in the

501 Alamut Visual software were able to detect the increase in the splicing probability

502 of the cryptic donor site as a result of variant chr8:g.87618576G>A. This

503 prompted us to test whether such splicing predictors can be used as preliminary

504    filters to identify candidate pathogenic variants for "gene by gene" validation

505    methods, circumventing the need for RNA-seq analysis of ROs. We found that

506    the more established approach represented by the Alamut Batch method of

507    combining the calculations of several splicing predictors that are designed to

508    identify known splicing motifs yielded too many candidates for gene by gene

509    validation. In contrast, the more recent approach of deep neuronal networks

510    algorithms, represented by SpliceAI, failed to assign high probability to the true

511    positive variant chr8:g.87618576G>A. Still, in cases where some prior knowledge

512    can help prioritize variants or highlight ones with lower than expected scores

513    these methods may yet be helpful. In cases where no prior knowledge can help

514    prioritize candidate variants, such as in patients where both pathogenic alleles

515    are non-coding, and especially in cases where a cell type relevant for functional

516    validation is not available, the approach established here is preferable.

# Appendices

## Differential gene expression analysis

519    As mentioned above, differential gene expression analysis was less informative

520    in the OGI-081 datasets. We compared gene expression levels from each of the

521    affected siblings to that of the unaffected sibling (Table S9). We found 401 genes

522    to be consistently down regulated, of which 27 were Y linked as expected given

523    that the two affected siblings are females while the unaffected sibling is a male.

524    We excluded these Y linked genes from further analysis. Tools are not currently

525    available to filter out non-Y linked genes that may be differentially expressed

526 between the sexes under normal conditions in ROs. Of the remaining 374 down

527 regulated genes, 29 also contained allelic variant pairs (Table S3). In addition,

528 we found 1120 genes to be consistently up regulated between the two affected

529 siblings and the unaffected sibling (Table S9). Of the up regulated genes, 15 also

530 contained allelic variant pairs (Table S3). None of the 44 differentially expressed

531 genes with allelic variant pairs are reported in RetNet as IRD genes.

## Supplemental Data

533 Supplemental data includes two figures and 11 tables.

## Acknowledgments

548    NINDS. We acknowledge the ENCODE Consortium and the Thomas Gingeras,

549    CSHL production lab for the use of skin RNA-seq samples ENCSR551NII,

550    ENCSR991HIR, ENCSR460YCS and ENCSR321PGV.

## 551    Declaration of Interests

552    The authors declare no competing interests.

## 553    Web Resources

554    ENCODE https://www.encodeproject.org/experiments/ENCSR551NII/

555    ExAC http://exac.broadinstitute.org/

556    gnomAD http://gnomad.broadinstitute.org/

557    GTEx Portal https://gtexportal.org/home/

558    Picard tools http://broadinstitute.github.io/picard/

559    RetNet  http://www.sph.uth.tmc.edu/RetNet/

## 560    Accession Numbers

561    The accession numbers for the RNAseq samples reported in this paper

562    (BioProject: PRJNA564377) are:

563    SRA:SRR10082823

564    SRA:SRR10082822

565    SRA:SRR10082821

26

566    SRA:SRR10082828

567    SRA:SRR10082829

568    SRA:SRR10082824

569    SRA:SRR10082830

570    SRA:SRR10082827

571    SRA:SRR10082826

572    SRA:SRR10082825

573    SRA:SRR10082820

574

# References

576    1. Sahel, J.A., Marazova, K., and Audo, I. (2015). Clinical characteristics and

577    current therapies for inherited retinal degenerations. Cold Spring Harb. Perspect.

578    Med.

579    2. Maguire, A.M., Simonelli, F., Pierce, E.A., Pugh, E.N., Mingozzi, F., Bennicelli,

580    J., Banfi, S., Marshall, K.A., Testa, F., Surace, E.M., et al. (2008). Safety and

581    efficacy of gene transfer for Leber's congenital amaurosis. N. Engl. J. Med. *358*,

582    2240–2248.

583    3. Hauswirth, W.W., Aleman, T.S., Kaushal, S., Cideciyan, A. V., Schwartz, S.B.,

584    Wang, L., Conlon, T.J., Boye, S.L., Flotte, T.R., Byrne, B.J., et al. (2008).

585    Treatment of Leber congenital amaurosis due to RPE65 mutations by ocular

586    subretinal injection of adeno-associated virus gene vector: Short-term results of a

587    phase I trial. Hum. Gene Ther. *19*, 979–990.

