



## Young genes have distinct gene structure, epigenetic profiles, and transcriptional regulation

Michael S Werner, Bogdan Sieriebriennikov, Neel Prabh, et al.

*Genome Res.* published online September 19, 2018  
Access the most recent version at doi:[10.1101/gr.234872.118](https://doi.org/10.1101/gr.234872.118)

---

<b>P&lt;P</b>	Published online September 19, 2018 in advance of the print journal.
<b>Accepted Manuscript</b>	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
<b>Open Access</b>	Freely available online through the <i>Genome Research</i> Open Access option.
<b>Creative Commons License</b>	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at <a href="http://creativecommons.org/licenses/by-nc/4.0/">http://creativecommons.org/licenses/by-nc/4.0/</a> .
<b>Email Alerting Service</b>	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or <a href="#">click here</a> .

---



---

To subscribe to *Genome Research* go to:  
<https://genome.cshlp.org/subscriptions>

---

Published by Cold Spring Harbor Laboratory Press

1  
2  
3 **Young genes have distinct gene structure, epigenetic profiles, and**  
4 **transcriptional regulation**  
5  
6  
7  
8  
9  
10  
11

12 Michael S. Werner<sup>1</sup>, Bogdan Sieriebriennikov<sup>1</sup>, Neel Prabh<sup>1</sup>, Tobias Loschko<sup>1</sup>, Christa  
13 Lanz<sup>1</sup> and Ralf J. Sommer<sup>1</sup>  
14

15 <sup>1</sup>Department of Evolutionary Biology, Max Planck Institute for Developmental Biology,  
16 72076 Tübingen, Germany

17 Corresponding author: [ralf.sommer@tuebingen.mpg.de](mailto:ralf.sommer@tuebingen.mpg.de)  
18  
19

20 **Keywords:** epigenetics, orphan genes, new genes, *Pristionchus pacificus*,  
21 *Caenorhabditis elegans*  
22  
23  
24  
25  
26  
27

28 **Abstract**

29 Species-specific, new, or ‘orphan’ genes account for 10-30% of eukaryotic genomes. Although  
30 initially considered to have limited function, an increasing number of orphan genes have been  
31 shown to provide important phenotypic innovation. How new genes acquire regulatory  
32 sequences for proper temporal and spatial expression is unknown. Orphan gene regulation may  
33 rely in part on origination in open chromatin adjacent to pre-existing promoters, although this  
34 has not yet been assessed by genome-wide analysis of chromatin states. Here we combine  
35 taxon-rich nematode phylogenies with Iso-Seq, RNA-seq, ChIP-seq, and ATAC-seq to identify  
36 the gene structure and epigenetic signature of orphan genes in the satellite model nematode  
37 *Pristionchus pacificus*. Consistent with previous findings we find young genes are shorter,  
38 contain fewer exons, and are on average less strongly expressed than older genes. However,  
39 the subset of orphan genes that are expressed exhibit distinct chromatin states from similarly  
40 expressed conserved genes. Orphan gene transcription is determined by a lack of repressive  
41 histone modifications, confirming long-held hypotheses that open chromatin is important for new  
42 gene formation. Yet orphan gene start sites more closely resemble enhancers defined by  
43 H3K4me1, H3K27ac and ATAC-seq peaks, in contrast to conserved genes that exhibit  
44 traditional promoters defined by H3K4me3 and H3K27ac. While the majority of orphan genes  
45 are located on chromosome arms that contain high recombination rates and repressive histone  
46 marks, strongly expressed orphan genes are more randomly distributed. Our results support a  
47 model of new gene origination by rare integration into open chromatin near enhancers.

48

49

50

51

52

53

## 54 **Introduction**

55 Gene regulation is a highly orchestrated process that includes transcription factor binding sites  
56 (TFBs), non-coding RNAs, histone modifications and chromatin structure (Voss and Hager  
57 2014). The identification and mechanism of these molecular factors have been revealed for  
58 several conserved gene networks leading towards a better understanding of development and  
59 disease. But how new genes, also referred to as orphan or taxon-restricted, acquire this  
60 complex architecture is unknown. For the increasing number of identified new genes that  
61 provide important biological function (Cai et al. 2008; Burki and Kaessmann 2004; Rosso et al.  
62 2008; Reinhardt et al. 2013; Chen et al. 2013b; 2010; Mayer et al. 2015; Santos et al. 2017), the  
63 evolutionary path from origin to integration into gene networks depends on their precise  
64 transcriptional regulation (Carelli et al. 2016). Yet in the majority of cases it is unclear how even  
65 the most fundamental *cis*-regulatory elements like promoter and termination sequences are  
66 obtained (Tautz and Domazet-Lošo 2011; Long et al. 2013). Orphan genes can originate *de*  
67 *novo* or by duplication, recombination or horizontal gene transfer into pre-existing regulatory  
68 architecture (Chen et al. 2012b; Betrán and Long 2003; Kaessmann et al. 2009; Kaessmann  
69 2010; McLysaght and Hurst 2016; Li et al. 2009). But the extent to which this occurs is limited  
70 by the potential to disrupt the genes already there (Vinckenbosch et al. 2006). In the few cases  
71 where integration has been observed, the presence of nearby regulatory sequences was largely  
72 detected by proximity, or sequence homology to known promoters, CpG islands or TFBs  
73 (Carvunis et al. 2012; Abrusán 2013; Ruiz-Orera et al. 2015; Li et al. 2016). Given these  
74 constraints, the contribution of pre-existing regulatory architecture to new gene transcription is  
75 still unknown, and with functional genomic information (e.g. chromatin states and enhancers)  
76 could potentially be much larger. Indeed, a recent analysis of mammalian ChIP-seq data sets  
77 found 51% of expressed mouse retrogenes (mRNAs that are reverse transcribed and inserted  
78 into the genome) exhibit robust H3K4 tri-methylation (Carelli et al. 2016), and transcription of the  
79 new gene *QQS* in *Arabidopsis thaliana* is inversely correlated with DNA methylation at 5'

80 transposable elements (Silveira et al. 2013), suggesting an important role for chromatin  
81 regulation in new gene transcription. We sought to employ the rich taxonomic resources of  
82 nematodes to first identify young and old genes, and then observe their regulatory architecture  
83 by several genome-wide approaches.

84         The diplogastrid nematode *Pristionchus pacificus* can be found in a necromenic  
85 relationship with beetles, but has been developed in the laboratory as a satellite model for  
86 comparative studies to *C. elegans* (Sommer and Streit 2011; Sommer and McGaughran  
87 2013)(Fig. 1A-D). More recent genetic analysis of dimorphic mouth-forms (Fig. 1E-G) has led to  
88 *P. pacificus* emerging as an important model system for phenotypic plasticity in its own right  
89 (Bento et al. 2010; Ragsdale et al. 2013; Seroby et al. 2016; Kieninger et al. 2016). In  
90 addition to the vast taxonomic diversity and corresponding genomes of other nematode species,  
91 the recent high quality chromosome-scale genome (Rödelsperger et al. 2017) and reverse  
92 genetic tools (Witte et al. 2015) in *P. pacificus* provide a robust framework for studying new  
93 genes (Prabh and Rödelsperger 2016; Baskaran et al. 2015). Here we probe the gene structure,  
94 expression, and regulatory architecture of *P. pacificus* evolutionary gene classes with long-read  
95 PacBio transcript sequencing (Iso-Seq), traditional high depth RNA-sequencing (RNA-seq),  
96 chromatin immunoprecipitation (ChIP-seq) of six histone post-translational modifications and  
97 assay for transposon-accessible chromatin (ATAC-seq). In addition to our findings, the data sets  
98 collected provide the first epigenomic map in *P. pacificus*, which is only the second  
99 comprehensive chromatin state annotation in nematodes, creating a resource for future  
100 functional and comparative studies.

101

102

103

104

105

## 106 **Results**

107

### 108 **Partitioning of *P. pacificus* genes into evolutionary classes**

109 The first *P. pacificus* draft genome published in 2008 (Dieterich et al. 2008) had a large number  
110 of genes with undetectable homology. Although the confidence in these gene predictions was  
111 initially low, every subsequent refinement of both the genome and gene annotation continually  
112 detected 20-40% of genes that appear as new, orphan, or taxon-restricted (Sinha et al. 2012;  
113 Baskaran et al. 2015; Prabh and Rödelsperger 2016; Baskaran and Rödelsperger 2015). Using  
114 our most recent chromosome-scale PacBio genome (Rödelsperger et al. 2017) and 24 other  
115 nematode species we re-evaluated the relative abundance of evolutionary gene classes (Fig.  
116 1H). We defined the most highly conserved genes as having 1:1 orthology with *C. elegans*  
117 (BLASTP e-value  $\leq 0.001$ ), which are estimated to have diverged from *P. pacificus* between 60-  
118 90 mya (Cutter 2008; Rota-Stabelli et al. 2013; Hedges et al. 2015). We also defined an  
119 intermediate conserved class as ‘homologs’ if they display homology to at least one gene in the  
120 other 24 nematode species (see Methods) – which could represent either relatively young  
121 genes, or old genes that have been lost. Finally, we define ‘orphan’ genes as having no  
122 homology to genes in the other 24 queried species. The resulting partition of genes  
123 approximates the ‘30% rule’ of new gene composition (Khalturin et al. 2009)(Fig. 1I). We then  
124 applied several genomic approaches to molecularly characterize each evolutionary gene class.

