



HAL
open science

An evolutionary model identifies the main evolutionary biases for the evolution of genome-replication profiles

Rossana Droghetti, Nicolas Agier, Gilles Fischer, Marco Gherardi, Marco Cosentino Lagomarsino

► **To cite this version:**

Rossana Droghetti, Nicolas Agier, Gilles Fischer, Marco Gherardi, Marco Cosentino Lagomarsino. An evolutionary model identifies the main evolutionary biases for the evolution of genome-replication profiles. *eLife*, 2021, 10, 10.7554/eLife.63542 . hal-03232379

HAL Id: hal-03232379

<https://hal.sorbonne-universite.fr/hal-03232379>

Submitted on 21 May 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **An evolutionary model identifies the main evolutionary biases for**
2 **the evolution of genome-replication profiles.**

3 Rossana Droghetti

4 *Università degli Studi di Milano, Via Festa del Perdono 7, Milan, Italy*

5 Nicolas Agier and Gilles Fischer

6 *Sorbonne Université, CNRS, Institut de Biologie Paris-Seine,*

7 *Laboratory of Computational and Quantitative Biology, Paris, France*

8 Marco Gherardi

9 *Dipartimento di Fisica, Università degli Studi di Milano, via Celoria 16, Milan, Italy and*

10 *INFN sezione di Milano, via Celoria 16, Milan, Italy*

11 Marco Cosentino Lagomarsino*

12 *IFOM Foundation, FIRC Institute for Molecular Oncology, via Adamello 16, Milan, Italy*

13 *Dipartimento di Fisica, Università degli Studi di Milano, via Celoria 16, Milan, Italy and*

14 *INFN sezione di Milano, via Celoria 16, Milan, Italy*

15 (Dated: May 11, 2021)

16 **Abstract**

17 Recent results comparing the temporal program of genome replication of yeast species belonging
18 to the *Lachancea* clade support the scenario that the evolution of replication timing program could
19 be mainly driven by correlated acquisition and loss events of active replication origins. Using
20 these results as a benchmark, we develop an evolutionary model defined as birth-death process for
21 replication origins, and use it to identify the evolutionary biases that shape the replication timing
22 profiles. Comparing different evolutionary models with data, we find that replication origin birth
23 and death events are mainly driven by two evolutionary pressures, the first imposes that events
24 leading to higher double-stall probability of replication forks are penalized, while the second makes
25 less efficient origins more prone to evolutionary loss. This analysis provides an empirically grounded
26 predictive framework for quantitative evolutionary studies of the replication timing program.

* marco.cosentino-lagomarsino@ifom.eu

27 I. INTRODUCTION

28 Eukaryotes, from yeast to mammals, rely on pre-defined “replication origins” along the
29 genome to initiate replication [1–4], but we still ignore most of the evolutionary principles
30 shaping the biological properties of these objects. Binding by initiation complexes defines
31 origins as discrete chromosomal loci, which are characterized by multiple layers of genomic
32 properties, including the necessary presence of autonomously replicating sequences, nucle-
33 osome depletion, and absence of transcription[5, 6]. Initiation at origins is stochastic, so
34 that different cells of the same population undergoing genome replication in S-phase will
35 typically initiate replication from different origins [7, 8].

36 Initiation from a single origin can be described by intrinsic rates and/or licensing
37 events [9]. Indeed, the genome-wide replication kinetics of a population of cells can be
38 accessed experimentally by different techniques [9–11]. Recent techniques also allow to
39 measure replication progression at the single-cell level [12, 13]. The estimation of key origin
40 parameters from data requires minimal mathematical models describing stochastic origin
41 initiation and fork progression [10, 14–16]. Typically, one can extract from the data origin
42 positions, as well as estimated origin-intrinsic characteristic firing times or rates. Knowledge
43 of origins positions and rates makes it possible to estimate the “efficiency” of an origin, i.e.
44 its probability of actively firing during S-phase, rather than being passively replicated.

45 Over evolution, a genome modifies its replication timing profile by “reprogramming”
46 origin positions and rates in order to maximize fitness, under the constraints of the possible
47 changes of these parameters that are physically and biologically accessible. Little is known
48 about this process, and finding basic rules that drive origin evolution is our main focus
49 here [17]. The main recognized constraint determining negative selection is due to replication
50 forks stalling between adjacent origins [18–20]. If two converging replication forks stall with
51 no origins in between them, it is generally agreed that replication cannot be rescued, and the
52 event leads to cell death. Such deadly “double stalls” can only happen with two converging
53 forks generated from consecutive origins. A pioneering study by Newman and coworkers [18]
54 used a combination of data analysis and mathematical models to understand the role of lethal
55 double stall events on origin placement. They found that the fork per-base stall probability
56 affects the distance between neighbor origins, and the optimal distance distribution tends
57 to a regular spacing, which is confirmed by experimental data. Thus, origin placement is far

58 from a uniform random distribution (which would translate into an exponential distribution
59 of neighbor origin distances). Instead, the regular lattice-like spacing that origin tend to
60 take is reminiscent of particles repelling each other.

61 Due to the streamlined genome and the experimental accessibility, yeasts are interesting
62 systems to study experimentally the evolution of replication programs. However, at the level
63 of the *Saccharomyces* genus, the replication program is highly conserved [21]. Hence, until
64 recently, no experimental account of the evolution of the replication program was available.
65 Our collaboration has recently produced data of this kind [22], by comparing replication
66 dynamics and origin usage of 10 distant *Lachancea* yeast species. This study highlights the
67 dominance of origin birth-death events (rather than e.g. chromosomal rearrangements) as
68 main evolutionary drive of the replication program changes, and characterizes the main prin-
69 ciples underlying origin birth-death events. Briefly, the fate of an origin strongly depends
70 on its neighbourhood, in particular the distance from neighbor origins and their efficiency.
71 Indeed, proximity to efficient origins correlates with weaker origin loss events. An evolu-
72 tionary bias against weak origins could be due to the fact that their presence is neutral
73 or even advantageous (e.g., in terms of reducing double stalls), but their advantage is not
74 sufficiently high for them to survive drift. These findings open the question of capturing the
75 relevant evolutionary biases acting on replication profiles in the framework of the empirical
76 birth-death evolutionary dynamics, for which the data set [22] provides an empirical testing
77 ground.

78 Here, we define a minimal evolutionary birth-death model for replication program evolu-
79 tion encompassing all the empirical observations made by Agier and coworkers [22], and we
80 use it to investigate the main evolutionary trade-offs that could explain the data.