588    4. Bainbridge, J.W.B., Smith, A.J., Barker, S.S., Robbie, S., Henderson, R.,

589    Balaggan, K., Viswanathan, A., Holder, G.E., Stockman, A., Tyler, N., et al.

590    (2008). Effect of gene therapy on visual function in Leber's congenital amaurosis.

591    N. Engl. J. Med. *358*, 2231–2239.

592    5. Kannabiran, C., and Mariappan, I. (2018). Therapeutic avenues for hereditary

593    forms of retinal blindness. J. Genet. *97*, 341–352.

594    6. Moore, N.A., Morral, N., Ciulla, T.A., and Bracha, P. (2018). Gene therapy for

595    inherited retinal and optic nerve degenerations. Expert Opin. Biol. Ther.

596    7. Scholl, H.P.N., Strauss, R.W., Singh, M.S., Dalkara, D., Roska, B., Picaud, S.,

597    and Sahel, J.A. (2016). Emerging therapies for inherited retinal degeneration.

598    Sci. Transl. Med.

599    8. Lek, M., Karczewski, K.J., Minikel, E. V., Samocha, K.E., Banks, E., Fennell,

600    T., O'Donnell-Luria, A.H., Ware, J.S., Hill, A.J., Cummings, B.B., et al. (2016).

601    Analysis of protein-coding genetic variation in 60,706 humans. Nature.

602    9. Karczewski, K.J., Francioli, L.C., Tiao, G., Cummings, B.B., Alföldi, J., Wang,

603    Q., Collins, R.L., Laricchia, K.M., Ganna, A., Birnbaum, D.P., et al. (2019).

604    Variation across 141,456 human exomes and genomes reveals the spectrum of

605    loss-of-function intolerance across human protein-coding genes. BioRxiv 531210.

606    10. Vaz-Drago, R., Custódio, N., and Carmo-Fonseca, M. (2017). Deep intronic

607    mutations and human disease. Hum. Genet. *136*, 1093–1111.

608    11. Zhang, F., and Lupski, J.R. (2015). Non-coding genetic variants in human

609     disease: Figure 1. Hum. Mol. Genet. *24*, R102–R110.

610     12. Khurana, E., Fu, Y., Chakravarty, D., Demichelis, F., Rubin, M.A., and

611     Gerstein, M. (2016). Role of non-coding sequence variants in cancer. Nat. Rev.

612     Genet. *17*, 93–108.

613     13. Gloss, B.S., and Dinger, M.E. (2018). Realizing the significance of noncoding

614     functionality in clinical genomics. Exp. Mol. Med.

615     14. Anna, A., and Monika, G. (2018). Splicing mutations in human genetic

616     disorders: examples, detection, and confirmation.

617     15. Ohno, K., Takeda, J.-I.I., and Masuda, A. (2018). Rules and tools to predict

618     the splicing effects of exonic and intronic mutations.

619     16. Cheng, J., Nguyen, T.Y.D., Cygan, K.J., Celik, M.H., Fairbrother, W.G.,

620     Avsec, Z., and Gagneur, J. (2019). MMSplice: modular modeling improves the

621     predictions of genetic variant effects on splicing. Genome Biol. *20*, 48.

622     17. Jagadeesh, K.A., Paggi, J.M., Ye, J.S., Stenson, P.D., Cooper, D.N.,

623     Bernstein, J.A., and Bejerano, G. (2018). S-CAP extends clinical-grade

624     pathogenicity prediction to genetic variants that affect RNA splicing. BioRxiv.

625     18. Jaganathan, K., Kyriazopoulou Panagiotopoulou, S., McRae, J.F., Darbandi,

626     S.F., Knowles, D., Li, Y.I., Kosmicki, J.A., Arbelaez, J., Cui, W., Schwartz, G.B.,

627     et al. (2019). Predicting Splicing from Primary Sequence with Deep Learning.

628     Cell *176*, 535-548.e24.

629     19. Evrony, G.D., Cordero, D.R., Shen, J., Partlow, J.N., Yu, T.W., Rodin, R.E.,

630     Hill, R.S., Coulter, M.E., Lam, A.T.N., Jayaraman, D., et al. (2017). Integrated

631    genome and transcriptome sequencing identifies a noncoding mutation in the

632    genome replication factor DONSON as the cause of microcephaly-micromelia

633    syndrome. Genome Res.

634    20. Cummings, B.B., Marshall, J.L., Tukiainen, T., Lek, M., Donkervoort, S.,

635    Foley, A.R., Bolduc, V., Waddell, L.B., Sandaradura, S.A., O'Grady, G.L., et al.

636    (2017). Improving genetic diagnosis in Mendelian disease with transcriptome

637    sequencing. Sci. Transl. Med.

