

1 **Functional filter for whole genome sequencing data**
2 **identifies HHT and stress-associated non-coding *SMAD4***
3 **polyadenylation site variants >5kb from coding DNA**

4
5 Sihao Xiao,^{1,2†*}, Zhentian Kai,³ Daniel Murphy,^{2,4} Dongyang Li,^{1,2} Dilip Patel,^{1,2} Adrianna Bielowka,^{1,2} Maria E.
6 Bernabeu-Herrero,^{1,2} Awatif Abdulmogith,^{1,2} Andrew D Mumford,⁵ Sarah Westbury,⁵ Micheala A Aldred,⁶ Neil
7 Vargesson,⁷ Mark J Caulfield,⁸ Genomics England Research Consortium,^{9‡} and Claire L Shovlin^{1,2,10*}

8
9 **Affiliations:**

10 ¹National Heart and Lung Institute, Imperial College London, London W12 ONN, UK.

11 ² National Institute for Health Research (NIHR) Imperial Biomedical Research Centre, London, W2 1NY
12 UK.

13 ³Topgen Biopharm Technology Co. Ltd; Shanghai, 201203, China.

14 ⁴Women's, Children's & Clinical Support (Pharmacy), Imperial College Healthcare NHS Trust, London,
15 W2 1NY, UK

16 ⁵School of Molecular and Cellular Medicine, University of Bristol, Bristol BS8 1QU, UK.

17 ⁶Division of Pulmonary, Critical Care, Sleep & Occupational Medicine, Indiana University School of
18 Medicine, Indianapolis, IN 46202 USA.

19 ⁷ School of Medicine, Medical Sciences and Nutrition, University of Aberdeen, Aberdeen AB24 3FX, UK.

20 ⁸William Harvey Research Institute, Queen Mary University of London, London E1 4NS, UK.

21 ⁹Genomics England, London EC1M 6BQ, UK,

22 ¹⁰Specialist Medicine, Imperial College Healthcare NHS Trust; London, W12 OHS, UK

23
24 *Correspondence c.shovlin@imperial.ac.uk (email to include both before and after publication – this
25 author will deal with manuscript submission). (include after publication sihao.xiao@bnc.ox.ac.uk)

26
27 † Current address: Big Data Institute, University of Oxford, Oxford, UK.

28 ‡ A full list of these authors is provided at the end of the main manuscript

29 **Abstract:** Despite whole genome sequencing (WGS), many single gene disorder cases remain
30 unsolved, impeding diagnosis and preventative care for people whose disease-causing variants
31 escape detection. Since early WGS data analytic steps prioritize protein-coding sequences, to
32 simultaneously prioritize variants in non-coding regions rich in transcribed and critical regulatory
33 sequences, we developed GROFFFY, an analytic tool which integrates coordinates for regions with
34 experimental evidence of functionality. Applied to WGS data from solved and unsolved hereditary
35 hemorrhagic telangiectasia (HHT) recruits to the 100,000 Genomes Project, GROFFFY-based
36 filtration reduced the mean number of variants per DNA from 4,867,167 to 21,486, without deleting
37 disease-causal variants. In three unsolved cases (two related), GROFFFY identified ultra-rare
38 deletions within the 3' untranslated region (UTR) of the proto-oncogene *SMAD4*, where germline
39 loss-of-function alleles cause combined HHT and colonic polyposis ([MIM: 175050](#)). Sited >5.4kb
40 distal to coding DNA, the deletions did not modify or generate microRNA binding sites, but instead
41 disrupted the sequence context of the final cleavage and polyadenylation site necessary for protein
42 production: By iFoldRNA, an AAUAAA-adjacent 16 nucleotide deletion brought the cleavage site
43 into inaccessible neighboring secondary structures, while a 4-nucleotide deletion unfolded the
44 downstream RNA polymerase II roadblock. *SMAD4* RNA expression differed to control-derived
45 RNA in resting and cycloheximide-stressed peripheral blood mononuclear cells. Patterns predicted
46 the mutational site for an unrelated HHT/polyposis- affected individual, where a complex insertion
47 was subsequently identified. In conclusion, we describe a functional rare variant type that impacts
48 regulatory systems based on RNA polyadenylation. Extension of coding sequence-focused gene
49 panels is required to capture these variants.

50

51

52 **INTRODUCTION**

53 Whole genome sequencing (WGS) is an established component of medical genetic and research
54 repertoires, but currently, the majority of its potential is unrealized. In any one individual, WGS
55 identifies millions of DNA variants compared to reference sequences. These are present in ~20,000
56 protein-coding genes, and also in much less understood regions of the genome that have diverse
57 functions including transcription into noncoding RNAs, participation in DNA chemical changes
58 that modify transcription, and binding to other nucleic acids or proteins.^{1,2}

59
60 Current WGS clinical foci are almost exclusively on a subgroup of protein-coding genes where
61 biological function is already known. In research spheres, in order to reduce the number of
62 variables per sample, interrogation of WGS data also commences with prioritization methods,
63 usually based on selection of specific genomic regions. Variants in the non-coding genome, while
64 not pre-depleted by the sequencing methodology, are effectively deleted in the early analytic stages
65 of variant prioritization. Importantly, application of these WGS methods leave large proportions
66 of ~~patients~~ individuals with hereditary conditions unsolved, without a genetic diagnosis.³

67
68 There is no accurate map of all functional genomic regions in human genomes, and it is difficult
69 to predict *a priori*, where all regulatory elements for a specific gene locus would be located. We
70 hypothesized however, that it would be possible to design a more efficient variant prioritization
71 method for WGS because markers of epigenetics and DNA-protein interactions have been applied
72 genome-wide by molecular laboratories, and an enormous body of biological experimental data
73 made publicly available. As a result, there now exist repositories of information indicating which
74 sections of DNA are more or less likely to have a functional role in at least one examined tissue.

75

76 We designed a ~~g~~enomic ~~r~~egions ~~o~~f ~~f~~unctionality ~~f~~ilter ~~f~~or priority (*GROFFFY*) based on published
77 experimental data particularly from ENCODE⁴⁻⁶ and performed validation and discovery analyses
78 in WGS data from ~~patients-individuals~~ recruited to the 100,000 Genomes Project.⁷ To accelerate
79 clinical impact, we focused discovery analyses on noncoding regions of a proto-oncogene
80 examined in diagnostic and screening gene panels. *SMAD4* is ubiquitously expressed and encodes
81 the common partner SMAD which regulates signaling by transforming growth factor (TGF)- β ,
82 bone morphogenetic protein (BMP), and activin ligands.^{8,9} As indicated by its function and earlier
83 gene names (*DPC4*-deleted in pancreatic cancer 4; *MADH4*-mothers against decapentaplegic), the
84 SMAD4 protein has major pathological and developmental roles.^{8,9} *SMAD4* is a target of cancer
85 genetic diagnostics because it is a driver gene for major cancers due to somatic loss,^{8,9} and because
86 germline heterozygous loss causes gastrointestinal polyposis (“juvenile polyposis syndrome”/JPS
87 [MIM: 174900]) where untreated hamartomatous polyps can undergo malignant transformation
88 leading to colon, gastric and other cancers.^{8,9} Heterozygous loss also causes TGF- β /BMP-related
89 vasculopathies including hereditary hemorrhagic telangiectasia (HHT) which usually results from
90 a loss-of-function variant in *ACVRL1* [MIM: 600376] or *ENG* [MIM: 187300].¹⁰ Where *SMAD4*
91 is identified (juvenile polyposis/hereditary hemorrhagic telangiectasia syndrome, JPHT [MIM:
92 175050]),^{10,11,12} ~~this~~ allows ~~patients-affected-individuals~~ to benefit from life-long polyposis and
93 aortopathy screening programs.⁹ Scientifically, *SMAD4* is of great interest because despite its
94 wide-ranging roles in development and disease, little is known of its regulation.^{8,9}

95

96 Here we report that this ~~new~~-WGS analytic approach identifies a ~~novel~~-type of functional DNA
97 variant uncaptured by usual clinical sequencing methodologies.

98

99 **MATERIALS AND METHODS**

100 The procedures followed were in accordance with the ethical standards of the responsible
101 committees on human experimentation (institutional and national) and proper informed consent
102 was obtained.

103

104 **Study Design**

105 The main elements of Study Design are outlined in [Figure 1](#). Following earlier recruitment of
106 individuals with hereditary hemorrhagic telangiectasia (HHT) to the 100,000 Genomes Project
107 (black arrows), the GROFFFY filter was designed as indicated in colored boxes, and applied to
108 the WGS data files within the 100,000 Genomes Project.

109

110 **Patient recruitment and sequencing**

111 The 100,000 Genomes Project was set up by the UK Department of Health and Social Security in
112 2013, to sequence whole genomes from National Health Service (NHS) patients. The study
113 received ethical approval from the Health Research Authority (HRA) Committee East England-
114 Cambridge South (REC Ref 14/EE/1112), and all participants provided written consent.
115 Anonymized raw sequencing data are available in the Genomics England Research
116 Environment.^{7,10} Separately, in a clinical diagnostic pipeline, Genomics England performed data
117 alignments and variant classifications fed back to recruiting clinicians.

118

119 The cohort recruited with hereditary hemorrhagic telangiectasia (HHT)^{10,11,12} were particularly
120 suited for *GROFFFY* methodological validation processes because a subset had not undergone
121 prior genetic testing, and because clinical pipelines were incomplete at the time of *GROFFFY*
122 analyses. This resulted in a validation dataset of 34 WGS sequences where clinical pipelines had
123 identified a causal variant in *ACVRL1*, *ENG* or *SMAD4*,¹²³⁻¹⁷⁸ and discovery dataset of 98 WGS

124 sequences where some DNAs were expected to have heterozygous loss-of-function variants in
125 these genes.

