

Guidelines for Human Gene Nomenclature

Elsbeth A. Bruford*^{1,2}, Bryony Braschi¹, Paul Denny¹, Tamsin E.M. Jones¹, Ruth L. Seal^{1,2},
Susan Tweedie¹

¹ HUGO Gene Nomenclature Committee (HGNC), European Molecular Biology Laboratory, European
Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK. ² Department of
Haematology, University of Cambridge School of Clinical Medicine, Cambridge Biomedical Campus, Cambridge
CB2 0AW, UK. * email: hgnc@genenames.org, elsbeth@genenames.org

Summary

Standardized gene naming is crucial for effective communication about genes, and as genomics becomes increasingly important in healthcare, the need for a consistent language for human genes becomes ever more vital. Here we present the current HUGO Gene Nomenclature Committee (HGNC) guidelines for naming not only protein-coding but also RNA genes and pseudogenes, and outline the changes in approach and ethos that have resulted from the discoveries of the last few decades.

Introduction

The first guidelines for human gene nomenclature were published in 1979¹, when the Human Gene Nomenclature Committee was originally established and charged with the authority to approve and implement standardized human gene symbols and names. In 1989 the Nomenclature Committee was placed under the auspices of the newly founded Human Genome Organisation (HUGO), becoming the HUGO Gene Nomenclature Committee (HGNC). Subsequent revisions to the nomenclature guidelines were published in 1987², 1995³, 1997⁴, and 2002⁵. In the intervening years the HGNC has published online updates to the guidelines to reflect the significant changes and increase in knowledge and data during this exciting

27 period in human genomics. Over 40,000 human loci have been named by the HGNC to date;
28 around half are protein-coding genes, and most resources now agree that there are around
29 19,000-20,000 protein-coding genes in the human genome, considerably lower than some
30 earlier estimates. As well as naming protein-coding genes, significant progress has been
31 made in different classes of RNA genes and pseudogenes. All approved human gene symbols
32 can be found in the online HGNC database (<https://www.genenames.org/>)⁶.

33 The philosophy of the HGNC used to be that "gene nomenclature should evolve with new
34 technology" and that symbol changes, if supported by most researchers working on a gene,
35 were considered if they reflected new functional information. Since the advent of clinical
36 genomics such changes have much wider impacts, and it is impossible to reach all clinicians,
37 patients, charities and other parties interested in genes. Therefore, the stability of gene
38 symbols, particularly those associated with disease, is now a key priority for the HGNC.
39 Nevertheless, novel information can still be encapsulated in the gene name without changing
40 the gene symbol.

41 As human gene symbols are also routinely transferred to homologous vertebrate genes,
42 including in our sister project the Vertebrate Gene Nomenclature Committee (VGNC), we now
43 avoid references to human-specific traits in nomenclature wherever possible.

44 We strongly advise researchers to contact us whenever they are considering naming a novel
45 gene, or renaming an existing gene or group of genes, of all locus types, not only for protein-
46 coding genes. It is not always possible to approve the symbol requested, but we strive to
47 work with researchers to find an acceptable alternative. Requesting an approved symbol
48 ensures that your published symbol is present in our and other biomedical databases. We
49 further encourage journal editors and reviewers to check that approved nomenclature is
50 being used and require that authors contact the HGNC prior to publication for any novel
51 symbols. Submitters should bear in mind that the HGNC is committed to minimal future
52 changes to gene symbols and that we do not take publication precedence into account when
53 approving nomenclature.

54 Readers should note that the following are guidelines and recommendations (Box 1), not
55 strict rules. We are aware of numerous exceptional legacy symbols and names that remain
56 approved. The HGNC considers the naming of each and every gene on a case-by-case basis,
57 and deviations from these guidelines may be made given sufficient evidence that the
58 nomenclature will ultimately aid communication and data retrieval.

59

60 **Gene naming**

61 For many years the HGNC has maintained the definition of a gene as "a DNA segment that
62 contributes to phenotype/function. In the absence of demonstrated function a gene may be
63 characterized by sequence, transcription or homology". As there is still no universally agreed
64 alternative we continue to use this definition.