81 II. RESULTS

82 Experimental data motivate an evolutionary model for origins turnover

83 This section presents a reanalysis of the experimental data from ref. [22]. We summarize
84 the main results of that study, and present additional considerations on the same data,
85 which motivate the evolutionary model framework used in the following.

86 Fig. 1 - Supplement 1 recapitulates the *Lachancea* clade phylogenetic tree used in the

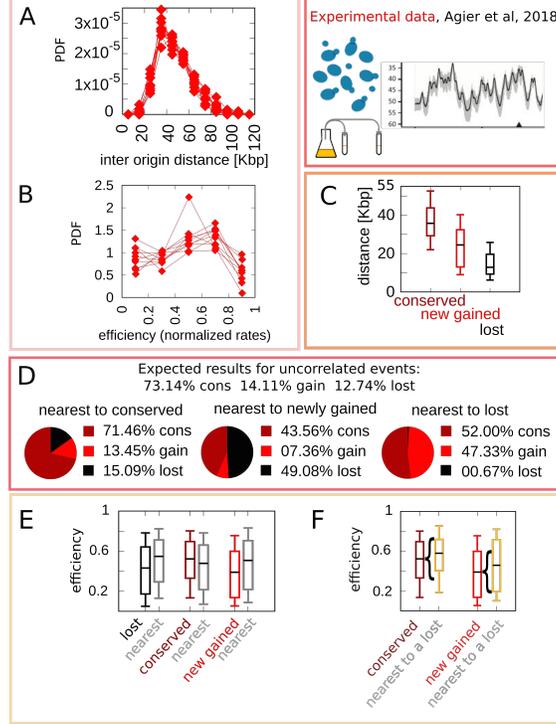


FIG. 1. Experimental data motivate an evolutionary model for replication origins turnover. **A:** Distribution of the distance between neighbor origins in ten *Lachancea* species, each histogram refers to a different species (data from ref. [22]), and all the plots show a marked peak around 35 Kbp. **B:** Distribution of the efficiency (calculated from a fit, using Eq. 4) for all origins in ten *Lachancea* yeast species [22]. **C:** From ref. [22], box plot of the distribution of the distance from the nearest origin split by evolutionary events, for conserved (dark red), newly gained (red) and lost origins (black), estimated comparing six sister species of the *Lachancea* clade [22]. **D:** Analysis of the origins that are nearest to conserved, newly gained and lost, compared to the expected result if events were uncorrelated [22]. **E:** Distribution of the efficiency of lost, conserved and newly gained origins (respectively in black, dark red and red) and their neighbors (grey). Note that the efficiency of lost origins is lower than average, while the efficiency of origins flanking a lost origin is higher. **F:** Box plot of efficiency of all conserved and newly gained origins compared to those flanking a lost origin, which tend to be more efficient. Braces indicate sub-sampling (the box plots on the right side are defined by a subset of points of the box plots on the left). Box plots show the median (bar), 25-75 (box), and 10-90 (whiskers) percentiles. The data in panel C, D, E and F refers to the six sister species of the *Lachancea* tree.

87 analysis. The evolution of the temporal program of genome replication can be quantified by
88 the divergence of the replication timing profiles across different species. Agier and coworkers
89 found that timing profiles diverge gradually with increasing evolutionary divergence between
90 species [22]. In principle, such divergence could be attributed to changes in the number,
91 placement, and biological properties of all origins. However, a careful analysis of correlations
92 (comparing the timing profiles and the activity of orthologous origins) shows that the main
93 driver of program differentiation across species is the acquisition and loss of active replication
94 origins. Specifically, the number of conserved origins decreases with increasing phylogenetic
95 distance between species, following the same trend as the conservation of the timing profiles.
96 This trend is the same in regions that are close to or away from breakpoints, pointing to
97 a secondary role of genome rearrangements. In addition, the authors of ref. [22] show that
98 the differences in the mere number of origins and the median difference in origin replication
99 timing between pairs of species are nearly constant with phylogenetic distance, leading to
100 exclude that origin reprogramming (rather than birth-death) plays a primary role in the
101 evolution of the timing program.

102 Any model for the evolution of the replication program must (i) reproduce the empiri-
103 cal distribution of the inter-origin distances, (ii) reproduce the empirical distribution of the
104 origin efficiencies, and (iii) account for the observed origin turnover dynamics. Previous
105 analyses [18, 22] have shown that origins are far from following a uniform distribution along
106 the genome. Fig. 1A shows that the inter-origin distance distribution robustly shows a uni-
107 modal shape across the ten *Lachancea* species studied in ref. [22]. Specifically, distributions
108 for each species show a marked peak around 35 Kbp. This peak corresponds to a typical
109 inter-origin distance, which is strikingly invariant across all *Lachancea* species. Fig. 1B shows
110 the distribution of the efficiencies, which is defined as the probability to actively fire during
111 the S phase, estimated for each origin in the *Lachancea* clade using Eq. 4 and a fit inferring
112 the firing rates of all origins assuming a standard nucleation-growth model (see Methods
113 and ref [16]). The single-species efficiency distributions show more variability across species
114 than the inter-origin distance distributions, but they are consistent with a common shape
115 and support.

116 As mentioned above, a key result of Agier and coworkers is the insight that the evolution
117 of the replication program is mainly shaped by the birth-death process of replication origins.
118 Fig. 1C-F recapitulate the main quantitative results that characterize this process. Note that

119 the analyses in Fig. 1C-F have been performed on the six sister species of *Lachancea* clade,
120 since the other species pairs are too distant to perform a reliable identification of conserved,
121 newly gained and lost origins [22].