126

127 **Generating GROFFFY**

128 Genomic coordinates of regions included in GROFFFY were generated from publicly available
129 databases using the Imperial College High Performing Computing service. Experimentally-
130 derived biological data were used in preference to computational predicted files with potential for
131 false negatives. Genomic coordinates were extracted from data aligned to GRCh38,¹⁹ and excluded
132 data originating in cancer cells. Regions being selected for are described in *Tables S1-S4* which
133 provide full details of coordinate derivation from transcribed loci and candidate regulatory element
134 (cRE) regions.^{6,18,19,20,21} Following merging of 18,828 bed files from 3,454 experiments,^{5,6}
135 sequences in the ENCODE blacklist²⁰ and RepeatMasker²¹ were excluded. In detail:

136

137 ***Genomic coordinates for candidate regulatory element (cRE) regions*** were generated using data
138 from the ENCODE Encyclopedia registry which includes data from both ENCODE^{4,5,24} and the
139 NIH Roadmap Epigenomics Consortium.^{25,22} We downloaded the call sets itemized in *Table S3* and
140 *Table S4* from the ENCODE portal.^{4,5} By merging DNA binding call sets²⁴ and representative
141 DNase hypersensitivity site (rDHS) call sets,^{26,23} a rough prediction of all cREs was made. Only
142 data aligned to GRCh38 were retained. Data generated from cancer cells were excluded as cancer
143 cells' genome are usually heavily modified and rearranged.^{27,24}

144 a) For DNA binding data, Histone ChIP-seq and Transcription Factor ChIP-seq which
145 target histones, transcription factors, chromatin remodelers, RNA polymerase complexes,
146 RNA binding proteins, cofactors, DNA replication proteins and DNA repair proteins, were

147 searched. We downloaded the call sets from the ENCODE portal⁴ as indicated in *Table S3*.
148 The 18,828 downloaded files (*Table S3*) were from 2,823 experiments performed in 9
149 different labs, and representing biosamples from 76 tissue types, 63 cell lines, 52 primary
150 cell types, 20 in-vitro differentiated cell types, and all life stages. Downloaded bed files
151 were divided into 10 subgroups with the first 9 subgroups containing 2,000 bed files and
152 the last subgroup containing 828 bed files. Bed files from each subgroup were merged
153 together using BEDOPS²⁵ and then merged bed files from each subgroup were joined
154 together last to obtain all DNA binding regions. The merged bed file for DNA-binding
155 regions was 49,941,695 Kb (i.e. greater than 15 genomes) before sorting, reflecting many
156 overlapping regions.

157 b) For rDHSs, DNA accessibility experiments were searched in the ENCODE
158 Encyclopedia database.⁵ We downloaded 4,409 files as indicated in *Table S4*. These were
159 from 631 experiments carried out in 2 different labs and represented biosamples from 69
160 cell lines, 76 primary cell types, 86 tissue types, and all life stages. Bed files were merged
161 together directly using BEDOPS (V2.4.26)²⁸⁻²⁵ to obtain all accessible DNA regions.

162 c) Region coordinates for CpG islands were downloaded from the UCSC database^{29,30-26,27}
163 using Rsync command tool. The downloaded file was in txt format and was converted to
164 bed files using awk function in the Linux³¹ system. The converted bed file was then sorted
165 by location using BEDOPS.^{28,25}

166

167 ***Genomic coordinates for transcribed loci*** were extracted as follows:

168 d) GENCODE human genome annotation version 31³⁰ for GRCh38⁴⁹ was downloaded
169 from the UCSC database^{31,32,26,27} using Rsync in the command line. All gene coordinates

170 were extracted by using awk function (including protein-coding gene and pseudogenes).
171 The downloaded file was in gff3 format and was converted to bed files using BEDOPS.^{28,25}
172 e) Long non-coding RNA annotations (lncRNA) were downloaded from both GENCODE
173 (release 31)³⁰ and Incipedia database (version 5.2).²⁰⁻¹⁸ LncRNA coordinates were merged
174 together using BEDOPS²⁸ BEDOPS²⁵ to obtain all possible lncRNA gene regions.

175 f) Micro RNA (miRNA) annotations were downloaded from miRbase (release 22.1).^{21,19}
176 The downloaded file was in gff3 format and was converted to bed files using BEDOPS.^{28,25}
177

178 ***Genomic coordinates of regions excluded*** were identified from the ENCODE blacklist which was
179 downloaded directly from the ENCODE²⁰ project website in bed format, and RepeatMasker³⁴
180 which was downloaded from the UCSC database^{31,32,26,27} in txt format. The Linux command awk
181 was used to grep out region coordinates, and BEDOPS²⁵ was used to sort the file.

182
183 The final size of the GROFFFY filter was 1,423,480,943bp approximating to 44.48% of GRCh38.

184
185 ***To assimilate for GROFFFY***, separately, the Ubuntu shell (version 16.04.2 LTS based on Linux
186 4.4.0-64-generic x86_64 system) was launched for the Genomics England Research
187 Environment,^{10,7} where the final genomics coordinates for GROFFFY were transferred. WGS
188 variant data were examined after analysis by the Illumina WGS Service Informatics pipeline.
189 This used Illumina Issac²⁹ and Starling for sequence alignment, and to identify variants. The
190 output files were in vcf files with .gz compression and decompressed using Gzip (Version 1.6).³⁷
191 Pandas module version 0.22.0³⁸⁻³⁰ in Python (version 3.6.5) was used to process vcf files.

192 Pegasus, the High-Performance computer cluster of Genomics England, was used to run
193 computationally intensive jobs, submitted to the Load Sharing Facility (LSF).

194
195 The Intersect function of Bedtools version 2.26.0³⁹⁻³¹ was then used to identify WGS variants from
196 vcf files that were in the GROFFFY bed file regions. Option -header was used to remove any
197 headers, and option -wa was used to ensure the output file format was the same as the input vcf
198 file. For the Intersect function, any intersection with GROFFFY was outputted to result files, even
199 if some part of the variation was outside of the filter region. Annotations of the WGS vcf files were
200 carried out using the Ensembl Variant Effect Predictor (VEP) version 96.3,⁴⁰⁻³² based on Perl
201 version 5.24,⁴¹ SAMtools version 1.5⁴²⁻⁵³³ specifically SAMtools HTSlib version 1.5,⁴²⁻³³ and a
202 list of options to optimize the process (*Table S2*). A Python script was written to produce 66 shell
203 scripts where each shell script contained 2 annotation jobs. R version 3.5.1,⁴³ within R studio
204 version 13.4.0⁴⁴ downloaded from the Comprehensive R Archive Network⁴⁵ was used to perform
205 statistical tests. Paired datasets were analysed using the non-parametric Mann Whitney (Wilcoxon
206 rank-sum) test, and multiple datasets by the Kruskal-Wallis rank-sum test with post-test Dunn's
207 multiple comparisons.

208
209 Regions from the ENCODE Blacklist,²²⁻²⁰ and RepeatMasker²³ were subtracted from regions
210 selected for using the 'difference' option in BEDOPS.²⁸⁻²⁵ All bed files were merged together to
211 obtain the selected genomic regions. Numeric data (*Table S5*) and Python scripts were approved
212 for export through the Research Environment AirLock under subproject RR42 (HHT-Gene-Stop,
213 *Table S6*).

214

215 **GROFFFY analysis of whole genome sequencing data**

216 As detailed in [Figure S1](#) and [Figure S2](#), stepwise filters excluded variants where general population
217 allele frequency exceeded 0.0002 in the 1000 Genome Project³⁴ or gnomAD³⁵ databases;
218 synonymous variants not in splice regions; all non HHT-causal variants in the Validation Set HHT
219 DNAs; and variants with a Combined Annotation-Dependent Depletion (CADD) score <10.^{48,36}
220 There was no *a priori* reason to follow any specific filter, for example, a CADD<10 does not
221 preclude such a variant being important, but our goal was to prioritise in the context of the current
222 question. In detail:

223 ***An autosomal dominant-specific disease application step*** was included as a high stringency
224 “white list” filter. For this, the annotated WGS files were retrieved for the 34 Validation Set DNAs
225 where a causative variant had already been identified in known HHT genes through clinical
226 pipelines.^{1312-1615,49,37} Variant information was collected through unique variant IDs consisting of
227 chromosome number, variant starting position, reference sequence and altered sequence (e.g.
228 chr1:111_C/TTT). To confirm that no two variants were represented by the same variant ID, the
229 full list was compared to a set where only unique values were stored, and shown to be identical.
230 The variant IDs were integrated in a white-list. Exclusion of these white-listed variants in other
231 patients was performed using the `isin` function of Pandas³⁰ module: Any variant in the white list
232 was deleted from the vcf files of the target set DNAs, and the number of variants after exclusion
233 was recorded and outputted to txt files.

234 ***For CADD score filtration and prioritization***, the plugin option of VEP was used to annotate
235 variants with CADD scores³⁶ in the enclosed Research Environment: databases for SNV
236 annotation (version 1.5) and small indel annotation (version 1.5) which were pre-installed in the
237 Research Environment were indexed. As the annotation of CADD score was quite slow, the

238 process was put towards the end of the analysis pipeline, so that there were fewer variants that
239 needed to be annotated. Prioritization by CADD score was performed by generating further
240 customized Python scripts. The CADD PHRED-scale score for all 9 billion SNVs and millions of
241 small indels was extracted from the information column. Variants absent from the CADD score
242 database were represented by an empty string by default and were replaced by number 999 instead.
243 Variants with a PHRED score less than 10 were removed so that both variants with top 10
244 percentiles deleteriousness and variants absent from the database were prioritized. The processed
245 files were stored in vcf format.

246

247 **Export of variant coordinates and bioinformatic analyses**

248 Following approval for export through AirLock (*Table S6*), variant genomic coordinates were
249 visualized in the UCSC Genome Browser.^{31,32,26,27} Endothelial expression of *SMAD4* was
250 examined in whole transcriptome data from primary human BOECs.⁵⁰⁻³⁸ Binary sequence
251 alignment map (bam) files aligned to GRCh38⁴⁹ were analyzed in Galaxy Version 2.4.1⁵¹⁻¹³⁹ and
252 the Integrated Genome Browser (IGB) 9.1.8.⁵²⁻⁴⁰

253

254 3'UTR alternative polyadenylation quantitative trait loci (3'aQTLs) from 46 tissues isolated from
255 467 individuals in GTEx⁴¹ were sourced through the 3'aQTL Atlas.⁵⁴⁻⁴² Genetic variants likely
256 affecting gene expression in GTEx V8⁴¹ data release were captured from UCSC^{31,32-26,27} CAVIAR
257 tracks, which define high confidence gene expression QTLs within 1MB of gene transcription start
258 sites (cis-eQTLs).^{31,32}

259

260 All variants were independently verified by Genomics England. Impact on microRNA binding
261 sites was examined through TargetScan Human Release 8.0⁵⁵⁻⁰⁴³ and miRDB.⁵⁶⁻⁴⁴ RNA structure
262 predictions were performed using iFoldRNA v2.0^{57,58-45,46} without restraints, and final models were
263 visualized using Mol* Viewer⁴⁷ via the Research Collaboratory for Structural Bioinformatics
264 Protein Data Bank server.⁶⁰⁻⁴⁸

265

266 **Clinical re-contact, correlations, and re-sampling**

267 Genomics England “Contact the Clinician” forms were submitted through the Research
268 Environment (Table S6) and clinicians who had recruited the participants were contacted and
269 joined the research team. Clinical correlations were performed through North Thames and South-
270 West NHS Genomic Medicine Service Alliances. The affected participants were contacted by their
271 clinicians and provided written consent for publication after reviewing the relevant sections of the
272 manuscript.