638    21. Lonsdale, J., Thomas, J., Salvatore, M., Phillips, R., Lo, E., Shad, S., Hasz,

639    R., Walters, G., Garcia, F., Young, N., et al. (2013). The Genotype-Tissue

640    Expression (GTEx) project. Nat. Genet.

641    22. Ratnapriya, R., Sosina, O.A., Starostik, M.R., Kwicklis, M., Kapphahn, R.J.,

642    Fritsche, L.G., Walton, A., Arvanitis, M., Gieser, L., Pietraszkiewicz, A., et al.

643    (2019). Retinal transcriptome and eQTL analyses identify genes associated with

644    age-related macular degeneration. Nat. Genet.

645    23. Aguet, F., Brown, A.A., Castel, S.E., Davis, J.R., He, Y., Jo, B., Mohammadi,

646    P., Park, Y.S., Parsana, P., Segrè, A. V., et al. (2017). Genetic effects on gene

647    expression across human tissues. Nature.

648    24. Danino, Y.M., Even, D., Ideses, D., and Juven-Gershon, T. (2015). The core

649    promoter: At the heart of gene expression. Biochim. Biophys. Acta - Gene Regul.

650    Mech. *1849*, 1116–1131.

651    25. Rickels, R., and Shilatifard, A. (2018). Enhancer Logic and Mechanics in

652    Development and Disease. Trends Cell Biol. *28*, 608–630.

653    26. Herzel, L., Ottoz, D.S.M., Alpert, T., and Neugebauer, K.M. (2017). Splicing

654    and transcription touch base: co-transcriptional spliceosome assembly and

655    function. Nat. Rev. Mol. Cell Biol. *18*, 637–650.

656    27. Ramanouskaya, T. V., and Grinev, V. V. (2017). The determinants of

657    alternative RNA splicing in human cells. Mol. Genet. Genomics *292*, 1175–1195.

658    28. Gonorazky, H., Liang, M., Cummings, B., Lek, M., Micallef, J., Hawkins, C.,

659    Basran, R., Cohn, R., Wilson, M.D., MacArthur, D., et al. (2016). RNAseq

660    analysis for the diagnosis of muscular dystrophy. Ann. Clin. Transl. Neurol. *3*,

661    55–60.

662    29. Gonorazky, H.D., Naumenko, S., Ramani, A.K., Nelakuditi, V., Mashouri, P.,

663    Wang, P., Kao, D., Ohri, K., Viththiyapaskaran, S., Tarnopolsky, M.A., et al.

664    (2019). Expanding the Boundaries of RNA Sequencing as a Diagnostic Tool for

665    Rare Mendelian Disease. Am. J. Hum. Genet. *104*, 466–483.

666    30. Capowski, E.E., Samimi, K., Mayerl, S.J., Phillips, M.J., Pinilla, I., Howden,

667    S.E., Saha, J., Jansen, A.D., Edwards, K.L., Jager, L.D., et al. (2019).

668    Reproducibility and staging of 3D human retinal organoids across multiple

669    pluripotent stem cell lines. Development *146*, dev171686.

670    31. Gonzalez-Cordero, A., Kruczek, K., Naeem, A., Fernando, M., Kloc, M.,

671    Ribeiro, J., Goh, D., Duran, Y., Blackford, S.J.I., Abelleira-Hervas, L., et al.

672    (2017). Recapitulation of Human Retinal Development from Human Pluripotent

673    Stem Cells Generates Transplantable Populations of Cone Photoreceptors. Stem

674    Cell Reports.

675  32. Meyer, J.S., Shearer, R.L., Capowski, E.E., Wright, L.S., Wallace, K.A.,

676  McMillan, E.L., Zhang, S.C., and Gamm, D.M. (2009). Modeling early retinal

677  development with human embryonic and induced pluripotent stem cells. Proc.

678  Natl. Acad. Sci. U. S. A.

679  33. Meyer, J.S., Howden, S.E., Wallace, K.A., Verhoeven, A.D., Wright, L.S.,

680  Capowski, E.E., Pinilla, I., Martin, J.M., Tian, S., Stewart, R., et al. (2011). Optic

681  vesicle-like structures derived from human pluripotent stem cells facilitate a

682  customized approach to retinal disease treatment. Stem Cells.

683  34. Nakano, T., Ando, S., Takata, N., Kawada, M., Muguruma, K., Sekiguchi, K.,

684  Saito, K., Yonemura, S., Eiraku, M., and Sasai, Y. (2012). Self-formation of optic

685  cups and storable stratified neural retina from human ESCs. Cell Stem Cell.