125

### 126 **Characterization of gene structure by long read RNA sequencing (Iso-Seq)**

127 We sought to improve the overall gene annotation in *P. pacificus*, and then characterize the  
128 genetic structures of each evolutionary gene class using PacBio Iso-Seq on mixed-  
129 developmental stage RNA (see Supplemental Methods, Supplemental Fig. S1A-C). After  
130 alignment we obtained 640,664 reads with a median insert size of 1,363 nucleotides  
131 (Supplemental Fig. S1D). Despite low read depth compared to conventional RNA-seq, our Iso-

132 Seq data covered 17,307 genes (68% of genes in the reference annotation 'El Paco')  
133 (Rödelsperger et al. 2017).

134 Relative to the current reference annotation, Iso-Seq identified a tighter distribution of  
135 gene lengths (Fig. 2A, median Iso-Seq=1,452 compared to median reference=1,599,  $p < 2.2 \times 10^{-16}$   
136 Wilcoxon rank sum test). This difference appears to be due to a more narrow distribution of  
137 exons, with 96.5% of Iso-Seq gene annotations containing between 1 to 20 exons, compared to  
138 85.7% for the reference annotation (Fig. 2B,  $p = < 2.2 \times 10^{-16}$  Wilcoxon rank sum test). The tighter  
139 distribution is also more consistent with the highly curated gene annotation of *C. elegans* in  
140 which 98.0% of genes contain between 1:20 exons (Supplemental Fig. S1E,F)(Michael and  
141 Manyuan 1999). This potential improvement in accuracy appears to result from fragmentation of  
142 excessively long gene predictions into distinct transcripts (see Supplemental Fig. S1G for  
143 example).

144 Long read Iso-Seq also provides more robust identification of isoforms, which are  
145 historically difficult to assemble from standard short-read RNA-sequencing (Conesa et al. 2016).  
146 Approximately half (50.6%) of expressed genes in *P. pacificus* exhibit greater than one isoform,  
147 and roughly a third (30.9%) exhibit greater than three isoforms (Supplemental Fig. S1H,I).  
148 However, some of these transcripts could be artifacts biased by incomplete coverage of 5' ends.  
149 Hence, we conservatively defined alternatively spliced isoforms as transcripts with the same  
150 start and stop coordinates, but differential exon inclusion or exclusion, intron retention, or  
151 differential splice site. Under this classification we observed 3,861 (24%) of expressed genes  
152 exhibit alternative splicing in *P. pacificus* (Fig. 2C), similar to the ~25% of genes estimated in *C.*  
153 *elegans* (Ramani et al. 2011). As an example we highlight gene *umm-s259-11.10-mRNA* where  
154 the majority of Iso-Seq reads (17/19) cover the entire transcript yielding eight isoforms, in stark  
155 contrast to standard short-read sequencing which rarely covers more than 3 exons per read  
156 (Fig. 2D). Collectively, a tighter distribution of transcript lengths and exon number, and diversity

157 of isoforms suggests that Iso-Seq improves the quality and quantity of gene annotation in *P.*  
158 *pacificus*.

159 Among evolutionary gene classes, most *C. elegans* 1:1 orthologs (88%), and  
160 approximately half of homologous and orphan genes (46% and 56%, respectively) exhibit Iso-  
161 Seq coverage, demonstrating that our Iso-Seq data is sensitive enough to detect thousands of  
162 transcripts from each evolutionary gene class (Fig. 2E). We also performed Iso-Seq on rRNA-  
163 depleted 'total RNA' (see Supplemental Methods) to assess whether young genes are un-, or  
164 under-polyadenylated, which is typical of non-coding RNAs (Derrien et al. 2012). We found a  
165 similar percent coverage from the direct and total RNA methods (Fig. 2E,F), and a consistent  
166 polyadenylation read bias for all gene classes (Fig. 2G-I). Hence most young genes, or at least  
167 transcribed young genes, are polyadenylated. As polyadenylation is an important component of  
168 transcriptional and translational regulation (Proudfoot 2011; Proudfoot et al. 2002), this argues  
169 that most young genes have retained, or already acquired, 3' termination and processing  
170 sequence architecture.

171 We then used our Iso-Seq annotation to characterize gene length and exon number  
172 between evolutionary gene classes. Consistent with other systems (Stein et al. 2018; Ruiz-  
173 Orera et al. 2015), we found a strong bias of *C. elegans* 1:1 orthologs to be longer and contain  
174 more exons than homologs, which in turn were longer and contained more exons than orphan  
175 genes (Fig. 2J,K,  $p = <2.2 \times 10^{-16}$  Wilcoxon rank sum test). The intermediate gene structure of  
176 intermediate conserved genes (homologs) is also consistent with a transitional evolutionary path  
177 between young and old genes proposed by Carvunis *et al.* in 2012 (Carvunis et al. 2012;  
178 Abrusán 2013; Neme and Tautz 2013). In the following sections we sought to characterize and  
179 compare the chromatin regulation of young versus old genes.

180

181

182

### 183 **The *P. pacificus* epigenome**

184 To identify regulatory regions and expression levels of orphan, homolog, and *C. elegans* 1:1  
185 orthologs we performed two to three replicates of ChIP-seq on nine histone modifications and  
186 two replicates of RNA-seq in *P. pacificus* adults, and two replicates of ATAC-seq on mixed-  
187 stage cultures (Supplemental Fig. S2, Supplemental Table 3, and Supplemental Methods). All  
188 data sets showed good correlations between biological replicates (Pearson correlation between  
189 0.70-0.93 for ChIP, 0.88 for ATAC-seq, and 0.98 for gene FPKMs in RNA-seq) (Supplemental  
190 Fig. S3). We identified enriched regions (i.e. peaks) for each replicate of ChIP and ATAC-seq  
191 using MACS2 (Zhang et al. 2008) (Supplemental Table 1, see Methods). H2bub, H3K9ac and  
192 H3K79me2 exhibited less than 50% peak reproducibility and were excluded from further  
193 analysis (Supplemental Fig. S4A). The majority of remaining samples exhibited >70% overlap  
194 between replicates, except for H3K9me3 (54% reproducibility). However, H3K9me3 is a broadly  
195 distributed histone modification that is challenging for peak-finding software (Wang et al. 2013),  
196 and although most H3K9me3 antibodies are of low specificity (Nishikori et al. 2012; Hattori et al.  
197 2013), they can nevertheless distinguish constitutive vs. facultative heterochromatin (Trojer and  
198 Reinberg 2007).

199 We also performed ATAC-seq for identifying regions of open chromatin (Buenrostro et  
200 al. 2013). Although the standard protocol led to reproducible peaks, initially we could not identify  
201 nucleosomal read density, perhaps suggesting a difficulty of obtaining higher resolution  
202 fragments from highly differentiated and heterogeneous cell populations. Yet the new *omni*  
203 ATAC method (Corces et al. 2017) yielded nucleosomal and sub-nucleosomal read densities  
204 (Supplemental Fig. 2E), which we used for subsequent analysis.

205 We clustered the six high-confidence histone marks and *omni* ATAC-seq data using a  
206 Hidden-Markov Model (ChromHMM) (Ernst and Kellis 2012) into eight chromatin states (Fig.  
207 3A). Each chromatin state is enriched in histone modifications that define specific functional  
208 domains, such as actively transcribed regions, heterochromatin, and regulatory loci. We

209 assigned putative chromatin state annotations based on established classifications (Fig.  
210 3B)(Ernst et al. 2011; Rada-Iglesias et al. 2011)(Supplemental Table 2). We find that, at least at  
211 the whole animal level, approximately half (57%) of the genome is repressed, approximately a  
212 fifth (16%) represents actively transcribed genes, and more than a quarter (27%) is regulatory  
213 (including 6,785 promoters, 13,648 active enhancers, and 3,853 'poised' enhancers) (Fig. 3C).

214 Next, we verified that histone marks are enriched at the center of promoters and active  
215 enhancer annotations, and performed a *de novo* motif search (Heinz et al. 2010) (Fig. 3D-G).  
216 Both promoters (30.6%) and enhancers (22.2%) were enriched in a recognition sequence for  
217 MBP1, a yeast transcriptional activator that controls cell cycle progression (Koch et al. 1993).  
218 There is weak homology (BlastP,  $e = 2 \times 10^{-4}$ ) to MBP1 in *P. pacificus* (UMM-S233-5.4-mRNA-1),  
219 and in the future it will be interesting to see if this gene is also involved in cell-cycle control.  
220 There were also notable differences between enhancers and promoters, including binding site  
221 matches to human homeobox, *Drosophila* GAGA, and eukaryotic GATA transcription factors,  
222 demonstrating the precision of promoter and enhancer annotations, and hinting at the existence  
223 of deeply conserved regulatory elements.

224 As expected, promoter annotations were strongly enriched at the 5' end of genes (Fig.  
225 3H). There was also another peak near the 3' end of genes. Enhancers were also enriched at  
226 both 5' and 3' ends, although they are more evenly distributed throughout gene bodies. The  
227 existence of promoter/enhancer elements at the 3' end of genes has been observed in other  
228 species, and while its function is still unclear, there are several reports supporting promoter-3'  
229 end chromatin looping to facilitate successive rounds of transcription and enforce directionality  
230 (Werner et al. 2017; Lainé et al. 2009; O'Sullivan et al. 2004; Grzechnik et al. 2014).