65 Ideally gene symbols are short, memorable and pronounceable, and most gene names are
66 long form descriptions of the symbol. Names should be brief, specific and convey something
67 about the character or function of the gene product(s), but not attempt to describe
68 everything known. Each gene is assigned only one symbol; the HGNC does not routinely
69 name isoforms (i.e. alternate transcripts or splice variants). This means no separate symbols
70 for protein-coding *or* non-coding RNA isoforms of a protein-coding locus or alternative
71 transcripts from a non-coding RNA locus (Box 2).

72 Where authors wish to use their own isoform notation, we advise stating clearly that this
73 notation denotes an isoform of a particular gene and quoting the HGNC symbol for that gene.

74 In exceptional circumstances, and following community demand, separate symbols have been
75 approved for gene segments in complex loci, i.e. the UGT1 locus, the clustered
76 protocadherins at 5q31 and the immunoglobulin and T cell receptor families. Putative
77 bicistronic loci may be assigned separate symbols to represent the distinct gene products. For
78 example, *PYURF*, "PIGY upstream reading frame" is encoded by the same transcript as *PIGY*,
79 "phosphatidylinositol glycan anchor biosynthesis class Y".

80 Table 1 summarizes key factors considered when assigning gene nomenclature. Additionally,
81 Supplementary table 1 lists characters recommended for specific usage in symbols,
82 Supplementary table 2 highlights specific conventions used in gene names, and
83 Supplementary tables 3 and 4 provide Greek-to-Latin alphabet conversions and single letter
84 amino acid symbols, respectively.

85

86 **Gene naming by biotype**

87 **Protein-coding genes**

88 We aim to name protein-coding genes based on a key normal function of the gene product.
89 Many protein-coding genes of known function are named in collaboration with internationally
90 recognized bodies composed of experts in a specific field. Where possible, related genes are
91 named using a common root symbol to enable grouping, typically based on sequence
92 homology, shared function or membership of protein complexes.

93 Gene group members should be designated by Arabic numerals placed immediately after the
94 root symbol, e.g. *KLF1*, *KLF2*, *KLF3*. More rarely single-letter suffixes may be used, e.g.
95 *LDHA*, *LDHB*, *LDHC*. Some large gene families may include a variety of number/letter
96 combinations to indicate subgroupings, e.g. *CYP1A1*, *CYP21A2*, *CYP51A1* (cytochrome P450
97 superfamily members).

98 For genes involved in specific immune processes, or encoding an enzyme, receptor or ion
99 channel, we consult with specialist nomenclature groups (see Supplementary Note). For other
100 major gene groups we consult a panel of advisors when naming new members and discussing
101 proposed nomenclature updates. A list of our specialist advisors is provided on our website⁶
102 and we welcome suggestions of new experts for specific gene groups.

103 In the absence of functional data, protein-coding genes may be named in the following ways:
104 (1) Based on recognized structural domains and motifs encoded by the gene (e.g. *ABHD1*
105 “abhydrolase domain containing 1”, *HEATR1* “HEAT repeat containing 1”). As these features
106 can provide insight into the character of the gene product, this type of symbol is commonly

107 retained even after the normal function of the gene product has been elucidated, though
108 further information may be added to the gene name; (2) Based on homologous genes within
109 the human genome. Where naming is based on characterized homologs, genes of unknown
110 function are given the next symbol within a designated series but with a different gene name
111 format, e.g. *CASTOR3*, "CASTOR family member 3" rather than "cytosolic arginine sensor for
112 mTORC1 subunit 3". The placeholder root symbol FAM ("family with sequence similarity") is
113 used when there is no information available for any of the homologous genes. Each
114 homologous family has a unique FAM number, e.g. *FAM3*, and each family member is
115 distinguished by a letter or letter and number, e.g. *FAM3A*, *FAM3C2P*. Note that this root can
116 be applied to both protein-coding and non-coding gene families; (3) Based on homologous
117 genes from another species. Where there is a 1-to-1 ortholog, the same/equivalent symbol
118 will be approved, e.g. human *CDC45* "cell division cycle 45" based on *S. cerevisiae CDC45*. A
119 unique number or letter suffix is added if there is more than one human homolog, e.g.
120 *UNC45A* and *UNC45B* are co-orthologs of *C. elegans unc-45*. Gene names are updated to be
121 appropriate for vertebrates, e.g. "unc-45 myosin chaperone" instead of "UNCoordinated 45";
122 (4) Based only on the presence of an open reading frame. Genes of unknown function that fit
123 none of the above criteria are designated by the chromosome of origin, the letters "orf" for
124 open reading frame (in lower case to prevent confusion between "O" and the numeral "0",
125 which may be part of the chromosome number) and a number in a series, e.g. *C3orf18*,
126 "chromosome 3 open reading frame 18". In cases where the coding potential of the locus is in
127 doubt we include the word "putative" in the name, e.g. "chromosome 18 putative open
128 reading frame 15".

