

## Title

**The high turnover of ribosome-associated transcripts from *de novo* ORFs produces gene-like characteristics available for *de novo* gene emergence in wild yeast populations**

Éléonore Durand<sup>1\*</sup>, Isabelle Gagnon-Arsenault<sup>1,2</sup>, Johan Hallin<sup>1,2</sup>, Isabelle Hatin<sup>3</sup>, Alexandre K Dubé<sup>1,2</sup>, Lou Nielly-Thibaut<sup>1</sup>, Olivier Namy<sup>3</sup> & Christian R Landry<sup>1,2</sup>

<sup>1</sup> Institut de Biologie Intégrative et des Systèmes, Département de Biologie, PROTEO, Centre de Recherche en Données Massives de l'Université Laval, Pavillon Charles-Eugène-Marchand, Université Laval, G1V 0A6 Québec, QC, Canada

<sup>2</sup> Département de biochimie, microbiologie et bio-informatique, Université Laval, G1V 0A6 Québec, QC, Canada

<sup>3</sup> Institute for Integrative Biology of the Cell (I2BC), CEA, CNRS, Université Paris-Sud, Université Paris-Saclay, 91190 Gif sur Yvette, France.

\* Current address : Université de Lille CNRS, UMR 8198-Evo-Eco-Paleo, Lille, France

**Correspondence to:** christian.landry@bio.ulaval.ca, eleonore.durand@univ-lille.fr

**Running title:** High turnover of translated *de novo* ORFs in yeast populations

**Keywords:** *De novo* gene birth, wild yeast populations, *Saccharomyces paradoxus*

## 1 **Abstract**

2 Little is known about the rate of emergence of genes *de novo*, how they spread in populations  
3 and what their initial properties are. We examined wild yeast (*Saccharomyces paradoxus*)  
4 populations to characterize the diversity and turnover of intergenic ORFs over short evolutionary  
5 time-scales. With ~34,000 intergenic ORFs per individual genome for a total of ~64,000  
6 orthogroups identified, we found *de novo* ORF formation to have a lower estimated turnover rate  
7 than gene duplication. Hundreds of intergenic ORFs show translation signatures similar to  
8 canonical genes. However, they have lower translation efficiency, which could reflect a  
9 mechanism to reduce their production cost or simply a lack of optimization. We experimentally  
10 confirmed the translation of many of these ORFs in laboratory conditions using a reporter assay.  
11 Translated intergenic ORFs tend to display low expression levels with sequence properties that  
12 generally are close to expectations based on intergenic sequences. However, some of the very  
13 recent translated intergenic ORFs, which appeared less than 110 Kya ago, already show gene-  
14 like characteristics, suggesting that the raw material for functional innovations could appear over  
15 short evolutionary time-scales.

16

## 17 **Introduction**

18 The emergence of new genes is a driving force for phenotypic evolution. New genes may arise  
19 from pre-existing gene structures through genome rearrangements leading to gene duplication,  
20 gene fusion or horizontal gene transfer, or *de novo* from previously non-coding regions (Chen et  
21 al. 2013). *De novo* gene birth was considered highly unlikely (Jacob 1977) up until the last  
22 decade when comparative genomics approaches shed light on the role of intergenic regions as a  
23 regular source of new genes (Tautz and Domazet-Lošo 2011; Landry et al. 2015; Schlotterer  
24 2015; McLysaght and Hurst 2016). Compared to other mechanisms, *de novo* gene origination is  
25 a source of complete innovation because genes emerge solely from mutations, not from the  
26 modification of preexisting genes, with preexisting functions (McLysaght and Hurst 2016).

27  
28 Non-coding regions need to go through three major steps to become gene-coding, the first two  
29 occurring in any order. i) The acquisition of an Open Reading Frame (ORF) by mutations  
30 conferring a gain of in-frame start and stop codons, and ii) the acquisition of regulatory sites to  
31 induce transcription and translation of the ORF. The third step is the retention of the expressed  
32 ORF and its selection because it encodes a less toxic or beneficial polypeptide (Schlotterer  
33 2015; Nielly-Thibault and Landry 2018). The subsequent maintenance of the structure by  
34 purifying selection will lead to the gene being shared among species, as we see for groups of  
35 homologous canonical genes. The birth of genes *de novo* could in theory be a frequent process  
36 since numerous ORFs in non-annotated regions are associated with ribosomes, indicating that  
37 they are likely translated and thus have the potential to produce *de novo* polypeptides, which are  
38 the raw material necessary for *de novo* gene birth (Ingolia et al. 2009; Wilson and Masel 2011;  
39 Carvunis et al. 2012; Ruiz-Orera et al. 2014; Lu et al. 2017; Vakirlis et al. 2017; Ruiz-Orera et al.  
40 2018). The different steps could be accelerated in some ways, depending on the genomic  
41 context. For instance, ORFs could emerge in long non-coding RNAs (lncRNAs) with relatively  
42 high pre-existing expression levels that reflect functions unrelated to the newly emerged ORF  
43 (Xie et al. 2012).

44  
45 Many putative *de novo* genes have been identified (McLysaght and Hurst 2016), but there is  
46 generally limited information about their translation and only few have been functionally  
47 characterized (Begun et al. 2006; Levine et al. 2006; Begun et al. 2007; Cai et al. 2008; Zhou et  
48 al. 2008; Knowles and McLysaght 2009; Li et al. 2010; Baalsrud et al. 2017). These young  
49 genes are generally small with a simple intron-exon structure, they are on average less  
50 expressed than canonical genes and they may diverge rapidly compared to older genes (Wolf et  
51 al. 2009; Tautz and Domazet-Lošo 2011). These properties make it challenging to differentiate  
52 *de novo* emerging genes from non-functional ORFs (McLysaght and Hurst 2016). The absence

53 of sequence similarities of a given gene with known genes in other species is not sufficient  
54 evidence for *de novo* origination, since it could also be due to rapid divergence between  
55 orthologs. This confusion resulted in spurious *de novo* origin annotations, especially over longer  
56 evolutionary time-scale (Gubala et al. 2017). One way to overcome the problem is to identify *de*  
57 *novo* genes and the corresponding orthologous non-coding sequences in closely related  
58 populations or species through synteny, which gives access to mutations occurring during the  
59 gene birth process rather than long after the appearance of the *de novo* genes (Begun et al.  
60 2006; Levine et al. 2006; Begun et al. 2007; Cai et al. 2008; Zhou et al. 2008; Knowles and  
61 McLysaght 2009; Li et al. 2010).

62  
63 The process of *de novo* gene birth has been framed under various hypotheses that consider the  
64 role of selection as acting at different time points. The continuum hypothesis involves a gradual  
65 change in characteristics from non-genic to genic and was used to explain patterns related to  
66 sizes of intergenic ORFs (Carvunis et al. 2012). The preadaptation hypothesis predicts extreme  
67 levels of gene-like characteristics in young *de novo* genes, as was observed for intrinsic  
68 structural disorder (Wilson et al. 2017). The two models both depend i) on the distribution of  
69 properties (non gene-like versus gene-like) of random polypeptides within intergenic regions and  
70 ii) whether these properties correlate with the probability that the peptides will have an adaptive  
71 potential. Examining the distribution of properties of novel polypeptides early after their  
72 emergence – before they potentially lose their initial properties – is therefore important to  
73 determine which one of the two models could be supported.

74  
75 Another question of interest is whether local composition along the genome can accelerate gene  
76 birth. The size of intergenic regions, their GC composition and the genomic context (e.g.  
77 spurious transcription) may affect the birth rate of *de novo* genes (Vakirlis et al. 2017; Nielly-  
78 Thibault and Landry 2018). A recent study on different yeast species found *de novo* genes to

79 preferentially emerge in GC-rich genomic regions, in recombination hotspot and near divergent  
80 promoters (Vakirlis et al. 2017). Another feature that may affect emergence, but also loss, of *de*  
81 *novo* genes is mutation rate; differences in mutation rate would affect the overall turnover of *de*  
82 *novo* genes. Finally, because turnover itself may covary with sequence base composition, the  
83 properties of *de novo* genes could also be biased towards specific properties (Nielly-Thibault  
84 and Landry 2018).

85  
86 Here we explore the contribution of intergenic genetic diversity in the emergence and retention  
87 of the raw material for *de novo* gene birth in wild *Saccharomyces paradoxus* populations. We  
88 focus on this yeast species because of its compact genome and close relatedness with the  
89 model species *Saccharomyces cerevisiae*. One advantage of *S. paradoxus* over *S. cerevisiae* is  
90 that the divergence of populations or lineages within species reflects natural events and not  
91 human domestication and human caused admixture since it has not been domesticated  
92 (Charron et al. 2014; Leducq et al. 2016). Most importantly, *S. paradoxus* harbors clearly defined  
93 lineages whose divergence times can be established and offers different levels of divergence  
94 that allow us to investigate recently emerged *de novo* genes. Finally, the use of natural  
95 populations may eventually allow for the connection between the evolution of *de novo* genes and  
96 key evolutionary processes such as adaptation and speciation, which have been intensively  
97 studied in *S. paradoxus* over the past few years (Charron et al. 2014; Naranjo et al. 2015;  
98 Leducq et al. 2016; Eberlein et al. 2017; Leducq et al. 2017; Weiss et al. 2018).

99  
100 Using this model, we characterized the repertoire and turnover of ORFs located in intergenic  
101 regions (named hereafter iORFs), as well as the associated putative *de novo* polypeptides using  
102 ribosome profiling, and examined how the properties of putative polypeptides covary with their  
103 age and expression, and how they compare with those of canonical genes.

104

## 105 **Results**

### 106 A large number of intergenic ORFs segregate in wild *S. paradoxus* populations

107 We first investigated the diversity and turnover of ORFs located in intergenic regions, which we  
108 named iORFs, and their characteristics in wild *S. paradoxus* strains (Supplementary  
109 information). Because eukaryotic genomes are pervasively transcribed (David et al. 2006;  
110 Pelechano et al. 2013), and lncRNAs may produce peptides (Ruiz-Orera et al. 2014), we initially  
111 assumed that any iORF could have the ability to be translated and thus, could contribute to the  
112 process of *de novo* gene birth. We used 24 *S. paradoxus* strains that are structured in three  
113 main lineages named *SpA*, *SpB* and *SpC* (Charron et al. 2014; Leducq et al. 2016) and two *S.*  
114 *cerevisiae* strains as outgroups (Fig. 1, see Fig. S1 for strain names). These lineages cover  
115 different levels of nucleotide divergence, ranging from ~ 13 % between *S. cerevisiae* and *S.*  
116 *paradoxus* to ~2.27 % between the *SpB* and *SpC* lineages (Kellis et al. 2003; Leducq et al.  
117 2016).

118  
119 We annotated iORFs as any first start and stop codons in the same reading frame not  
120 overlapping with known features, and with no minimum size (Carvunis et al. 2012; Sieber et al.  
121 2018). We then measured the conservation of iORFs between strains using a conservative  
122 approach (Fig. S1, see Methods and Supplementary information). To understand how the iORF  
123 repertoire changes over a short evolutionary time scale, we also estimated the age of iORFs and  
124 their turnover using ancestral sequence reconstruction (see Methods).

125  
126 We identified between 34,216 and 34,503 iORFs per *S. paradoxus* strain, for a total of 64,225  
127 orthogroups annotated in at least one strain (Table 1 and Supplemental Table S1). We observed  
128 that the iORF repertoire of yeast populations is the result of frequent gains, losses, and size  
129 changes (Supplementary information). 56 % of the most ancient iORFs (detected at N2, Fig 1)

130 are still segregating within *S. paradoxus*, showing the role of wild populations as a reservoir of  
131 iORFs that can be used to address the dynamics of early *de novo* gene evolution.

132

### 133 Hundreds of intergenic ORFs show signatures of active translation

134 We performed ribosome profiling to identify iORFs that are potentially translated and that thus  
135 possibly produce polypeptides. Only iORFs with a minimum size of 60 bp were considered for  
136 this analysis. Among them, 12 that displayed a significant blast hit when searched in the  
137 proteomes of 417 species, including 237 fungi, were removed for the downstream analysis (see  
138 Methods). The final set examined consisted of 19,689 iORFs. We prepared ribosome profiling  
139 sequencing libraries for four strains, one belonging to each lineage or species: YPS128 (*S.*  
140 *cerevisiae*), YPS744 (*SpA*), MSH-604 (*SpB*) and MSH-587-1 (*SpC*), in two biological replicates.  
141 All strains were grown in synthetic oak exudate (SOE) medium (Murphy et al. 2006) to  
142 approximate the natural conditions in which *de novo* genes could emerge in wild yeast strains.