231 To verify that our chromatin states correlate with a dynamically regulated gene we  
232 looked at *Ppa-pax-3*, which our lab has shown is expressed in early juvenile stages but is  
233 repressed during development (Yi and Sommer 2007). Indeed, we found *Ppa-pax-3* is in a large  
234 H3K27me3-repressed domain in adults (Fig. 3I). However, we also noticed two putative

235 enhancers after the first exon and 3' end, perhaps suggesting preparation for activation in  
236 developing embryos. Collectively, these data represent the first genome-wide annotation of  
237 chromatin regulation in *P. pacificus*, and to our knowledge represents only the second  
238 comprehensive data set in nematodes.

239

#### 240 **Chromatin regulation corresponds to gene expression**

241 We extended the previous single gene example to genome-wide high-depth RNA-seq and  
242 binned the adult transcriptome into four expression categories (Fig. 4A), then assessed the  
243 chromatin states of each. As predicted, gene bodies (exons and introns) of the highest  
244 expressed categories (groups 1 and 2) exhibited enrichment in chromatin states designated as  
245 'transcriptional transition' and 'elongation'. Conversely, repressive chromatin states were  
246 virtually absent from genes in the top two categories. In contrast, the lowest two expressed  
247 categories (groups 3 and 4) exhibited proportionally greater enrichment in repressive chromatin  
248 states and decreased enrichment in transcriptional transition and elongation states (Fig. 4B).  
249 While there was a minor enrichment in promoter chromatin states at 5' ends and 5' UTRs  
250 between high vs. low expressed categories, there was a larger difference in repressive  
251 chromatin states, especially at the 5' ends. There was also an increase in enhancer enrichment  
252 at 5' ends and 5' UTRs in the low expressed categories, perhaps reflecting a 'poised' chromatin  
253 state that is reactive to environmental influence. While promoters and enhancers exhibit a  
254 relatively small portion of the genome (15.6%), they comprise the majority of intergenic regions  
255 (Fig. 4B), hinting at a large and mostly unexplored regulatory circuitry in the compact nematode  
256 genome.

257

#### 258 **Chromatin regulation of evolutionary gene classes**

259 Next we assessed the chromatin states of evolutionary gene classes. *C. elegans* 1:1 orthologs  
260 resembled the highest expression categories (groups 1 and 2), while conserved and orphan

261 genes more closely resembled the lower expression categories (groups 3 and 4) (Fig. 4C-E).  
262 These histone patterns reflect the higher expression of *C. elegans* 1:1 orthologs compared to  
263 less conserved gene classes (Fig. 4F). Nevertheless, we noticed a significant number of orphan  
264 and homologous expressed gene outliers (Fig. 4F), and wondered if their chromatin signatures  
265 resembled that of expressed *C. elegans* 1:1 orthologs. Here, we found differences. Specifically,  
266 strongly expressed (group 1 and 2) orphan and homologous genes, which represent only 9.3%  
267 and 12.8% of their respective categories, broadly resembled the general chromatin state pattern  
268 of their classes except for having reduced repressive histone marks (Fig. 4G-I). Second,  
269 chromatin states 3 and 4 representing transcriptional transition and elongation are more highly  
270 represented in *C. elegans* 1:1 orthologs compared to expressed orphan and conserved genes.  
271 Third, *C. elegans* 1:1 orthologs exhibit little to no signature of active enhancers (chromatin state  
272 1) at their 5' ends or 5' UTRs, which are instead dominated by the promoter chromatin state  
273 consisting of H3K4me3 and H3K27ac. However, expressed homologous and orphan genes  
274 exhibited both promoter and enhancer enrichment at their 5' ends and 5' UTRs, and orphan  
275 genes in particular exhibited greater enrichment in enhancer vs. promoter chromatin states.

276 To investigate this difference more closely we examined the distribution of histone ChIP-  
277 seq and ATAC-seq around the 5' ends of each evolutionary class. Whereas expressed *C.*  
278 *elegans* 1:1 ortholog TSSs are dominated by H3K4me3 and H3K27ac, expressed orphan and  
279 homologous genes exhibit comparatively stronger enrichment of H3K4me1 and ATAC-seq (Fig.  
280 4J). Specifically, *C. elegans* 1:1 orthologs exhibit a 5' H3K4me3/H3K4me1 ratio of 10.1,  
281 compared to 2.4 for homologs, and 1.4 for orphan genes. Furthermore, while 54% of expressed  
282 *C. elegans* 1:1 ortholog 5' ends are within 1 kb of an annotated promoter, only 27% of  
283 expressed homologous gene and 21% of expressed orphan genes are in similar proximity to  
284 promoters ( $p < 2.2 \times 10^{-16}$  Wilcoxon rank sum test)(Fig. 5A). Conversely, 46% of expressed  
285 homologous and orphan gene TSSs are in 1 kb of active and poised enhancers, compared to  
286 33% of expressed *C. elegans* 1:1 ortholog TSSs ( $p < 2.2 \times 10^{-16}$  Wilcoxon rank sum test for both

287 comparisons)(Fig. 5B). Importantly, the expression of group 1 and 2 orphan and homologous  
288 genes are actually higher than *C. elegans* 1:1 orthologs ( $p < 2.2 \times 10^{-16}$  Wilcoxon rank sum  
289 test)(Supplemental Fig. S5), demonstrating that their chromatin architecture is independent of  
290 the general correlation with expression. There are two key points from these results. First, the  
291 transcription of young genes appears to depend on the absence of repressive heterochromatin,  
292 demonstrating a widely held but unconfirmed theory for the requirement of open chromatin in  
293 new gene expression. Second, the 5' end of expressed young genes, especially orphan genes,  
294 resemble enhancers rather than canonical promoters.

295

### 296 **Genomic position affects chromatin regulation of evolutionary gene classes**

297 Finally, we analyzed the general pattern of chromatin marks and evolutionary gene classes at  
298 the chromosomal level (Fig. 6). We observed strong patterns of activating marks in the center of  
299 autosomes II-V, and the repressive mark H3K27me3 on the chromosome arms. Conversely the  
300 X Chromosome was highly enriched in both repressive marks H3K27me3 and H3K9me3. These  
301 patterns have been observed in *C. elegans* (Liu et al. 2011), and correspond to general patterns  
302 in nematodes of dense clusters of conserved genes and low recombination rates in the center of  
303 chromosomes, and species-specific genes and high recombination rates in the chromosome  
304 arms (Coghlan 2005). However, Chromosome I in *P. pacificus* is an exception to other  
305 autosomes, where we observed two bands of activating marks instead of one. Recent analysis  
306 suggests that roughly half of *P. pacificus* Chromosome I is homologous to *C. elegans*  
307 Chromosome X, and the other half is homologous to Chromosome V (Rödelsperger et al. 2017).  
308 The *P. pacificus* chromosome pattern was viewed as ancestral because this organization is also  
309 found in the distantly related nematode *Bursaphelenchelus xylophilus*. However, the bipartite  
310 presence of histone modifications and conserved genes hints at an ancient chromosomal fusion  
311 from an even earlier origin, or frequent and repeated chromosomal fission and fusion events.

312           The chromosome-scale distribution of evolutionary gene classes was consistent with  
313 histone modification patterns. Specifically, *C. elegans* 1:1 orthologs, which are strongly  
314 expressed, are enriched in the active histone-mark chromosome centers. Conversely, the lower  
315 expressed homologous and orphan genes are enriched in the chromosome arms, which contain  
316 higher recombination rates and a greater density of repressive histone marks. However, these  
317 patterns are lost when controlling for expression. Highly transcribed (group 1 and 2) orphan and  
318 homologous genes were more randomly distributed throughout chromosomes, if not slightly  
319 biased to be closer towards the centers (Fig. 6). This genome-wide perspective also supports  
320 the model that location into open chromatin is a critical factor for new gene origination, or at  
321 least transcription, of new genes.

322

## 323 **Discussion**

324 Here we combine the first chromatin state analysis in *P. pacificus* with taxon-rich nematode  
325 phylogenies to analyze the transcriptional regulation of young genes. We identified eight  
326 chromatin states that partition the genome into varying levels of repression, transcription, or  
327 regulatory elements. Expressed young genes were found in open chromatin states, supporting a  
328 widely held model of new gene origination. But to our surprise, young gene 5' ends are more  
329 similar to enhancers than traditional promoters. We also analyzed young gene transcript  
330 structure by long read Iso-Seq, which revealed a unique signature for each evolutionary gene  
331 class. Finally, a bipartite pattern of active histone marks in Chromosome I provides molecular  
332 evidence of an ancient chromosomal fusion event ~180 mya. The ability to probe over 20,000  
333 high confidence promoters and enhancers will be a valuable resource for future mechanistic  
334 studies, especially when combined with the powerful array of genetic, phylogenetic, and  
335 ecological tools recently available to *P. pacificus*.