273

274 Two of these participants also consented to further blood samples together with 3 healthy
275 volunteers and a further unsolved patient-individual recruited to the 100,000 Genomes Project with
276 a JPHT *SMAD4*^{+/+} [MIM: 175050] clinical phenotype. This study was approved by the East of
277 Scotland Research Ethics Service (EoSRES: 16/ES/0095), and the 6 participants provided written
278 informed consent. Using methods we have developed to perform experimental treatments on
279 human cells while resuspended in endogenous plasma,⁶¹⁻⁶³⁴⁹⁻⁵¹ peripheral blood mononuclear cells
280 (PBMCs, ‘monocytes’) were prepared using BD Vacutainer® CPT™ tubes (Bunzl Healthcare,

281 Coalville, UK) according to manufacturer's instructions with minor modifications. As detailed
282 further in the Appendix, these were to provide comparative resources of cells in stressed and
283 unstressed states, where alternate transcripts/exon region use might be impacted by modified
284 efficiency of final AAUAAA cleavage and polyadenylation due to the 3'UTR variants.

285

286 Briefly, immediately after venesection, the blood was gently remixed by inverting 8-10 times, and
287 centrifuged within 2 hours of collection for 30 minutes at 1600 relative centrifugal force (RCF) at
288 room temperature. The PBMC-containing buffy coat and plasma were collected by pipetting from
289 above the gel layer, transferred to a single 50ml tube for each donor, and gently inverted to
290 resuspend. After PBMC resuspension in plasma, for each donor, equal volumes were distributed
291 to separate experimental treatment tubes, prewarmed at 37°C for 10 minutes, then subjected to 4
292 different treatment conditions for 1 hour including control at 37°C, and low temperature in a 32°C
293 waterbath for 1hr to mimic the stress incurred at the threshold between mild and moderate
294 hypothermia.⁶⁴⁻⁵² Additional stresses previously optimized in our laboratory^{61-63,49-51,65,66,53,54} were
295 inhibition of translation by cycloheximide 100µg/ml (cycloheximide inhibits eukaryotic
296 translation elongation by mechanisms include binding to the 60S ribosomal subunit E-site^{67,68,55,56})
297 and a clinically-relevant mild reactive oxygen species (ROS) stress using ferric citrate
298 10µmol/L.^{66,69,54,57} After 1hr, all tubes were centrifuged at 520 RCF at room temperature for 15
299 minutes. Cell pellets were lysed in Tri reagent (Cambridge Bioscience Ltd, Cambridge, UK)
300 before distribution to replicate tubes for paired rRNA-depleted and polyA-selected RNA
301 sequencing library generation.

302

303 **RNA Sequencing and Differential Expression Analyses**

304 RNA extraction and quality control for 96 samples was performed by Genewiz (Leipzig,
305 Germany). For RNASeq library preparations, 48 samples were polyA selected for polyadenylated
306 RNA enrichment, and 48 paired samples underwent ribosomal (r)RNA depletion. RNA was
307 fragmented and random primed for first and second strand cDNA synthesis, end repair, 5'
308 phosphorylation, dA-tailing, adaptor ligation, PCR enrichment and Illumina HiSeq sequencing
309 using paired-end 150bp reads (Genewiz, Leipzig, Germany). Sequenced reads were trimmed using
310 Trimmomatic v.0.36,⁷⁰⁻⁵⁸ aligned to Homo sapiens GRCh38⁴⁹ using STAR aligner v2.5.2b, and
311 unique gene reads that fell within exon regions counted using Subread package v1.5.2 (Genewiz,
312 Leipzig, Germany).

313

314 Blinded to the types of donors and treatments, Genewiz performed differential gene expression
315 analyses using DESeq2,⁷⁴⁻⁵⁹ and differential exon expression using ~~DEXSeq~~⁷²-~~DEXSeq~~⁶⁰ to
316 identify differentially spliced genes by testing for significant differences in read counts on exon
317 regions (and junctions) of the genes. In DEXSeq,⁷²⁻⁶⁰ read counts are normalised by size factors:
318 Contributions to the average are weighted by the reciprocal of an estimate of their sampling
319 variance, and the expected variance used to derive weights for the "balanced" coefficients reported
320 as estimates for the strengths of differential exon usage and DEXSeq plotting, that are of similar
321 magnitude to the original read counts.⁷²⁻⁶⁰ The output indicates alternative transcript isoform
322 regulation, noting individual exon region assignment is reliable as long as only a small fraction of
323 counting regions (bins) in the gene is called significant.⁷²⁻⁶⁰

324

325 Noting control variability in initial DESeq2 analyses (Figure S3), the least variable of human
326 transcripts (the 25 genes with GINI Coefficients (GCs) <0.15 in diverse cells^{61,62}) were used to

327 evaluate individual library quality (Figure S4), and subsequently employed for DESeq2
328 normalisation. For these normalisations, the intra-assay coefficient of variation (CV,
329 $100 * \text{standard deviation (SD)} / \text{mean}$)⁷⁵⁻⁶³ was calculated for replicate pairs using alignment per gene
330 adjusted for total read counts per library, and analyses restricted to libraries where >50% of GINI
331 genes had a CV < 10% ('met CV10'). Three rRNA depletion datasets failed this quality control.
332 The remaining datasets were ~~DeSeq2~~⁷¹-~~DeSeq2~~⁵⁹ normalised using the GINI^{73,74}-^{61,62} genes as
333 housekeepers: For this, the ratio of alignment counts for each selected housekeeper gene in each
334 dataset to the geometric mean of that gene was calculated across the remaining 45 datasets from
335 rRNA-depleted libraries. The median value of these ratios in each library was used to generate the
336 'size factor' to scale that library's alignments (Table S7).

337

338 **Statistics**

339 Descriptive statistical analyses were performed using Python, and STATA v17.0 (Statacorp,
340 College Station, Texas). Comparative statistics of the number of variants before and after
341 filtration was performed using Mann Whitney two-group comparisons. RNASeq expression was
342 analysed in STATA v17.0 (Statacorp, College Station, Texas) and GraphPad Prism 9 (GraphPad
343 Software, San Diego, CA), compared using Kruskal Wallis and Dunn's post test applied for
344 selected pairwise comparisons.

345

346

347 **RESULTS**

348

349 **GROFFFY defines biologically validated regions of functionality**

350 By only including regions where biological experiments have generated evidence in favor of
351 functional roles, GROFFFY essentially excludes biologically less important regions of the
352 genome. Nevertheless, the GROFFFY filter region based on positive selection of transcribed loci
353 and candidate regulatory element (cREs), and masking of repetitive regions, included 44.4% of
354 the human genome. A heatmap at 500kb resolution is provided in [Figure 2A](#). A more detailed
355 view of GROFFFY is provided in [Figure 2B](#).

356
357

358 **GROFFFY substantially reduces the number of DNA variants per DNA**

359 The scale of the bioinformatics challenge was emphasized by the pre-filtration number of DNA
360 variants per individual which ranged from 4,786,039 to 5,070,340 (mean 4,867,167). Applying
361 GROFFFY as a first filter reduced the mean number of variants by 2,812,015 ([Figure 3A](#), [Figure](#)
362 [3B](#)). Restricting to rare variants with population allele frequencies $<2 \times 10^{-4}$ ⁴⁶⁻³⁴ removed means of
363 2,476,589⁴⁶⁻³⁴ and 2,483,377⁴⁷⁻³⁵ variants/DNA according to database ([Figure 3A](#), [Figure 3B](#)).
364 After removing variants with a [CADD](#)⁴⁸-[CADD](#)³⁶ score <10 , the mean number of unique, rare and
365 impactful DNA variants per genome was 21,486 ([Figure 3C](#)).

366

367 GROFFFY did not delete key variants, as shown by the validation dataset: All already-known
368 pathogenic variants in the unfiltered dataset were retained post filtration ([Figure 4A](#)). Further, in
369 the discovery set of 98 whole genomes, for *ACVRL1* and *ENG*, the majority of identified ~~novel~~
370 variants clustered to the exons and flanking regions sequenced in clinical diagnostics ([Figure 4B](#)).

371

372 **Hot spot of rare deletion variants in the distal *SMAD4* 3' untranslated region**

373 No coding *SMAD4* variants were identified in the discovery dataset (Figure 4B). We focused on a
374 hot spot of 3 deletion variants in the 3' untranslated region (UTR) of *SMAD4* (Figure 4B). There
375 were two unique variants, one of which was identified in both affected members of a single family.
376 The variants deleted nucleotides 5,519 and 5,649bp distal to the *SMAD4* stop codon, and did not
377 affect any microRNA binding sites.^{55,5643,44} The wild-type sequences were consistently expressed in
378 human primary blood outgrowth endothelial cells (BOECs) derived from donors with normal
379 ~~*SMAD4*⁵⁰~~-~~*SMAD4*³⁸~~ (Figure 5A). General population common variant data also supported the
380 importance of the region: while the 3' UTR did not contain any expression quantitative trait loci
381 (eQTLs)⁵³⁴¹ (Figure 5B), the variants were within the only kilobase of the 3'UTR to contain 3'
382 alternate polyadenylation QTLs (3' aQTLs,⁵⁴⁻⁴² Figure 5C).

383

384 **Variants delete nucleotides near the final *SMAD4* alternate polyadenylation site**

385 The *SMAD4* UTR used by all coding transcripts contains 7 alternate polyadenylation site (PAS)⁷⁷⁻⁶⁵
386 AAUAAA hexamers. These cluster in two proximal groups of 3, before a single final AAUAAA
387 hexamer at chr18:51,083,977 (Figure 5A). This final hexamer lay immediately proximal to the two
388 deletion variants, and as expected,⁷⁷⁻⁶⁵ was flanked by an upstream AU-rich element suited to binding
389 of proteins in the cleavage and polyadenylation (CPA) complex, and downstream repeat elements
390 predicting intermolecular interactions in single stranded RNA that would generate secondary
391 structures to block the progress of RNA polymerase II (Figure 6A). One GROFFFY-filtered variant
392 deleted the 16 nucleotides sited +3 to +18 from the PAS hexamer with 5 further single nucleotide
393 substitutions, and the second deleted 4 nucleotides in the downstream repetitive element region
394 (Figure 6B, Table S8).