122 Fig. 1C shows a box plot of the distance from the nearest origin for all the conserved (dark
123 red), newly gained (red) and lost (black) origins. Lost replication origins tend to be closer
124 to their neighbors, much more so than newly gained or conserved origins. This observation
125 reveals that the distance of an origin from its nearest neighbor is correlated to the loss rate
126 of the same origin over evolution. This is an essential feature that any evolutionary model
127 of this process must take into account [18, 22]. More in detail, Fig. 1D further quantifies the
128 correlation between gain and loss events of neighboring origins, by comparing the fraction of
129 observed events of loss, gain, or conservation, given the state of the nearest origin (conserved,
130 lost, or gained). The distribution of event types for origins that are nearest neighbors of a
131 newly gained origin deviates significantly from the null expectation of random uncorrelated
132 events (i.e., in a simple scenario where the fractions of conserved, newly gained, and lost
133 origins are fixed to the empirical values, and birth and death events of neighboring origins are
134 independent). The same non-null behavior is observed for origins that are nearest to a lost
135 origin, with the roles of gain and loss events exchanged. In summary, successive birth/death
136 or death/birth events happen more frequently in the same genomic location than expected
137 by chance. Beyond such a spatial correlation along the chromosomal coordinate, the analysis
138 illustrates that birth and death events are correlated in time as well (in fact, the analyzed
139 evolutionary events took place in the terminal branches of the phylogenetic tree, and thus
140 they must have been close in term of evolutionary time).

141 Finally, Fig. 1E and 1F show that origins lying near loci where origins were recently lost
142 are typically in the high-efficiency range of the distribution, and that lost origins tend to
143 be less efficient than conserved origins. Fig. 1E compares the distribution of the efficiency
144 of lost, conserved and newly gained origins with the distribution of efficiency of the nearest
145 origins. The efficiency of origins neighboring a loss event is higher than average, while the
146 efficiency of lost origins is lower than average. These results clearly support the influence
147 of origin efficiency on origin death events. This is confirmed by Fig. 1F, which shows
148 the distribution of efficiency of all conserved and newly gained origins. For both classes,
149 considering only those origins that are nearest neighbors to a recently lost origin yields an
150 increase in the efficiency.

151 Different mechanisms could lead to the correlations described above. Overall, it is clear
152 that origin strength is somehow “coupled” to birth-death events. For example, conserved
153 origins may become more efficient after the loss of neighbor origins, or the birth of new highly
154 efficient origins could facilitate the loss of neighbors, or losing an origin could expedite the
155 acquisition of a new origin nearby. Overall, these results reveal that the origin birth-death
156 process is following some specific “rules” that involve both inter-origin distances and origin
157 efficiency.

158 Note that the results of Fig. 1E might appear to be incompatible with Fig. 1D, but they
159 are not. Fig. 1E shows that the efficiency of newly gained origins is lower than average, and
160 Fig. 1D shows that the majority of origins that are nearest to a locus with a recent loss
161 event are newly gained. The apparent contradiction arises from Fig. 1E, which shows that
162 the average efficiency of origins close to a lost one is higher than average. This inconsistency
163 is resolved by the analysis shown in Fig. 1F, which shows that origins appearing close to
164 recently lost ones are among the most efficient.

165 **A birth-death model including evolutionary bias from inter-origin replication fork** 166 **double stalling recapitulates the main features of replication origin turnover**

167 The joint stalling of two replication forks in the same inter-origin region along the genome
168 is a well-characterized fatal event that may occur during S-phase. The frequency of this event
169 in a clonal population clearly affects fitness. A previous modeling study [18] focusing on yeast
170 demonstrated that, in order to minimize the probability of a double stall anywhere along the
171 chromosome, origins must be placed in the most ordered spatial configuration, namely all
172 the consecutive origins must be equidistant from each other. However, the previous study
173 did not incorporate this principle into an evolutionary dynamics of origin turnover. Thus,
174 the important question arises of whether the tendency to avoid double stalls is related to
175 origin gain and loss. To address this question, we defined a birth-death model, rooted in
176 the experimental observations discussed in the previous section. This “double-stall aversion
177 model”, described in detail below, biases the turnover of replication origins in such a way
178 that events (in particular birth events) leading to a decreasing double-stall probability are
179 promoted, because they increase the fitness of the cell.

180 In the double-stall aversion model, the extent to which the acquisition of a new origin

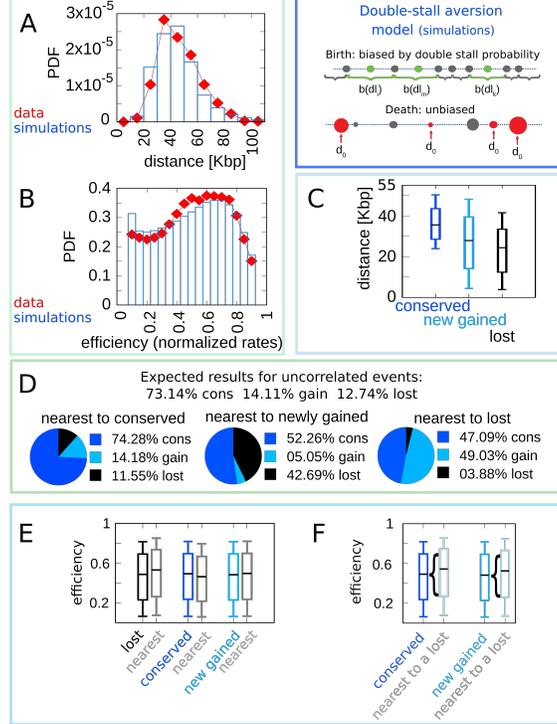


FIG. 2. The double-stall aversion model reproduces origin turnover and distributions but fails to capture correlations between origin turnover and origin strength. The plots show the simulations of the best-fitting double stall aversion model compared with empirical data. **A:** Inter-origin distance distribution in simulated species (blue bars) compared to the empirical distribution for the ten *Lachancea* species (red diamonds). **B:** Origin efficiency distribution in simulated (blue bars) *vs* empirical species (red diamonds). **C:** Box plot of the distance from the nearest origin split by evolutionary events, i.e. for conserved (dark blue), newly gained (blue) and lost origins (black), for simulated species. **D:** Fraction of origins that are nearest to conserved, newly gained and lost, for simulated species, compared to the expected result for uncorrelated events. **E:** Box plot of efficiency of lost, conserved and newly gained origins (respectively in black, dark blue and blue) and their neighbors (grey), in simulated species. The six distributions show very little variation. **F:** The efficiency of all conserved and newly gained origins compared to the ones flanking a lost origin. Braces indicate sub-sampling. Box plots show the median (bar), 25-75 (box), and 10-90 (whiskers) percentiles. Simulation parameters (see methods): $\gamma = 2.4$, overall birth and death rate $\bar{b} = 13.6 \text{Mbp}^{-1}t^{-1}$, $\bar{d} = 0.61t^{-1}$ and firing-rate resampling rate $R = 0.92t^{-1}$, where t is measured by protein-sequence divergence. The panels A and B are generated using data from approximately 320.000 simulated origins, while panels C, D, E and F are built using data from about 60.000 birth and death events and 240.000 conservation events.