143

144 Typically, a ribosome profiling density pattern is characterized by a strong initiation peak located  
145 at the start codon followed by a trinucleotide periodicity at each codon of protein-coding ORFs.  
146 We used this feature to identify a set of iORFs that are the most likely to be translated and we  
147 compared their translation intensity with that of annotated genes. We first detected peaks of  
148 initiation sites around start codons. As expected, the number of ribosome profiling reads located  
149 at this position is on average lower for iORFs than for annotated genes (Fig. 2A). However, we  
150 observed an overlap between the two read density distributions, illustrating a similar read density  
151 between highly expressed iORFs and lowly expressed genes. We observed an initiation peak for  
152 73.9 to 87.9 % of standard annotated genes depending on the strain, and for 1.4 to 6.9 % of  
153 iORFs (Table 3 and Fig. 2B). Detected peaks were classified using three levels of precision and  
154 intensity: 'p1' for less precise peaks (+/- 1nt relative to the first base of the start codon), 'p2' for  
155 precise peaks (detected at the exact first base of the start codon) and 'p3' for precise peaks with

156 strong initiation signals characterized here by the highest read density in the ORF (see  
157 Methods). Among all iORFs with a detected initiation peak, 30, 35 and 34% respectively belong  
158 to p1, p2 and p3. A comparable repartition (Chi-square test, p-value= 0.59) was observed for  
159 annotated genes with 24, 40 and 36% for each precision group, showing that the precision levels  
160 used in our analysis are reliable.

161  
162 We measured codon periodicity, which is characterized by an enrichment of reads at the first  
163 nucleotide of each codon in the first 50 nt excluding the start codon. As for the start codon  
164 region, the number of ribosome profiling reads is lower for iORFs compared to known genes  
165 (Fig. 2C). Among the features with a detected initiation peak, 91.8 to 94.8% of genes and 29.4 to  
166 41 % of iORFs show a significant codon periodicity per strain (Table 3 and Fig. 2D). The number  
167 of detected translation signals is lower in the *SpB* strain, which is most likely due to a lower  
168 number of reads obtained for this strain and the use of raw read density in this analysis (see  
169 Methods). iORFs with an initiation peak and a significant periodicity in at least one strain were  
170 considered as significantly translated and labeled tORFs, whereas iORFs with no significant  
171 translation signatures were labeled ntORFs. We performed a metagene analysis on annotated  
172 genes and tORFs, which revealed a similar ribosome profiling read density pattern between low  
173 expressed genes and tORFs, and confirmed a distinct codon periodicity with significant  
174 translation signature for tORFs (Fig. 2E and S2). The resulting tORF set contains 418  
175 orthogroups with sizes ranging from 60 to 369 nucleotides. They are represented in all age and  
176 conservation categories, which suggests a continuous emergence of potentially translated ORFs  
177 along the phylogeny (Fig. 2F).

178  
179 We compared our results with those resulting from an alternative method (RiboTaper) that is  
180 based on the quantification of the in-frame nucleotide periodicity to detect *de novo* translated  
181 ORFs (Calviello et al. 2016). Among the 418 tORFs detected in our primary analysis, 170 were



182 also annotated *de novo* with RiboTaper (Fig. S3). Additionally, we detected 373 translated ORFs  
183 private to the RiboTaper method. Compared with tORFs private to our methods, tORFs private  
184 to RiboTaper are also characterized by an overall clear trinucleotidic periodicity but they display  
185 on average weaker initiation peaks as well as an overall lower ribosome profiling read coverage  
186 in the first 50 bp (Fig. S3). Below, we describe the results of the analysis performed with the set  
187 of 418 tORFs detected using our method, and which have been confirmed using the subset of  
188 170 tORFs detected by both methods, as well as with the set of 525 tORFs detected using  
189 RiboTaper (see Supplementary information and Supplemental Table S2). tORFs represent a  
190 small fraction (~2 %) of the 19,689 iORF orthogroups longer than 60 nt. This percentage may be  
191 a conservative estimate because the detection depends on the chosen method the conditions  
192 examined, the filters and the ribosome profiling sequencing depth. However, the number of  
193 tORFs is consistent with pervasive transcription measurements in *S. cerevisiae*, with several  
194 hundreds of transcripts detected in non-annotated genomic regions (David et al. 2006). Overall,  
195 for a genome of about 5,000 genes, the roughly 400 *de novo* tORFs which may produce *de novo*  
196 polypeptides could be an important contribution to the proteome diversity of these natural  
197 populations.

#### 198 Translational buffering acts on intergenic ORFs

199 We compared the expression levels of ancient and recent tORFs with that of known genes to  
200 examine if *de novo* polypeptides display gene-like expression levels. Because the *de novo* gene  
201 birth process under the continuum hypothesis involves an increase of tORF size, we also  
202 compared tORF properties while controlling for size ranges per age. The overlap between the  
203 size distributions of tORFs and genes is at the extremes of both distributions and the number of  
204 tORFs is not large enough to generalize the overall properties of longer tORFs with those of  
205 smaller genes (Fig. 3A).

207

208 We investigated translation and transcription levels using ribosome profiling and total RNA  
209 sequencing. We estimated translation efficiency (TE) per gene and tORF as the ratio of  
210 ribosome profiling footprints (RPFs) over total mRNA normalized read counts. TE values  
211 increase with the number of translating ribosomes per molecule of mRNA, illustrating a more  
212 effective translation per mRNA unit (Ingolia et al. 2009). Note that RPF and total RNA coverages  
213 were calculated on the first 60 nt for both genes and tORFs to reduce the bias introduced by the  
214 high number of reads at the initiation codon compared to the rest of the sequences, which tends  
215 to increase TE estimates in short tORFs compared to longer genes. After this correction, TE  
216 values remain significantly correlated with gene size but the effect is small and should not  
217 interfere in our analysis (Fig. 3E).

218  
219 As expected for intergenic regions, tORFs were less transcribed and translated than genes  
220 (Wilcoxon test, p-values  $< 2.2 \times 10^{-16}$ , Fig. 3A-B). We also observed, on average, a significantly  
221 lower TE (Wilcoxon test, p-value =  $< 2.2 \times 10^{-16}$ , Fig. 3C) for most of tORFs compared to genes,  
222 suggesting that tORFs are less actively translated than genes, even when considering the same  
223 size ranges (Fig. 3C). We note, however, that the longest tORF size range category contains  
224 only one tORF (tORF\_102655), which displays a much higher TE value compared with tORFs  
225 from all other size ranges.

226  
227 More generally, the most highly transcribed tORFs display a more reduced TE compared to  
228 genes (Fig. 3D, ANCOVA, p-value  $< 2.2 \times 10^{-16}$ ). This buffering effect was confirmed when  
229 considering the 525 tORFs detected using RiboTaper, as well as in the subset of 170 tORFs  
230 detected using both methods (Fig. S4). A potential consequence of this post-transcriptional  
231 buffering is a reduction of polypeptides translated per molecule of mRNA. The buffering of highly  
232 transcribed tORFs may be due to a rapid selection to reduce the production of toxic polypeptides  
233 or may simply be a consequence of a recent increase in transcription without a change in

234 features that would increase translation rate (e.g. codon usage). The buffering effect is similar  
235 among tORFs of different ages, with no significant pairwise differences between slopes (data not  
236 shown), which supports the hypothesis of no selection for or against translation. Although, on  
237 average, tORFs have lower expression levels and TE values than genes, we noted a significant  
238 overlap between expression levels and TEs in the two sets, which means that some tORFs have  
239 gene-like expression levels and translation efficiencies.

240

#### 241 Translated intergenic polypeptides display a high variability for gene-like traits

242 A recent study suggested that selection favors pre-adapted *de novo* young genes with high  
243 levels of intrinsic protein structural disorder (ISD). They showed that young *de novo* genes were  
244 more disordered than old genes, whereas random polypeptides in intergenic regions were on  
245 average less disordered (Wilson et al. 2017). This would suggest that young polypeptides with  
246 an adaptive potential would already be biased in terms of structural properties compared to the  
247 neutral expectations based on random sequences. We used our data on within species diversity  
248 to examine whether such features are indeed present among tORFs. We examined the  
249 properties of predicted polypeptides as a function of emergence timing in order to follow  
250 evolution before or at the early beginning of the action of selection. We compared the level of  
251 intrinsic disorder, GC-content and genetic diversity (based on SNPs density) in tORFs as a  
252 function of age and with the properties of annotated known genes. We noted that these  
253 properties were confirmed when considering the 525 tORFs detected using RiboTaper, as well  
254 as in the subset of 170 tORFs detected using both methods (Fig. S5). On average, protein  
255 disorder and GC-content are lower in tORFs than in canonical genes regardless of the tORFs  
256 age (Wilcoxon test, p-values < 0.001, Fig. 4B-C). This pattern was confirmed for most of tORFs  
257 and genes sharing the same size range of 45-100 amino acid long (Fig. 4B-C).

258

259 We examined if SNP density variation along the genome could influence tORF turnover.  
260 Regardless of their ages, tORFs are located in regions displaying a higher SNP density  
261 compared to genes, which is consistent with stronger purifying selection on canonical genes  
262 (Fig. 4D). Moreover, younger tORFs, appearing along the terminal branches, tend to be in  
263 regions with higher SNP rates compared to older ones at N2, even when considering the same  
264 size ranges (Fig. 4D). This may be due to mutation rate variation or differences in evolutionary  
265 constraints acting on tORF in an age specific manner. Older tORFs are not preferentially located  
266 at the proximity of genes where selection may be stronger (Fig. 4G), suggesting that the lower  
267 diversity observed at N2 is mainly due to a lower mutation rate. These observations suggest that  
268 younger tORFs are more likely to occur in rapidly evolving sequences with higher mutation rates.  
269 We performed a multivariate analysis to look for polypeptides with extreme values for multiple  
270 traits as an indicator of their functional potential. We observed a subset of tORFs sharing all  
271 characteristics that are typically considered to be gene-like in both more ancient or recent tORFs  
272 (Fig. 4F). Among them, tORF\_102655, which is the only representative of the longest tORF size  
273 range on Fig. 3 and 4, is characterized by multiple gene-like characteristics with extreme intrinsic  
274 disorder, GC%, SNP rate and TE values (Fig. 3 and Fig. 4). This tORF, acquired along the *SpC*  
275 terminal branch and fixed in all strains of the *SpC* lineage, might be recruited by natural selection  
276 if gene-like characteristics increase its functional potential. Sequences are too similar between  
277 strains to test for purifying selection individually on each tORF. Instead, we estimated the  
278 likelihood of the global dN/dS ratio for two merged set of tORFs, containing ancient tORFs  
279 conserved in all *S. paradoxus* strains (set 1) or tORFs appearing at N1 and conserved between  
280 the *SpB* and *SpC* lineages (set 2). Both sets seem to evolve neutrally without significant  
281 purifying selection (NS p-values). Altogether, tORFs do not display significant purifying selection,  
282 but it appears that as a neutral pool, they provide raw material with gene-like characteristics for  
283 selection to act.

284

285 Some intergenic translated ORFs display strong expression changes between lineages in SOE  
286 conditions

287  
288 Our analysis has so far revealed that natural populations are provided with a regular supply of  
289 *de novo* putative polypeptides in intergenic regions (Table 2) at a rate sufficient to provide  
290 lineages that diverged less than 500,000 years ago with different gene contents. We looked for  
291 lineage-specific emerging putative polypeptides among tORFs based on significant differences  
292 of ribosome profiling coverage between each pair of strains (see Methods). Note that a  
293 translation gain or increase may be due to an iORF gain, a transcription/translation increase, or  
294 both. 33 tORFs display a significant lineage-specific expression increase, with 20, 5 and 8 tORFs  
295 in *SpA*, *SpB* and *SpC* respectively (Fig. 5 and S6). Among them, 24 are lineage-specific, and 16  
296 of those were acquired along terminal branches, like the *SpB*-specific tORF\_70680 (Fig. 5).  
297 Nearly 70 % of strong lineage-specific expression patterns are correlated with the presence of  
298 the tORF in one lineage only. This suggests that iORF turnover (gain and loss of start and stop  
299 codons) mostly explain translation differences and not a lineage expression increase in a region  
300 already containing a conserved iORF for instance. Three tORFs are more expressed in both *SpB*  
301 and *SpC* strains compared to *SpA* and *Scer*, suggesting an event occurring along branch b2  
302 (Fig. 1A and S6). We also detected older expression gain/increase events in *S. paradoxus*  
303 relative to *S. cerevisiae* for 9 tORFs, for instance tORF\_69174 (Fig. 5 and S6).