336           The origin and subsequent regulation of orphan genes is a widely debated topic that has  
337 garnered several theoretical models. Among these are that orphan genes can be transcribed by

338 integrating into open chromatin, or near gene promoters, effectively hijacking their regulatory  
339 sequences and thereby mitigating the need to evolve them *de novo* (Long et al. 2013;  
340 Kaessmann 2010; Kaessmann et al. 2009; Tautz and Domazet-Lošo 2011; McLysaght and  
341 Hurst 2016; Chen et al. 2012b; 2012a). This model was supported by analyzing the position of  
342 transcribed retrogenes in the human genome (Vinckenbosch et al. 2006), however it was also  
343 demonstrated that such integration is often deleterious to the host genes. Indeed a recent  
344 analysis found only ~14% of mammalian retrogenes utilized pre-existing promoters (Carelli et al.  
345 2016). Nevertheless, to date very few tests of these predictions have been performed beyond  
346 retrogenes, and the identify of *cis*-regulatory elements has traditionally been inferred through  
347 spatial proximity to genes or known regulatory sequences like TFBs and CpG islands (Betrán  
348 and Long 2003; Chen et al. 2012b; Ni et al. 2012; Carvunis et al. 2012; Li et al. 2016; Ruiz-  
349 Orera et al. 2015). The phylogenetic diversity of nematode genomes and our recent  
350 chromosome-scale genome (Rödelsperger et al. 2017) allowed us to query all orphan genes in  
351 *P. pacificus*, including but not limited to retrogenes. By applying ChIP-seq and ATAC-seq we  
352 could then interrogate the functional *P. pacificus* genome, including *cis* and *trans* enhancers,  
353 and up to eight different chromatin states. The data presented herein support a model of orphan  
354 gene integration into open chromatin near enhancers preferentially over promoters.

355         Enhancers were originally thought of as inactive DNA elements that harbor TFBs  
356 (Wasylyk 1988), however over the past decade research from several labs has shown  
357 enhancers exhibit bi-directional transcription by RNA polymerase (Andersson et al. 2014; Chen  
358 et al. 2013a; Lam et al. 2014). Now there is a growing consensus that enhancers and promoters  
359 are similar regulatory elements (Andersson et al. 2015; Kim and Shiekhattar 2015), but  
360 promoters have evolved additional sequences to enforce directionality (Grzechnik et al. 2014) or  
361 increased expression. Indeed, promoters can even function as enhancers for other genes  
362 (Engreitz et al. 2016). Under this paradigm, we propose a model whereby a new gene that  
363 originates near an enhancer and is adaptive, will eventually acquire more sophisticated

364 regulatory architecture – thereby transitioning the enhancer into a promoter (Fig. 5C). This  
365 model is complimentary to ‘proto-promoters’ proposed by the Kaessman lab for 8-9% of rat  
366 expressed retrogenes which have H3K4me3 enrichment in rats but H3K4me1 enrichment in  
367 syntenic non-expressed regions in mouse (Carelli et al. 2016). However, here we show that  
368 proximity of new genes near enhancers correlates with and presumably drives their  
369 transcription, and extend this argument beyond retrogenes as a general feature of new genes. If  
370 these transcripts prove functional, then selection can convert the enhancer to a promoter. This  
371 model potentially solves two problems faced by new genes, (1) expression via origination near  
372 enhancers, and (2) introduction near enhancers, especially *trans* enhancers, as opposed to  
373 promoters, does not require exchange or competition with the pre-existing gene landscape.  
374 Nevertheless, we caution that such interpretation is speculative at this point, and examining and  
375 then experimentally manipulating H3K4me1/3 at syntenic loci in closely related strains and  
376 species is necessary to test these hypotheses.

377         In principle, this model could operate regardless of the method of new gene origination  
378 (*de novo*, duplication and divergence, or retrotransposition). Transcripts from enhancers or  
379 lncRNA promoters generally exhibit less splicing, 3' processing and polyadenylation relative to  
380 protein coding genes (Derrien et al. 2012), and are often digested by the nuclear exosome  
381 (Schlackow et al. 2017). In *de novo* gene evolution, mutations that recruit sequence specific  
382 splicing factors or 3' processing factors such as CPSF73 could stabilize enhancer transcripts  
383 allowing for their translation and potential functionalization. In some cases a *de novo* gene that  
384 is acting as a functional ncRNA, referred to as ‘moonlighting’, could lead to greater expression  
385 and a greater window of time to accrue such mutations (Jalali et al. 2016). New genes formed  
386 by duplication and insertion, or retrotransposition near an enhancer could similarly be  
387 transcribed, but without the parental regulatory architecture. In this new genomic context the  
388 gene will likely be expressed in different developmental stages or tissues, possibly providing  
389 new functions. While misregulation of genes often coincides with deleterious effects and

390 disease (Lee and Young 2013), in this case the parental gene is still maintained and operating  
391 under normal control, while the copied gene is freed, within limits (Geiler-Samerotte et al. 2011),  
392 to explore neofunctionalization.

393         Compared to the current reference annotation ('El Paco'), our Iso-Seq annotation  
394 identified shorter genes with fewer exons. The distribution was more similar to *C. elegans* gene  
395 structures. Nevertheless, we note that there are still substantially more genes in *P. pacificus*  
396 greater than 10 exons compared to *C. elegans* (Supplemental Fig. S1F), arguing that further  
397 refinement is still required (although an evolutionary divergence in gene length is formally  
398 possible). We also explored the genetic structure of young vs. old genes. Orphan genes  
399 displayed the shortest gene lengths and fewest exons, and *C. elegans* 1:1 orthologs were the  
400 longest and contained the most exons. The result that homologs appear to be intermediate in  
401 length and exon number is consistent with a transitional path between old and young genes  
402 (Abrusán 2013; Carvunis et al. 2012). But whether this indicates divergence from old genes, or  
403 *de novo* evolution from young genes is unknown, and likely reflects examples of both. Iso-Seq  
404 also identified that almost a quarter (24%) of expressed genes in *P. pacificus* have multiple  
405 isoforms. Although alternative splicing is often due to stochastic transcriptional errors (Pickrell et  
406 al. 2010), there are important examples of alternative isoforms that affect diverse biological  
407 processes (Baralle and Giudice 2017). Whether the multiplicity of transcripts observed here are  
408 differentially expressed during development or in different environmental conditions, and  
409 ultimately if they are functional, will be the focus of future experiments.

410         Iso-Seq of polyadenylated transcripts and rRNA-depleted total RNA demonstrated that  
411 most young genes are polyadenylated. In mammals noncoding RNAs are un-, or  
412 underpolyadenylated (Derrien et al. 2012), arguing that most new genes in *P. pacificus*  
413 represent coding transcripts. However, retrogene transcripts that contain their polyadenylation  
414 'scar' in the genome may be transcribed directly with a polyA tail, and thus appear as  
415 polyadenylated regardless of whether they have been pseudogenized or not. Nevertheless, our

416 interpretation that they are mostly coding is consistent with a previous investigation of orphan  
417 genes in *P. pacificus* that found appreciable peptide coverage from mass spectrometry and  
418 evidence of negative selection (Prabh and Rödelsperger 2016). Beyond orphan genes,  
419 comparing polyadenylated and total RNA Iso-Seq data sets should also be valuable for  
420 investigating gene structures of long non-coding RNAs (lncRNAs), including antisense ncRNA  
421 that have been shown to affect phenotypic plasticity in *P. pacificus* (Seroby et al. 2016).  
422 Genome-wide, we found most young genes are present in chromosome arms where  
423 recombination and repressive chromatin in nematodes is the highest (Coghlan 2005; The *C.*  
424 *elegans* Sequencing Consortium 1998). However, the approximately 10% of young genes that  
425 are highly transcribed (expression groups 1 and 2) were more randomly distributed. Thus, while  
426 recombination in the arms appears to be a furnace for new gene generation, most of these  
427 genes are repressed (expression groups 3 and 4), and have a higher barrier for  
428 functionalization. This pattern highlights several unresolved questions. In particular, does the  
429 presence in open chromatin reflect rare recombination events or *de novo* origination? Further,  
430 are these transcribed new genes ‘born’ into open chromatin and serve as a template for  
431 evolution, or have they already acquired nascent function and their presence in open chromatin  
432 is a result of translocation to increase their expression? Additional functional genomic  
433 comparisons and synteny analysis may shed light on these questions.

434 At chromosome scale resolution, we observe a double-banding of active histone marks  
435 on Chromosome I, in contrast to all other autosomes in both *P. pacificus* and *C. elegans* (Liu et  
436 al. 2011). Based on previous phylogeny and synteny analysis (Rödelsperger et al. 2017), we  
437 propose this pattern is a remnant from a fusion event that occurred prior to the split between  
438 *Diplogasterida* and *Tylenchida*, estimated at ~ 180 mya (Hedges et al. 2015; Cutter 2008). Then  
439 more recently, this portion broke off in the *Caenorhabditis* lineage. This interpretation  
440 parsimoniously explains the long-standing conundrum that Chromosome V in *C. elegans* has  
441 unusual chromosome ‘arm-like’ characteristics, including relatively high recombination rates and

442 a low density of conserved genes (Barnes et al. 1995; The *C. elegans* Sequencing Consortium  
443 1998; Parkinson et al. 2004). In essence, it looks like a chromosome arm because it is, or more  
444 precisely, was, prior to breaking off. If true, the remarkable stability of histone mark patterns  
445 suggests that chromatin organization per se could serve as a molecular fossil of past genomic  
446 re-arrangements. Perhaps probing chromatin structure in conjunction with recombination rates  
447 could provide a historical record of genome evolution in other nematodes or organisms.