129 Historically, genes of unknown function identified by the Human cDNA project at the Kazusa
130 DNA Research Institute¹⁰ have been named using the *KIAA#* identifiers assigned by this
131 project.

132 **Pseudogenes**

133 We define a pseudogene as a sequence that is incapable of producing a functional protein
134 product but has a high level of homology to a functional gene. In general, we only name
135 pseudogenes that retain homology to a significant proportion of the functional ancestral gene.

136 The majority of pseudogenes are processed and named based on a specific parent gene, e.g.
137 *DPP3P1*, "DPP3 pseudogene 1". Such pseudogene numbering is usually species-specific and
138 hence orthology cannot be inferred from identical pseudogene symbols in different species.

139 Pseudogenes that retain most of the coding sequence compared to other family members
140 (and are usually unprocessed) are named as a new family member with a "P" suffix, e.g.
141 *CBWD4P*, "COBW domain containing 4, pseudogene". This naming format is also used for
142 genes that are pseudogenized relative to their functional ortholog in another species, e.g.
143 *ADAM24P*, "ADAM metallopeptidase domain 24, pseudogene" is the pseudogenized ortholog
144 of mouse *Adam24*. Note, rarely such pseudogenes do not include the "P" if the symbol is well
145 established, e.g. *UOX*, "urate oxidase (pseudogene)".

146 A small number of genes are currently pseudogenized in the reference genome, but known to
147 have coding alleles segregating in the population. Such loci are given the locus type "protein-
148 coding" and indicated by including "(gene/pseudogene)" at the end of the gene name, e.g.
149 *CASP12*, "caspase 12 (gene/pseudogene)".

150 **Non-coding RNA genes**

151 We name non-coding RNA (ncRNA) genes according to their RNA type, please see our recent
152 review¹¹. For small RNAs where an expert resource exists, we follow their naming schema,
153 e.g. miRBase¹² for microRNAs and the genomic tRNA database (GtRNAdb)¹³ for tRNAs. Other
154 classes of ncRNA such as small nuclear RNAs are named in collaboration with specialist
155 advisors.

156 For long non-coding RNAs (lncRNAs), wherever possible we name these based on a key
157 function or characteristic of the encoded RNA. Where functional information is not available, a
158 systematic nomenclature is applied, see Figure 1.

159 **Readthrough transcripts**

160 Readthrough transcripts are normally produced from adjacent loci and include coding and/or
161 non-coding parts of two (or more) genes. The HGNC only names readthrough transcripts
162 that are consistently annotated by both the RefSeq annotators at NCBI¹⁴ and the GENCODE

163 annotators at Ensembl¹⁵. These transcripts have the locus type “readthrough transcript” and
164 are symbolized using the two (or more) symbols from the parent genes, separated by a
165 hyphen, e.g. *INS-IGF2*, and the name “[symbol] readthrough”, e.g. “INS-IGF2 readthrough”.
166 The name may also include additional information about the potential coding status of the
167 transcript, such as “(NMD candidate)”.

168 **Gene segments**

169 For specific complex loci the HGNC assigns symbols to individual gene segments, solely based
170 on community request. Examples of this are the immunoglobulins and T-cell receptors, the
171 UGT1 locus and clustered protocadherins.