304  
305 Several tORFs show significant translation using a reporter assay

306 Finally, we selected the 45 tORFs displaying significant translation changes described above to  
307 test for translation using a reporter gene. We chose to cover ancient and recent polypeptide gain  
308 events (i.e. lineage-specific or older events). We used a mutated dihydrofolate reductase gene  
309 (DHFR) as a reporter enzyme to fuse with the tORFs (Tarassov et al. 2008; Freschi et al. 2013).  
310 This enzyme confers resistance to methotrexate (MTX) when expressed at significant levels. We

311 integrated the DHFR coding sequence that excludes the start codon in fusion at the 3' end of the  
312 candidate tORFs in *SpA*, *SpB* and *SpC* genetic backgrounds. We fused the DHFR in the same  
313 reading frame as the tORF (construction tORF\_DHFR\_in\_frame) to test for translation controlled  
314 by the native tORF promoter and most likely translation initiation codon (Fig. 6). We also fused  
315 the DHFR with the tORFs in a different reading frame as a negative control  
316 (tORF\_DHFR\_out\_of\_frame). We then tested the translation of the constructs using cell growth  
317 assay on a medium supplemented with MTX and on a medium supplemented with DMSO as a  
318 control (Fig. 6) (Methods). We also fused the DHFR with 12 canonical genes as positive  
319 controls.

320  
321 We found support for the translation of 26 of the 45 tORFs in at least one strain (Fig 6 and Fig.  
322 S8) and 6 tORFs with a translation signal potentially from a different reading frame, where out of  
323 frame fusion cells grew better on selective medium than in frame fusions (Fig. 6, Fig. S7 and Fig.  
324 S8). Interestingly, four of these tORFs have overlapping iORFs in different reading frames,  
325 which suggests that they could be translated instead of the tORF we were focusing on  
326 (tORF\_230326, tORF\_80553, tORF\_102655 and tORF\_70680, see Fig. S5 and S8). Eleven of  
327 the remaining tORFs display no translation signals and 8 had growth differences in the control  
328 conditions so we could not conservatively detect an effect (Fig. S8). Note that among the  
329 translated tORFs detected using this approach, 13 were identified only by our custom method for  
330 ribosome profiling data.

331 We next asked if translation was conserved between conditions and strains. We compared the  
332 translation of tORFs between the three strains and with the translation pattern observed with  
333 ribosome profiling data for the strains in which the DHFR constructs were successful in all three  
334 backgrounds. We succeed in transforming five tORFs in all lineages (*SpA*, *SpB* and *SpC*), with  
335 translation signals that were consistent with our expression criteria (see Methods). However, we  
336 observed that the expression patterns of the tORFs are likely specific to an environment, for

337 instance in SOE medium, tORF\_7665 was found to be translated in the *SpC* strain, whereas on  
338 the MTX medium, the translation was found only in the *SpB* strain (Fig. S9). Some translation  
339 signals were also conserved between strains and conditions, for instance for tORF\_14438,  
340 which is translated in all three strains in both conditions. These results confirmed the translation  
341 detected by ribosomal profiling and indicate that the transcription and translation of tORFs could  
342 be highly condition specific, at least for the two conditions considered here (note that the DHFR  
343 assay requires a very specific condition). However, the two methods measure slightly different  
344 parameters, for instance steady state protein abundance for the DHFR assay and steady state  
345 mRNA/ribosome association for the RPF data, which could also contribute to the difference in  
346 signals.

347

## 348 **Discussion**

349 To better understand the early stages of *de novo* gene birth, we characterized the properties and  
350 turnover of recently evolving iORFs and their putative peptides over short evolutionary time-  
351 scales using closely related wild yeast populations. The number of iORFs identified almost  
352 doubled when considering within species diversity, which illustrates the possible role of  
353 intergenic diversity and the high turnover in providing molecular innovation. Note that we likely  
354 underestimate the total number of iORFs segregating in *S. paradoxus* genomes because of our  
355 conservative approach to identify a set of unambiguous iORF orthogroups in which we excluded  
356 regions too highly divergent that resulted in poor alignments. We focused on ORFs strictly  
357 located in intergenic regions but it is important to note that they represent only a subset of non-  
358 coding ORFs. Indeed, a recent study has shown that >65% of *de novo* genes arose from  
359 transcript isoforms of ancient genes in *Saccharomyces sensus stricto* (Lu et al. 2017). ORFs  
360 overlapping known genes (in a different reading frame or in the opposite strand) and  
361 pseudogenes may also provide an unneglectable source of ORFs and could be an important

362 contribution to the proteome diversity in wild populations (Ji et al. 2015; Lu et al. 2017; Casola  
363 2018).

364  
365 The repertoire of iORFs within *S. paradoxus* came from ancient iORFs that are still segregating  
366 within *S. paradoxus*, and is regularly supplied with *de novo* iORFs gains. The turnover and  
367 retention of iORFs appear at least partly guided by mutation rate variation affecting the number  
368 of gains and losses, or by size changes with some larger changes. In addition, longer iORFs  
369 were more likely to be submitted to size changes, because of the longer mutational target  
370 between the start and stop codons. The iORF turnover rate is lower than the rate of gene  
371 duplication or gene loss estimated in yeast (not considering whole genome duplication, (Lynch et  
372 al. 2008)) but is high enough to provide closely related lineages with distinct sets of novel ORFs  
373 with coding potential.

374 Among the ~20,000 iORF orthogroups of 60 nt and longer, a small fraction (~2%) showed  
375 translation signatures similar to expressed canonical genes in the single condition we tested.  
376 Among the 418 tORFS detected using our custom methods, 40% (n=170) were confirmed with  
377 the use of another tool (RiboTaper). The detection of translated non canonical ORFs particularly  
378 varies depending on the methods, and may lead to only a subset of shared annotated translated  
379 ORFs detected, probably due to their generally low translation levels (Xiao et al. 2018). Here, we  
380 observed that the use of different methods to detect translation may favor tORFs with different  
381 characteristics. For instance, the analysis performed with RiboTaper showed that this tool has  
382 more power to detect translation signals on less expressed tORFs, with small initiation peaks.  
383 Because we focused on intergenic regions, we gave more importance to translation initiation  
384 signals. However, our analysis on expression and sequence properties were robust to  
385 translation detection methods.

386 We observed a stronger post-transcriptional buffering in the tORFs with the highest transcription,  
387 reflecting either selection against translation or a lack of selection for optimal translation. This



388 buffering was observed with the use of another ribosome profiling sequencing dataset in *S.*  
389 *cerevisiae* (Fig. S10, McManus et al. 2014). The buffering effect was previously observed in  
390 interspecies yeast hybrids, especially for genes that show transcriptional divergence, and was  
391 hypothesized to be the result of stabilizing selection on the amount of proteins produced  
392 (McManus et al. 2014). In our case, the post-transcription buffering effect is similar between  
393 older and younger tORFs, suggesting that selection has instead not been acting or has been too  
394 weak to affect this feature.

395  
396 Consistent with a model in which most tORFs are neutral, the corresponding *de novo*  
397 polypeptide properties are on average close to the expectation for random sequences. However,  
398 the diversity is large enough that some tORFs have gene-like properties, suggesting a small set  
399 of neutrally evolving polypeptides with a potential for new functions. iORF translation signatures  
400 (tORFs) were detected for both ancient and recent iORFs and are represented in all  
401 conservation groups. This illustrates that there are regular gains and losses of tORFs along the  
402 phylogeny. The overall absence of purifying selection acting on tORFs suggests a neutral  
403 evolution of most intergenic polypeptides, as observed in rodents (Ruiz-Orera et al. 2018). A  
404 study recently found that the expression of random sequences are likely to have an effect on  
405 fitness (Neme et al. 2017). By analogy with the fitness effect distribution of new mutations, which  
406 are characterized by a large number of mutations of neutral or small effect and few mutations of  
407 large effect (Bataillon and Bailey 2014), we hypothesized that only a small fraction of tORFs  
408 appearing from random mutations could provide an adaptive advantage strong enough to  
409 display a purifying selection signature early after birth. Given this, the resemblance of tORFs to  
410 random sequences does not entirely preclude any potential molecular function or fitness effect.

411  
412 Recently emerging tORFs along terminal branches are more frequent in regions with a higher  
413 SNP density, whereas older tORFs tend to be located in slowly evolving regions. This

414 observation suggests variable turnover rates depending on the local mutation rate. Regions with  
415 low mutation rates could act as a reservoir of ancient tORFs segregating in the population for a  
416 longer time before being lost. On the other hand, mutation hotspots may allow rapid testing of  
417 many molecular combinations, which could be advantageous in a changing environment. Most  
418 tORFs have a subset of gene-like characteristics, implying that they would require limited  
419 refinement by natural selection to acquire new functions. They belong to ancient and recent  
420 tORF gain events, suggesting that gene-like characteristics may be conserved over longer  
421 evolutionary time scales. These properties could be available immediately for selection to act or  
422 when populations are exposed to a changing environment. In addition, even if for a subset of  
423 tORFs, the properties are getting closer to gene properties, changes are generally small. This  
424 suggests that if they are retained by drift or selection, they provide the raw material to gradually  
425 evolve as in the continuum hypothesis (Carvunis et al. 2012). We identified a recently emerging  
426 tORF that had several gene-like characteristics, suggesting that it is pre-adapted to be  
427 biochemically functional. This example illustrates that the birth of a *de novo* polypeptide may be  
428 immediately accompanied with larger gains of gene-like properties, as in the pre-adapted  
429 hypothesis (Wilson et al. 2017).

430

## 431 **Material and methods**

### 432 Characterization of the intergenic ORFs diversity

433 We investigated intergenic ORF (iORF) diversity in wild *Saccharomyces paradoxus* populations,  
434 which are structured in three main lineages named *SpA*, *SpB* and *SpC* (Charron et al. 2014;  
435 Leducq et al. 2016). The wild *S. cerevisiae* strain YPS128 was used in our experiments and the  
436 reference S288C (version R64-2-1) was added in our analysis for the functional annotation.

437

### 438 *Genome assemblies*

439 New genomes assemblies were performed using high-coverage sequencing data from five, ten  
440 and nine North American strains belonging to lineages *SpA*, *SpB* and *SpC* respectively 1 (Fig.  
441 S1) (Leducq et al. 2016) using IDBA\_UD (Peng et al. 2012). For strain YPS128, raw reads were  
442 kindly provided by J. Schacherer from the 1002 Yeast Genomes project (Peter et al. 2018). We  
443 used the default option for IDBA-UD parameters: a minimum k-mer size of 20 and maximum k-  
444 mer size of 100, with 20 increments in each iteration. Scaffolds were then ordered and  
445 orientated along a reference genome using ABACAS (Assefa et al. 2009), using the `-p` nucmer  
446 parameter. *S. paradoxus* and *S. cerevisiae* scaffolds were respectively aligned along the  
447 reference genome of the CBS432 (Scannell et al. 2011) and S288C (version R64-2-1 from the  
448 Saccharomyces Genome Database (<https://www.yeastgenome.org/>)) strains. Unused scaffolds  
449 in the ordering and longer than 200 bp were also conserved in the dataset for further analysis.

450  
451 *Identification of homologous intergenic regions*  
452 We detected homologous intergenic regions using synteny. Genes were predicted using  
453 Augustus (Stanke et al. 2008) with the complete gene model for the species parameter  
454 “*saccharomyces\_cerevisiae\_S288C*”. Orthologs were annotated using a reciprocal best hit  
455 (RBH) approach implemented in SynChro (Drillon et al. 2014) against the reference S288C  
456 (version R64-2-1) using a delta parameter of 3. We used RBH gene pairs provided by SynChro  
457 and the Clustering methods implemented in Silixx (Miele et al. 2011) to identify conserved  
458 orthologs among the 26 genomes. We selected orthologs conserved among all strains and with  
459 a conserved order to extract orthologous microsyntenic genomic regions  $\geq 100$  nt between each  
460 pair of genes (Fig. S1).