448

## 449 **Methods**

### 450 **Evolutionary Gene Classification**

451 Nematode phylogenies were schematically drawn using data downloaded and analyzed from  
452 Holterman et al. (Holterman et al. 2017) and Van Megen et al. (van Megen et al. 2009).  
453 Evolutionary gene classes were defined in a tiered process. First, we defined conserved genes  
454 by blastp and tblastn analysis; we performed all pairwise searches of *P. pacificus* proteins as  
455 query against 24 nematode protein sets as target, and the proteins of each of the 24 nematodes  
456 as query against *P. pacificus* proteins as target. One blastp hit ( $e < 10^{-3}$ ) in any of these 48  
457 comparisons, or one tblastn hit ( $e < 10^{-5}$ ) using *P. pacificus* proteins as query against any of the  
458 24 nematode species genomes was enough to classify a gene in *P. pacificus* as conserved. Any  
459 gene that did not fit these criteria was defined as a *P. pacificus* 'Orphan gene'. Within the  
460 conserved gene class, we then defined 1:1 orthologs as having the best reciprocal blastp hit ( $e$   
461  $\leq 10^{-3}$ ) between *C. elegans* and *P. pacificus* (sorted by e value, and raw scores were used to  
462 break ties). Conserved genes that were not in this 1:1 ortholog class, but were previously  
463 identified by homology in at least one of the 24 nematodes species, were defined as  
464 'Homologous genes'. We kept the e value cutoff relatively 'high' because of the large  
465 phylogenetic distance between *C. elegans* and *P. pacificus*, and hence more conservative with  
466 respect to our orphan gene lists.

### 467 **Nematode synchronization and collection**

468 *P. pacificus* (PS312) cultures for ChIP-, ATAC-, and RNA-seq were synchronized with bleach  
469 and grown on agar to young adults following (Werner et al. 2017, see Supplemental Methods).  
470 Worm pellets were flash frozen until processing. For Iso-Seq we used mixed-developmental  
471 stage (egg, J2, and J4/young adult) RNA at equimolar ratios.

#### 472 **Native histone ChIP-seq**

473 Native (non-cross-linked) chromatin immunoprecipitation (ChIP) of histone post-translational  
474 modifications was performed by combining nematode nuclear isolation (Steiner et al. 2012) with  
475 native ChIP (Brand et al. 2008). Co-precipitated DNA was PCR-amplified and converted to  
476 Illumina libraries using the TruSeq Nano kit (Illumina) and sequenced on a HiSeq 3000. See  
477 Supplemental Methods for detailed protocol.

#### 478 **ATAC-seq**

479 *omni* ATAC-seq was performed on mixed stage purified nuclei following the Corces et al. 2017  
480 protocol, with a few modifications (see Supplementary Methods) and sequenced on an Illumina  
481 HiSeq 3000.

#### 482 **Iso-Seq**

483 RNA was extracted from different developmental time-points separately using TRIzol Reagent  
484 (Invitrogen), then after the quality control equal amounts of RNA from different time points were  
485 pooled. cDNA synthesis of 'direct' Iso-Seq was performed directly using SMARTer PCR cDNA  
486 Synthesis Kit (Clontech Laboratories), while 'total RNA' was first rRNA-depleted with Ribo-Zero  
487 rRNA Removal Kit (Human/Mouse/Rat) (Illumina), then *in vitro* polyadenylated with Poly(A)  
488 Polymerase (New England Biolabs). Direct and rRNA-depleted cDNA were converted into  
489 SMRTbell libraries following the guidelines provided by Pacific Biosciences. SMRT Link  
490 software version 4.0.0 (Pacific Biosciences) was used to convert subreads to circular consensus  
491 sequences and identify full-length non-chimeric reads, which were mapped to the El Paco  
492 genome using GMAP (Wu et al. 2005). See Supplemental Methods for detailed protocol.

#### 493 **RNA-seq**

494 Whole animal young adult (64-68 hour post-bleaching) frozen pellets were freeze-thawed 3x in  
495 TRIzol Reagent (Invitrogen) before purification, and converted to sequencing libraries with the  
496 NEBNext Ultra Directional RNA-seq for Illumina kit, and sequenced on a HiSeq 3000. See  
497 Supplemental Methods for detailed protocol.

#### 498 **Bioinformatics**

499 All sequencing data was mapped to the El Paco genome (Rödelsperger et al. 2017) using  
500 GMAP (Wu et al. 2005) for Iso-Seq, Bowtie 2 (Ben Langmead and Salzberg 2012) for ChIP and  
501 ATAC-seq, and HISAT2 (Kim et al. 2015) for RNA-seq. Peaks were obtained by MACS2  
502 (Zhang et al. 2008), and only samples containing 50% overlap between replicates were kept.  
503 Overlapping peaks were merged using BEDTools (Quinlan and Hall 2010)(Supplemental Table  
504 1). Coverage plots were calculated using BEDTools or HOMER (Heinz et al. 2010) with merged  
505 replicate files, and plotted in R (R Core Team 2016). Chromatin states were obtained with  
506 ChromHMM (Ernst and Kellis 2012) using merged replicate input files. Distance to nearest  
507 promoter or enhancers were performed with BEDTools. See Supplemental Methods for detailed  
508 procedure of all Bioinformatic steps.

509

510

#### 511 **Data Access**

512  
513 Raw and processed data sets from this study have been submitted to the European Nucleotide  
514 Archive (ENA; <https://www.ebi.ac.uk/ena>) under accession number PRJEB24584.

515

#### 516 **Acknowledgements**

517 We are thankful to current members of the Sommer laboratory and Dr. Talia Karasov for  
518 thoughtful critique of experiments, results, and interpretations. This study was funded by the  
519 Max Planck Society.

520

521 **Author contributions**

522 M.S.W. and R.J.S. conceived and designed all experiments. M.S.W. conducted ChIP- and  
523 ATAC-seq with assistance from T.L.; M.S.W, B.S. and C.L. performed Iso-Seq; N.P. conducted  
524 phylogenetic analysis and prepared evolutionary gene category data sets; M.S.W. performed all  
525 bioinformatic analysis. M.S.W. wrote the manuscript with assistance from R.J.S.

526

527 **Disclosure Declaration**

528 The authors declare no competing conflicts of interest.

529

530 **Figure 1. Comparison of *Pristionchus pacificus* and *Caenorhabditis elegans* and**  
531 **phylogenetic relationship. (A-B). *P. pacificus* is often found in a necromenic relationship with**  
532 **insect hosts, preferentially scarab beetles, in the dormant dauer state. When the beetle dies,**  
533 **worms exit the dauer stage to feed on bacteria that bloom on the decomposing carcass. (C-D)**  
534 ***C. elegans*, the classic nematode model organism is often found in leaf detritus and rotting**  
535 **fruits. Ex. rotting apple photo taken by M.S.W. (E-G) *P. pacificus* has become an important**  
536 **model for developmental (phenotypic) plasticity. Adults can adopt (E) a narrow mouth form with**  
537 **one tooth (stenostomatous, St) that makes them strict bacterial feeders. However, the ‘boom-**  
538 **and-bust’ life cycle creates significant competition for resources, and under crowded conditions**  
539 **adults can develop an alternative mouth form (F) with a wider buccal cavity and an extra tooth**  
540 **(eurystomatous, Eu) that allows them to prey on other nematodes. (G) Shown here is a**  
541 **eurystomatous *P. pacificus* preying on a *C. elegans* larva. (H) Schematic phylogeny of**  
542 **nematodes, generated based on the publications of Holterman et. al 2017, and Van Megen et**  
543 **al. 2009. (I) Break-down of *P. pacificus* genes by evolutionary category: one to one orthology**  
544 **with *C. elegans* (*C. elegans* 1:1) is the most conserved, followed by genes sharing homology**  
545 **with at least one gene from the 24 other nematodes (homologous), and finally genes that are**

546 only found in *Pristionchus* (orphan). All categories were defined by BLASTP homology (e-value  
547  $\leq 0.001$ , see Methods).

548

549

550 **Figure 2. Long read RNA sequencing (Iso-Seq) improves gene annotation, identifies**

551 **alternative splicing and can distinguish different evolutionary gene classes by gene**

552 **structure.** (A) Density distribution of cDNA gene lengths between the ‘El Paco’ reference (grey)

553 and Iso-Seq annotation (black). The Iso-Seq annotation was derived from guided assembly

554 using Stringtie (Pertea et al. 2016; see Methods), and plots were created using the Density

555 function in R. (B) Density distribution of exons per gene between ‘El Paco’ reference and Iso-

556 Seq annotations. Method and color scheme are similar to (A). (C) Alternatively spliced isoforms,

557 defined as having multiple detected isoforms with the same start and stop coordinates. White

558 column represents genes containing isoforms that have the same exon-intron structure but

559 different splice-sites, and red columns represent genes containing isoforms with different

560 numbers of exons due to intron retention or exon inclusion/exclusion. (D) Example locus of Iso-

561 Seq reads compared to standard short read RNA-seq. Also shown are Iso-Seq assembled

562 isoforms compared to the single reference gene *umm-S259-11.10-mRNA-1*, made in Integrated

563 Genome Viewer (IGV). (E,F) Percent coverage of evolutionary gene classes by Iso-Seq with

564 either the ‘direct’ method (E) or rRNA-depleted total RNA (F). (G-I) Iso-Seq coverage per gene

565 of each evolutionary class in direct (y-axis) compared to total (x-axis). Coverage was

566 determined by BEDTools, and median ratios of direct/total RNA are presented. Dotted line of

567 slope = 1, y intercept = 0 represents equal coverage between methods. (J,K) Similar density

568 distributions of cDNA length and exon number as in (A,B) but for the three evolutionary gene

569 classes.