181 changes the probability of a double stall P_i^{DS} depends on the length l_i of the inter-origin
 182 region where the event occurred. This probability is therefore coordinate-dependent, and
 183 can be derived by a procedure similar to the one carried out in [18] (see more details in
 184 Methods),

$$P_i^{\text{DS}} = 1 - (1 + \pi l_i) \exp(-\pi l_i) , \quad (1)$$

185 where l_i is the length of the genome region between the $(i + 1)$ -th and the i -th origin and π
 186 is the mean per-nucleotide fork stall rate; we use the value from ref. [18], $\pi = 5 \times 10^{-8}$ per
 187 nucleotide. Note that the double stall probability is completely independent from the origin
 188 firing rates and efficiency, and depends only on the distance between the origins.

189 In our simulations of the model (see Methods for a more detailed explanation), the genome
 190 was represented as a vector of origins, identified by the position and the firing rate. The
 191 model is a discrete-time Markov chain, and for the double-stall aversion variant the chain is
 192 specified by the following update rules,

- 193 • In each inter-origin region, the origin birth rate is biased by the value of the double-
 194 stall probability in that region. Specifically, the origin birth rate (per unit time) in
 195 the region i , of length l_i between the i -th and $(i + 1)$ -th origin is given by

$$b_i = N\bar{b}(P_i^{\text{DS}})^\gamma l_i , \quad (2)$$

196 where P_i^{DS} is the (constant) double stall probability density in region i (Eq.1), \bar{b} is
 197 the birth rate (per Mbp and per unit time) extracted from experimental data (see
 198 Methods), and γ is a positive real parameter that controls the strength of the bias. N
 199 is a normalization factor added to match the empirical birth rate \bar{b} . Newborn origins
 200 are placed in the middle of the inter-origin region i .

- 201 • Death (i.e., loss of origins) is unbiased, and occurs at random origins with rate \bar{d}
 202 (estimated from experimental data, see Methods), regardless of their efficiency or
 203 their neighbor's efficiency.

204 The justification for the assumption that newborn origins are placed at midpoints in the
 205 model ultimately comes from data (Fig. 1 - Supplement 2) where a strong bias in this
 206 direction is found. Relaxing this assumptions has consequences on the distance distribution
 207 and leads to poorer-performing models. We interpret this bias as the result of a faster

208 (hence undetectable in our data) evolutionary process that counter-selects origins far from
209 midpoints.

210 Firing rates in the model evolve by reshuffling of the empirical firing rate distribution,
211 with a time scale that is set empirically (see Methods and Figures 1 - Supplement 3 and 1 -
212 Supplement 4). On shorter time scales, firing rate changes are likely more gradual, making
213 firing-rate evolution similar to a diffusion process. However, such changes are not quantifiable
214 in our data set, which would leave the model with many extra parameters (a firing rate
215 diffusion constant and bounds to set the empirical distributions) that are very difficult to
216 estimate. Additionally, the firing-rate distributions of the conserved (thus older) origins and
217 of newborn (younger) ones are quite similar (Fig. 1 - Supplement 3B), and this condition is
218 not generally met under a simple diffusive process.

219 Fig. 2 shows the simulation results of the model with best-fitting parameter values (see
220 Methods and Fig. 2 for other parameter values). Fig. 2A and B, show that the double-stall
221 aversion model reproduces the two main “structural” features of yeast genome, namely the
222 inter-origin distance distribution and the origin efficiency distribution. Additionally, Fig. 2C
223 and D show that the same model reproduces the observed correlations between the inter-
224 origin distance and origin birth-death events, as well as the correlation between birth-death
225 events and nature of the neighbor origins observed in the data (conserved, newly gained, or
226 lost).

227 **The double stall hypothesis alone fails to capture correlations of origin turnover** 228 **with efficiency**

229 In spite of the good performance of the double-stall aversion model in explaining the
230 empirical marginal distributions, we find that it fails to reproduce the observed correlations
231 between the efficiency of an origin and the recent history of the nearest ones. Fig. 2E shows
232 very faint variations in efficiency of origins that are nearest neighbors to origins of different
233 evolutionary fate. In particular, the observed huge divergence in efficiency between lost
234 origins and their neighbors is absent in the model simulations. Note that Fig. 4 and Fig. 2F
235 show that in the double-stall aversion model origins nearest to a loss event are slightly more
236 efficient than average. This trend is due to the fact that after an origin is lost, its neighbours
237 are subject to lower interference, and automatically become more efficient. However, Fig. 4

238 shows that this null trend is too weak to explain the experimental data. These considerations
239 indicate that a model without a direct mechanism linking the efficiency of an origin to the
240 birth-death events of its neighbors cannot reproduce the data.

241 **Double-stall aversion and interference between proximate origins explain the cor-**
242 **related evolution of origin presence and efficiency**

243 Based on the above considerations, we defined a joint model that takes into account both
244 the evolutionary pressure given by the double-stall probability and the direct effect of origin
245 efficiency on birth-death events.

246 Specifically, this model is defined as follows.

- 247 • The birth process is the same as in the double-stall aversion model described above: the
248 birth rate is biased by the double-stall probability in each inter-origin region [eq. (1)],
249 and newborn origin are placed in the middle of the region.
- 250 • Death of an origin is biased by its efficiency: less efficient origins are more easily lost.
251 Specifically, the death rate (per unit time) for the i -th origin is

$$d_i = N\bar{d}\exp(-\beta \text{eff}_i), \quad (3)$$

252 where eff_i is the efficiency of the i -th origin, Eq. (4), \bar{d} is the mean death rate extracted
253 from experimental data (see Methods). The positive parameter β tunes the interaction
254 strength: the larger β , the steeper the dependence of d_i on eff_i . The normalizing factor
255 N is chosen so as to match the empirical total death rate.

256 We note that the bias parameters β and γ are not inferred based on branch data, but on
257 distributions of extant species (see Methods).

258 Fig. 3 gathers plots of the structural features (distribution of inter-origin distances and
259 efficiencies, Fig. 3A-B) and the evolutionary correlations involving efficiency, evolutionary
260 fate, distance to nearest neighbor, and fate of nearest neighbor (Fig. 3C-D-E-F). Overall, the
261 joint model reproduces all the observations considered here regarding the layout of origins
262 and their evolutionary dynamics, indicating that the experimental data can be rationalized
263 by a fitness function that includes both the detrimental effects of non replicated regions and
264 the evolutionary cost of maintaining inefficient replication origins.