395

396 **The deletion variants disrupt RNA secondary structures required for cleavage and**
397 **polyadenylation**

398 iFoldRNA secondary structures^{57,58,45,46} visualized using Mol* ~~Viewer~~⁵⁹-~~Viewer~~⁴⁷ via the Research
399 Collaboratory for Structural Bioinformatics Protein Data Bank server,⁶⁰~~48~~ indicated that both
400 deletion variants disturbed secondary structures that substantially altered the sequence context for
401 CPA activity. In wildtype sequence, the AAUAAA hexamer was in a near-linear conformation
402 with stacked pyrimidine and purine rings evident on magnified views (Figure 7Ai, Movie S1).
403 Strikingly, with the neighboring complex deletion variant, the AAUAAA nucleotides acquired
404 new inter-molecular interactions, lost the stacked alignment of bases, and were incorporated into
405 inaccessible secondary structures (Figure 7Aii, Movie S2). In contrast, the second variant which
406 deleted 4 nucleotides 134bp downstream of the AAUAAA hexamer, disrupted and unfolded the
407 downstream structured region expected to be the major RNA polymerase II
408 ~~roadblock~~⁷⁷~~roadblock~~⁶⁵(Figure 7B).

409
410 **Clinical correlation**

411 All three ~~patients~~~~individuals~~ with Variants 1 and 2 had clinically-confirmed HHT.^{44,10-18,17,79}~~67~~
412 After identification of the *SMAD4* variants, recruiting clinicians also reported *SMAD4*-compatible
413 clinical phenotypes: The first-degree relatives with Variant 2 had no other identified cause to HHT.
414 They each experienced daily nosebleeds, had classical HHT telangiectasia, and one had pulmonary
415 arteriovenous malformations requiring treatment, and hemihypertrophy (left-right axis defect).
416 Gastrointestinal and aortopathy screening had not been considered. The ~~patient~~~~individual~~ with
417 Variant 1 did have a missense variant in *ACVRL1*, though in addition to severe nosebleeds needing
418 blood transfusion and intravenous iron, classical HHT telangiectasia and pulmonary arteriovenous

419 malformations, they had multiple colonic and rectal polyps requiring excision over a 6 year period
420 of observation.

421

422 **Peripheral blood mononuclear cell *SMAD4* RNA Expression**

423 As described in the Data Supplement, peripheral blood mononuclear cells (PBMCs) were isolated
424 from affected individuals with the 3 *SMAD4* variants and controls, and cultured in conditions
425 predicted to modify 3'UTR use, before RNA sequencing. DESeq2⁵⁹ analyses of PolyA-selected
426 RNAs indicated that *SMAD4* polyadenylated transcripts increased after a 1hr hypothermic stress,
427 and this was also seen in 2 individuals with the *SMAD4* variants (*Figure S3*). However, for the
428 rRNA-depleted libraries representing “total” RNA, variability between control samples assessed
429 by initial DESeq2 analyses high (*Figure S3*). This reduced after normalising with low GINI^{73,74}
430 ^{61,62}coefficient genes (*Figure S4*).

431

432 Whether normalised to read counts per library, or GINI^{73,74}_{61,62}genes, total *SMAD4* RNA
433 expression was lower in the “inaccessible AAUAAA” Variant 1 donor than 3 controls (*Figure*
434 *8Ai*). Decrements were also apparent in untreated PBMC exon regions by DEXSEQ⁶⁰ (*Figure*
435 *8Bi*). In controls, *SMAD4* transcript expression was modified following 1hr cycloheximide
436 100µg/mL, with lower use of exon region (ER)60 containing the AAUAAA site and variants,
437 consistent with shorter 3'UTRs after stress (*Figure S5, Figure S6*). Despite this, ER60 use was
438 further reduced in the Variant 1 donor after CHX (*Figure 8Ci, Figure 8Di*) with increased use of
439 penultimate exon regions ER52-55 (*Figure 8Di, Figure 78Ei*), supporting different 1hr changes in
440 RNA splicing on stress.

441

442

443 Total *SMAD4* RNA was higher in the “roadblock unfolding” Variant 2 donor than 3 normal
444 controls across all conditions (Figure 8Aii). Although exon region use was similar to controls in
445 untreated PBMCs (Figure 8Bii, Figure 8Cii), after 1hr cycloheximide, compared to controls there
446 was higher use of regions corresponding to two of the 3’ aQTLs (Figure 8Dii, Figure 8Eii). We
447 concluded that the contrasting overall expression patterns were consistent with the opposing
448 predictions following RNA modelling of Variants 1 and 2; that Variant 2 data also supported
449 different 1hr changes in RNA splicing after CHX stress, but that precise transcript changes would
450 need to be the subject of future RNA studies.

451

452

453 **Validation of positive control variant.**

454 The third donor had been recruited as a *SMAD4* positive control due to HHT-JP syndrome (colonic
455 and gastric polyposis; HHT nosebleeds; HHT mucocutaneous telangiectasia, pulmonary AVMs
456 treated by embolization, and antecedent HHT-JP family history). However, no *SMAD4* variant had
457 been identified by clinical service panel testing, the 100,000 Genomes Project clinical pipelines,
458 or by GROFFFY. Total PBMC *SMAD4* expression levels were lower than controls (Figure 8Aiii),
459 and similar to Variant 1 (Figure 8Ai/iii) with additional similarities to Variants 1 and 2 post
460 cycloheximide (Figure 8E). A new team member blinded to the findings and project, was invited
461 to examine the raw *SMAD4* WGS data in the binary alignment map (bam) file, and identified a
462 single exonic variant in the donor’s DNA (Figure S7). This was sited between Variants 1 and 2 in
463 the 3’UTR, with the complex insertion/rearrangement separated from Variant 1 by only two bases
464 (Figure S7).

465

466

467 **DISCUSSION**

468 We have presented and validated a system that synthesizes biologically-generated signals of
469 function in order to filter out variants in DNA regions with no such evidence of functionality. This
470 generically applicable method was highly effective in reducing the number of WGS variants from
471 almost 5 million per individual to an average of ~21,000. Critically, the method retained
472 pathogenic variants already known in a validation dataset, and identified ultra-rare, disease-
473 associated variants in the distal *SMAD4* 3' UTR. These variants disrupted RNA secondary
474 structures required for cleavage and polyadenylation, and subsequent RNASeq and clinical
475 correlations supported *SMAD4* etiology.

476

477 Study strengths include the development and application of an unbiased, genome-wide method
478 with no prior assumptions. Of other variant filtration methods already used in WGS, most depend
479 on union and intersection rules of existing annotation tracks. The candidate cis-regulatory elements
480 file produced by ENCODE has been particularly favored with its specific predictions of each
481 possible CRE position and size. By using the raw biological data providing broader areas for
482 inclusion, GROFFFY may better suit the purpose of a first pass filter for definition of variants
483 worthy of further study, than computational predicted files with potential false negatives.
484 Simultaneous evaluation of WGS data from nearly 100 individuals with a similar clinical
485 phenotype enabled resource direction to unstudied non-coding sequences where multiple rare, high
486 impact variants were identified. Study strength was further augmented by replicate RNASeq
487 expression data from primary human endothelial cells, the cell type responsible for the *SMAD4*
488 clinical phenotype (HHT) where causal loss-of-function variants were being sought, and Genomics

489 England clinician contact pathways that identified *SMAD4*-specific phenotypes after the draft
490 manuscript was approved for submission. This also enabled recontact, leading to evidence from
491 sequenced individuals' PBMCs that support perturbations in *SMAD4* RNA expression. We do not
492 expect the PBMC responses to be a complete model of the variant effects in all pathobiological
493 contexts, but they are presented in order to provide functional evidence of- molecular impact.
494 Additionally, extensive open-source datasets and code enabled exploration of common human
495 variation responsible for *SMAD4* QTLs containing exons where expression was impacted by the
496 identified variants, while the variants themselves highlighted an emerging field in biology that has
497 had limited recognition in medicine.

498

499 A potential study weakness, the presented discovery elements that focus on a single gene, can be
500 justified because of the immediate pathway to translational impact. In addition to somatic cancer
501 genetic diagnostics, early diagnosis of a germline heterozygous *SMAD4* loss-of-function allele
502 offers proven methods to save lives and emergency healthcare resources by institution of
503 gastrointestinal (from adolescence) and aortic screens,⁹ in addition to standard HHT screening and
504 pre-symptomatic interventions.^{10,11,12} It is not possible to perform further segregation analyses in
505 these families as all known affected relatives were in the antecedent generations and deceased.
506 Two of these ultra-rare variants have been detected previously (rs1599209874, absent in gnomAD,
507 TOMMO MAF of 0.00006; rs1375437193, gnomAD MAF of 0.000071): Since the phenotypes
508 are late onset, it is very likely that such variants could be identified in members of the population
509 who did not yet have a clinical diagnosis. Thus, while detailed mechanistic dissection can be the
510 subject of future work, we suggest the presented data support immediate extension of the *SMAD4*
511 regions included in biological and virtual gene panels for patients with HHT, juvenile polyposis

512 and cancer to include the 3'UTR sequences flanking the final AAUAAA hexamer. For the HHT
513 patients harboring the identified variants, there seems sufficient evidence for them to be considered
514 as “likely *SMAD4* HHT” for at least one round of endoscopic and echocardiographic surveillance,
515 while further functional studies are pending For other HHT patients where conventional screening
516 of HHT genes has not identified a causal variant, the possibility of undetected *SMAD4* variation
517 can be considered.