686  35. Phillips, M.J., Wallace, K.A., Dickerson, S.J., Miller, M.J., Verhoeven, A.D.,

687  Martin, J.M., Wright, L.S., Shen, W., Capowski, E.E., Percin, E.F., et al. (2012).

688  Blood-derived human iPS cells generate optic vesicle-like structures with the

689  capacity to form retinal laminae and develop synapses. Invest. Ophthalmol. Vis.

690  Sci.

691  36. Reichman, S., Terray, A., Slembrouck, A., Nanteau, C., Orieux, G., Habeler,

692  W., Nandrot, E.F., Sahel, J.A., Monville, C., and Goureaua, O. (2014). From

693  confluent human iPS cells to self-forming neural retina and retinal pigmented

694  epithelium. Proc. Natl. Acad. Sci. U. S. A.

695  37. Wahlin, K.J., Maruotti, J.A., Sripathi, S.R., Ball, J., Angueyra, J.M., Kim, C.,

696  Grebe, R., Li, W., Jones, B.W., and Zack, D.J. (2017). Photoreceptor Outer

697    Segment-like Structures in Long-Term 3D Retinas from Human Pluripotent Stem

698    Cells. Sci. Rep.

699    38. Zhong, X., Gutierrez, C., Xue, T., Hampton, C., Vergara, M.N., Cao, L.H.,

700    Peters, A., Park, T.S., Zambidis, E.T., Meyer, J.S., et al. (2014). Generation of

701    three-dimensional retinal tissue with functional photoreceptors from human

702    iPSCs. Nat. Commun.

703    39. Cora, V., Haderspeck, J., Antkowiak, L., Mattheus, U., Neckel, P.H., Mack,

704    A.F., Bolz, S., Ueffing, M., Pashkovskaia, N., Achberger, K., et al. (2019). A

705    Cleared View on Retinal Organoids. Cells *8*, 391.

706    40. Kim, S., Lowe, A., Dharmat, R., Lee, S., Owen, L.A., Wang, J., Shakoor, A.,

707    Li, Y., Morgan, D.J., Hejazi, A.A., et al. (2019). Generation, transcriptome

708    profiling, and functional validation of cone-rich human retinal organoids. Proc.

709    Natl. Acad. Sci. *116*, 10824–10833.

710    41. Capowski, E.E., Samimi, K., Mayerl, S.J., Phillips, M.J., Pinilla, I., Howden,

711    S.E., Saha, J., Jansen, A.D., Edwards, K.L., Jager, L.D., et al. (2018).

712    Reproducibility and staging of 3D human retinal organoids across multiple

713    pluripotent stem cell lines. Development *146*, dev171686.

714    42. Consugar, M.B., Navarro-Gomez, D., Place, E.M., Bujakowska, K.M., Sousa,

715    M.E., Fonseca-Kelly, Z.D., Taub, D.G., Janessian, M., Wang, D.Y., Au, E.D., et

716    al. (2015). Panel-based genetic diagnostic testing for inherited eye diseases is

717    highly accurate and reproducible, and more sensitive for variant detection, than

718    exome sequencing. Genet. Med. *17*, 253–261.

719   43. Falk, M.J., Zhang, Q., Nakamaru-Ogiso, E., Kannabiran, C., Fonseca-Kelly,

720   Z., Chakarova, C., Audo, I., MacKay, D.S., Zeitz, C., Borman, A.D., et al. (2012).

721   NMNAT1 mutations cause Leber congenital amaurosis. Nat. Genet.

722   44. Hardwick, S.A., Chen, W.Y., Wong, T., Deveson, I.W., Blackburn, J.,

723   Andersen, S.B., Nielsen, L.K., Mattick, J.S., and Mercer, T.R. (2016). Spliced

724   synthetic genes as internal controls in RNA sequencing experiments. Nat.

725   Methods *13*, 792–798.

726   45. McLaren, W., Gil, L., Hunt, S.E., Riat, H.S., Ritchie, G.R.S., Thormann, A.,

727   Flicek, P., and Cunningham, F. (2016). The Ensembl Variant Effect Predictor.

728   Genome Biol.

729   46. Handsaker, R.E., Van Doren, V., Berman, J.R., Genovese, G., Kashin, S.,

730   Boettger, L.M., and Mccarroll, S.A. (2015). Large multiallelic copy number

731   variations in humans. Nat. Genet.

732   47. Ewels, P., Magnusson, M., Lundin, S., and Käller, M. (2016). MultiQC:

733   summarize analysis results for multiple tools and samples in a single report.