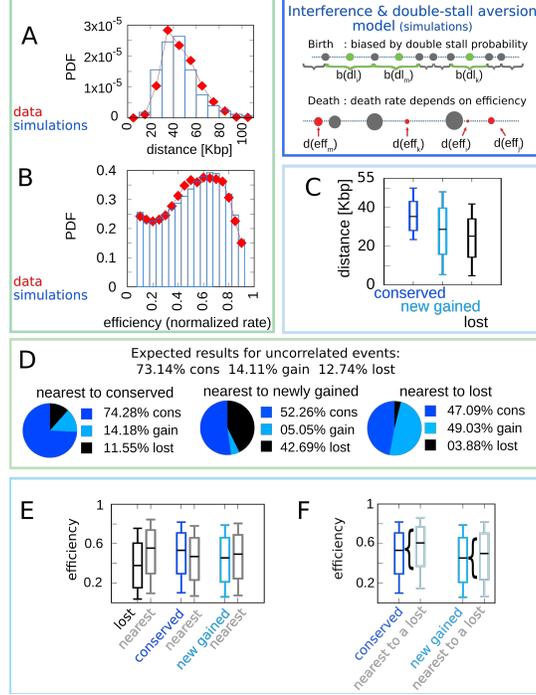


FIG. 3. **A model where both fork stalling and interference affect fitness explains the correlations between origins evolutionary events.** Result of the joint model best-fitting simulation compared with empirical data. **A:** Inter-origin distance distribution in simulated species (blue bars) *vs* empirical distribution for the ten *Lachancea* species (red diamonds). **B:** Origin efficiency distribution in simulated (blue bars) *vs* empirical species (red diamonds). The agreement between simulation and experimental data shows that this joint evolutionary model reproduces the typical structural features of a yeast genome. **C:** Box plot of the distance from the nearest origin split by evolutionary events, i.e. for conserved (dark blue), newly gained (blue) and lost origins (black), for simulated species. **D:** Fraction of origins that are nearest to conserved, newly gained and lost, for simulated species, compared to the expected result for uncorrelated events. **E:** Box plot of efficiency of lost, conserved and newly gained origins (respectively in black, dark blue and blue) and their neighbors (grey), in simulated species. **F:** The efficiency of all conserved and newly gained origins compared to the ones flanking a lost origin. Braces indicate sub-sampling. Box plots show the median (bar), 25-75 (box), and 10-90 (whiskers) percentiles. Panels D - F show that the model correctly reproduces the correlation between origin birth-death events over evolution and efficiency of the nearest origin. Simulation parameters (see Methods): $\gamma = 2.2$, $\beta = 1.9$, overall birth and death rate $\bar{b} = 13.6Mbp^{-1}t^{-1}$, $\bar{d} = 0.61t^{-1}$ and rate of origin firing-rate reshuffling $R = 0.92t^{-1}$, where t is measured by protein-sequence divergence. The panels A and B show data from approximately 600.000 simulated origins, while panels C, D, E and F data from about 100.000 birth and death events and 500.000 conservation events.

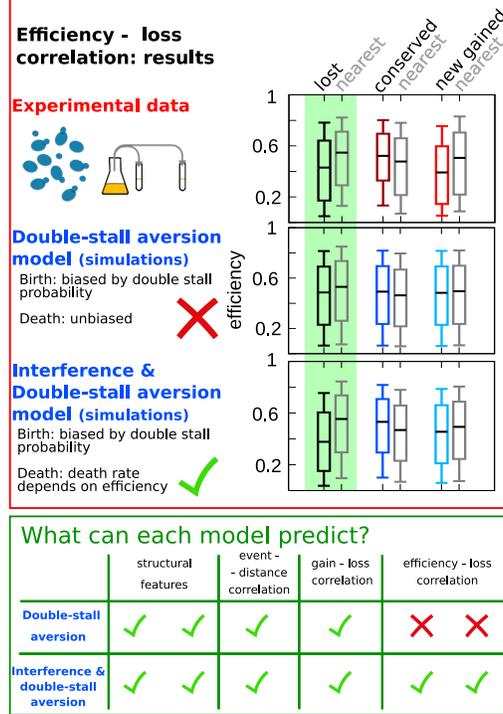


FIG. 4. **Comparison of model predictions for the correlations of origin birth-death events.** The plots in the red upper box compare efficiency distributions of the best-fitting simulation of the two different models (bottom and central panels) with experimental data (top panel). Comparison of the box plot of efficiency of lost, conserved and newly gained origins (red for the data, blue for the models) shows better agreement of the joint efficiency/double-stall aversion model (bottom panel) with the experimental data. Hence, the joint model reproduces well the correlation between evolutionary birth-death events of origins and efficiency of the nearest origin, while the double-stall aversion model fails. Box plots show the median (bar), 25-75 (box), and 10-90 (whiskers) percentiles. Simulation parameters for the joint model (see Methods): $\gamma = 2.2$, $\beta = 1.9$, and for the double-stall aversion one: $\gamma = 2.4$. General parameters: overall birth and death rate $\bar{b} = 13.6Mbp^{-1}t^{-1}$, $\bar{d} = 0.61t^{-1}$ and rate of origin firing-rate reshuffling $R = 0.92t^{-1}$, where t is measured by protein-sequence divergence. In the green lower box we compare the predictive power of the two models for each of the tested feature of the experimental data. The box highlights that both the double stall-aversion model and the joint efficiency - double stall model are able to reproduce the structural features of the genome. Also the correlation between events - distance from the nearest and event - event of the nearest are correctly predicted by both models. The important difference between the two proposed models is found for the correlation between evolutionary events and origin efficiency, which is predicted and can be explained solely by the joint model.

265 In particular, the coupling between the efficiency of an origin and the death rate of its
266 neighbors, through the probability of passive replication, reproduces the empirical correla-
267 tions shown in Fig. 1. Figure 4 summarizes this crucial point of comparison between the
268 joint efficiency/double-stall aversion model and the pure double-stall aversion case. The
269 three plots compare efficiency distributions of lost, conserved and newly gained origins (red
270 for the data, blue for the models) with those of their neighbors (grey). Comparison of these
271 plots shows that only the joint model reproduces the differences in efficiency of lost origins
272 and their neighbors.