461  
462 *Ancestral reconstructions of intergenic sequences*  
463 We reconstructed ancestral genomic sequences of intergenic regions. Because the divergence  
464 between strains belonging to the same lineage is low, we chose one strain per lineage to

465 estimate the ancestral intergenic sequences at each divergence node between lineages (Fig. S1  
466 and 1A), that is YPS128 (*S. cerevisiae*), YPS744 (*SpA*), MSH-604 (*SpB*) and MSH-587-1 (*SpC*).  
467 The ancestral sequence reconstruction was done using Historian (Holmes 2017), which allows  
468 the reconstruction of ancestral indels in addition to nucleotide sequences. Note that indel  
469 reconstruction is essential here to not introduce artefactual frameshifts in ancestral iORFs, see  
470 below, which depends on the conservation of the same reading frame between the start and the  
471 stop codon. Historian was run with a Jukes-Cantor model and using a phylogenetic tree inferred  
472 from aligned intergenic sequences by PhyML version 3.0 (Guindon et al. 2010) with the Smart  
473 Model Selection (Lefort et al. 2017) and YPS128 as outgroup.

474  
475 *iORF annotation and conservation level*  
476 Orthologous regions identified between each pair of conserved genes in contemporary strains  
477 and their ancestral sequence reconstructions were aligned using Muscle (Edgar 2004) with  
478 default parameters. Intergenic regions with a global alignment of less than 50% of identity  
479 among strains (including gaps) were removed. We annotated iORFs defined as any sequence  
480 between canonical start and stop codons, in the same reading frame and with a minimum size of  
481 three codons, using a custom Python script. Because we are working on homologous aligned  
482 regions, the presence-absence pattern does not suffer from limitation alignment bias occurring  
483 when we are working with short sequences. We extracted a presence/absence matrix based on  
484 the exact conservation of the start and the stop codon in the same reading frame (Fig. S1). iORF  
485 aligned coordinates were then converted to genomic coordinates on the respective genomes of  
486 each strain, and removed if there was any overlap with a known feature annotation, such as  
487 rRNA, a tRNA, a ncRNA, a snoRNA, non-conserved genes and pseudogenes annotated on the  
488 reference S288C (version R64-2-1 <https://www.yeastgenome.org/>). Additional masking was  
489 performed by removing iORFs i) located in a region with more than 0.6 % of sequence identity  
490 with *S. cerevisiae* ncRNA or gene (including pseudogenes and excluding dubious ORFs) from

491 the reference genome, or *Saccharomyces kudriavzevii* and *Saccharomyces eubayanus* genes  
492 (Zerbino et al. 2018), ii) in a low complexity region identified with repeat masker  
493 (<http://www.repeatmasker.org/>) and iii) when local alignments of iORFs +/- 300 bp displayed less  
494 than 60% of identity (including gaps). If an iORF overlapped a masked region detected in only  
495 one strain, it was removed for all the other strains in order to not introduce presence-absence  
496 patterns due to strain specific masking.

497 iORFs that do not overlap a known feature were then classified according to the conservation  
498 level: 1) conserved in both species, 2) specific and conserved within *S. paradoxus*, 3) fixed  
499 within lineages and divergent among, 4) specific and fixed in one lineage, 4) polymorphic in at  
500 least one lineage (Fig. S1).

501  
502 For iORFs with a minimum size of 60 nt, we also performed a sequence similarity search against  
503 the proteome of NCBI RefSeq database (O'Leary et al. 2016) for 417 species in the reference  
504 RefSeq category and the representative fungi RefSeq category (containing 237 fungi species).  
505 iORFs with a significant hit (e-value <  $10^{-3}$ ) were removed to exclude any risks of having an  
506 ancient pseudogene. Among the 19,701 iORFs tested, only 12 displayed a significant hit,  
507 illustrating the stringency of our thresholds for the iORF annotation and filtering above.

508  
509 *Evolutionary history of iORFs*

510 Gain and loss events were inferred by comparing presence/absence patterns between ancestral  
511 nodes and actual iORFs. Because the ancestral reconstruction was done using one strain per  
512 lineage (see above), polymorphic iORFs absent in all the considered strains have been removed  
513 from this analysis. iORFs with no detected ancestral homologs were considered as appearing on  
514 terminal branches. We estimated the rate of iORF gain/substitution on each branch as the  
515 number of iORF gain divided by the number of substitution (*i.e* branch length × sequence size)  
516 and calculated the mean of the four branches. The iORF gain rate per cell per division was

517 estimated by calculating the number of expected substitution per cell per division (from the  
518 substitution rate estimated at  $0.33 \times 10^{-9}$  per site per cell division by Lynch et al. (2008), multiplied  
519 by the iORF gain rate per substitution.

520 The evolution of iORF sizes was inferred by connecting iORFs with their ancestral homologs  
521 along the phylogeny if they shared the same start and/or stop position on aligned intergenic  
522 sequences. iORF sizes of two connected iORFs may be conserved if there are no changes, an  
523 increase or a decrease if there are connected only by the same start or stop position because  
524 the position of the other extremity of the iORFs changed.

525

#### 526 Ribosome profiling and mRNA sequencing libraries

527 Ribosome profiling and mRNA sequencing experiments were conducted with the strains  
528 YPS128 (*S. cerevisiae*) (Sniegowski et al. 2002) and YPS744 (*S. paradoxus*), MSH604 (*S.*  
529 *paradoxus*) and MSH587 (*S. paradoxus*) belonging respectively to groups *SpA*, *SpB* and *SpC*  
530 according to Leducq et al. (2016). We prepared two replicates per strain and library type. The  
531 protocol is described in supplementary methods. Briefly, strains were grown in SOE (Synthetic  
532 Oak Exudate) medium (Murphy et al. 2006). Ribosome profiling footprints were purified using the  
533 protocol described in Baudin-Baillieu et al. (2016) with modifications (see supplementary  
534 methods). The rRNA was depleted in purified ribosome footprints and total mRNA samples using  
535 the Ribo-Zero Gold rRNA Removal Kit for yeast (Illumina) according to the manufacturer's  
536 instructions. Ribosome profiling and total mRNA libraries were constructed using the TruSeq  
537 Ribo Profile kit for yeast (illumina), using manufacturer's instructions starting from fragmentation  
538 and end repair step. Libraries were sequenced with Illumina HiSeq 2500 at The Genome  
539 Quebec Innovation Center (Montreal, Canada).

540

#### 541 Detection of translated iORFs

542 Both total RNA and ribosome profiling sequencing libraries were processed using the same  
543 procedure. Raw sequences were trimmed of 3' adapters using CUTADAPT (Martin 2011). For  
544 RPF data, reads with lengths of 27–33 nucleotides were retained for further analysis as this size  
545 is most likely to represent footprinted fragments. For mRNA, reads with lengths of 27–40  
546 nucleotides were retained. Adapter trimmed reads were aligned to the respective genome of  
547 each sample using Bowtie version 1.1.2 (Langmead et al. 2009) with parameters –best –  
548 chunkmbs 500.

549  
550 We used ribosome profiling reads to identify translated iORFs using a custom method. This  
551 analysis was performed on iORFs longer or equal to 60 nucleotides to detect translation  
552 signatures and codon periodicity on at least 20 codons. Annotated iORFs may be overlapping  
553 because of the three possible reading frames for each strand. Ribosomal speed differences  
554 during translation cause an accumulation of ribosome footprints at specific positions within a  
555 gene (Ingolia 2016). We used ribosome profiling read density, which is typically characterized by  
556 a strong initiation peak located at the start codon followed by a codon periodicity at each codon,  
557 to detect the translated iORF among overlapping ones. For each strain, we performed a  
558 metagene analysis at the start codon region of iORFs and annotated conserved genes to detect  
559 the p-site offset for each read length between 28 and 33 nt. Because the ribosome profiling  
560 density pattern is stronger in highly translated regions, metagene analyses were done using the  
561 two replicates of each strain pooled in one coverage file. Ribosome footprints were mapped to  
562 their 5' ends, and the distance between the largest peak upstream of the start codon and the  
563 start codon itself is taken to be the P-site offset per read length. When comparing annotated  
564 genes and iORFs, we obtained similar P-site offset estimates per read length, which were used  
565 for next analysis. We then extracted the aligned read densities, subtracted by the P-offset  
566 estimates, per iORF or gene for next analyses. Metagene analyses were performed using the

567 metagene, psite and get\_count\_vectors scripts from the Plastid package (Dunn and Weissman  
568 2016), metagene figures were done using R scripts (R Core Team 2013).

569  
570 We identified translation initiation signals from ribosome profiling per base read densities, by  
571 detecting peaks at the start codon using a custom R script. We defined three precision levels of  
572 peak initiation: 'p3' if the highest peak is located at the first nucleotide of the start codon, 'p2'  
573 there is a peak at the first position of the start codon and 'p1' if there is a peak at the first position  
574 of the start codon +/- 1 nucleotide. A minimum of five reads was required for peak detection.  
575 Read phasing was estimated by counting the number of aligned reads at the first, second or the  
576 third position for all codons, excluding the first one, of the considered iORF or gene, to test for a  
577 significant deviation from expected ratio with no periodicity, that is 1/3 of each, with a binomial  
578 test. We applied an FDR correction for multiple testing. A minimum of 15 reads was required for  
579 phasing detection.

580 iORF families or genes with an initiation peak and a significant periodicity, *i.e.* a FDR corrected  
581 p-value < 0.05, in at least one strain were considered as translated and named tORFs.

582 We detected translation signature using the RiboTaper software (Calviello et al. 2016). We used  
583 read lengths for which we obtained the best in frame phasing with annotated genes according to  
584 quality check plots provided by RiboTaper, and which are 30-31 nt for *SpA*, 30-32 for *SpB* and  
585 31-32 for *SpC*, and a P-offset of 13.

586  
587 Differential expression analysis

588 Reads were strand-specifically mapped to tORFs and conserved genes using the coverageBed  
589 command from the bedTools package version 2.26.0 (Quinlan and Hall 2010), with parameter -s  
590 (Supplemental Table S3). We then examined significant tORF expression changes between  
591 strains. The differential expression analysis was performed using DESeq2 (Love et al. 2014).  
592 Significant differences were identified using 5% FDR and 2-fold magnitude. We identified lineage



593 specific expression increase when the expression of the tORFs in the considered lineage was  
594 significantly more expressed than the others strains in all pairwise comparisons. For *SpB-SpC*  
595 increase, we selected tORFs when *SpB* and *SpC* strains were both more expressed than  
596 YPS128 and *SpA*, and *S. paradoxus* increase when all *S. paradoxus* lineages were more  
597 expressed than YPS128.

598  
599 For the visualization of tORF coverages (Fig. 5 and Fig. S6), we extracted the per base  
600 coverage on the same strand using the `genomecov` command from the `bedTools` package  
601 version 2.26.0 (Quinlan and Hall 2010). The normalization was performed by dividing the  
602 perbase coverage of each library with the size factors estimated with DESeq2 (Love et al. 2014).

603  
604 Strain construction for in vivo translation confirmation

605 45 tORFs along with 12 canonical genes (Supplemental Table S5) were tagged with a modified  
606 full-length DHFR — a marker that gives resistance to methotrexate (Tarasov et al. 2008) — in  
607 frame and out of frame (as a control). The tORFs were chosen due to their strong translation  
608 signature differences between lineages as found by the differential expression analysis with  
609 ribosome profiling. If the tORF is indeed expressed, in-frame DHFR-tagged strains should grow  
610 in medium supplemented with methotrexate. This complements the ribosome profiling as an *in*  
611 *vivo* confirmation of tORF expression.

612 DHFR along with a HPH resistance module (on a pAG32-DHFR1,2-3 (synthesized by Synbio  
613 Tech, New Jersey, USA)) were PCR amplified (Kapa Hifi DNA polymerase – Kapa Biosystems  
614 Inc., Wilmington, USA) using primers that, at each end, added homology regions flanking the  
615 stop codon of the tORF of interest (Supplemental Table S4). Forward primers were flush with the  
616 stop codon for the in frame integration, and -2bp for the out of frame one (figure 6A). To fuse the  
617 DHFR with the tORFs, 8  $\mu$ l of the PCR products were then used for transformations in *SpA*

618 (*YPS744*), *SpB* (MSH604) and *SpC* (MSH587-1) (only *SpC* for the canonical genes) according  
619 to the method described in (Bleuven et al. 2018).