734   Bioinformatics *32*, 3047–3048.

735   48. Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S.,

736   Batut, P., Chaisson, M., and Gingeras, T.R. (2013). STAR: ultrafast universal

737   RNA-seq aligner. Bioinformatics *29*, 15–21.

738   49. Liao, Y., Smyth, G.K., and Shi, W. (2014). featureCounts: an efficient general

739   purpose program for assigning sequence reads to genomic features.

740   Bioinformatics *30*, 923–930.

741   50. Wong, T., Deveson, I.W., Hardwick, S.A., and Mercer, T.R. (2017).

742   ANAQUIN: a software toolkit for the analysis of spike-in controls for next

743   generation sequencing. Bioinformatics *33*, 1723–1724.

744   51. Wu, W., Zong, J., Wei, N., Cheng, J., Zhou, X., Cheng, Y., Chen, D., Guo, Q.,

745   Zhang, B., and Feng, Y. (2018). CASH: a constructing comprehensive splice site

746   method for detecting alternative splicing events. Brief. Bioinform. *19*, 905–917.

747   52. Vaquero-Garcia, J., Barrera, A., Gazzara, M.R., González-Vallinas, J.,

748   Lahens, N.F., Hogenesch, J.B., Lynch, K.W., and Barash, Y. (2016). A new view

749   of transcriptome complexity and regulation through the lens of local splicing

750   variations. Elife *5*, e11752.

751   53. Robinson, M.D., McCarthy, D.J., and Smyth, G.K. (2010). edgeR: a

752   Bioconductor package for differential expression analysis of digital gene

753   expression data. Bioinformatics *26*, 139–140.

754   54. Wang, L., Wang, S., and Li, W. (2012). RSeQC: quality control of RNA-seq

755   experiments. Bioinformatics *28*, 2184–2185.

756   55. Bray, N.L., Pimentel, H., Melsted, P., and Pachter, L. (2016). Near-optimal

757   probabilistic RNA-seq quantification. Nat. Biotechnol. *34*, 525–527.

758   56. Farkas, M.H., Grant, G.R., White, J. a, Sousa, M.E., Consugar, M.B., Pierce,

759   E. a, and Genomics, B. (2013). Transcriptome analyses of the human retina

760   identify unprecedented transcript diversity and 3.5 Mb of novel transcribed

761   sequence via significant alternative splicing and novel genes. BMC Genomics *14*,

762   486.

763    57. Mossman, S.S., Bronstein, A.M., Gresty, M.A., Kendall, B., and Rudge, P.

764    (1990). Convergence nystagmus associated with Arnold-Chiari malformation.

765    Arch. Neurol. *47*, 357–359.

766    58. Milhorat, T.H., Chou, M.W., Trinidad, E.M., Kula, R.W., Mandell, M., Wolpert,

767    C., and Speer, M.C. (1999). Chiari I malformation redefined: clinical and

768    radiographic findings for 364 symptomatic patients. Neurosurgery *44*, 1005–

769    1017.

770    59. Urbizu, A., Toma, C., Poca, M.A., Sahuquillo, J., Cuenca-León, E., Cormand,

771    B., and Macaya, A. (2013). Chiari malformation type I: a case-control association

772    study of 58 developmental genes. PLoS One *8*, e57241.

773    60. Shaikh, A., and Ghasia, F. (2015). Neuro-ophthalmology of type Chiari

774    malformation. Expert Rev Ophtalmol *10*, 351–357.

775    61. Ray, C., Nagy, L., and Mobley, J. (2014). Familial aggregation of chiari

776    malformation: presentation, pedigree, and review of the literature. Turk.

777    Neurosurg. *26*, 315–320.

778    62. Yuan, X.X., Li, Y., Sha, S.F., Sun, W.X., Qiu, Y., Liu, Z., Zhu, W.G., and Zhu,

779    Z.Z. (2017). [Genetic analysis of posterior cranial fossa morphology in families of

780    Chiari malformation type Ⅰ]. Zhonghua Yi Xue Za Zhi *97*, 1140–1144.

781    63. Sarnat, H.B. (2018). Cerebellar networks and neuropathology of cerebellar

782    developmental disorders (Elsevier).

783    64. Merello, E., Tattini, L., Magi, A., Accogli, A., Piatelli, G., Pavanello, M.,

784    Tortora, D., Cama, A., Kibar, Z., Capra, V., et al. (2017). Exome sequencing of

785   two Italian pedigrees with non-isolated Chiari malformation type I reveals

786   candidate genes for cranio-facial development. Eur. J. Hum. Genet. *25*, 952–959.