273 In order to show that the stall-aversion and interference model has better quantitative
274 agreement with the data, we also performed a simplified likelihood ratio analysis. The full
275 likelihood of the model is complex, but we have defined “partial” likelihoods for the joint and
276 the double stall aversion model just taking into account the marginal probabilities shown as
277 box plots in Fig. 4 and Fig. 4 - Supplement 1 (see Methods). Fig. 1 shows that the joint
278 model performs better for all the four chosen features. In our view, the qualitative difference
279 shown in Fig. 4 may be taken as a stronger argument in favor of the combined model, in
280 the sense that, beyond any quantitative agreement relying on parameters, the additional
281 ingredient of a coupling between origin birth-death dynamics and origin rates is needed to
282 explain the data.

283 **The joint efficiency / double-stall aversion model correctly predicts origin family** 284 **divergence**

285 Having established that the joint model is required to reproduce observations on single
286 lineages, we turned to its predictions on observations that require knowledge of the whole
287 phylogenetic tree, such as origin evolutionary families, defined as sets of orthologous ori-
288 gins [22].

289 We thus set up a simulation of the model on a cladogenetic structure, fixed by the observed
290 structure of the *Lachancea* clade phylogenetic tree (see Methods for the simulations details).
291 The output of each run in such simulations are nine different simulated genomes whose
292 lineages are interconnected in the same way as the empirical species, and each branch follows
293 the empirical divergence. We stress that these simulations just include intersecting lineages
294 whose branched structure corresponds precisely to the lineages of the empirical tree. The

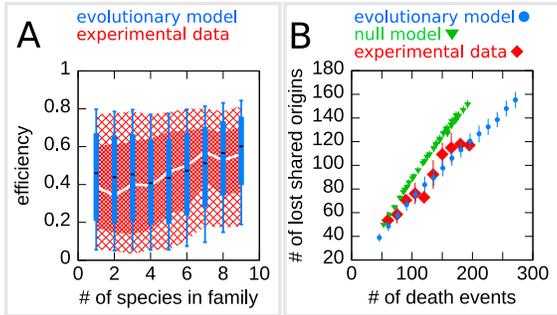


FIG. 5. **The efficiency/double-stall aversion model predicts origin divergence.** The plots compare predictions of the evolutionary model on the extent of origin divergence (simulations of the *Lachancea* phylogenetic tree) with empirical data. **A:** Box plot of origins efficiency distributions split by family size. The plot compares origin families (sets of orthologous origins) in the nine *Lachancea* species (white line and red shaded areas) and in simulated species (blue boxes, for 100 simulation runs). Medians are shown as white line for data, black bar for simulation, 25-75 percentiles as shaded area for data, box for simulation, and 10-90 percentiles as coarse shaded area for data, whiskers for simulation. **B:** Origin divergence measured by the number of origins in the common ancestor that were lost in a pair of species, plotted as a function of total origin loss events. The plot compares model simulations (blue circles, 100 simulation runs), the experimental data (red squares) and a null model that shuffles the empirical birth - death events in each branch (green triangles, 1000 simulation runs). Error bars are standard deviations on y -axis values. Simulation parameters (for the evolutionary model, see Methods): $\gamma = 2.2$, $\beta = 1.9$, overall birth and death rate $\bar{b} = 13.6\text{Mbp}^{-1}t^{-1}$, $\bar{d} = 0.61t^{-1}$ and rate of origin firing-rate reshuffling $R = 0.92t^{-1}$, where t is measured by protein-sequence divergence.

295 phylogenetic structure does not emerge from the simulation, as our model does not describe
 296 speciation. The model for the tree can simulate nine species, all the species except for *L.*
 297 *kluyveri*, as this species was used as outgroup for the computation of the length of the tree
 298 branches [22]. We have repeated all the analyses on these simulations, and verified that
 299 all the previous results hold, Fig. 5 - Supplement 1. We then turned to other independent
 300 predictions of the joint model, which could be compared to measurements in ref. [22].

301 Fig. 5A reports the dynamics of origin families. As reported in ref. [22], origins that
 302 belong to larger evolutionary families tend to have a higher efficiency compared to origins
 303 in smaller families, which is possibly due to the fact that, on average, high efficiency origins

304 tend to survive longer. Note however that there is no deterministic relation between family
305 size and origin age because the relationship between these two is determined by the structure
306 of the phylogenetic tree. Indeed, two families of the same size may have roots in different
307 points of the tree, and thus the origins belonging to them may have very different ages. Thus,
308 the prediction of the relation between origin efficiency and origin-family size is not trivial.
309 Fig. 5A shows the results for the origin efficiency for families of varying size, comparing the
310 experimental data and 100 different runs of the simulation.

311 As a second step, we have considered the model prediction for the divergence of the shared
312 origins in two species descending from a common ancestor. Specifically, we asked whether
313 the number of origin death events occurring in two branches of the tree could justify the
314 number of common origins in the two species. Indeed, whenever in a pair of species the
315 number of shared origins is lower than the number of origins belonging to their common
316 ancestor, this discrepancy must be due to the evolutionary loss events. These events are
317 predicted by our model to be correlated in diverging species, due to the common ancestry
318 and the coupling of loss events to origin efficiency and distance. This correlation should
319 lower the number of shared origins losses compared to a null expectation where loss events
320 are not correlated. Fig. 5B shows that the model correctly predicts the divergence in the
321 number of shared origins lost during evolution, without any parameter adjustment. We also
322 verified that, as expected, a null evolutionary model is not able to reproduce this feature.
323 The null model fixes in each branch of the simulated tree the same number of birth and
324 death events that are present in the corresponding branch of *Lachancea* tree, but these
325 events occur uniformly along the genome. The difference between the null model and the
326 evolutionary model predictions shown in Fig. 5B is a consequence of correlated origins losses
327 due to the common genome structure, in terms of origins positions and efficiencies, that each
328 pair of species inherit from their common ancestor.

329 We note that birth and death rate are inferred as global parameters, ignoring correlations.
330 Despite this, Fig. 5B shows that the model reproduces the higher correlation in birth and
331 death events in closer-related branches than in distant branches as a consequence of the
332 common positions and firing rates of the origins in the ancestor.