172 **Genomic regions**

173 The HGNC previously named genomic regions referenced in the literature, such as *XIC*, “X
174 chromosome inactivation center”, and gene clusters were assigned symbols suffixed with the
175 “@” character, e.g. *HOXA@*, “homeobox A cluster”. We no longer routinely provide symbols
176 for genomic regions but some, such as those for fragile sites, have been retained where they
177 have been used in publications and this information would otherwise be lost.

178

179 **Genes only found within subsets of the population**

180 Historically, the HGNC has only approved symbols for genes that are on the human reference
181 genome. Rare exceptions have been made when requested by particular communities, e.g.
182 structural variants within the HLA and KIR gene families, both of which have dedicated
183 nomenclature committees. Future naming of structural variants will be restricted to those on
184 alternate loci that have been incorporated into the human reference genome by the Genome
185 Reference Consortium (GRC, <https://www.ncbi.nlm.nih.gov/grc>). The underscore character is
186 reserved for genes annotated on alternate reference loci, e.g. *GTF2H2C_2* is a second copy of
187 *GTF2H2C* on a 5q13.2 alternate reference locus; *APOBEC3A_B* is a deletion hybrid on a
188 22q13 alternate reference locus that includes exons from both the *APOBEC3A* and *APOBEC3B*
189 parent genes.

190

191 **Status**

192 All HGNC gene records have a status: the vast majority are "approved", but when new
193 evidence shows that a previously named gene is no longer considered to be real the entry
194 changes to the status "entry withdrawn". Wherever possible we avoid reusing symbols from
195 "entry withdrawn" records, as this can cause considerable confusion.

196

197 **Naming across vertebrates**

198 We recommend that orthologous genes across vertebrate (and where appropriate, non-
199 vertebrate) species should have the same gene symbol.

200 **The Vertebrate Gene Nomenclature Committee**

201 The Vertebrate Gene Nomenclature Committee (VGNC, <https://vertebrate.genenames.org/>) is
202 an extension of the HGNC responsible for assigning standardized names to genes in
203 vertebrate species that currently lack a nomenclature committee. The VGNC coordinates with
204 the five established existing vertebrate nomenclature committees, MGNC (mouse)¹⁶, RGNC
205 (rat, <https://rgd.mcw.edu/nomen/nomen.shtml>), CGNC (chicken)¹⁷, XNC (Xenopus frog)¹⁸
206 and ZNC (zebrafish)¹⁹, to ensure vertebrate genes are named in line with their human
207 homologs.

208 Orthologs of human *C#orf#* genes are assigned the human symbol with the other species
209 chromosome number as a prefix and an H denoting human. Therefore, as the ortholog of
210 human *C1orf100* is on cow chromosome 16, the cow symbol is *C16H1orf100* with the
211 corresponding gene name "chromosome 16 C1orf100 homolog".

212 Gene families with a complex evolutionary history should ideally be named with the help of
213 an expert in the field, as has already been implemented for the olfactory receptor²⁰ and
214 cytochrome P450 gene families.

215 **Species designation**

216 To distinguish the species of origin for homologous genes with the same gene symbol, we
217 recommend citing the NCBI taxonomy ID²¹, as well as either the current name or the
218 GenBank common name, e.g. Taxonomy ID: 9598 and either *Pan troglodytes* or chimpanzee.

219

220 **Nomenclature updates**

221 While we are committed to minimizing symbol changes some updates will still be appropriate.
222 All requests for change are considered on a case-by-case basis and often involve community
223 consultation. We anticipate most future changes will fall into one of the following categories.

224 **Symbol updates for placeholders**

225 FAMs, C#orfs and KIAAs are regarded as placeholder symbols and updated with structure
226 and/or function-based designations whenever possible. However, where specific placeholder
227 symbols have become entrenched in the literature, we may make exceptions and retain the
228 placeholder, while updating the gene name, e.g. *FAM20B* has been retained with the updated
229 gene name FAM20B glycosaminoglycan xylosylkinase.