787   65. Solis-Moruno, M., de Manuel, M., Hernandez-Rodriguez, J., Fontsere, C.,

788   Gomara-Castaño, A., Valsera-Naranjo, C., Crailsheim, D., Navarro, A., Llorente,

789   M., Riera, L., et al. (2017). Potential damaging mutation in LRP5 from genome

790   sequencing of the first reported chimpanzee with the Chiari malformation. Sci.

791   Rep. *7*, 15224.

792   66. Wiszniewski, W., Lewis, R.A., and Lupski, J.R. (2007). Achromatopsia: the

793   CNGB3 p.T383fsX mutation results from a founder effect and is responsible for

794   the visual phenotype in the original report of uniparental disomy 14. Hum. Genet.

795   *121*, 433–439.

796   67. Nishiguchi, K.M., Sandberg, M.A., Gorji, N., Berson, E.L., and Dryja, T.P.

797   (2005). Cone cGMP-gated channel mutations and clinical findings in patients with

798   achromatopsia, macular degeneration, and other hereditary cone diseases. Hum.

799   Mutat. *25*, 248–258.

800   68. Kohl, S. (2000). Mutations in the CNGB3 gene encoding the beta-subunit of

801   the cone photoreceptor cGMP-gated channel are responsible for achromatopsia

802   (ACHM3) linked to chromosome 8q21. Hum. Mol. Genet. *9*, 2107–2116.

803   69. Maguire, J., McKibbin, M., Khan, K., Kohl, S., Ali, M., and McKeefry, D.

804   (2018). CNGB3 mutations cause severe rod dysfunction. Ophthalmic Genet. *39*,

805   108–114.

806   70. Mayer, A.K., Van Cauwenbergh, C., Rother, C., Baumann, B., Reuter, P., De

807    Baere, E., Wissinger, B., and Kohl, S. (2017). CNGB3 mutation spectrum

808    including copy number variations in 552 achromatopsia patients. Hum. Mutat. *38*,

809    1579–1591.

810    71. Nishiguchi, K.M., and Rivolta, C. (2012). Genes associated with retinitis

811    pigmentosa and allied diseases are frequently mutated in the general population.

812    PLoS One *7*,.

813    72. Baralle, F.E., and Giudice, J. (2017). Alternative splicing as a regulator of

814    development and tissue identity. Nat. Rev. Mol. Cell Biol. *18*, 437–451.

815    73. Bush, S.J., Chen, L., Tovar-Corona, J.M., and Urrutia, A.O. (2017).

816    Alternative splicing and the evolution of phenotypic novelty. Philos. Trans. R.

817    Soc. Lond. B. Biol. Sci. *372*,.

818    74. ENCODE Project, Bernstein, B.E., Birney, E., Dunham, I., Green, E.D.,

819    Gunter, C., and Snyder, M. (2012). An integrated encyclopedia of DNA elements

820    in the human genome. Nature.

821    75. Davis, C.A., Hitz, B.C., Sloan, C.A., Chan, E.T., Davidson, J.M., Gabdank, I.,

822    Hilton, J.A., Jain, K., Baymuradov, U.K., Narayanan, A.K., et al. (2018). The

823    Encyclopedia of DNA elements (ENCODE): data portal update. Nucleic Acids

824    Res. *46*, D794–D801.

825    76. Bönigk, W., Altenhofen, W., Müller, F., Dose, A., Illing, M., Molday, R.S., and

826    Kaupp, U.B. (1993). Rod and cone photoreceptor cells express distinct genes for

827    cGMP-gated channels. Neuron *10*, 865–877.

828    77. Matveev, A. V, Quiambao, A.B., Browning Fitzgerald, J., and Ding, X.-Q.

829    (2008). Native cone photoreceptor cyclic nucleotide-gated channel is a

830    heterotetrameric complex comprising both CNGA3 and CNGB3: a study using

831    the cone-dominant retina of Nrl-/- mice. J. Neurochem. *106*, 2042–2055.

832    78. Shapiro, M.B., and Senapathy, P. (1987). RNA splice junctions of different

833    classes of eukaryotes: sequence statistics and functional implications in gene

834    expression. Nucleic Acids Res. *15*, 7155–7174.

835    79. Yeo, G., and Burge, C.B. (2004). Maximum entropy modeling of short

836    sequence motifs with applications to RNA splicing signals. J. Comput. Biol. *11*,

837    377–394.