570

571

572 **Figure 3. The epigenome of *Pristionchus pacificus*.** (A) Chromatin states determined  
573 through a hidden markov model (ChromHMM) clustered by histone modifications and ATAC-  
574 seq, normalized by coverage. Darker blue represents greater enrichment. (B) Candidate  
575 annotation of each chromatin state according to ENCODE/modENCODE data sets (Ernst et al.  
576 2011; The Roadmap Epigenomics Consortium 2015). Repressive chromatin states are divided  
577 into three categories according to standard definitions of constitutive (repressed 3) and  
578 facultative (repressed 1 and 2) heterochromatin. Poised enhancers are defined according to  
579 previous annotations of loci containing H3K27me3 and DNase sensitivity. (C) Genome-wide  
580 distribution of chromatin states, and further clustering into three categories: repressive,  
581 transcribed, or regulatory. (D) Heatmap of indicated histone modifications for promoter  
582 chromatin states, where each line represents a single 6 kb locus centered on the promoter.  
583 Heatmap matrices were generated in HOMER, clustered from highest to lowest enrichment and  
584 plotted in R. (E) Position weight matrices of *de novo* sequence motifs in promoters, queried  
585 using HOMER. Table also includes percent of promoters containing motif, *p* value and matches  
586 to known transcription factors. (F,G) Similar to (D,E) but for enhancer chromatin states. (H)  
587 Average density plots of promoter (light blue) and enhancer (dark blue) locations relative to  
588 gene bodies, extended 5 kb in each direction from their 5' and 3' ends. Density values  
589 measured using HOMER and plotted in excel. (I) Epigenomic data of histone modification ChIP-  
590 seq, ATAC-seq and RNA-seq surrounding the *Ppa-pax3* gene. Input is included as a reference,  
591 and chromatin state annotations are included on the bottom matching the colors in (C). ChIP  
592 and ATAC-seq coverage is autoscaled per sample, and RNA-seq coverage is in log-scale.

593

594 **Figure 4. Chromatin states correlate with expression, but expressed young genes exhibit**  
595 **distinct profiles.** (A) Average expression (fpkm) from two biological replicates of RNA-seq,  
596 plotted for each gene from highest to lowest along the x-axis. Expression categories were  
597 binned according to approximate inflection points. (B) Chromatin state enrichment of each

598 expression category broken down by genetic element (i.e. TSS, UTRs, exons and introns). (C-  
599 E) Similar to (B) but for each evolutionary gene class. (F) Expression of each evolutionary gene  
600 class determined from average RNA-seq fpkms. Stars indicate p value < 0.05, Welch's two-  
601 tailed test. (G-I) Similar to (B-E), but only for highly expressed (group 1 and 2) genes belonging  
602 to each category. (J) Normalized average densities of H3K4me3, H3K4me1, H3K27ac and  
603 ATAC-seq over a 7 kb window centered at 5' ends. Densities were measured in HOMER and  
604 normalized to the highest and lowest values in each gene class.

605

606 **Figure 5. Distance of promoters and enhancers to evolutionary gene classes.** (A) Distance  
607 cumulative frequency distribution of the nearest promoter, or (B) enhancer (active and poised)  
608 to transcription start sites (TSSs) from each evolutionary gene category. (C) Model of new gene  
609 transcriptional regulation. Enhancers exhibit bi-directional transcription, which can lead to new  
610 gene transcription if it originates near an enhancer, either by *de novo* formation or  
611 duplication/insertion. If the new gene provides useful function, selection will occur not only on  
612 protein function, but also the gene structure leading to more exons, and on regulatory elements  
613 to provide more temporal or spatial control, and more or less transcription. Ultimately, evolution  
614 on enhancer sequences will convert it to a traditional promoter.

615

616 **Figure 6: Chromosome-wide distribution of histone modifications reveals distinct**  
617 **patterns for evolutionary gene classes and a double-band pattern on Chr1.** Genome-wide  
618 patterns of histone modifications from ChIP-seq and ATAC-seq presented as a heatmap with  
619 increasing abundance from white to blue, and white to red for RNA-seq (normalized by depth).  
620 Also plotted are gene densities of each evolutionary class binned by expressed (group 1 and 2)  
621 or transcriptionally repressed (group 3 and 4) for each class.

622

623

624

625 **References**

- 626 Abrusán G. 2013. Integration of new genes into cellular networks, and their structural  
627 maturation. *Genetics* **195**: 1407–1417.
- 628 Andersson R, Gebhard C, Miguel-Escalada I, Hoof I, Bornholdt J, Boyd M, Chen Y, Zhao X,  
629 Schmidl C, Suzuki T, et al. 2014. An atlas of active enhancers across human cell types and  
630 tissues. *Nature* **507**: 455–461.
- 631 Andersson R, Sandelin A, Danko CG. 2015. A unified architecture of transcriptional regulatory  
632 elements. *Trends in Genetics* **31**: 426–433.
- 633 Baralle FE, Giudice J. 2017. Alternative splicing as a regulator of development and tissue  
634 identity. *Nature Reviews Molecular Cell Biology* **18**: 437–451.
- 635 Barnes TM, Kohara Y, Coulson A, Hekimi S. 1995. Meiotic recombination, noncoding DNA and  
636 genomic organization in *Caenorhabditis elegans*. *Genetics* **141**: 159–179.
- 637 Baskaran P, Rödelsperger C. 2015. Microevolution of Duplications and Deletions and Their  
638 Impact on Gene Expression in the Nematode *Pristionchus pacificus* ed. D. Dupuy. *PLoS*  
639 *ONE* **10**: e0131136.
- 640 Baskaran P, Rödelsperger C, Prabh N, Serobyann V, Markov GV, Hirsekorn A, Dieterich C.  
641 2015. Ancient gene duplications have shaped developmental stage-specific expression in  
642 *Pristionchus pacificus*. *BMC Evol Biol* **15**: 185.
- 643 Ben Langmead, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**:  
644 357–359.
- 645 Bento G, Ogawa A, Sommer RJ. 2010. Co-option of the hormone-signalling module dafachronic  
646 acid-DAF-12 in nematode evolution. *Nature* **466**: 494–497.
- 647 Betrán E, Long M. 2003. Dntf-2r, a Young *Drosophila* Retroposed Gene With Specific Male  
648 Expression Under Positive Darwinian Selection. *Genetics* **164**: 977–988.
- 649 Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. 2013. Transposition of native  
650 chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding  
651 proteins and nucleosome position. *Nat Methods* **10**: 1213–1218.
- 652 Burki F, Kaessmann H. 2004. Birth and adaptive evolution of a hominoid gene that supports  
653 high neurotransmitter flux. *Nat Genet* **36**: 1061–1063.
- 654 *C. elegans* Sequencing Consortium. 1998. Genome sequence of the nematode *C. elegans*: a  
655 platform for investigating biology. *Science* **282**: 2012–2018.
- 656 Cai J, Zhao R, Jiang H, Wang W. 2008. De Novo Origination of a New Protein-Coding Gene in  
657 *Saccharomyces cerevisiae*. *Genetics* **179**: 487–496.
- 658 Carelli FN, Hayakawa T, Go Y, Imai H, Warnefors M, Kaessmann H. 2016. The life history of