230 **Replacing underused and problematic nomenclature**

231 We may consider updating symbols that have been rarely/never published, are not suitable
232 for transfer to other vertebrates, and/or have been widely used but could cause significant
233 problems. Examples are shown in Box 3.

234

235 **Gene symbol usage**

236 The HGNC endorses the use of italics to denote genes, alleles and RNAs to distinguish them
237 from proteins.

238 We advise that authors quote the approved gene symbol at least once in the abstract of any
239 publication. Every gene with an approved symbol also has a unique HGNC ID in the format
240 HGNC:number (e.g. gene symbol *BRAF*, HGNC ID HGNC:1097). While we aim to minimize
241 symbol changes some updates are inevitable and sometimes an approved symbol can be
242 used to denote a different gene in the literature; therefore we advise quoting the HGNC ID
243 for each gene to avoid ambiguity. HGNC IDs are associated with the gene sequence and do
244 not change unless the gene structure undergoes extreme alteration (i.e. merged with another
245 locus or split into multiple loci). This ensures effective and reliable tracking of data regardless
246 of any nomenclature changes.

247

248 **Acknowledgements**

249 *Many thanks to all current and former members of the HGNC team, in particular the late*
250 *Professor Sue Povey who was HGNC's PI from 1996-2007, our specialist advisors and*
251 *advisory board members past and present. The HGNC relies heavily on the expertise and*
252 *feedback of researchers and are grateful for all input we receive. The HGNC is currently*
253 *funded by the National Human Genome Research Institute (NHGRI) grant U24HG003345 (to*
254 *EAB) and Wellcome Trust grant 208349/Z/17/Z (to EAB).*

255

256 **Contributions**

257 EAB directed and obtained funding for the project. EAB, RLS and ST wrote the original draft.
258 EAB, RLS, ST, BB and TEMJ revised the manuscript. TEMJ constructed figure 1. All authors
259 (EAB, BB, PD, TEMJ, RLS, ST) contributed to, and commented on, the manuscript prior to
260 submission and contributed to the development of the current nomenclature guidelines.

261

262 **Competing Interests**

263 The authors declare no competing interests.

264

265 **References**

- 266 1. Shows, T.B. *et al.* International system for human gene nomenclature (ISGN, 1979).
267 *Cytogenet Cell Genet.* **25**, 96-116 (1979).
- 268 2. Shows, T.B. *et al.* Guidelines for human gene nomenclature. An international system for
269 human gene nomenclature (ISGN, 1987). *Cytogenet Cell Genet.* **46**, 11-28 (1987).
- 270 3. McAlpine, P. Genetic nomenclature guide. Human. *Trends Genet.* Mar, 39-42 (1995).
- 271 4. White, J.A. *et al.* Guidelines for human gene nomenclature. *Genomics* **45**, 468-471 (1997).
- 272 5. Wain, H.M., Bruford, E.A., Lovering, R.C., Lush, M.J., Wright, M.W., Povey, S. Guidelines
273 for Human Gene Nomenclature. *Genomics* **79**, 464-470 (2002).
- 274 6. Braschi, B. *et al.* Genenames.org: the HGNC and VGNC resources in 2019. *Nucleic Acids*
275 *Res.* **47**, D786-D792 (2019).
- 276 7. den Dunnen, J.T. Describing Sequence Variants Using HGVS Nomenclature. *Methods Mol*
277 *Biol.* **1492**, 243-251 (2017).
- 278 8. Mayer, J., Blomberg, J., Seal, R.L. A revised nomenclature for transcribed human
279 endogenous retroviral loci. *Mob DNA* **2**, 7 (2011).
- 280 9. Amberger, J.S., Bocchini, C.A., Scott, A.F., Hamosh, A. OMIM.org: leveraging knowledge
281 across phenotype-gene relationships. *Nucleic Acids Res.* **47**, D1038-D1043 (2019).
- 282 10. Nagase, T., Koga, H., Ohara, O. Kazusa mammalian cDNA resources: towards functional
283 characterization of KIAA gene products. *Brief Funct Genomic Proteomic.* **5**, 4-7 (2006).
- 284 11. Seal, R.S. *et al.* A guide to naming human non-coding RNA genes. *EMBO J.* **39(6)**,
285 e103777 (2020).
- 286 12. Kozomara, A., Birgaoanu, M., Griffiths-Jones, S. miRBase: from microRNA sequences to
287 function. *Nucleic Acids Res.* **47**, D155-D162 (2019).