620 Successful transformations were confirmed by growth on YPD + 250 µg/ml hygromycin B (HYG)  
621 + 100 µg/ml Nourseothricin (NAT) and by PCR amplification of the region containing the tORF  
622 tagged with DHFR.

623

#### 624 Phenotyping of DHFR-tagged strains

625 Transformed strains were incubated at 30°C in 2ml 96-deepwell plates containing 1ml of liquid  
626 YPD+HYG+NAT for 24h. From there, different 96-arrays were made and the strains were printed  
627 onto solid YPD+HYG+NAT plates (omnitrays) using a robotic platform (BM5-SC1, S&P Robotics  
628 Inc.) with appropriate pin tools (96, 384 and 1536). Plates were incubated two days at 30°C. The  
629 solid media 96-arrays were pinned into 384-arrays and then, into the 1536-array with which the  
630 phenotyping was done. The final 1536-plate was then replicated into the same format on a  
631 second YPD+HYG+NAT plate to get more uniformly sized colonies. Plates were incubated two  
632 days at 30°C between each steps. All strains were present in five or six replicates. To avoid  
633 positional effects of the plate borders, the two outer rows and columns were filled with a control  
634 strain (BY4743 LSM8-DHFR[1,2]/CDC39-DHFR[3]).

635

636 To test for methotrexate resistance, all strains were then transferred to DMSO (control) and MTX  
637 DHFR PCA media (0.67% yeast nitrogen base without amino acids and without ammonium  
638 sulfate, 2% glucose, 2.5% noble agar, drop-out without adenine, methionine and lysine, and 200  
639 µg/mL methotrexate (MTX) diluted in DMSO (or only DMSO in the control medium)). Plates were  
640 incubated at 30°C for four days, after which a second round of MTX selection was performed.  
641 Plates were incubated at 30°C for another four days. Images were taken with an EOS Rebel T5i  
642 camera (Canon) every two hours during the entire course of the experiment. Incubation and  
643 imaging was performed in a splmager custom platform (S&P Robotics Inc.).

644  
645 Images were processed using the `gitter.batch` function in the R package Gitter (Wagih, Parts  
646 2014 – Version 1.1.1). The last image of each experiment was used as a reference image to  
647 ensure accurate identification of colonies at early timepoints. The size after 60 hours of growth  
648 (the 30th image) was extracted and the median was calculated for the replicates, these values  
649 are the base for figure 6B (Supplemental Table S6). In-frame and out of frame strains were  
650 phenotyped together on the same plate to alleviate batch effects. Translation was detected i)  
651 when we observed colony size differences between in-frame and out of frame constructions on  
652 MTX medium with a student t-test ( $p$ -value  $< 0.05$ ), and ii) if both positive controls display colony  
653 sizes of more than 1000 and with similar growth for both controls.

654  
655 Some of the observed results were confirmed by measuring cell growth in a spot-dilution assay.  
656 Briefly, precultures of cells expressing DHFR fused to tORFs of interest were adjusted to an  
657 OD<sub>600</sub>/mL of 1 in water. 5-fold serial dilutions were performed and 6  $\mu$ L of each dilution were  
658 spotted on DMSO and MTX DHFR PCA media. Plates were incubated for five days at 30°C and  
659 imaged each day with an EOS Rebel T3i camera (Canon).

660  
661 Expression and sequence properties  
662 Normalized read counts for ribosome profiling and total mRNA samples were extracted with  
663 DESeq2 software (Love et al. 2014) and we calculated the mean of the two replicates per library  
664 type. Translation efficiency (TE) was calculated as the ratio of RPF over total mRNA normalized  
665 read counts on the first 60 nt. We excluded tORFs and genes with less than 10 total RNA reads  
666 in the first 60 nt for the TE calculation. Slope differences between genes and tORFs were tested  
667 with an ANCOVA. We confirmed the buffering effect on tORFs annotated in the *S. cerevisiae*  
668 reference strain S288C with ribosome profiling and RNA sequencing data obtained in (McManus  
669 et al. 2014) (Fig. S10).

670  
671 The intrinsic disorder was calculated for genes and intergenic tORFs using IUPRED (Dosztanyi  
672 et al. 2005). The SNP rate was calculated for each syntenic intergenic region by dividing the  
673 total number of intergenic SNPs in *S. paradoxus* alignments, by the total number of nucleotides  
674 in the region, as in Agier and Fischer (2012) study for intergenic sequences. We used the  
675 *codeml* program from the PAML package version 4.7 (Yang 2007) to estimate the likelihood of  
676 the dN/dS ratio, using the same procedure as employed by Carvunis et al. (2012) with codon  
677 model 0.  
678 All analyses were conducted and figures were created using python and R (R Core Team 2013).  
679

680 **Data access**

681 High-throughput sequencing data generated in this study have been submitted to the NCBI  
682 BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject>) under accession number  
683 PRJNA400476. Assemblies and annotations are available at  
684 [https://landrylab.ibis.ulaval.ca/?page\\_id=2211](https://landrylab.ibis.ulaval.ca/?page_id=2211).

685

686 **Acknowledgments**

687 We thank G. Charron and the IBIS sequencing platform (B. Boyle) for technical help and A.R  
688 Carvunis, R. Dandage and the reviewers for comments on the manuscript. This project was  
689 funded by a FRQNT Team grant to C.R.L and Xavier Roucou and NSERC discovery grant to  
690 C.R.L. C.R.L. holds the Canada Research Chair in Evolutionary Cell and Systems Biology.

691

692 **Author contributions**

693 E.D and C.R.L conceived the project. E.D, O.N, I.H, and I.G.A designed ribosome profiling  
694 experiments. E.D, I.G.A and I.H performed ribosome profiling experiments. A.K.D, J.H, I.G.A and  
695 C.R.L designed and performed functional validation experiments. E.D performed the  
696 bioinformatics analyses with helpful advices from L.N.T, C.R.L and O.N. E.D wrote the  
697 manuscript with revisions from all authors.

698

699 **Disclosure declaration**

700 The authors have no conflict of interest to declare.

701

## 702 References

- 703
- 704 Agier N, Fischer G. 2012. The mutational profile of the yeast genome is shaped by  
705 replication. *Mol Biol Evol* **29**: 905-913.
- 706 Assefa S, Keane TM, Otto TD, Newbold C, Berriman M. 2009. ABACAS: algorithm-based  
707 automatic contiguation of assembled sequences. *Bioinformatics* **25**: 1968-1969.
- 708 Baalsrud HT, Torresen OK, Hongro Solbakken M, Salzburger W, Hanel R, Jakobsen KS,  
709 Jentoft S. 2017. De novo gene evolution of antifreeze glycoproteins in codfishes  
710 revealed by whole genome sequence data. *Mol Biol Evol*  
711 doi:10.1093/molbev/msx311.
- 712 Bataillon T, Bailey SF. 2014. Effects of new mutations on fitness: insights from models and  
713 data. *Ann N Y Acad Sci* **1320**: 76-92.
- 714 Baudin-Baillieu A, Hatin I, Legendre R, Namy O. 2016. Translation Analysis at the Genome  
715 Scale by Ribosome Profiling. *Methods Mol Biol* **1361**: 105-124.
- 716 Begun DJ, Lindfors HA, Kern AD, Jones CD. 2007. Evidence for de novo evolution of testis-  
717 expressed genes in the *Drosophila yakuba*/*Drosophila erecta* clade. *Genetics* **176**:  
718 1131-1137.
- 719 Begun DJ, Lindfors HA, Thompson ME, Holloway AK. 2006. Recently evolved genes  
720 identified from *Drosophila yakuba* and *D. erecta* accessory gland expressed  
721 sequence tags. *Genetics* **172**: 1675-1681.
- 722 Bleuven C, Dubé AK, Nguyen GQ, Gagnon-Arsenault I, Martin H, Landry CR. 2018. A  
723 collection of barcoded natural isolates of *Saccharomyces paradoxus* to study  
724 microbial evolutionary ecology. *MicrobiologyOpen* doi:DOI:10.1002/mbo3.773.
- 725 Cai J, Zhao R, Jiang H, Wang W. 2008. De novo origination of a new protein-coding gene in  
726 *Saccharomyces cerevisiae*. *Genetics* **179**: 487-496.
- 727 Calviello L, Mukherjee N, Wyler E, Zauber H, Hirsekorn A, Selbach M, Landthaler M,  
728 Obermayer B, Ohler U. 2016. Detecting actively translated open reading frames in  
729 ribosome profiling data. *Nat Methods* **13**: 165-170.
- 730 Carvunis AR, Rolland T, Wapinski I, Calderwood MA, Yildirim MA, Simonis N, Charlotiaux B,  
731 Hidalgo CA, Barbette J, Santhanam B et al. 2012. Proto-genes and de novo gene birth.  
732 *Nature* **487**: 370-374.
- 733 Casola C. 2018. From De Novo to "De Nono": The Majority of Novel Protein-Coding Genes  
734 Identified with Phylostratigraphy Are Old Genes or Recent Duplicates. *Genome Biol*  
735 *Evol* **10**: 2906-2918.
- 736 Charron G, Leducq JB, Landry CR. 2014. Chromosomal variation segregates within incipient  
737 species and correlates with reproductive isolation. *Mol Ecol* **23**: 4362-4372.
- 738 Chen S, Krinsky BH, Long M. 2013. New genes as drivers of phenotypic evolution. *Nat Rev*  
739 *Genet* **14**: 645-660.
- 740 David L, Huber W, Granovskaia M, Toedling J, Palm CJ, Bofkin L, Jones T, Davis RW,  
741 Steinmetz LM. 2006. A high-resolution map of transcription in the yeast genome.  
742 *Proc Natl Acad Sci U S A* **103**: 5320-5325.
- 743 Dosztanyi Z, Csizmok V, Tompa P, Simon I. 2005. IUPred: web server for the prediction of  
744 intrinsically unstructured regions of proteins based on estimated energy content.  
745 *Bioinformatics* **21**: 3433-3434.
- 746 Drillon G, Carbone A, Fischer G. 2014. SynChro: a fast and easy tool to reconstruct and  
747 visualize synteny blocks along eukaryotic chromosomes. *PLoS One* **9**: e92621.

- 748 Dunn JG, Weissman JS. 2016. Plastid: nucleotide-resolution analysis of next-generation  
749 sequencing and genomics data. *BMC Genomics* **17**: 958.
- 750 Eberlein C, Nielly-Thibault L, Maaroufi H, Dube AK, Leducq JB, Charron G, Landry CR. 2017.  
751 The Rapid Evolution of an Ohnolog Contributes to the Ecological Specialization of  
752 Incipient Yeast Species. *Mol Biol Evol* **34**: 2173-2186.
- 753 Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high  
754 throughput. *Nucleic Acids Res* **32**: 1792-1797.
- 755 Freschi L, Torres-Quiroz F, Dube AK, Landry CR. 2013. qPCA: a scalable assay to measure  
756 the perturbation of protein-protein interactions in living cells. *Mol Biosyst* **9**: 36-43.
- 757 Gubala AM, Schmitz JF, Kearns MJ, Vinh TT, Bornberg-Bauer E, Wolfner MF, Findlay GD.  
758 2017. The Goddard and Saturn Genes Are Essential for Drosophila Male Fertility and  
759 May Have Arisen De Novo. *Mol Biol Evol* **34**: 1066-1082.
- 760 Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. 2010. New algorithms  
761 and methods to estimate maximum-likelihood phylogenies: assessing the  
762 performance of PhyML 3.0. *Syst Biol* **59**: 307-321.
- 763 Holmes IH. 2017. Historian: accurate reconstruction of ancestral sequences and  
764 evolutionary rates. *Bioinformatics* **33**: 1227-1229.
- 765 Ingolia NT. 2016. Ribosome Footprint Profiling of Translation throughout the Genome. *Cell*  
766 **165**: 22-33.
- 767 Ingolia NT, Ghaemmaghani S, Newman JR, Weissman JS. 2009. Genome-wide analysis in  
768 vivo of translation with nucleotide resolution using ribosome profiling. *Science* **324**:  
769 218-223.
- 770 Jacob F. 1977. Evolution and tinkering. *Science* **196**: 1161-1166.
- 771 Ji Z, Song R, Regev A, Struhl K. 2015. Many lncRNAs, 5'UTRs, and pseudogenes are  
772 translated and some are likely to express functional proteins. *Elife* **4**: e08890.
- 773 Kellis M, Patterson N, Endrizzi M, Birren B, Lander ES. 2003. Sequencing and comparison of  
774 yeast species to identify genes and regulatory elements. *Nature* **423**: 241-254.
- 775 Knowles DG, McLysaght A. 2009. Recent de novo origin of human protein-coding genes.  
776 *Genome Res* **19**: 1752-1759.
- 777 Landry CR, Zhong X, Nielly-Thibault L, Roucou X. 2015. Found in translation: functions and  
778 evolution of a recently discovered alternative proteome. *Curr Opin Struct Biol* **32**: 74-  
779 80.
- 780 Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-efficient  
781 alignment of short DNA sequences to the human genome. *Genome Biol* **10**: R25.
- 782 Leducq JB, Henault M, Charron G, Nielly-Thibault L, Terrat Y, Fiumera HL, Shapiro BJ,  
783 Landry CR. 2017. Mitochondrial Recombination and Introgression during Speciation  
784 by Hybridization. *Mol Biol Evol* **34**: 1947-1959.
- 785 Leducq JB, Nielly-Thibault L, Charron G, Eberlein C, Verta JP, Samani P, Sylvester K,  
786 Hittinger CT, Bell G, Landry CR. 2016. Speciation driven by hybridization and  
787 chromosomal plasticity in a wild yeast. *Nat Microbiol* **1**: 15003.
- 788 Lefort V, Longueville JE, Gascuel O. 2017. SMS: Smart Model Selection in PhyML. *Mol Biol*  
789 *Evol* **34**: 2422-2424.
- 790 Levine MT, Jones CD, Kern AD, Lindfors HA, Begun DJ. 2006. Novel genes derived from  
791 noncoding DNA in Drosophila melanogaster are frequently X-linked and exhibit  
792 testis-biased expression. *Proc Natl Acad Sci U S A* **103**: 9935-9939.