- 659 retrocopies illuminates the evolution of new mammalian genes. *Genome Res* **26**: 301–314.
- 660 Carvunis A-R, Rolland T, Wapinski I, Calderwood MA, Yildirim MA, Simonis N, Charlotheaux B,  
661 Hidalgo CA, Barbette J, Santhanam B, et al. 2012. Proto-genes and *de novo* gene birth.  
662 *Nature* **487**: 370.
- 663 Chen RAJ, Down TA, Stempor P, Chen QB, Egelhofer TA, Hillier LW, Jeffers TE, Ahringer J.  
664 2013a. The landscape of RNA polymerase II transcription initiation in *C. elegans* reveals  
665 promoter and enhancer architectures. *Genome Res* **23**: 1339–1347.
- 666 Chen S, Krinsky BH, Long M. 2013b. New genes as drivers of phenotypic evolution. *Nature*  
667 *Reviews Genetics* **14**: 645–660.
- 668 Chen S, Ni X, Krinsky BH, Zhang YE, Vibranovski MD, White KP, Long M. 2012a. Reshaping of  
669 global gene expression networks and sex-biased gene expression by integration of a  
670 young gene. *EMBO J* **31**: 2798–2809.
- 671 Chen S, Spletter M, Ni X, White KP, Luo L, Long M. 2012b. Frequent Recent Origination of  
672 Brain Genes Shaped the Evolution of Foraging Behavior in *Drosophila*. *Cell Rep* **1**: 118–  
673 132.
- 674 Chen S, Zhang YE, Long M. 2010. New Genes in *Drosophila* Quickly Become Essential.  
675 *Science* **330**: 1682–1685.
- 676 Conesa A, Madrigal P, Tarazona S, Gomez-Cabrero D, Cervera A, McPherson A, Szcześniak  
677 MW, Gaffney DJ, Elo LL, Zhang X, et al. 2016. A survey of best practices for RNA-seq data  
678 analysis. *Genome Biol* **17**: 13.
- 679 Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A,  
680 Heravi-Moussavi A, Kheradpour P, Zhang Z, Wang J, et al. 2015. Integrative analysis of  
681 111 reference human epigenomes. *Nature* **518**: 317–330.
- 682 Coghlan, A. 2005. Nematode genome evolution. WormBook, ed. The *C. elegans* Research  
683 Community, WormBook, doi/10.1895/wormbook.1.15.1, <http://www.wormbook.org>.
- 684 Corces MR, Trevino AE, Hamilton EG, Greenside PG, Sinnott-Armstrong NA, Vesuna S,  
685 Satpathy AT, Rubin AJ, Montine KS, Wu B, et al. 2017. An improved ATAC-seq protocol  
686 reduces background and enables interrogation of frozen tissues. *Nat Methods* **14**: 959–962.
- 687 Cutter AD. 2008. Divergence times in *Caenorhabditis* and *Drosophila* inferred from direct  
688 estimates of the neutral mutation rate. *Mol Biol Evol* **25**: 778–786.
- 689 Derrien T, Johnson R, Bussoti G, Tanzer A, Djebali S, Tilgner H, Guernec G, Martin D, Merkel  
690 A, Knowles DG, et al. 2012. The GENCODE v7 catalog of human long noncoding RNAs:  
691 analysis of their gene structure, evolution, and expression. *Genome Res* **22**: 1775–1789.
- 692 Dieterich C, Clifton SW, Schuster LN, Chinwalla A, Delehaunty K, Dinkelacker I, Fulton L, Fulton  
693 R, Godfrey J, Minx P, et al. 2008. The *Pristionchus pacificus* genome provides a unique  
694 perspective on nematode lifestyle and parasitism. *Nat Genet* **40**: 1193–1198.
- 695 Engreitz JM, Haines JE, Perez EM, Munson G, Chen J, Kane M, McDonel PE, Guttman M,

- 696 Lander ES. 2016. Local regulation of gene expression by lncRNA promoters, transcription  
697 and splicing. *Nature* **539**: 452–455.
- 698 Ernst J, Kellis M. 2012. ChromHMM: automating chromatin-state discovery and  
699 characterization. *Nat Methods* **9**: 215–216.
- 700 Ernst J, Kheradpour P, Mikkelsen TS, Shores N, Ward LD, Epstein CB, Zhang X, Wang L,  
701 Issner R, Coyne M, et al. 2011. Mapping and analysis of chromatin state dynamics in nine  
702 human cell types. *Nature* **473**: 43–49.
- 703 Geiler-Samerotte KA, Dion MF, Budnik BA, Wang SM, Hartl DL, Drummond DA. 2011.  
704 Misfolded proteins impose a dosage-dependent fitness cost and trigger a cytosolic unfolded  
705 protein response in yeast. *Proc Natl Acad Sci USA* **108**: 680–685.
- 706 Grzechnik P, Tan-Wong SM, Proudfoot NJ. 2014. Terminate and make a loop: regulation of  
707 transcriptional directionality. *Trends Biochem Sci* **39**: 319–327.
- 708 Hattori T, Taft JM, Swist KM, Luo H, Witt H, Slattery M, Koide A, Ruthenburg AJ, Krajewski K,  
709 Strahl BD, et al. 2013. Recombinant antibodies to histone post-translational modifications.  
710 *Nat Methods* **10**: 992–995.
- 711 Hedges SB, Marin J, Suleski M, Paymer M, Kumar S. 2015. Tree of life reveals clock-like  
712 speciation and diversification. *Mol Biol Evol* **32**: 835–845.
- 713 Heinz S, Benner C, Spann N, Bertolino E, Lin YC, Laslo P, Cheng JX, Murre C, Singh H, Glass  
714 CK. 2010. Simple combinations of lineage-determining transcription factors prime cis-  
715 regulatory elements required for macrophage and B cell identities. *Mol Cell* **38**: 576–589.
- 716 Holterman M, Karegar A, Mooijman P, van Megen H, van den Elsen S, Vervoort MTW, Quist  
717 CW, Karssen G, Decraemer W, Opperman CH, et al. 2017. Disparate gain and loss of  
718 parasitic abilities among nematode lineages ed. W.O. Wong. *PLoS ONE* **12**: e0185445.
- 719 Jalali S, Gandhi S, Scaria V. 2016. Navigating the dynamic landscape of long noncoding RNA  
720 and protein-coding gene annotations in GENCODE. *Hum Genomics* **10**: 35.
- 721 Kaessmann H. 2010. Origins, evolution, and phenotypic impact of new genes. *Genome Res* **20**:  
722 1313–1326.
- 723 Kaessmann H, Vinckenbosch N, Long M. 2009. RNA-based gene duplication: mechanistic and  
724 evolutionary insights. *Nature Reviews Genetics* **10**: 19–31.
- 725 Khalturin K, Hemmrich G, Fraune S, Augustin R, Bosch TCG. 2009. More than just orphans: are  
726 taxonomically-restricted genes important in evolution? *Trends in Genetics* **25**: 404–413.
- 727 Kieninger MR, Ivers NA, Rödelsperger C, Markov GV, Sommer RJ, Ragsdale EJ. 2016. The  
728 Nuclear Hormone Receptor NHR-40 Acts Downstream of the Sulfatase EUD-1 as Part of a  
729 Developmental Plasticity Switch in *Pristionchus*. *Curr Biol*.
- 730 Kim D, Langmead B, Salzberg SL. 2015. HISAT: a fast spliced aligner with low memory  
731 requirements. *Nat Methods* **12**: 357–360.

- 732 Kim T-K, Shiekhattar R. 2015. Architectural and Functional Commonalities between Enhancers  
733 and Promoters. *Cell* **162**: 948–959.
- 734 Koch C, Moll T, Neuberg M, Ahorn H, Nasmyth K. 1993. A role for the transcription factors  
735 Mbp1 and Swi4 in progression from G1 to S phase. *Science* **261**: 1551–1557.
- 736 Lainé J-P, Singh BN, Krishnamurthy S, Hampsey M. 2009. A physiological role for gene loops in  
737 yeast. *Genes & Development* **23**: 2604–2609.
- 738 Lam MTY, Li W, Rosenfeld MG, Glass CK. 2014. Enhancer RNAs and regulated transcriptional  
739 programs. *Trends Biochem Sci* **39**: 170–182.
- 740 Lee TI, Young RA. 2013. Transcriptional Regulation and Its Misregulation in Disease. *Cell* **152**:  
741 1237–1251.
- 742 Li L, Foster CM, Gan Q, Nettleton D, James MG, Myers AM, Wurtele ES. 2009. Identification of  
743 the novel protein QQS as a component of the starch metabolic network in Arabidopsis  
744 leaves. *The Plant Journal* **58**: 485–498.
- 745 Li Z-W, Chen X, Wu Q, Hagmann J, Han T-S, Zou Y-P, Ge S, Guo Y-L. 2016. On the Origin of  
746 De Novo Genes in Arabidopsis thaliana Populations. *Genome Biol Evol* **8**: 2190–2202.
- 747 Liu T, Rechtsteiner A, Egelhofer TA, Vielle A, Latorre I, Cheung M-S, Ercan S, Ikegami K,  
748 Jensen M, Kolasinska-Zwierz P, et al. 2011. Broad chromosomal domains of histone  
749 modification patterns in *C. elegans*. *Genome Res* **21**: 227–236.
- 750 Long M, VanKuren NW, Chen S, Vibranovski MD. 2013. New Gene Evolution: Little Did We  
751 Know. *Annu Rev Genet* **47**: 307–333.
- 752 Mayer MG, Rödelberger C, Witte H, Riebesell M, Sommer RJ. 2015. The Orphan Gene  
753 dauerless Regulates Dauer Development and Intraspecific Competition in Nematodes by  
754 Copy Number Variation ed. S.K. Kim. *PLoS Genet* **11**: e1005146.
- 755 McLysaght A, Hurst LD. 2016. Open questions in the study of *de novo* genes: what, how and  
756 why. *Nature Reviews Genetics* **17**: 567–578.
- 757 Michael D, Manyuan L. 1999. Intron—exon structures of eukaryotic model organisms. *Nucleic  
758 Acids Res* **27**: 3219–3228.
- 759 Neme R, Tautz D. 2013. Phylogenetic patterns of emergence of new genes support a model of  
760 frequent de novo evolution. *BMC Genomics* **2012 13:1 14**: 117.
- 761 Ni X, Zhang YE, Nègre N, Chen S, Long M, White KP. 2012. Adaptive Evolution and the Birth of  
762 CTCF Binding Sites in the Drosophila Genome ed. H.S. Malik. *PLOS Biol* **10**: e1001420.
- 763 Nishikori S, Hattori T, Fuchs SM, Yasui N, Wojcik J, Koide A, Strahl BD, Koide S. 2012. Broad  
764 Ranges of Affinity and Specificity of Anti-Histone Antibodies Revealed by a Quantitative  
765 Peptide Immunoprecipitation Assay. *Journal of Molecular Biology* **424**: 391–399.
- 766 O'Sullivan JM, Tan-Wong SM, Morillon A, Lee B, Coles J, Mellor J, Proudfoot NJ. 2004. Gene  
767 loops juxtapose promoters and terminators in yeast. *Nat Genet* **36**: 1014–1018.