838    80. Reese, M.G., Eeckman, F.H., Kulp, D., and Haussler, D. (1997). Improved

839    Splice Site Detection in Genie. J. Comput. Biol. *4*, 311–323.

840    81. Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with

841    Bowtie 2. Nat. Methods *9*, 357–359.

842    82. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth,

843    G., Abecasis, G., Durbin, R., and 1000 Genome Project Data Processing

844    Subgroup, 1000 Genome Project Data Processing (2009). The Sequence

845    Alignment/Map format and SAMtools. Bioinformatics *25*, 2078–2079.

846    83. Cosset, E., Locatelli, M., Marteyn, A., Lescuyer, P., Dall Antonia, F., Mor,

847    F.M., Preynat-Seauve, O., Stoppini, L., and Tieng, V. (2019). Human Neural

848    Organoids for Studying Brain Cancer and Neurodegenerative Diseases. J. Vis.

849    Exp.

850    84. Gopalakrishnan, J. (2019). The Emergence of Stem Cell-Based Brain

851    Organoids: Trends and Challenges. Bioessays *41*, e1900011.

852    85. Low, J.H., Li, P., Chew, E.G.Y., Zhou, B., Suzuki, K., Zhang, T., Lian, M.M.,

853    Liu, M., Aizawa, E., Rodriguez Esteban, C., et al. (2019). Generation of Human

854    PSC-Derived Kidney Organoids with Patterned Nephron Segments and a De

855    Novo Vascular Network. Cell Stem Cell.

856    86. Prior, N., Inacio, P., and Huch, M. (2019). Liver organoids: from basic

857    research to therapeutic applications. Gut.

858    87. Torresi, J., Tran, B.M., Christiansen, D., Earnest-Silveira, L., Schwab,

859    R.H.M., and Vincan, E. (2019). HBV-related hepatocarcinogenesis: the role of

860    signalling pathways and innovative ex vivo research models. BMC Cancer *19*,

861    707.

862    88. Strikoudis, A., Cieślak, A., Loffredo, L., Chen, Y.W., Patel, N., Saqi, A.,

863    Lederer, D.J., and Snoeck, H.W. (2019). Modeling of Fibrotic Lung Disease

864    Using 3D Organoids Derived from Human Pluripotent Stem Cells. Cell Rep.

865

866    ## Figure Legends

867    **Figure 1: Family OGI-081 variant segregation scheme and retinal**

868    **phenotypes.** A. The OGI-081 pedigree with variant segregation scheme. B.

869    Fundus (upper panel) and OCT image (lower panel) for the OGI-081-197 at age

870    8, area of retinal degeneration is indicated by the red bar.

871    **Figure 2: RO differentiation.** A. Schema of the differentiation process and a

872    light microscopy image of a typical RO. Arrow head indicating photoreceptors,

873    scale bar = 100 microns.  B-G. Immunocytochemistry on cryosections of ROs.

874    B&E. NR2E3 staining of rod nuclei (Green). C&F. Mature cones show staining of

875    cone opsins in the cone photoreceptor outer segments (Red). C. S opsin. F. ML

876    opsin. D&G. Overlay of rod and cone staining. All cones are stained with ARR3

877    (Purple). Scale bars = 20 microns.

878    **Figure 3: Comparison of IRD gene expression and splice junctions.** Human

879    Normal Retina (HNR; N=3, gray), RO from the unaffected sibling (N=5, blue),

880    Skin-Sun Exposed (SSE; N=473, green) and Whole Blood (WB; N=407, red).

881    A&B Average TPM values of IRD genes. A. IRD genes are sorted by their

882    expression in HNR overlaid with RO, SSE or WB. B. Violin plot. C&D Number of

883    splice junctions detected by MAJIQ. C. All annotated genes. D. IRD genes.

884    **Figure 4: Alternative splicing in the *NCALD* and *CNGB3* genes.** A. Venn

885    diagram of genes found to have alternative splicing events in OGI-081

886    comparison of affected vs. unaffected siblings and genes found to have

887    segregating allelic pairs (green). Alternative splicing analysis was conducted by

888    MAJIQ (blue) and CASH (red). B. Collapsed diagram of exons (black boxes) from

889    all isoforms of the *NCALD* and *CNGB3* genes. DNA variants (red); MAJIQ (blue)

890    and CASH (yellow) alternative splicing events (E). Events detected by MAJIQ are

891    depicted as split reads arches. The event range detected by CASH is depicted by

892    the left (L) and right (R) borders. Genomic locations of variants, junctions and

893    event borders are given in table S6.

894    **Figure 5: Aberrant splicing of *CNGB3* in the affected vs. unaffected**

895    **siblings.** A. Sashimi plot presenting RNA-seq results showing a cryptic exon

896    spliced into the isoform as a result of the intronic variant chr8:g.87618576G>A.