- 793 Li CY, Zhang Y, Wang Z, Zhang Y, Cao C, Zhang PW, Lu SJ, Li XM, Yu Q, Zheng X et al. 2010. A  
794 human-specific de novo protein-coding gene associated with human brain functions.  
795 *PLoS Comput Biol* **6**: e1000734.
- 796 Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for  
797 RNA-seq data with DESeq2. *Genome Biol* **15**: 550.
- 798 Lu TC, Leu JY, Lin WC. 2017. A Comprehensive Analysis of Transcript-Supported De Novo  
799 Genes in *Saccharomyces sensu stricto* Yeasts. *Mol Biol Evol* **34**: 2823-2838.
- 800 Lynch M, Sung W, Morris K, Coffey N, Landry CR, Dopman EB, Dickinson WJ, Okamoto K,  
801 Kulkarni S, Hartl DL et al. 2008. A genome-wide view of the spectrum of spontaneous  
802 mutations in yeast. *Proc Natl Acad Sci U S A* **105**: 9272-9277.
- 803 Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing  
804 reads. *EMBnetjournal* **17**: 10-12.
- 805 McLysaght A, Hurst LD. 2016. Open questions in the study of de novo genes: what, how and  
806 why. *Nat Rev Genet* **17**: 567-578.
- 807 McManus CJ, May GE, Speakman P, Shteyman A. 2014. Ribosome profiling reveals post-  
808 transcriptional buffering of divergent gene expression in yeast. *Genome Res* **24**: 422-  
809 430.
- 810 Miele V, Penel S, Duret L. 2011. Ultra-fast sequence clustering from similarity networks  
811 with SiLiX. *BMC Bioinformatics* **12**: 116.
- 812 Murphy HA, Kuehne HA, Francis CA, Sniegowski PD. 2006. Mate choice assays and mating  
813 propensity differences in natural yeast populations. *Biol Lett* **2**: 553-556.
- 814 Naranjo S, Smith JD, Artieri CG, Zhang M, Zhou Y, Palmer ME, Fraser HB. 2015. Dissecting  
815 the Genetic Basis of a Complex cis-Regulatory Adaptation. *PLoS Genet* **11**: e1005751.
- 816 Neme R, Amador C, Yildirim B, McConnell E, Tautz D. 2017. Random sequences are an  
817 abundant source of bioactive RNAs or peptides. *Nat Ecol Evol* **1**: 0217.
- 818 Nielly-Thibault L, Landry CR. 2018. Differences between the de novo proteome and its non-  
819 functional precursor can result from neutral constraints on its birth process, not  
820 necessarily from natural selection alone. *bioRxiv*: doi: 10.1101/289330.
- 821 O'Leary NA, Wright MW, Brister JR, Ciuffo S, Haddad D, McVeigh R, Rajput B, Robbertse B,  
822 Smith-White B, Ako-Adjei D et al. 2016. Reference sequence (RefSeq) database at  
823 NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids*  
824 *Res* **44**: D733-745.
- 825 Pelechano V, Wei W, Steinmetz LM. 2013. Extensive transcriptional heterogeneity revealed  
826 by isoform profiling. *Nature* **497**: 127-131.
- 827 Peng Y, Leung HC, Yiu SM, Chin FY. 2012. IDBA-UD: a de novo assembler for single-cell and  
828 metagenomic sequencing data with highly uneven depth. *Bioinformatics* **28**: 1420-  
829 1428.
- 830 Peter J, De Chiara M, Friedrich A, Yue JX, Pflieger D, Bergstrom A, Sigwalt A, Barre B, Freil K,  
831 Llored A et al. 2018. Genome evolution across 1,011 *Saccharomyces cerevisiae*  
832 isolates. *Nature* **556**: 339-344.
- 833 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic  
834 features. *Bioinformatics* **26**: 841-842.
- 835 R Core Team. 2013. R: A language and environment for statistical computing. *R Foundation*  
836 *for Statistical Computing*.
- 837 Ruiz-Orera J, Messeguer X, Subirana JA, Alba MM. 2014. Long non-coding RNAs as a source  
838 of new peptides. *Elife* **3**: e03523.



- 839 Ruiz-Orera J, Verdaguer-Grau P, Villanueva-Canas JL, Messeguer X, Alba MM. 2018.  
840 Translation of neutrally evolving peptides provides a basis for de novo gene  
841 evolution. *Nat Ecol Evol* **2**: 890-896.
- 842 Scannell DR, Zill OA, Rokas A, Payen C, Dunham MJ, Eisen MB, Rine J, Johnston M, Hittinger  
843 CT. 2011. The Awesome Power of Yeast Evolutionary Genetics: New Genome  
844 Sequences and Strain Resources for the *Saccharomyces sensu stricto* Genus. *G3*  
845 (*Bethesda*) **1**: 11-25.
- 846 Schlotterer C. 2015. Genes from scratch--the evolutionary fate of de novo genes. *Trends*  
847 *Genet* **31**: 215-219.
- 848 Sieber P, Platzer M, Schuster S. 2018. The Definition of Open Reading Frame Revisited.  
849 *Trends Genet* **34**: 167-170.
- 850 Sniegowski PD, Dombrowski PG, Fingerman E. 2002. *Saccharomyces cerevisiae* and  
851 *Saccharomyces paradoxus* coexist in a natural woodland site in North America and  
852 display different levels of reproductive isolation from European conspecifics. *FEMS*  
853 *Yeast Res* **1**: 299-306.
- 854 Stanke M, Diekhans M, Baertsch R, Haussler D. 2008. Using native and syntenically mapped  
855 cDNA alignments to improve de novo gene finding. *Bioinformatics* **24**: 637-644.
- 856 Tarassov K, Messier V, Landry CR, Radinovic S, Serna Molina MM, Shames I, Malitskaya Y,  
857 Vogel J, Bussey H, Michnick SW. 2008. An in vivo map of the yeast protein  
858 interactome. *Science* **320**: 1465-1470.
- 859 Tautz D, Domazet-Lošo T. 2011. The evolutionary origin of orphan genes. *Nat Rev Genet* **12**:  
860 692-702.
- 861 Vakirlis NN, Hebert AS, Opulente DA, Achaz G, Hittinger CT, Fischer G, Coon JJ, Lafontaine I.  
862 2017. A molecular portrait of de novo genes in yeasts. *Mol Biol Evol*  
863 doi:10.1093/molbev/msx315.
- 864 Weiss CV, Roop JI, Hackley RK, Chuong JN, Grigoriev IV, Arkin AP, Skerker JM, Brem RB.  
865 2018. Genetic dissection of interspecific differences in yeast thermotolerance. *Nat*  
866 *Genet* doi:10.1038/s41588-018-0243-4.
- 867 Wilson BA, Foy SG, Neme R, Masel J. 2017. Young Genes are Highly Disordered as Predicted  
868 by the Preadaptation Hypothesis of De Novo Gene Birth. *Nat Ecol Evol* **1**: 0146-0146.
- 869 Wolf YI, Novichkov PS, Karev GP, Koonin EV, Lipman DJ. 2009. The universal distribution of  
870 evolutionary rates of genes and distinct characteristics of eukaryotic genes of  
871 different apparent ages. *Proc Natl Acad Sci U S A* **106**: 7273-7280.
- 872 Xiao Z, Huang R, Xing X, Chen Y, Deng H, Yang X. 2018. De novo annotation and  
873 characterization of the translome with ribosome profiling data. *Nucleic Acids Res*  
874 **46**: e61.
- 875 Xie C, Zhang YE, Chen JY, Liu CJ, Zhou WZ, Li Y, Zhang M, Zhang R, Wei L, Li CY. 2012.  
876 Hominoid-specific de novo protein-coding genes originating from long non-coding  
877 RNAs. *PLoS Genet* **8**: e1002942.
- 878 Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* **24**:  
879 1586-1591.
- 880 Zerbino DR, Achuthan P, Akanni W, Amode MR, Barrell D, Bhai J, Billis K, Cummins C, Gall A,  
881 Giron CG et al. 2018. Ensembl 2018. *Nucleic Acids Res* **46**: D754-D761.
- 882 Zhou Q, Zhang G, Zhang Y, Xu S, Zhao R, Zhan Z, Li X, Ding Y, Yang S, Wang W. 2008. On the  
883 origin of new genes in *Drosophila*. *Genome Res* **18**: 1446-1455.
- 884

885

**Table 1. Number of iORFs per conservation group**

Conservation group	iORF family numbers	Proportion (%)
Conserved	3,961	6
Spar	9,315	15
Div	12,750	20
Spe group	22,740	35
Pol	15,459	24
<b>Total</b>	<b>64,225</b>	

**Table 2. Estimated age of iORFs in *S. paradoxus* lineages**

Age (Node or branch) <sup>1</sup>	Total	Numbers > or equal to 60nt <sup>2</sup>	Numbers with translation signature <sup>2</sup>
N2	34,092	8,336	221
N1	6,782	2,664	56
b1 ( <i>SpA</i> )	8,454	3,608	73
b3 ( <i>SpB</i> )	6,860	2,948	13
b4 ( <i>SpC</i> )	5,324	2,235	48
Total without redundancy <sup>2</sup>	61,243	19,689	418

<sup>1</sup> N1 and N2 refer to phylogenetic nodes (see Fig. 1A). b1, b3 and b4 are terminal branches, these categories refer to iORFs absent in ancestral sequences (based on the conservation of the start and stop position in the same reading frame). iORFs present in none of the strains used for reconstruction analysis were removed (see Methods).

<sup>2</sup> The 12 iORFs with significant blastp hits against reference proteomes (see results and Methods) were removed.

886  
887  
888

**Table 3. Detection of translated genes or iORFs**

Strain	Genes peak	Genes phasing <sup>1</sup>	iORFs peak	iORFs phasing <sup>1</sup>
YPS128 ( <i>S. cer</i> )	4,095 (85.7%)	3,874 (94.6%)	83 (6.9%)	29 (34.9%)
YPS744 ( <i>SpA</i> )	4,190 (87.7%)	3,846 (91.8%)	643 (6.7%)	188 (29.4%)
MSH-604 ( <i>SpB</i> )	3,531 (73.9%)	3,287(93.1%)	139 (1.4%)	57 (41.0%)
MSH-587-1 ( <i>SpC</i> )	4,203 (87.9%)	3,985 (94.8%)	472 (4.9%)	190 (40.5%)
Total (without redundancy if shared between strains)	4,573	4,443	1,151	418

<sup>1</sup> Number of iORFs or genes with a significant trinucleotide periodicity in ribosome profiling data among those with an initiation peak

889  
890

891 **Figure legends**

892  
893 **Figure 1. A large pool of iORFs segregate within and among *S. paradoxus* lineages. A)**  
894 Phylogenetic tree of strains used for the reconstruction of ancestral intergenic sequences. Node  
895 and branch names are indicated in orange and grey respectively. **B)** Scheme of the iORFs  
896 annotation procedure (see Methods and Figure S1 for a complete description). Pairs of genes  
897 annotated as syntenic were used to align intergenic genomic regions in which iORFs were  
898 characterized. **C)** Number of annotated iORFs per age group, corresponding to the oldest node  
899 in which they were detected. 'Term' refers to iORFs appearing on terminal branches and being  
900 absent in ancestral reconstructions. iORFs are colored according to their conservation group  
901 (see Methods and Fig. S1): conserved (cons), *S. paradoxus* (Spar) specific and fixed, divergent  
902 (Div), divergent group-specific (DivG) and polymorphic (Pol). iORFs detected only in ancestral  
903 sequences are shown in gray.