- 288 13. Chan, P.P., Lowe, T.M. GtRNAdb 2.0: an expanded database of transfer RNA genes
289 identified in complete and draft genomes. *Nucleic Acids Res.* **44**, D184-D189 (2016).
- 290 14. O'Leary, N.A. *et al.* Reference sequence (RefSeq) database at NCBI: current status,
291 taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733-D745 (2016).
- 292 15. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes.
293 *Nucleic Acids Res.* **47**, D766-D773 (2019).
- 294 16. Maltais, L.J., Blake, J.A., Eppig, J.T., Davisson, M.T. Rules and guidelines for mouse gene
295 nomenclature: a condensed version. *Genomics* **45**, 471-476 (1997).
- 296 17. Burt, D.W. *et al.* The Chicken Gene Nomenclature Committee Report. *BMC Genomics* **10**,
297 S5 (2009).
- 298 18. James-Zorn, C. *et al.* Xenbase: Core features, data acquisition, and data processing.
299 *Genesis* **53**, 486-497 (2015).
- 300 19. Ruzicka, L. *et al.* The Zebrafish Information Network: new support for non-coding genes,
301 richer Gene Ontology annotations and the Alliance of Genome Resources. *Nucleic Acids Res.*
302 **47**, D867-D873 (2019).
- 303 20. Olender, T. *et al.* A unified nomenclature for vertebrate olfactory receptors. *BMC Evol.*
304 *Biol.* **20**, 42 (2020).
- 305 21. Federhen, S. The NCBI Taxonomy Database. *Nucleic Acids Res.* **40**, D136-D143 (2012).
- 306

307 Figure 1: The HGNC has a systematic process for naming long non-coding (lnc)RNA genes. In the
 308 absence of suitable published information, lncRNA genes are named based on genomic context.
 309 Abbreviations are as follows: HG: host gene, PC: protein-coding, DT: divergent transcript (used for
 310 lncRNA genes that share a promoter with a PC gene), IT: intronic transcript, OT: overlapping transcript,
 311 AS: antisense RNA, LINC: long intergenic non-protein coding RNA.

312

313 **Table 1**

Symbols	Names
Must be unique within a given genome	Should be brief and specific
Must not be offensive or pejorative (ideally in any language)	
Must not use superscripts or subscripts or punctuation*	Should minimize punctuation; commas, hyphens and parentheses are included for clarity**
Must only contain uppercase*** Latin letters and Arabic numerals	Must be written in American English
Must start with a letter	Must start with a lowercase letter (unless starting with an eponymous term or capitalised abbreviation)
Should not include "G" for gene, "H" for human, Roman numerals or Greek letters	Should not include the words "gene" or "human"
Should not spell proper names or common words or match commonly used abbreviations	Should start with the same letter as the symbol (to facilitate alphabetical listing and grouping)
Should avoid duplicating symbols in other species (unless orthologous)	Should not reference: any species, taxa, tissue specificity, molecular weight, chromosomal location, human-specific features and phenotypes, familial terms

Some letters or combinations of letters are used in a symbol to give a specific meaning, and their use for other meanings should be avoided where possible (see supplementary table 1).

Descriptive modifiers usually follow the main part of the name, to enable the use of a common root symbol for a gene group, e.g. **ACADM** "acyl-CoA dehydrogenase **medium chain**" and **ACADS** "acyl-CoA dehydrogenase **short chain**".

314 * see Supplementary Table 2 for punctuation exceptions in symbols

315 ** exceptions on punctuation are made for enzyme names

316 *** sole exception of C#orfs

317

318 **Box 1:**

319 **A summary of the guidelines is:**

320 1. Each gene is assigned a unique symbol, HGNC ID and descriptive name.

321 2. Symbols contain only uppercase Latin letters and Arabic numerals.