518

519 Alternate polyadenylation has not been explored to date for *SMAD4*, or for other heritable diseases
520 beyond triplet expansion neurodegenerative diseases^{65,81-69} Long 3'UTRs with their abundance of
521 regulatory motifs provide greater opportunity for regulatory control than short 3' UTRs, while
522 switching between alternate polyadenylation sites to provide shorter or longer 3'UTRs is
523 increasingly recognized to modify protein translation, for example differentially transporting
524 mRNAs to condensates which can result in translation repression or enrichment in specified
525 cellular regions or states.⁷⁷⁻⁶⁵ Our data suggest this will be important for regulation of *SMAD4*, a
526 ubiquitous and essential protein with diverse functions,^{8,9} where ~7kb of 3' UTR is transcribed at
527 high levels in coding and non-coding transcripts (Figure 5, Figure 8, Fig. S5, Fig. S6). Recent data
528 highlight that polyadenylation sites differ in strength: weaker proximal CPA sites are used in genes
529 with cell type-specific transcription, (requiring transcriptional enhancers to strengthen CPA
530 activity), while distal and single PAS sites are strongest to ensure mature mRNAs are produced.⁸²
531 ⁷⁰As recently reviewed,⁷⁷⁻⁶⁵ cleavage and polyadenylation occurs while RNA polymerase II (Pol
532 II) is transcribing a gene, and is regulated by Pol II elongation dynamics. Pol II pausing
533 immediately downstream to a final AAUAAA hexamer CPA cleavage site is necessary in order to
534 enable CPA complex assembly and co-transcriptional addition of the “poly-A tail” that is essential

535 for mRNA generation and subsequent protein translation (Figure 6). If at the final polyadenylation
536 site, the full cycle of polymerase pausing, CPA complex binding and cleavage/polyadenylation is
537 impaired, different sites and efficiency of polyadenylation would modify function. Our current
538 data examining 1hr stress responses when the cell has to rely predominantly on reuse of existing
539 RNA transcripts, highlight further mechanisms to explore. These include 3' UTR variant impacts
540 on alternate splice site selection, and maintenance of polyadenylated transcripts that may be less
541 successfully achieved in the setting of stress conditions necessitating rapid changes (Figure S3A).
542 The potential to facilitate future development of 3'UTR therapeutics is augmented given repetitive
543 regions of pol II "roadblocks" provide fertile and previously hidden substrates for impactful human
544 DNA variation,

545
546 In conclusion, we present and validate a filter that reduces the overwhelming number of variants
547 identified by WGS, while retaining functional genome variation of importance to patients.
548 Exposure of non-coding variants in the top 10 percentile of deleteriousness, and clusters in
549 unexplored genomic regions, enhances the near-term value of WGS. The GROFFFY filter enabled
550 identification of rare *SMAD4* variants that disrupt the final site for RNA cleavage and
551 polyadenylation, necessary for protein production. However, the full extent to which rare stress
552 impact, functional alternate polyadenylation site (SIFAPS) variants contribute to diseases will
553 only be exposed if untranslated sequences spanning the sites are included in virtual and physical
554 diagnostic gene panels. Wider use of WGS, and inclusion of 3'aQTL UTR regions in exome-
555 based sequencing are recommended to capture relevant disease-specific variants.

556
557

558

559 **Web Resources section**

560 **Genome Reference Consortium Human Build 38:**

561 https://www.ncbi.nlm.nih.gov/datasets/genome/GCF_000001405.40/

562 **GENCODE Human Genome Release 31:** <https://www.gencodegenes.org/human/>

563 **LINUX:** <https://opensource.com/resources/linux>

564 **Online Mendelian Inheritance in Man:** <https://www.omim.org/>

565 **Perl 5.24:** <https://docs.activestate.com/activeperl/5.24/get/relnotes/>

566 **R: The R Project for Statistical Computing:** <https://www.r-project.org/>

567 **RStudio:** <https://www.rstudio.com/>

568 **RepeatMasker:** <http://www.repeatmasker.org>

569 **The Comprehensive R Archive Network:** <https://cran.r-project.org/>

570 **The National Genomics Research and Healthcare Knowledgebase v5 (2019) Genomics England.**
571 [doi:10.6084/m9.figshare.4530893.v5.](https://doi.org/10.6084/m9.figshare.4530893.v5)

572

573 **Appendix**

574 Supplementary Data include a single file of Supplementary Methods, 7 Supplementary Figures,
575 and 6 Supplementary Tables; 2 Supplementary Tables provided as separate resource files; and 2
576 Movies.

577

578

579 **Declaration of interests**

580 The authors declare no competing interests

581

582 **Acknowledgments:**

583 This research was made possible through access to the data and findings generated by the 100,000
584 Genomes Project. The work was cofounded by the National Institute for Health Research Imperial
585 Biomedical Research Centre, the D’Almeida Charitable Trust, and Imperial College Healthcare
586 NHS Trust. AA was supported by Prince Sultan Military Medical City, Saudi Arabia. MAA was
587 supported by the National Institutes of Health (grant R35HL140019). The 100,000 Genomes
588 Project is managed by Genomics England Limited (a wholly owned company of the Department
589 of Health and Social Care). The 100,000 Genomes Project uses data provided by patients and
590 collected by the National Health Service as part of their care and support. We thank the National
591 Health Service staff of the UK Genomic Medicine Centres and the participants for their willing
592 participation; the Genomics England Clinical Research Interface team, specifically Susan Walker,
593 for separately reviewing bam file variant sequences; Charlotte Bevan, Michael Hubank and
594 Santiago Vernia for helpful discussions and manuscript review; and our academic and public
595 partners within the NIHR Imperial BRC’s Social Genetic and Environmental Determinants of
596 Health (SGE) theme. We specifically thank the presented families for confirmation of their clinical
597 phenotypes and consent to share in this manuscript. The views expressed are those of the authors
598 and not necessarily those of funders, the NHS, the NIHR, or the Department of Health and Social
599 Care.

600

601

602 **Author contributions:**

603 Conceptualization: SX, CLS

604 Methodology: SX, ZK, DP, AB, MEB, AA, MAA, NV, MJC, GERC, CLS

605 Investigation: SX, ZK, ADM, SW, CLS

606 Visualization: SX, CLS

607 Funding acquisition: SX, AA, ADM, SW, MAA, MJC, CLS

608 Project administration: GERC, CLS

609 Supervision: DP, MEB, MAA, CLS

610 Writing – original draft: CLS

611 Writing – review & editing: SX, ZK, DP, AB, MEB, AA, ADM, SW, MAA, NV, MJC,

612 GERC, CLS

613

614 SX devised and generated the GROFFFY approach, devised all scripts to generate GROFFFY,
615 and generated all GROFFFY numeric data, Figures 1, 2 and 3, Figures S1 and S2, and Tables S1,
616 S2, S3, S4 S5, and S6. ZK advised on Linux and script generation. DM interrogated Donor 3
617 bam files. DL assisted in PBMC cultures. DP, AB, MBH, and MAA performed BOEC cultures
618 and RNA sequencing. AA designed primers for validations. AM contributed to patient
619 recruitment. SW contributed to clinical correlations. NV advised on *SMAD4* regulation. GERC
620 performed all sequencing. MJC contributed to specific project set up at Genomics England. CLS
621 recruited patients and performed clinical correlations; devised concepts and advised on
622 GROFFFY approaches; devised and performed PBMC cultures; devised and performed in-house
623 endothelial and PBMC RNASeq and variant level data analyses; generated Figures 4, 5, 6, 7 and

624 8, Figures S3, S4, S5, S6 and S7, Tables S6, S7 and S8, and wrote the manuscript. All authors
625 have reviewed and approved the final manuscript.

626

627 **Data and code availability:**

628 The publicly available file accession numbers used to generate the code are provided in full
629 within the Data Supplement and have been submitted to the NCBI BioProject database
630 (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number PRJNA596860, referencing
631 the WGS data source under accession number SAMN13640532. Primary data from the 100,000
632 Genomes Project, which are held in a secure Research Environment, are available to registered
633 users. Please see [https://www.genomicsengland.co.uk/about-gecip/for-gecip-members/data-and-](https://www.genomicsengland.co.uk/about-gecip/for-gecip-members/data-and-data-access)
634 [data-access](https://www.genomicsengland.co.uk/about-gecip/for-gecip-members/data-and-data-access) for further information.

635

636

637

638 **References and Notes**

639 1 Ransohoff, J.D., Wei, Y., and Khavari, P.A. (2018) The functions and unique features of
640 long intergenic non-coding RNA. *Nat. Rev. Mol. Cell. Biol* 19, 143—157.

641 2 Marchal, C., Sima, J., and Gilbert, D.M. (2019). Control of DNA replication timing in the
642 3D genome. *Nat. Rev. Mol. Cell. Biol.* 20(12), 721—737.

643 3 Halley, M.C., Ashley, E.A., and Tabor, H.K. (2022). Supporting undiagnosed participants
644 when clinical genomics studies end. *Nat. Genet.* 54, 1063—1065

645 4 Sloan, C.A., Chan, E.T., Davidson, J.M., Malladi, V.S., Strattan, J.S., Hitz, B.C., Gabdank,
646 I., Narayanan, A.K., Ho, M., Lee, B.T. et al. (2016). ENCODE data at the ENCODE portal.
647 *Nucleic Acids Res.* 44(D1), D726—32

648 5 Davis, C.A., Hitz, B.C., Sloan, C.A., Chan, E.T., Davidson, J.M., Gabdank, I., Hilton, J.A.,
649 Jain, K., Baymuradov, U.K., Narayanan, A.K. et al. (2018) The Encyclopedia of DNA elements
650 (ENCODE): data portal update. *Nucleic Acids Res.* 46(D1), D794—D801, accessed for data
651 25.07.2019

652 6 Harrow, J., Denoeud, F., Frankish, A., Reymond, A., Chen, C.K., Chrast, J., Lagarde, J.,
653 Gilbert, J.G., Storey, R., Swarbreck, D. et al. (2006). GENCODE: producing a reference
654 annotation for ENCODE. *Genome Biol.* 7 Suppl 1 S4,1—9.

655 7 The National Genomics Research and Healthcare Knowledgebase v5, Genomics England.
656 doi:10.6084/m9.figshare.4530893.v5. 2019

657 8 Hayashi, A., Hong J. and Iacobuzio-Donahue C.A. (2021). The pancreatic cancer genome
658 revisited. *Nat Rev Gastroenterol Hepatol.* 18(7), 469—481.

659 9 Larsen Haidle, MacFarland, J.S.P. and Howe, J.R. (2022) Juvenile Polyposis Syndrome.
660 2003 May 13 [updated 2022 Feb 3]. In: Adam M.P., Everman D.B., Mirzaa, G.M. editors.
661 GeneReviews® [Internet]. Seattle (WA): University of Washington, Seattle; 1993—2022.