904  
905 **Figure 2. A fraction of the iORFs display translation signatures similar to genes. A)**  
906 Distribution of the ribosome profiling read counts for genes (grey) and iORFs (purple) at the start  
907 codon position. **B)** Number of genes (Gen) or iORFs with a detected initiation peak at the start  
908 codon position. Peaks are colored according to the precision of the detection (see Methods),  
909 from the most precise (p3) to the least precise (p1). Genes and iORFs with no peaks detected  
910 are shown in green (p0). **C)** Distribution of the ribosome profiling read counts in the first 50 nt of  
911 iORFs excluding the start codon **D)** Proportions of genes or iORFs with a significant in frame  
912 codon periodicity (read phasing in blue) among genes and iORFs with a detected initiation peak.  
913 Genes and iORFs with no detected phasing are shown in green. **E)** Metagene analysis for  
914 significantly translated highly (HE, left) or lowly (LE, middle) expressed genes (grey), and  
915 intergenic translated ORFs (tORFs) (purple, right). The mean of the 5' read counts is plotted  
916 along the position relative to the start codon for significantly translated genes or tORFs. The  
917 lines of the matrix indicate the normalized coverage of genes or tORFs with significant  
918 translation signatures, with one feature per line. **A-E)** Results for the *SpC* strain MSH-587-1 are  
919 shown (see Fig. S2 for *SpA* and *SpB* results). **F)** Total Number of tORFs per conservation group  
920 per age detected in all sequenced strains.

921  
922 **Figure 3. Putative intergenic polypeptides are less efficiently translated compared to**  
923 **genes. A-C)** Ribosome profiling (RPF), total RNA and translation efficiency (TE) - read counts in  
924 the first 60 nt, normalized to correct for library size differences in  $\log_2$  - are displayed for genes  
925 (Gen) and tORFs depending on their ages (N2, N1 and Term). Significant differences in pairwise  
926 comparisons are displayed above each plot (Wilcoxon test, \*\*\* for p-values < 0.001, \*\* for p-  
927 values < 0.01 and \* for p-values < 0.05). Mean estimates per size range are colored by green  
928 intensities (from pale for low values to dark green for high values) below. Numbers per size  
929 range and age are indicated below the graph. **D)** RPF plotted as a function of total RNA for  
930 tORFs in purple, or genes in grey. **E)** TE plotted as a function of tORF or gene sizes (number of  
931 amino acid residues in  $\log_2$ ). Regression lines are plotted for significant Spearman correlations  
932 (p-values < 0.05). Expression levels were calculated using the mean of the two replicates.

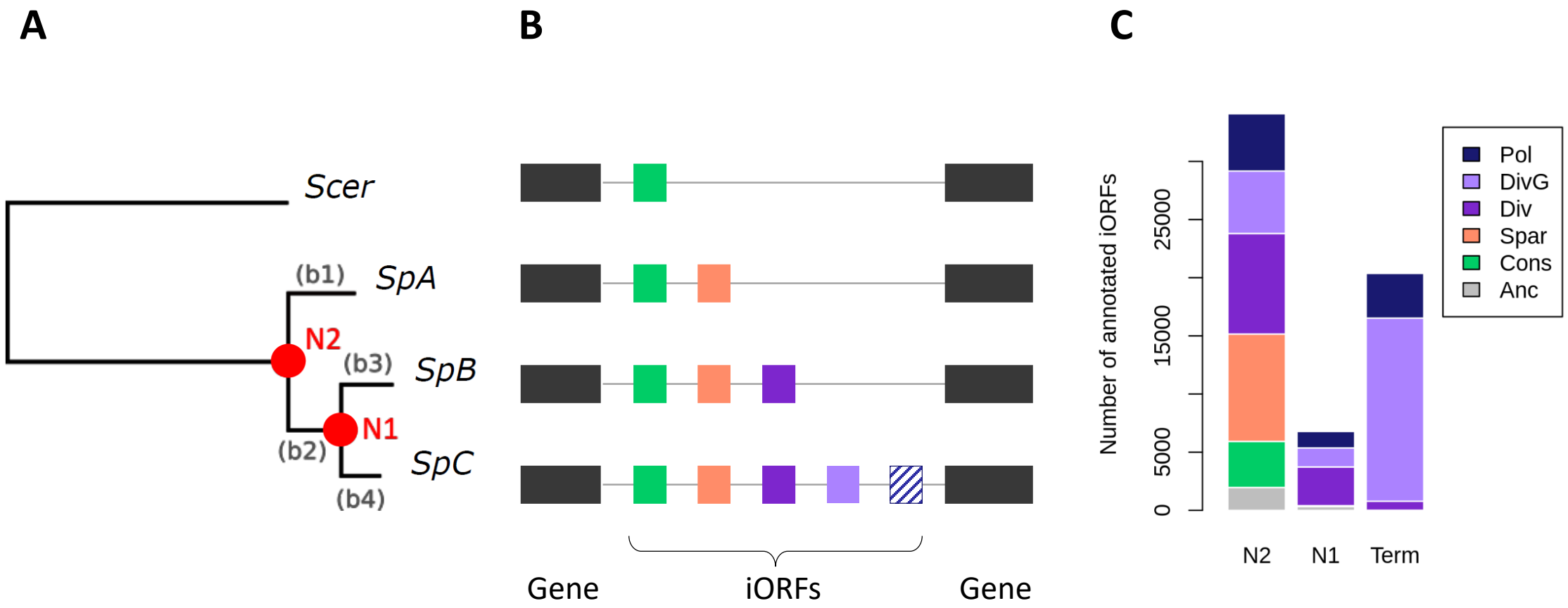
933  
934 **Figure 4. Age-dependent characteristics of intergenic polypeptides. A-E)** Sizes ( $\log_2$   
935 number of residues), mean disorder (ISD), GC %, SNP density and distance to the closest gene

936 are displayed for genes and tORFs as a function of their ages (N2, N1 and Term). Pairwise  
937 significant differences are displayed above each plot (Wilcoxon test, \*\*\* for p-values < 0.001, \*\*  
938 for p-values < 0.01 and \* for p-values <0.05). Mean estimates per size ranges are colored with  
939 green intensities (from pale for low values to dark green high values) below. **F)** Principal  
940 component analysis using the number of residues (SIZE in log<sub>2</sub>), ribosome profiling (RPF), total  
941 RNA (TOT) and translation efficiency (TE) (as read counts in the first 60 nt normalized to correct  
942 for library size differences and in log<sub>2</sub>), intrinsic disorder (ISD), the GC% and SNP density  
943 (SNP). tORFs are colored as a function of their ages. The two first axis explain 33 and 20 % of  
944 the variation (total 53 %).

945  
946 **Figure 5. A continuous emergence of putative polypeptides in *S. paradoxus*.** Normalized  
947 RPF read coverage for a selection of lineage specific (or group specific) tORFs per strain. RPF  
948 read coverages are displayed for replicate 1 and 2 with a blue or pink area respectively. The  
949 positions of all iORFs (including ntORFs and tORFs) in the genomic area are drawn below each  
950 plot. The tORF of interest is labeled with a yellow dot and is plotted in black. iORFs overlapping  
951 the iORF of interest are plotted in black when they are in the same reading frame, and in grey  
952 when they are in a different reading frame as the selected tORF.

953  
954 **Figure 6. DHFR tagging confirms expression of tORFs.** **A)** Conceptual figure of the  
955 approach, 45 tORFs were tagged with a full-length DHFR, in frame or out of frame in *SpA*, *SpB*  
956 and *SpC*, then phenotyped by time-resolved imaging and spot-dilution assays. **B)** Log<sub>2</sub> colony  
957 sizes of strains tagged with DHFR in frame (y-axis) and out of frame (x-axis). The colony size is  
958 taken after 60 hours of growth (shown as a red vertical line in panel A) on medium  
959 supplemented with methotrexate. Colors represent the different strains, the CTRL strains are  
960 tagged in canonical genes, these constructs were made in the *SpC* strain. Dashed line:  $y = x$ . **C)**  
961 Spot-dilution assays further confirm expression of the tORFs, and shows differential expression  
962 of tORF\_159125 and tORF\_153359. 10-fold dilutions go from top to bottom. **B, C)** For the  
963 corresponding controls in medium not supplemented with methotrexate, see Fig. S7.