DISCUSSION AND CONCLUSIONS

Overall, this study provides a framework to study replication-program evolution driven by replication-origin birth-death events, and demonstrates that both fork stalling and efficiency shape the adaptive evolution of replication programs. The model framework is predictive and falsifiable and it can be used to formulate predictions on the phylogenetic tree. In future studies, it would be interesting to explore the predictions for the evolutionary dynamics under perturbations, such as evolution under increased replication stress or conditions where fork stalling becomes more frequent. Additionally, the framework can be used to discover specific trends, such as different evolutionary dynamics of specific genomic regions (subtelomeres [23], regions containing repeats, etc. [24, 25]), role of genome spatial organization [26], and correlated firing of nearby origins.

A general question concerns the predictive value of the model proposed here on out-of-sample data. Fig. 5 shows that fit-independent predictions apply across the tree. Importantly, the model is based on simple global parameters, and not fine tuned on local features of the tree. To underline this point, we verified that a model fit using only the subtree between LADA and LAWA yielded similar parameters. Clearly, we cannot exclude that the values of the birth and death rate, and also the bias parameters γ and β could be *Lachancea*-specific, while we speculate that the conclusions on the relevant evolutionary mechanisms might apply more generally.

The previous approach by Newman and coworkers [18] described the evolution of origin distance as an optimization process that minimizes double fork-stall events, without attempting to characterize explicitly the evolutionary dynamics. Such approaches are limited compared to the framework presented here, because they can predict only the origin-distance distribution, and they do not allow any prediction regarding origin and replication-program evolution along lineages and across phylogenetic trees. In accordance with the results of Newman *et al.*, we confirm that double-stall events are a primary driver of the evolution of replication programs, and we frame this finding into the empirically measured birth-death evolutionary dynamics of replication origins. Additionally, we show that next to fork-stall events, origin efficiency plays an important role into shaping the evolutionary landscape seen by a replication timing profile.

What could be the mechanisms coupling efficiency to origin birth death? The actual pro-

364 cess of origin death could be nearly neutral [27], as low-efficiency origins, are - by definition -
365 rarely used, and unused origins, over evolutionary times are more prone to decay in sequence,
366 and consequently in firing-rate until they disappear. Equally, a new-born origin close to a
367 very strong one (which would make the new-born origin relatively inefficient) could be used
368 rarely. This would make this origin relatively less likely to establish over evolutionary times
369 compared to an isolated new-born origin. However, rarely used origins could be essential
370 in situations of stress (and in particular they could resolve double-stall events). Finally,
371 a fitness cost for maintaining too many origins might set up an overall negative selection
372 preventing a global increase in origin number [16, 28, 29].

373 III. MATERIALS AND METHODS

374 Data

375 The experimental data used in this work come from ref. [22]. In particular, we made use
376 of the data regarding the replication origins. For each origin in each of the ten *Lachancea*
377 species, this dataset includes the chromosome coordinate and firing rate, and the inferred
378 birth and death events occurred in the branches of the phylogenetic tree shown in Fig. 1 -
379 Supplement 1. Focusing on the terminal branches of the tree and on the extant replication
380 origins, this study defines three categories of origins: (i) “conserved” origins (which survived
381 from the last ancestor) (ii) “newly gained” origins gained in the last branch of the phylo-
382 genetic tree, (iii) “lost” origins, which were present in the last ancestor species and are not
383 present in the terminal branch. Properties of the lost origins (e.g. position and firing rate)
384 are inferred from the projection of the corresponding ones on the closest species, keeping into
385 account synteny. Since the synteny map is less precise in distant species, the information on
386 the origins events is only available for the six sister species in the tree, which belong to the
387 three closest species pairs, highlighted with the red shaded area in Fig. 1 - Supplement 1.

388 Computation of the efficiency

389 Origin efficiency was defined as the probability of actively firing during S phase (or,
390 equivalently, the probability of not being passively replicated by forks coming from nearby

origins). In practice we computed it by the following formula

$$\text{eff}_i = (1 - P_{i,i-1})(1 - P_{i,i+1}) , \quad (4)$$

where $P_{i,i+1}$ and $P_{i,i-1}$ are the probabilities for the i -th origin of passive replication respectively from the $(i + 1)$ -th and $(i - 1)$ -th origins. Note that this efficiency formula Eq. 4, is an approximation that only takes into account the possibility to be passively replicated by neighbor origins, neglecting the influence of other nearby origins. Following ref [22], for computing the efficiency we assumed that the origin firing process has constant rate [16], and we thus obtain the following closed expressions for the probabilities of passive replication

$$P_{i,i+1} = \frac{\lambda'_{i+1}}{\lambda'_{i+1} + \lambda'_i} \exp \left[-\lambda'_i \frac{|x_{i+1} - x_i|}{v} \right] , \quad (5)$$

and

$$P_{i,i-1} = \frac{\lambda'_{i-1}}{\lambda'_{i-1} + \lambda'_i} \exp \left[-\lambda'_i \frac{|x_{i-1} - x_i|}{v} \right] . \quad (6)$$

In the above equations, v is the typical velocity of replication forks, x_i is the i -th origin chromosome coordinate, and λ'_i is the i -th origin firing rate divided by the mean firing rate of the species the origin belong to. The raw firing rates in the data are affected by the different physiology of the nine *Lachancea* species in the experimental growth conditions (which were the same for all the species). In order to reduce these differences, we normalized the rates by their average for each given species. For this reason, we did not make use of the origin efficiency data already present in [22].

Computation of the double-stall probability

The probability P_i^{DS} that two converging forks stall is easily computed in the limit where the stall probability per base-pair is small and the number of base-pairs is large. Under these assumptions, stalling is a Poisson process with rate (per base-pair) π . P_i^{DS} can be written in terms of the probability $P^{\text{S}}(x)$ that a single fork stalls after replicating x nucleotides,

$$P_i^{\text{DS}} = \int_0^{l_i} dx \int_0^{l_i-x} dy P^{\text{S}}(x) P^{\text{S}}(y) , \quad (7)$$

where l_i is the length (number of base-pairs) of the i -th inter-origin region. Imagine two converging replication forks starting from origins i and $i + 1$: the two integration variables x and y represent the number of base-pairs that each fork replicates before stalling. By using

414 the Poisson-process result $P^S(x) = \pi \exp(-\pi x)$ and performing the integration, one obtains
415 the result in Eq. (1).