662 ~~10—The National Genomics Research and Healthcare Knowledgebase v5, Genomics England.~~
663 ~~doi:10.6084/m9.figshare.4530893.v5. 2019~~

664 110 Shovlin, C.L., Buscarini, E., Sabbà, C., Mager, H.J., Kjeldsen, A.D., Pagella, F., Sure, U.,
665 Ugolini, S., Topping, P.M., Suppressa, P. et al. (2022). The European Rare Disease Network for
666 HHT Frameworks for management of hereditary haemorrhagic telangiectasia in general and
667 speciality care. *Eur J Med Genet.* 65(1):104370.

668 [1211](#) Faughnan, M.E., Mager, J.J., Hetts, S.W., Palda, V.A., Lang-Robertson, K., Buscarini, E.,
669 Deslandres, E., Kasthuri, R.S., Lausman, A., Poetker, et al. (2020). Second International
670 Guidelines for the Diagnosis and Management of Hereditary Hemorrhagic Telangiectasia. *Ann*
671 *Intern Med.* 173(12):989—1001.

672 [1312](#) Clarke, J.M., Alikian, M., Xiao, S., Kasperaviciute, D., Thomas, E., Turbin, I., Olupona,
673 K., Cifra, E., Curetean, E., Ferguson, T. et al. (2020) Low grade mosaicism in hereditary
674 haemorrhagic telangiectasia identified by bidirectional whole genome sequencing reads through
675 the 100,000 Genomes Project clinical diagnostic pipeline. *J Med Genet.* 57(12):859—862.

676 [1413](#) Balachandar, S., Graves, T.J., Shimonty, A., Kerr, K., Kilner, J., Xiao, S., Slade, R., Sroya,
677 M., Alikian, M., Curetean, E., et al. (2022) Identification and validation of a novel pathogenic
678 variant in GDF2 (BMP9) responsible for hereditary hemorrhagic telangiectasia and pulmonary
679 arteriovenous malformations. *Am J Med Genet A.* 188(3):959—964.

680 [1514](#) Joyce, K.E., Onabanjo, E., Brownlow, S., Nur, F., Olupona, K., Fakayode, K., Sroya, M.,
681 Thomas, G.A., Ferguson, T., Redhead, J. et al. (2022) Whole genome sequences discriminate
682 hereditary hemorrhagic telangiectasia phenotypes by non-HHT deleterious DNA variation. *Blood*
683 *Adv.* 6(13):3956—3969.

684 [1615](#) Shovlin, C.L., Almaghlouth, F., Alsafi, A., Coote, N.C., Rennie, C.R., Wallace, G.M.F.,
685 Govani, F.S., and Genomics England Research Consortium. (2023) Updates on diagnostic criteria
686 for hereditary haemorrhagic telangiectasia in the light of whole genome sequencing of “Gene
687 Negative” individuals recruited to the 100,000 Genomes Project. *J Med Genet* In Press

688 [1716](#) Sharma L., Almaghlouth F., Mckernan, H., Springett, J., Tighe, H.C., Genomics England
689 Research Consortium and Shovlin C.L. (2023) Iron deficiency responses and integrated

690 compensations in patients according to hereditary haemorrhagic telangiectasia ACVRL1, ENG
691 and SMAD4 genotypes. *Haematologica*. In Press

692 ~~1817~~ Shovlin, C.L., Simeoni, I., Downes, K., Frazer, Z.C., Megy, K., Bernabeu-Herrero, M.E.,
693 Shurr, A., Brimley, J., Patel, D., Kell, L.,- et al. (2020) Mutational and phenotypic characterization
694 of hereditary hemorrhagic telangiectasia. *Blood*. 136(17):1907—1918.

695 ~~19~~ ~~Genome Reference Consortium Human Build 38.~~ ~~Available at~~
696 ~~<https://www.ncbi.nlm.nih.gov/assembly/?term=GRCh38&cmd=DetailsSearch>~~

697 ~~2018~~ Volders, P.J., Anckaert, J., Verheggen, K., Nuytens, J., Martens, L., Mestdagh, P., and
698 Vandesompele, J. (2019). LNCipedia 5: towards a reference set of human long non-coding RNAs.
699 *Nucleic Acids Res.*;47(D1):D135-D139.

700 ~~2119~~ Kozomara, A. Birgaoanu M., Griffiths-Jones S. (2019). miRBase: from microRNA
701 sequences to function. *Nucleic Acids Res.* 47(D1):D155—62

702 ~~2220~~ Amemiya, H.M., Kundaje, A. and Boyle A.P. (2019) The ENCODE Blacklist:
703 Identification of Problematic Regions of the Genome. *Sci Rep.* 9(1):9354

704 ~~2321~~ ~~Smit, A.F.A., Hubley, R., Green, P. (2013-2015) RepeatMasker Open-4.0. Institute for~~
705 ~~Systems Biology RepeatMasker available at: <http://www.repeatmasker.org>, accessed for data~~
706 ~~25.09.2019~~

707 ~~24~~ ~~Encode Experiment search. Available at~~

708 ~~2522~~ Bernstein, B.E., Stamatoyannopoulos, J.A., Costello, J.F., Ren, B., Milosavljevic, A.,
709 Meissner, A., Kellis, M., Marra, M.A., Beaudet, A.L., Ecker, J.R et al. (2010). The NIH roadmap
710 epigenomics mapping consortium. *Nature Biotechnol* 28;1045—8

711 [2623](#) Mundade, R., Ozer, H.G., Wei, H., Prabhu, L., and Lu, T. (2014). Role of ChIP-seq in the
712 discovery of transcription factor binding sites, differential gene regulation mechanism, epigenetic
713 marks and beyond. *Cell Cycle* 13(18):2847—52.

714 [2724](#) Koch, L. (2017) Cancer genetics: A 3D view of genome rearrangements. *Nat Rev Genet*
715 18(8):456.

716 [2825](#) Neph, S., Kuehn, M.S., Reynolds, A.P., Haugen, E., Thurman, R.E., Johnson, A.K., Rynes,
717 E., Maurano, M.T., Vierstra, J., Thomas, S. et al. (2012) BEDOPS: high-performance genomic
718 feature operations. *Bioinformatics* 28(14):1919—20

719 [26](#) Kent, W.J., Sugnet, C.W., Furey, T.S., Roskin, K.M., Pringle, T.H., Zahler, A.M., and
720 Haussler, D. (2002). The human genome browser at UCSC. *Genome Res.* 12(6):996—1006.

721 [3227](#) Nassar, L.R., Barber, G.P., Benet-Pagès, A., Casper, J., Clawson, H., Diekhans, M.,
722 Fischer, C., Gonzalez, J.N., Hinrichs, A.S., Lee, B.T., et al. (2023). The UCSC Genome Browser
723 database: 2023 update. *Nucl. Acid Res.* 51(D1):D1188—D1195.

724 [28](#) Raczy, C., Petrovski, R., Saunders, C.T., Chorny, I., Kruglyak, S., Margulies, E.H.,
725 Chuang, H.Y., Källberg, M., Kumar, S.A., Liao, A., et al. (2013) Isaac: ultra-fast whole-genome
726 secondary analysis on Illumina sequencing platforms. *Bioinformatics* 29(16):2041—3

727 [3629](#) Raczy, C., Petrovski, R., Saunders, C. T., Chorny, I., Kruglyak, S., Margulies, E. H.,
728 Chuang, H. Y., Källberg, M., Kumar, S. A., Liao, A., Little, K. M., Strömberg, M. P., and Tanner,
729 S. W. (2013). Isaac: ultra-fast whole-genome secondary analysis on Illumina sequencing
730 platforms. *Bioinformatics* (Oxford, England), 29(16), 2041–2043

731 [30](#) McKinney W. (2010). *Data Structures for Statistical Computing in Python*: 51—6.
732 Available from: <http://conference.scipy.org/proceedings/scipy2010/mckinney.html>

733 [3931](#) Quinlan, A. R., and Hall, I. M (2010) BEDTools: a flexible suite of utilities for comparing
734 genomic features. *Bioinformatics*; 26(6):841–2.

735 [4032](#) Ensembl Variant Effect Predictor (VEP), available at
736 <http://grch37.ensembl.org/info/docs/tools/vep/index.html>

737 [33](#) Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis,
738 G., Durbin, R., ~~&~~[and](#) 1000 Genome Project Data Processing Subgroup. (2009) The Sequence
739 Alignment/Map format and SAMtools. *Bioinformatics* 25(16):2078—9.

740 [34](#) Zheng-Bradley, X., Streeter, I., Fairley, S., Richardson, D., Clarke, L., Flicek, P., and 1000
741 Genomes Project Consortium. (2017) 1000 Genomes Project Consortium, Alignment of 1000
742 Genomes Project reads to reference assembly GRCh38. *Gigascience*. 6(7):1—8

743 [4735](#) Karczewski, K.J., Francioli, L.C., Tiao, G., Cummings, B.B., Alföldi, J., Wang, Q.,
744 Collins, R.L., Laricchia, K.M., Ganna, A., Birnbaum, D.P. et al. (2020) The mutational constraint
745 spectrum quantified from variation in 141,456 humans. *Nature*. 581(7809):434—443.

746 [4836](#) Rentzsch, P., Witten, D., Cooper, G.M., Shendure, J., and Kircher, M. (2019) CADD:
747 predicting the deleteriousness of variants throughout the human genome. *Nucleic Acids Res*.
748 47(D1):D886—D894.

749 [4937](#) Richards, S., Aziz, N., Bale, S., Bick, D., Das, S., Gastier-Foster, J., Grody, W. W., Hegde,
750 M., Lyon, E., Spector, E., et al (2015). Standards and guidelines for the interpretation of sequence
751 variants: a joint consensus recommendation of the American College of Medical Genetics and
752 Genomics and the Association for Molecular Pathology. *Genet Med*. 17(5):405—24

753 [5038](#) Bernabeu-Herrero, M.E., Patel, D., Bielowka, A., Chaves Guerrero, P., Marciniak S.J.,
754 Nosedá, M., Aldred, M.A., Shovlin, C.L. Heterozygous transcriptional signatures unmask variable
755 premature termination codon (PTC) burden alongside pathway-specific adaptations in blood

756 outgrowth endothelial cells from patients with nonsense DNA variants causing hereditary
757 hemorrhagic telangiectasia. BioRxiv 471269 v2. 27th June 2023. [doi:
758 https://doi.org/10.1101/2021.12.05.471269](https://doi.org/10.1101/2021.12.05.471269)

759 [5139](#) Afgan, E., Baker, D., Batut, B., van den Beek, M., Bouvier, D., Cech, M., Chilton, J.,
760 Clements, D., Coraor, N., Grüning, B.A., et al. (2018) The Galaxy platform for accessible,
761 reproducible and collaborative biomedical analyses: 2018 update, *Nucleic Acids Res.* 46;W1:
762 W537—W544

763 [5240](#) Freese, N.H., Norris, D.C. and Loraine, A.E. (2016) Integrated genome browser: visual
764 analytics platform for genomics. *Bioinformatics.* 32(14):2089—95

765 [5341](#) GTEx Consortium. (2013). The Genotype-Tissue Expression (GTEx) project. *Nat Genet.*
766 45(6):580—5.