964  
965  
966



**Figure 1**

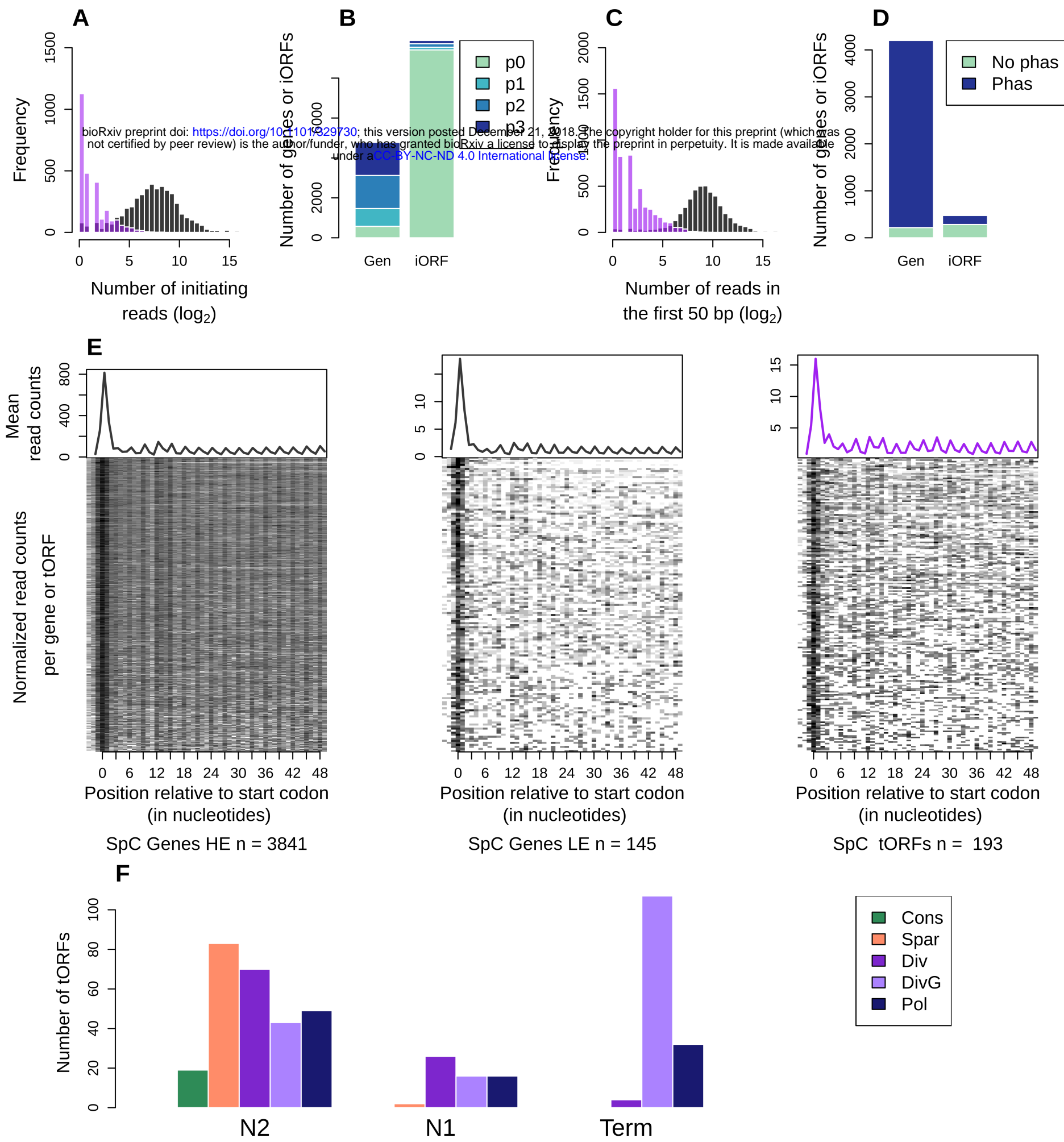


Figure 2



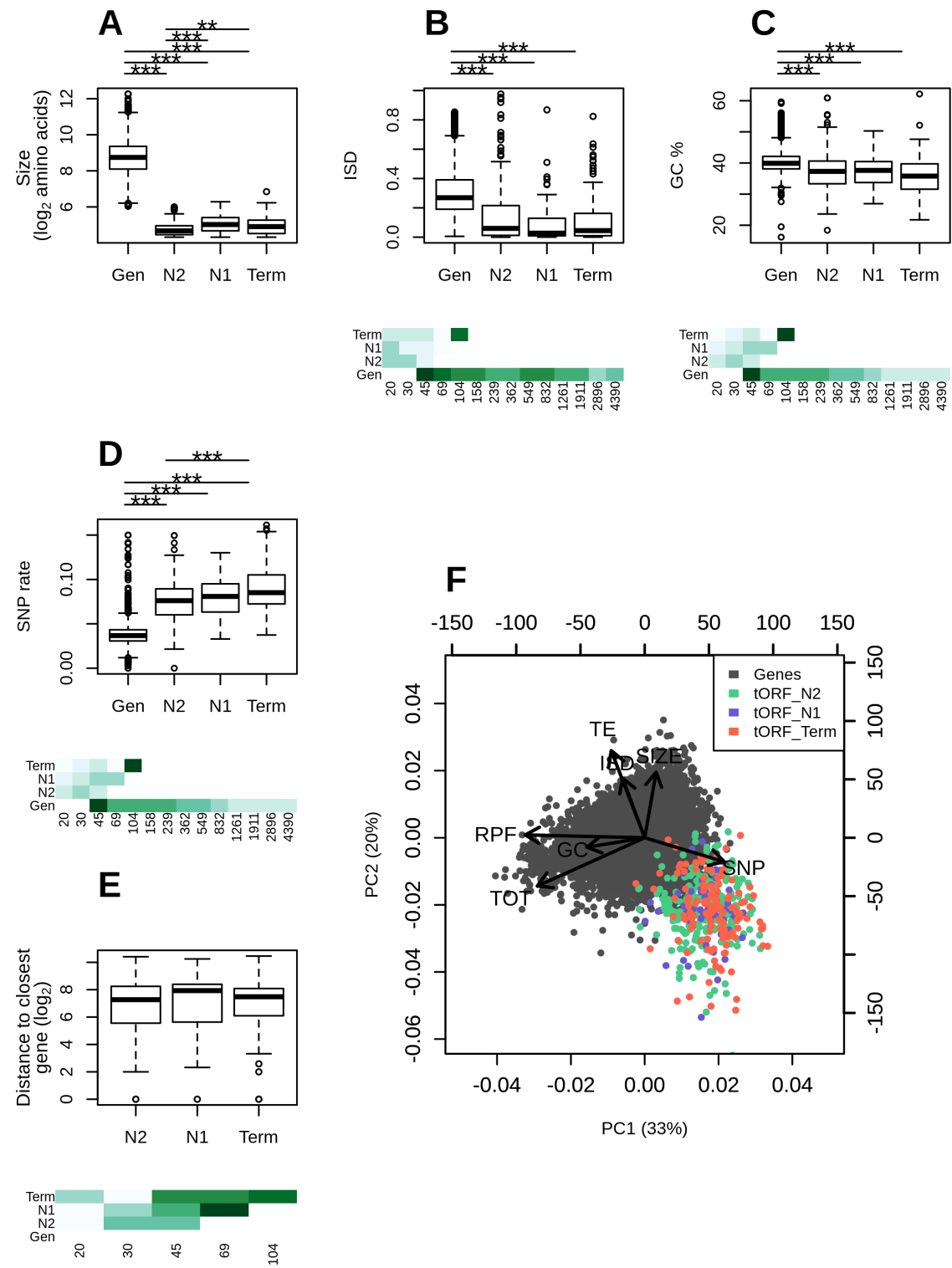


Figure 4



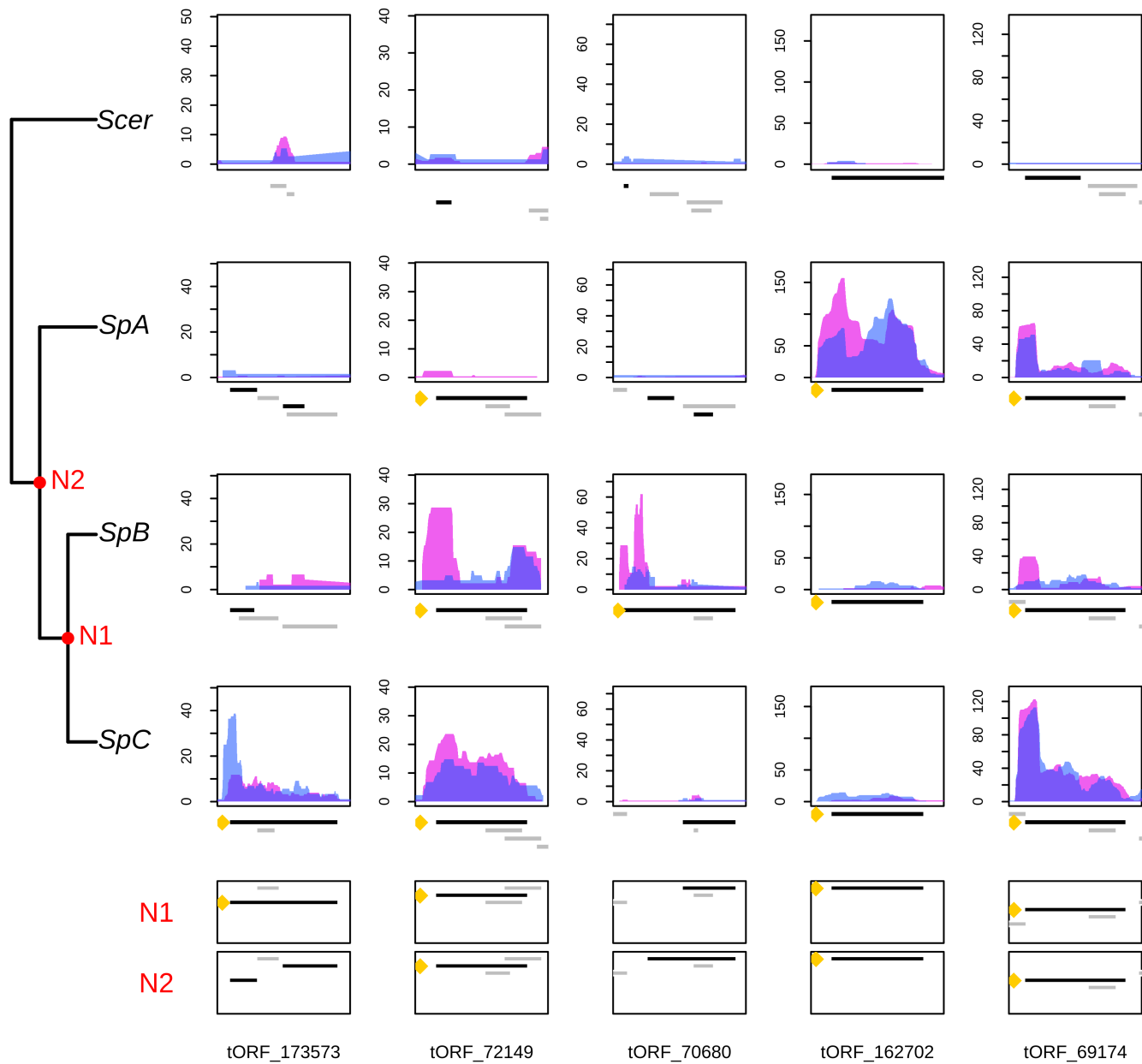


Figure 5

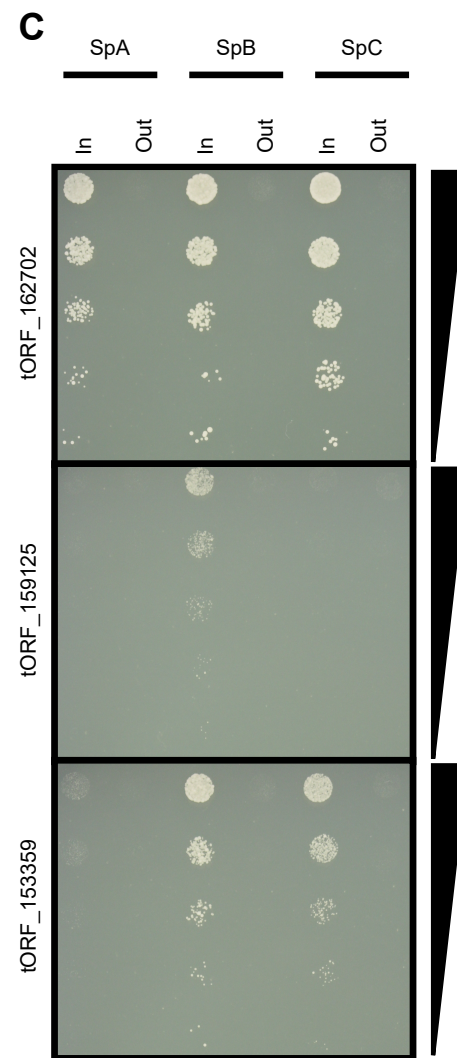
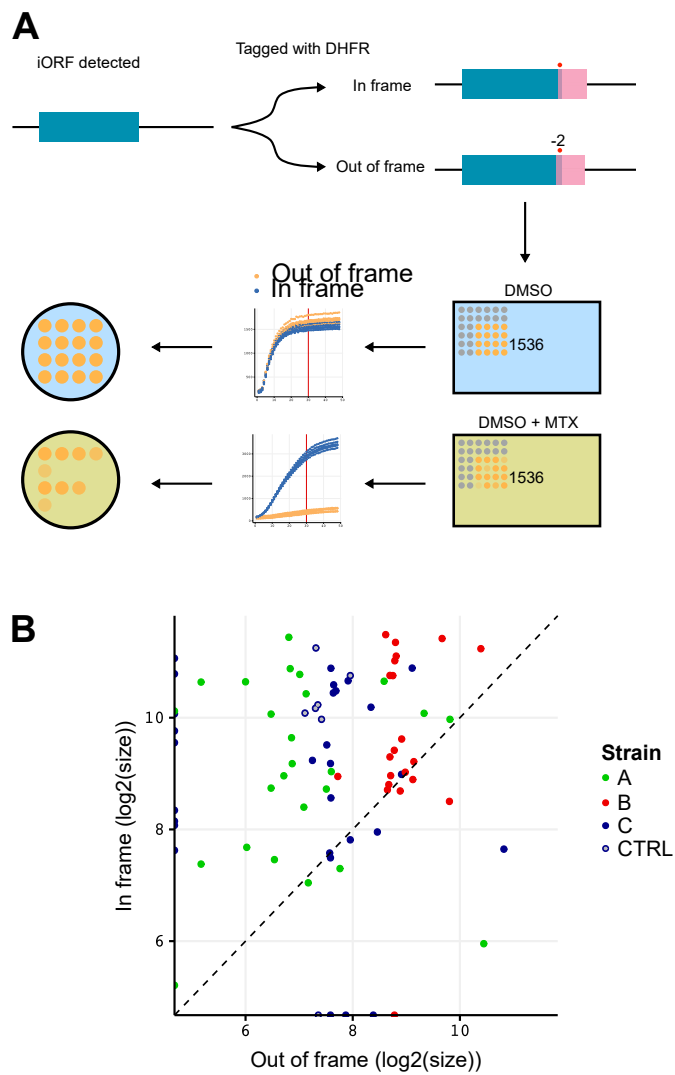


Figure 6