322 3. Symbols should not be the same as commonly used abbreviations

323 4. Nomenclature should not contain reference to any species or "G" for gene.

324 5. Nomenclature should not be offensive or pejorative.

325

326 **Box 2:**

327 HGNC does not provide official nomenclature for the following:

328 **sequence variant nomenclature**, which is the responsibility of the Human Genome

329 Variation Society (HGVS)⁷. They provide recommendations for defining variations found in

330 DNA, RNA and protein sequences, and endorse the use of HGNC gene symbols within their

331 notation.

332 **products of gene translocations or fusions:** we are not aware of official naming
333 guidelines for these. SYMBOL1-SYMBOL2 is widely used, but we use this format for
334 readthrough transcripts (see section 2.4) and hence would specifically *not* recommend this
335 for translocations or fusions. We recommend the format SYMBOL1/SYMBOL2, which has
336 been used in some publications, e.g. *BCR/ABL1*.

337 **protein nomenclature:** we have no authority over naming proteins, but coordinate closely
338 with specialist groups who name specific subsets of proteins, such as the Enzyme
339 Commission. The recently devised International Protein Nomenclature Guidelines
340 (https://www.ncbi.nlm.nih.gov/genome/doc/internatprot_nomenguide/) were written with
341 the involvement of the HGNC, and in agreement with these guidelines we recommend that
342 “protein and gene symbols should use the same abbreviation”. We further advise that
343 proteins are referenced using non-italicised gene symbols, to distinguish them from genes.

344 **nomenclature for regulatory genomic elements** such as promoters, enhancers and
345 transcription factor binding sites. We also do not provide nomenclature for transposable
346 element insertions in the human genome. Protein-coding and long non-coding RNA genes
347 that fit the criteria outlined in Mayer *et al.*⁸ may be named as ERV-derived genes, but ERV
348 insertions will not be named.

349 **nomenclature for human loci associated with clinical phenotypes and complex**
350 **traits.** While HGNC historically named these, this activity has been taken over by OMIM⁹. All
351 HGNC entries with the locus type “phenotype only” now have the status “entry withdrawn”.
352 Note that some uncharacterized genes shown to be causative for a specific phenotype
353 adopted the phenotype symbol and name. Where these phenotypic symbols have become
354 entrenched in the literature, we aim to update the corresponding gene names to reflect an
355 aspect of the normal function of the gene and its products, e.g. TSC1, “tuberous sclerosis 1”
356 is now TSC1, “TSC complex subunit 1”.

357

358 **Box 3:**

359 Scenarios that may merit a symbol change include:

- 360 • adoption of a more appropriate/popular alias, e.g. *RNASEN* was updated to
361 *DROSHA* (drosha ribonuclease III) due to overwhelming community usage.
- 362 • domain or motif-based nomenclature, e.g. *TMEM206* (transmembrane protein
363 206) is now *PACC1* (proton activated chloride channel 1).
- 364 • phenotype/disease-based nomenclature, e.g. *CASC4* (cancer susceptibility
365 candidate 4) was renamed *GOLM2* (golgi membrane protein 2), consistent with its
366 paralog *GOLM1*.
- 367 • location-based nomenclature, e.g. *TWISTNB* (TWIST neighbour) is now *POLR1F*
368 (RNA polymerase I subunit F).
- 369 • pejorative symbols, e.g. *DOPEY1* was renamed to *DOP1A* (DOP1 leucine zipper like
370 protein A).
- 371 • misleading/incorrect nomenclature, e.g. *OTX3* was initially named erroneously as
372 an OTX family member and has been renamed as *DMBX1*
373 (diencephalon/mesencephalon homeobox 1).
- 374 • symbols that affect data handling and retrieval, e.g. all symbols that auto-
375 converted to dates in Microsoft Excel have been changed (*SEPT1* is now *SEPTIN1*;
376 *MARCH1* is now *MARCHF1* etc); tRNA synthetase symbols that were also common
377 words have been changed (*WARS* is now *WARS1*, *CARS* is now *CARS1*, etc.).

378

379

380

381