897    The cryptic exon is only present in the affected sibling (lower panel, red) and not

898    in the unaffected sibling (upper panel, blue). The splice junction between exon

899    14b and exon 15 is not represented by split reads in the Sashimi plot due to an

900    alignment error (FigureS1). B. RT-PCR using primers from the canonical exon14

901    and exon 15. All three siblings express the normal size isoform lacking exon 14b

902    (lower band). A larger abnormal band containing exon14b (upper band) is

903    present in the two affected siblings Af1 and Af2 but not in the unaffected sibling

904    (Un). Negative controls lacking RNA template in the RT reaction (NC1) and NC1

905    used as template for PCR amplification (NC2). Sanger sequencing of the larger

906    band confirming the inclusion of exon 14b. C. Schematic representation of the

907    protein domains in W.T CNGB3 and the two mutant alleles found in the affected

908    siblings of OGI-081.

909    **Figure 6: Mislocalization of the CNGB3 truncated proteins.**

910    Immunocytochemical analysis of day 262 ROs from the heterozygous parent

911    OGI-081-200 (A-C) and an affected sibling OGI-081-197 (D-F). In the

912    heterozygote, both ML opsins (red) and CNGB3 (green) are localized to the

913    photoreceptor outer segments whereas in the affected sibling, CNGB3 localizes

914    to the photoreceptor inner segments. An exemplary photoreceptor outer segment

915    is indicated by the white brackets. Nuclei are counterstained with DAPI (blue).

916    Scale bars = 20 micron

917

918 **Table 1: Sample level alignment report and QC summary.** A high average

919 unique percentage alignment rate is reported.

| SampleName | RawReads[a] | HQReads[b] | % UniqAligned[c] |
|---|---|---|---|
| OGI-081-197-1 | 159,976,749 | 129,671,633 | 87.73 |
| OGI-081-197-2 | 256,399,765 | 234,701,223 | 90.52 |
| OGI-081-197-3 | 166,237,226 | 155,442,261 | 88.88 |
| OGI-081-198-1 | 167,327,080 | 145,058,230 | 88.66 |
| OGI-081-198-2 | 197,037,935 | 179,026,100 | 89.13 |
| OGI-081-198-3 | 191,485,904 | 177,663,416 | 90.15 |
| OGI-081-340-1 | 175,199,604 | 156,566,678 | 91.01 |
| OGI-081-340-2 | 227,928,431 | 216,526,830 | 89.53 |
| OGI-081-340-3 | 137,096,260 | 129,192,437 | 90.47 |
| OGI-081-340-4 | 151,572,341 | 140,982,711 | 90.59 |
| OGI-081-340-5 | 164,908,758 | 154,550,198 | 99.66 |

920 [a] Raw reads count (RawReads), [b] Filtered high quality reads (HQReads), [c]

921 Percent of uniquely aligned reads (% UniqAligned).

922 **Table 2: Splicing junctions involved with the inclusion of exon14b.**

| | Exon 14 Donor site | Exon 14b Acceptor site | Exon 14b Donor site | | Exon 15 Acceptor site | |
|---|---|---|---|---|---|---|
| | Annotated | W.T | W.T | Variant | Annotated | W.T |
| NM_019098.4[a] | c.1662 | c.1663-2238 | c.1663-2137 | c. 1663-2137C>T | c.1663 | c.1663-54 |
| SSF[b] | 95.3 | 83.3 | 67.8 | 73.7 | 76.7 | 90.4 |
| MaxEntScan[c] | 10.5 | 4.8 | 0 | 4.8 | 6.1 | 8.7 |
| NNSplice[d] | 1 | 0.5 | 0 | 1 | 0.1 | 0.8 |

923 [a] The NM_019098.4 isoform of CNGB3 is used to define the cDNA coordinates. [b-]

924 [d] Splicing prediction scores are given from the three algorithms, SSF,

925 MaxEntScan and NNSplice.

44

A.



B.

A.

**A.**

RO      SSE      WB

IRD Genes

**B.** IRD Genes

TPM

HNR   RO   SSE   WB

**C.** All Genes

#Splice Junctions

HNR   RO   SSE

**D.** IRD Genes

#Splice Junctions

HNR   RO   SSE

A.

B.

*NCALD*

V2.1 V1 E1L ES1 ES1 V2.2 ES2 ES2 E1R V2.3 V2.4

*CNGB3*

ES1 ES1 V1.1 V1.2

E1L ES2 ES2 V2 E2L E1R ES3 E2R

14b