767 [5442](#) Li, L., Huang, K. L., Gao, Y., Cui, Y., Wang, G., Elrod, N. D., Li, Y., Chen, Y. E., Ji, P.,
768 Peng, F., et al. (2021) An atlas of alternative polyadenylation quantitative trait loci contributing to
769 complex trait and disease heritability. *Nat Genet.* 53(7):994-1005.

770 [5543](#) McGeary, S. E., Lin, K. S., Shi, C. Y., Pham, T. M., Bisaria, N., Kelley, G. M., and Bartel,
771 D. P. (2019) The biochemical basis of microRNA targeting efficacy. *Science.*
772 366(6472):eaav1741

773 [5644](#) Chen, Y. Wang X. (2020) miRDB: an online database for prediction of functional
774 microRNA targets. *Nucleic Acids Res;* 48(D1):D127—D131.

775 [5745](#) Ding, F., Sharma, S., Chalasani, P., Demidov, V.V., Broude, N.E., and Dokholyan, N.V
776 (2008). Large scale simulations of 3D RNA folding by discrete molecular dynamics: From
777 structure prediction to folding mechanisms. *RNA* 14:1164—1173.

778 [5846](#) Krokhotin, A., Houlihan K., and Dokholyan N.V. (2015). iFoldRNA v2: folding RNA with
779 constraints. *Bioinformatics* 31 2891—2893.

780 [5947](#) Sehnal, D., Bittrich, S., Deshpande, M., Svobodová, R., Berka, K., Bazgier, V., Velankar,
781 S., Burley, S.K., Koča, J., and Rose, A.S. (2021). Mol* Viewer: modern web app for 3D
782 visualization and analysis of large biomolecular structures. *Nucleic Acids Res.* 49(W1): W431—
783 W437.

784 [6048](#) Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov,
785 I.N., and Bourne, P. E (2000). The Protein Data Bank. *Nucleic Acids Res* 28:235-242, accessed at
786 <https://www.rcsb.org/3d-view>.

787 [6149](#) Shi, C. (2012) Developing a new mutational analysis system using hereditary haemorrhagic
788 telangiectasia as a genetic model. Imperial College London MSc Thesis (supervisor CL Shovlin).

789 [6250](#) Li, Y. (2013) Gene mutations and transcripts in hereditary haemorrhagic telangiectasia
790 (HHT). Imperial College London MSc Thesis (supervisor CL Shovlin)..

791 [6351](#) Shurr, A.Y.L., Maurer, C., Turbin, I.G., Bernabeu-Herrero, M.E., Aldred, M., Patel, D.,
792 and Shovlin, C.L. (2019) Addressing the problem of variants of uncertain significance in genetic
793 diagnosis of vascular pulmonary disease: a role for transcript expression in blood monocytes?
794 *Thorax* 74:A152

795 [6452](#) Duong, H. and Patel, G. (2022). Hypothermia. In: *StatPearls Treasure Island (FL):*
796 *StatPearls*.

797 [6553](#) Govani, F.S., Giess, A., Mollet, I.G., Begbie, M.E., Jones, M.D., Game, L., and Shovlin,
798 C. L. (2013) Directional next-generation RNA sequencing and examination of premature
799 termination codon mutations in endoglin/hereditary haemorrhagic telangiectasia. *Mol Syndromol.*
800 4(4):184—96.

801 [6654](#) Mollet, I.G., Patel, D., Govani, F.S., Giess, A., Paschalaki, K., Periyasamy, M., Lidington,
802 E.C., Mason, J.C., Jones, M.D., Game, L., et al. (2016) Low dose iron treatments induce a DNA
803 damage response in human endothelial cells within minutes. *PLoS One*. 11(2):e0147990.

804 [6755](#) Schneider-Poetsch, T., Ju, J., Eyler, D.E., Dang, Y., Bhat, S., Merrick, W.C., Green, R.,
805 Shen, B., and Liu, J.O. (2010) Inhibition of eukaryotic translation elongation by cycloheximide
806 and lactimidomycin. *Nat Chem Biol*. 6(3):209—217.

807 [6856](#) Shen, L., Su, Z., Yang, K., Wu, C., Becker, T., Bell-Pedersen, D., Zhang, J., and Sachs,
808 M.S (2021) Structure of the translating *Neurospora* ribosome arrested by cycloheximide. *Proc Natl*
809 *Acad Sci U S A*. 118(48):e2111862118.

810 [6957](#) Kartikasari, A.E., Georgiou, N.A., Visseren, F.L., van Kats-Renaud, H., van Asbeck, B.S.,
811 and Marx, J.J. (2006) Endothelial activation and induction of monocyte adhesion by
812 nontransferrin-bound iron present in human sera. *FASEB J*. 20(2):353—5.

813 [7058](#) Bolger, A.M. Lohse, M. and Usadel, B. (2014) Trimmomatic: a flexible trimmer for
814 Illumina sequence data. *Bioinformatics* 30(15):2114—20. v0.36

815 [7159](#) Anders S., and Huber, W. (2010) Differential expression analysis for sequence count data.
816 *Genome Biol* 11(10):R106

817 [7260](#) Anders, S., Reyes, A., and Huber, W. (2012) Detecting differential usage of exons from
818 RNA-seq data. *Genome Res*. 22(10):2008—17

819 [7361](#) Wright Muelas, M., Mughal, F., O'Hagan,S., Day, P.J. and Kell, D.B.(2019) The role and
820 robustness of the Gini coefficient as an unbiased tool for the selection of Gini genes for normalising
821 expression profiling data. *Sci Rep*.9(1):17960.

822 [7462](#) O'Hagan,S., Wright Muelas, M., Day, P.J., Lundberg, E., and Kell, D.B. (2018) GeneGini:
823 Assessment via the Gini Coefficient of Reference "Housekeeping" Genes and Diverse Human
824 Transporter Expression Profiles. *Cell Syst.* 6(2):230—244.e1

825 [7563](#) Reed, G.F. Lynn, F. and Meade, B.D. (2002) Use of coefficient of variation in assessing
826 variability of quantitative assays. *Clin Diagn Lab Immunol.* 9(6):1235—9

827 [7664](#) Lek, M., Karczewski, K.J., Minikel, E.V., Samocha, K.E., Banks, E., Fennell, T.,
828 O'Donnell-Luria, A.H., Ware, J.S., Hill, A.J., Cummings, B.B. et al. (2016). Analysis of protein-
829 coding genetic variation in 60,706 humans. *Nature.* 536(7616):285—91

830 [7765](#) Mitschka, S. and Mayr, C. (2022). Context-specific regulation and function of mRNA
831 alternative polyadenylation. *Nat Rev Mol Cell Biol.* 7:1—18.

832 [7866](#) Sayers, E.W., Bolton, E.E., Brister, J.R., Canese, K., Chan, J., Comeau, D.C., Farrell, C.M.,
833 Feldgarden, M., Fine, A.M., Funk, K., et al. (2023). Database resources of the National Center for
834 Biotechnology Information in 2023. *Nucleic Acids Res.* 51(D1):D29—D38.

835 [7967](#) Shovlin, C.L. (2010). Hereditary haemorrhagic telangiectasia: pathophysiology, diagnosis
836 and treatment. *Blood Rev.* 24(6):203—19.

837 [8068](#) Evans, C., Hardin, J., and Stoebel, D. M. (2018) Selecting between-sample RNA-Seq
838 normalization methods from the perspective of their assumptions. *Brief Bioinform.* 19(5):776—
839 792

840 [8169](#) Li, Y., Li, J., Wang, J., Zhang, S., Giles, K., Prakash, T.P., Rigo, F., Napierala, J.S., and
841 Napierala, M. (2022) Premature transcription termination at the expanded GAA repeats and
842 aberrant alternative polyadenylation contributes to the Frataxin transcriptional deficit in
843 Friedreich's ataxia. *Hum Mol Genet.* 31(20):3539—3557.

844 [8270](#) Kwon, B., Fansler, M. M., Patel, N. D., Lee, J., Ma, W., and Mayr, C. (2022) Enhancers
845 regulate 3' end processing activity to control expression of alternative 3'UTR isoforms. Nat
846 Commun. 13(1):2709.

847

848 **The Genomics England Research Consortium Members comprised on 8th May 2022:**
849 Ambrose, J. C. 1 ; Arumugam, P.1 ; Bevers, R.1 ; Bleda, M. 1 ; Boardman-Pretty, F. 1,2 ;
850 Boustred, C. R. 1 ; Brittain, H.1 ; Brown, M.A.; Caulfield, M. J.1,2 ; Chan, G. C. 1 ; Giess A. 1;
851 Griffin, J. N. ; Hamblin, A.1; Henderson, S.1,2; Hubbard, T. J. P. 1 ; Jackson, R. 1 ; Jones, L. J.
852 1,2; Kasperaviciute, D. 1,2 ; Kayikci, M. 1 ; Kousathanas, A. 1; Lahnstein, L. 1 ; Lakey, A.;
853 Leigh, S. E. A. 1 ; Leong, I. U. S. 1 ; Lopez, F. J. 1 ; Maleady-Crowe, F. 1 ; McEntagart, M.1;
854 Minneci F. 1 ; Mitchell, J. 1 ; Moutsianas, L. 1,2 ; Mueller, M. 1,2 ; Murugaesu, N. 1; Need, A.
855 C. 1,2 ; O'Donovan P. 1; Odhams, C. A. 1 ; Patch, C. 1,2 ; Perez-Gil, D. 1 ; Pereira, M. B.1 ;
856 Pullinger, J. 1 ; Rahim, T. 1 ; Rendon, A. 1 ; Rogers, T. 1 ; Savage, K. 1 ; Sawant, K. 1; Scott, R.
857 H. 1 ; Siddiq, A. 1 ; Sieghart, A. 1 ; Smith, S. C. 1; Sosinsky, A. 1,2 ; Stuckey, A. 1 ; Tanguy M.
858 1 ; Taylor Tavares, A. L.1; Thomas, E. R. A. 1,2 ; Thompson, S. R. 1 ; Tucci, A. 1,2 ; Welland,
859 M. J. 1 ; Williams, E. 1 ; Witkowska, K. 1,2 ; Wood, S. M. 1,2; Zarowiecki, M. 1 .