- 768 Parkinson J, Mitreva M, Whitton C, Thomson M, Daub J, Martin J, Schmid R, Hall N, Barrell B,  
769 Waterston RH, et al. 2004. A transcriptomic analysis of the phylum Nematoda. *Nat Genet*  
770 **36**: 1259–1267.
- 771 Pertea M, Kim D, Pertea GM, Leek JT, Salzberg SL. 2016. Transcript-level expression analysis  
772 of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat Protoc* 11: 1650–1667.
- 773 Pickrell JK, Pai AA, Gilad Y, Pritchard JK. 2010. Noisy Splicing Drives mRNA Isoform Diversity  
774 in Human Cells ed. E.T. Dermitzakis. *PLoS Genet* **6**: e1001236.
- 775 Prabh N, Rödelberger C. 2016. Are orphan genes protein-coding, prediction artifacts, or non-  
776 coding RNAs? *BMC Bioinformatics* 2016 **17**: 226.
- 777 Proudfoot NJ. 2011. Ending the message: poly(A) signals then and now. *Genes & Development*  
778 **25**: 1770–1782.
- 779 Proudfoot NJ, Furger A, Dye MJ. 2002. Integrating mRNA Processing with Transcription. *Cell*  
780 **108**: 501–512.
- 781 Rada-Iglesias A, Bajpai R, Swigut T, Brugmann SA, Flynn RA, Wysocka J. 2011. A unique  
782 chromatin signature uncovers early developmental enhancers in humans. *Nature* **470**: 279–  
783 283.
- 784 Ragsdale EJ, Müller MR, Rödelberger C, Sommer RJ. 2013. A Developmental Switch Coupled  
785 to the Evolution of Plasticity Acts through a Sulfatase. **155**: 922–933.
- 786 Ramani AK, Calarco JA, Pan Q, Mavandadi S, Wang Y, Nelson AC, Lee LJ, Morris Q,  
787 Blencowe BJ, Zhen M, et al. 2011. Genome-wide analysis of alternative splicing in  
788 *Caenorhabditis elegans*. *Genome Res* **21**: 342–348.
- 789 Reinhardt JA, Wanjiru BM, Brant AT, Saelao P, Begun DJ, Jones CD. 2013. De Novo ORFs in  
790 *Drosophila* Are Important to Organismal Fitness and Evolved Rapidly from Previously Non-  
791 coding Sequences ed. E. Betrán. *PLoS Genet* **9**: e1003860.
- 792 Rosso L, Marques AC, Weier M, Lambert N, Lambot M-A, Vanderhaeghen P, Kaessmann H.  
793 2008. Birth and Rapid Subcellular Adaptation of a Hominoid-Specific CDC14 Protein ed. M.  
794 Long. *PLOS Biol* **6**: e140.
- 795 Rota-Stabelli O, Daley AC, Pisani D. 2013. Molecular timetrees reveal a Cambrian colonization  
796 of land and a new scenario for ecdysozoan evolution. *Curr Biol* **23**: 392–398.
- 797 Rödelberger C, Meyer JM, Prabh N, Lanz C, Bemm F, Sommer RJ. 2017. Single-Molecule  
798 Sequencing Reveals the Chromosome-Scale Genomic Architecture of the Nematode Model  
799 Organism *Pristionchus pacificus*. *Cell Rep* **21**: 834–844.
- 800 Ruiz-Orera J, Hernandez-Rodriguez J, Chiva C, Sabidó E, Kondova I, Bontrop R, Marques-  
801 Bonet T, Albà MM. 2015. Origins of De Novo Genes in Human and Chimpanzee ed. J.  
802 Noonan. *PLoS Genet* **11**: e1005721.
- 803 Santos ME, Le Bouquin A, Crumière AJJ, Khila A. 2017. Taxon-restricted genes at the origin of  
804 a novel trait allowing access to a new environment. *Science* **358**: 386–390.

- 805 Schlackow M, Nojima T, Gomes T, Dhir A, Carmo-Fonseca M, Proudfoot NJ. 2017. Distinctive  
806 Patterns of Transcription and RNA Processing for Human lincRNAs. *Mol Cell* **65**: 25–38.
- 807 Serobyán V, Xiao H, Namdeo S, Rödelšperger C, Sieriebriennikov B, Witte H, Röseler W,  
808 Sommer RJ. 2016. Chromatin remodelling and antisense-mediated up-regulation of the  
809 developmental switch gene *eud-1* control predatory feeding plasticity. *Nat Commun* **7**:  
810 12337.
- 811 Silveira AB, Trontin C, Cortijo S, Barau J, Del Bem LEV, Loudet O, Colot V, Vincentz M. 2013.  
812 Extensive Natural Epigenetic Variation at a De Novo Originated Gene ed. M.D.  
813 Purugganan. *PLoS Genet* **9**: e1003437.
- 814 Sinha A, Sommer RJ, Dieterich C. 2012. Divergent gene expression in the conserved dauer  
815 stage of the nematodes *Pristionchus pacificus* and *Caenorhabditis elegans*. *BMC Genomics*  
816 *2012 13:1* **13**: 254.
- 817 Sommer RJ, McGaughan A. 2013. The nematode *Pristionchus pacificus* as a model system for  
818 integrative studies in evolutionary biology. *Molecular Ecology* **22**: 2380–2393.
- 819 Sommer RJ, Streit A. 2011. Comparative Genetics and Genomics of Nematodes: Genome  
820 Structure, Development, and Lifestyle. [http://dxdoi.org/101146/annurev-genet-110410-](http://dxdoi.org/101146/annurev-genet-110410-132417)  
821 [132417](http://dxdoi.org/101146/annurev-genet-110410-132417).
- 822 Stein JC, Yu Y, Copetti D, Zwickl DJ, Zhang L, Zhang C, Chougule K, Gao D, Iwata A,  
823 Goicoechea JL, et al. 2018. Genomes of 13 domesticated and wild rice relatives highlight  
824 genetic conservation, turnover and innovation across the genus *Oryza*. *Nat Genet* **50**: 285–  
825 296.
- 826 Steiner FA, Talbert PB, Kasinathan S, Deal RB, Henikoff S. 2012. Cell-type-specific nuclei  
827 purification from whole animals for genome-wide expression and chromatin profiling.  
828 *Genome Res* **22**: 766–777.
- 829 Tautz D, Domazet-Lošo T. 2011. The evolutionary origin of orphan genes. *Nature Reviews*  
830 *Genetics* **12**: 692–702.
- 831 Trojer P, Reinberg D. 2007. Facultative Heterochromatin: Is There a Distinctive Molecular  
832 Signature? *Mol Cell* **28**: 1–13.
- 833 van Megen H, van den Elsen S, Holterman M, Karssen G, Mooyman P, Bongers T, Holovachov  
834 O, Bakker J, Helder J. 2009. A phylogenetic tree of nematodes based on about 1200 full-  
835 length small subunit ribosomal DNA sequences. *Nematology* **11**: 927–950.
- 836 Vinckenbosch N, Dupanloup I, Kaessmann H. 2006. Evolutionary fate of retroposed gene  
837 copies in the human genome. *Proc Natl Acad Sci USA* **103**: 3220–3225.
- 838 Voss TC, Hager GL. 2014. Dynamic regulation of transcriptional states by chromatin and  
839 transcription factors. *Nature Reviews Genetics* **15**: 69–81.
- 840 Wang J, Lunyak VV, Jordan IK. 2013. BroadPeak: a novel algorithm for identifying broad peaks  
841 in diffuse ChIP-seq datasets. *Bioinformatics* **29**: 492–493.

- 842 Wasylyk B. 1988. Enhancers and transcription factors in the control of gene expression.  
843 *Biochimica et Biophysica Acta (BBA) - Gene Structure and Expression* **951**: 17–35.
- 844 Werner MS, Sullivan MA, Shah RN, Nadadur RD, Grzybowski AT, Galat V, Moskowitz IP,  
845 Ruthenburg AJ. 2017. Chromatin-enriched lncRNAs can act as cell-type specific activators  
846 of proximal gene transcription. *Nat Struct Mol Biol* **24**: 596–603.
- 847 Witte H, Moreno E, Rödelsperger C, Kim J, Kim J-S, Streit A, Sommer RJ. 2015. Gene  
848 inactivation using the CRISPR/Cas9 system in the nematode *Pristionchus pacificus*. *Dev*  
849 *Genes Evol* **225**: 55–62.
- 850 Yi B, Sommer RJ. 2007. The pax-3 gene is involved in vulva formation in *Pristionchus pacificus*  
851 and is a target of the Hox gene *lin-39*. *Development* **134**: 3111–3119.
- 852 Zhang Y, Liu T, Meyer CA, Eeckhoute J, Johnson DS, Bernstein BE, Nusbaum C, Myers RM,  
853 Brown M, Li W, et al. 2008. Model-based analysis of ChIP-Seq (MACS). *Genome Biol* **9**:  
854 R137.
- 855