860

861 1. Genomics England, London, UK

862 2. William Harvey Research Institute, Queen Mary University of London, London, EC1M 6BQ,

863 UK.

864

865

866 **FIGURE LEGENDS**

867

868

869 **Fig. 1 GROFFFY Study Protocol:** Flow chart illustrating sequence of stages described in the
870 text and online Data Supplement

871

872

873 **Figure 2: GROFFFY and the human genome.** A) Heatmap displaying GROFFFY

874 categorization of Genome Reference Consortium Human [GRCh] Build 38.⁴⁹ The heatmap maps
875 61,799 data points at 500kb resolution, and heights represent the percentage of each 500kb region

876 included in GROFFFY. B) Higher resolution image from a randomly chosen region of the genome

877 (on chromosome 3: chr3:25,597,986-2-25,783,443). The top 4 tracks illustrate sources, from top:

878 GENCODE⁶ gene annotations, CpG islands,^{31,32,26,27} long non-coding RNAs^{18,30} and miRNAs.^{24,19}

879 The lowest track illustrates the final filter. Note that this filter contains both intra and intergenic
880 regions for the region, and that the raw data were not subjected to any processed annotation tracks.

881

882 **Figure 3: Application of GROFFFY to whole genome sequences.** A) Serial application of

883 GROFFFY; allele frequency filters based on frequencies in the 1000 Genomes (1000g)⁴⁶⁻³⁴ or
884 gnomAD;⁴⁷⁻³⁵ synonymous (Synon.) filter, and white-listed filter (see Methods and *Tables S1-S4*

885 for further details). Where error bars are not visible at the illustrated scale, exact numeric data are

886 provided in *Table S5*. B) Number of variants remaining per DNA after applying each comparator

887 filter set. C) Number, site and type of DNA variants present in 98 human whole genomes before

888 and after application of GROFFFY and other filters, scaled in one dimension (black bars) and two

889 dimensions (blue circles). CADD, combined annotation dependent depletion score where >10
890 represents a variant in the top 10% of deleteriousness.^{48,36} Irrespective of other filters applied,
891 GROFFFY, and its individual components, significantly reduced the number of variants compared
892 to the other tested filter sets (*Figure S1, Figure S2*).

893

894

895

896

897 **Figure 4: GROFFFY variant-level validation:** Comparison of HHT gene variants in the validation
898 and discovery datasets before and after application of GROFFFY. **A)** Validation dataset: **i)** Total
899 number of variants with indicated CADD scores (note logarithmic scale pre filtration versus linear
900 scale post filtration). **ii)** Non-pathogenic variants by CADD score categories: Molecular subtype are
901 indicated in the key. **iii)** Pathogenic variants by CADD score categories, and molecular subtype as in
902 key. Note identical plots in **iii)** pre and post filtration because all pathogenic variants were still present
903 post filtration. **B)** Discovery dataset: Location of GROFFFY-captured variants in the major HHT
904 genes¹⁸⁷ -**i)** *ACVRL1*, **ii)** *ENG*, **iii)** *SMAD4*. The cartoons include screenshots of GRCh38 from the
905 University of California Santa Cruz (UCSC) Genome Browser,^{31,32} and major transcripts. Red inverted
906 triangles indicate location of variants after application of all filters (dark red for coding/splice regions,
907 bright red for non-coding regions).

908

909

910

911

912 **Figure 5: Expression of *SMAD4* in primary human cells**

913 **A)** Endothelial *SMAD4* total RNA expression: RNASeq data from 8 different cultures of blood
914 outgrowth endothelial cells (BOECs) with normal *SMAD4* sequence.⁵⁰⁻³⁸ The consistent peaks
915 sharply define exon boundaries in GRCh38.⁴⁹ Sites of start and stop codons, unique filtered variants
916 (red triangles), and the seven alternate cleavage and polyadenylation sites at c.3121, c.3487, c.3791,
917 c.5186, c.5452, c.5615 and c.7709 are also highlighted. The cartoon below links the RNASeq
918 expression by grey dotted lines to the main (upper) and alternate *SMAD4* RefSeq⁷⁸-RefSeq⁶⁶ splice
919 isoforms that share the final UTR-containing exon, with exons to scale. Blue: coding, grey: non-coding
920 regions. Note although isoform 6 shares the majority of nucleotides with isoforms 1,2 and 3, it does
921 not share the same ribosomal reading frame, and contains a unique penultimate exon with stop codons
922 in all 3 reading frames that enhance fidelity as a non-coding transcript.

923 **B)** Number of general population *SMAD4* expression QTLs (eQTLs⁵⁴) per interval of DNA flanking
924 the TGA stop codon, as listed by USCSC CAVIAR tracks^{31,32,26,27} for data from the Genotype Tissue
925 Expression project (GTEx).⁵³⁻⁴¹ The graphs are centered on the *SMAD4* natural stop codon site
926 (vertical red arrow), with relevant gene loci indicated to scale horizontally above graphs. **i)** Overview
927 of *SMAD4* locus and flanking regions at 10kb intervals. **ii)** Magnified view of penultimate and final
928 exons at 1kb intervals.

929 **C)** Number of general population 3'UTR alternative polyadenylation QTLs (~~3'aQTLs~~⁵⁴3'aQTLs⁴²)
930 per kilobase of DNA flanking the TGA stop codon, as determined in GTEx.⁵³⁻⁴¹ **i)** Overview of
931 *SMAD4* locus and flanking regions at 10kb intervals. **ii)** Magnified view of penultimate and final
932 exons at 1kb intervals.

933

934

935 **Figure 6: Schematic of *SMAD4* 3' UTR variants in the context of RNA function.**

936 **A)** Color-coded nucleotides 7561_7920 of the *SMAD4* main coding transcript NM_005359. These
937 span the final AAUAAA hexamer (red bar) and include the upstream AU-rich (blue/green) region,
938 downstream repetitive elements, and sites of the two variants. Deleted residues are indicated by
939 black bars, missense substitutions as black triangles. For FASTA format sequences, see [Table S8](#).

940 **B)** Variant 1 (chr18:51083986
941 CTTAACGCGCGTGCACGCGCGCGCACACA>CAACGCGCGTGCACGCG and Variant
942 2 (chr18:51084116 AACT>A) in detail. The AAUAAA hexamer is shown by DNA sequence
943 (AATAAA) and highlighted by a pink box, deleted residues by red underline (wildtype) or vertical
944 red line (variants), and missense substitutions as red stars. [Table S8](#) provides further sequence
945 details.

946

947

948 **Figure 7: Replicate iFoldRNA structures:** iFoldRNA^{57,58-45,46} simulations as visualized in Mol*

949 Viewer.^{59,6047,48} **A) Variant 1: i)** Three representative simulations of the wildtype 150 nucleotides
950 spanning the AAUAAA hexamer (light green), selected as those best illustrating the 3 dimensional
951 relationships in two dimensions. The site of the 32 nucleotide deletion/insertion is highlighted in
952 yellow within the upper structure where they are best demarcated. The lower structures provide
953 further simulations highlighting in brighter green, the accessibility of the near-linear AAUAAA
954 sequences. A camera spin is provided in [Movie S1](#). **ii)** Magnified view of five separate simulations
955 of Variant 1 sequence with AAUAAA hexamer site highlighted in purple. All 5 simulations were
956 consistent and showed the AAUAAA hexamer now inaccessible, incorporated into a secondary
957 structure. A camera spin is provided in [Movie S2](#). **B) Variant 2: i/ii)** Two representative

958 simulations of wildtype sequence. The site of the 4 nucleotides deleted in the variant are
959 highlighted in red in panoramic (i) and magnified (ii) views. (iii) Five simulations of the variant
960 sequence.

961

962 **Figure 8: SMAD4 RNA expression in ribosomal (r)RNA-depleted libraries from 3 controls**

963 **compared to affected donors from the 3 separate HHT families: (i) Variant 1, (ii) Variant 2, (iii)**

964 **an unsolved clinical SMAD4 positive control. For preceding methodological data on the rRNA**

965 **depleted libraries, see [Figures S3-S6](#). RNA from peripheral blood mononuclear cells (PBMCs)**

966 **cultured with and without 3 different 1hr stresses ([Figure S3](#)): SMAD4 splice site changes met**

967 **DEXSeq2⁵⁹ significance after cycloheximide (CHX). A) Total SMAD4 RNA from control (grey,**

968 **N=23) and *patient*-SMAD4 variant-affected donors (red: (i) Variant 1 N=7; (ii) Variant 2 N=8; (iii)**

969 **Donor 3, N=7) following DESeq2 normalisation^{59,80-68} using GINI housekeeper genes.^{73,74,61,62}**

970 **Note contrast between (i)/(ii) (Variant 1 and 2 donors), but similarities between (i)/(iii) (Variant 1 and**

971 **Donor 3). B-E) DEXSEQ⁶⁰ splicing patterns across 61 exon regions in 22 SMAD4 exons. B)**

972 **Exon region (ER) use in untreated PBMCs by donor, plotting data from the individual patients**

973 **(red) and the same 3 controls (black). Exons are colour-coded to highlight 3' aQTL loci.⁵⁴⁻⁴² C)**

974 **Use of ER60, the variant-containing 3'UTR region in untreated and CHX-treated PBMCs: note**

975 **again (i)/(iii) similarities. D) The ratio of exon region use between CHX-treated and untreated**

976 **PBMCs, plotted as in B). E) The ratios in the final 8 exons (ER40-ER61) containing all ERs**

977 **differentially used by the Variant 1 and 2 donors after CHX. Each graph is annotated with the**

978 **genomic DNA origins and kb markers (upper bar); sites of the 3' aQTLs (*, see [Figure 5C](#)); and**

979 **variant outlier values (red) that were not accompanied by increased polyadenylated transcripts**

980 **([Figure S3](#)).**

981