416 **Evolutionary model**

417 We defined origin birth-death models incorporating different evolutionary biases. In these
418 models, the genome is described as a one-dimensional circle with discrete origin location x_i ,
419 where the length of the genome is equal to the average genome length in *Lachancea* clade
420 (10.7Mbp). We made use of a circular genome in order to avoid border effects. In the
421 model, the set of origins change over evolution by three basic (stochastic) processes, birth of
422 an origin in a certain genome region, origin death and change of origins firing rate. We have
423 verified that choosing linear chromosome does not alter significantly our findings, although
424 it affects the distances between origins close to chromosome ends (Fig. 3 - Supplement 1).

425 Overall origin birth/death rates were estimated from the data as follows. To estimate the
426 overall birth rate \bar{b} we considered, for all the terminal branches of the phylogenetic tree, the
427 number of birth events N_b , the genome length of the corresponding species L and the length
428 of the tree branch T , and divided N_b by LT . Then we averaged over all terminal branches.
429 To estimate the overall death rate \bar{d} , we followed a similar approach, taking the number
430 of death events N_d in the terminal branches, the length of the branch T and the number
431 of origins in the corresponding species n_{ori} , then computing $N_d T^{-1} n_{ori}^{-1}$ for all the terminal
432 branches and averaging these values. The final results for overall birth and death rates from
433 the origin birth death events across the *Lachancea* clade are $\bar{b} = 13.5627Mbp^{-1}t^{-1}$ and
434 $\bar{d} = 0.612287t^{-1}$.

435 We verified that the assumption of constant rates was consistent with the the empirical
436 variability of the numbers of birth and death events per unit time along different branches of
437 the tree, by comparing simulations with data. Fig. 5 - Supplement 2 shows that simulations
438 and empirical data present similar spreading,

439 The process by which origin firing rates change over evolution was described as stochastic,
440 with every origin having a fixed probability per unit time of changing its firing rate, given by
441 $R = 0.92t^{-1}$, a value fixed from experimental data (see appendix and Fig. 1 - Supplement 4).
442 When a firing rate changes, it is resampled from the distribution of all the empirical normal-
443 ized firing rates, computed using the data in [22] (see appendix and Fig. 1 - Supplement 3

444 for more details).

445 Simulations

446 *Code Availability.* The code used to run the simulations, together with instructions to
447 run it, was shared as a repository on Mendeley data, and is available at the url <https://data.mendeley.com/datasets/vg3r5355bj/2>.
448 *Algorithm.* The prediction of the dif-
449 ferent evolutionary models were derived numerically, making use of custom simulations
450 written in C++, which implement the origin birth-death dynamics as a Gillespie algo-
451 rithm [30]. Every model variant was required to reproduce the experimental overall rates,
452 $\bar{b} = 13.5627Mbp^{-1}t^{-1}$ for origin birth, $\bar{d} = 0.612287t^{-1}$ for origin death and $R = 0.92t^{-1}$
453 for firing rate change. We simulated the three processes defining the model as follows. (i)
454 The birth process has a common definition for the stall aversion and joint model. The al-
455 gorithm first tests each subsequent inter-origin region, calculates the birth probability from
456 Eq. 2 and stores the results. Subsequently, it computes the normalization factor N , in order
457 to match the empirical birth rate per nucleotide \bar{b} . Finally, it samples all the inter-origin
458 regions drawing birth events from the computed birth probability (Eq. 2). New origins are
459 placed the mid points of the tested intervals. (ii) The death process is different for the
460 stall-aversion model (unbiased) and the joint model (related to the origin efficiency). In the
461 joint model, the algorithm first calculates the death rate for each origin using Eq. 3 and
462 stores the results. Subsequently, it computes the normalization factor N , in order to match
463 the empirical mean death rate \bar{d} . Finally, it samples all origin drawing death events from the
464 computed death probability. For the unbiased process (stall-aversion model) the dynamics
465 is identical, but all the origins have the same death rate \bar{d} , so that the algorithm can skip
466 the calculation of N . (iii) The process updating origin firing rates over evolutionary times
467 is common to all model variants. The probability of update per origin per unit time is R .
468 Origins are sampled for each time step and assigned a new rate uniformly extracted from
469 the empirical distribution of all normalized firing rates with probability Rdt .

470 During the simulation the genome configuration (chromosome position, firing rate, effi-
471 ciency for each origin) is known at each time step, which matches the empirical time (tree-
472 branch length, measured by protein-sequence divergence). For simulating single lineages, we
473 started with a collection of 50 origins, with positions and firing rate uniformly drawn from

474 all the possible ones. Rapidly, the inter-origin distances distribution, the efficiency one and
475 the number of origins reach a steady state (for the number of origins, set by the balance
476 of birth and death rate, and characterized by approximately 225 origins). Configurations,
477 including birth-death events, were printed at regular time intervals after steady state is
478 reached. The time interval between prints is chosen to be equal to the average length of the
479 *Lachancea* phylogenetic tree terminal branches, in order to compare single-lineage simula-
480 tions with empirical data. For simulations on a phylogenetic tree, after one species reaches
481 the steady state, it is used as a root. To reproduce the empirical branching structure of the
482 tree, we run the simulation, one for each branch of the phylogenetic tree, each time starting
483 from the species at the previous branching point, for a period that matches the length of the
484 branch. If the simulated branch is terminal then the configuration corresponds to one of the
485 empirical species, otherwise it corresponds to a “branching-point species” and it can be used
486 as starting point for other simulations. Each simulation run gives nine different simulated
487 species with the same cladogenetic structure as the empirical species (Fig. 1 - Supplement 1).

488 *Fitting procedure.* The biased birth-death processes in the simulations rely on some param-
489 eters to tune the strength of the bias, these are the only parameters to fix by a fit, since
490 all the other parameter values are fixed empirically. In the joint model there are two free
491 parameters, γ and β that tune respectively the strength of the bias on the origin birth and
492 on the origin death process. For a discrete set of parameter pairs spanning realistic intervals
493 we run hundred different simulations, each starting with a randomized genome. Considering
494 the simulated species for all the pairs of parameter values, we quantify the discrepancy with
495 experimental data by evaluating the L1 distance of the normalized histogram of efficiency
496 and inter-origin distances. This quantity is a number between 0 and 2, 0 if the histograms
497 perfectly overlap and 2 if they have completely different supports. For each pair of param-
498 eters the analysis gives two values of discrepancy. We choose the value of γ (the parameter
499 that tunes the bias on the birth rate based on double-stall aversion) by taking the smaller
500 discrepancy from the inter-origin distances distribution. For the value of β (which tunes the
501 interference bias on the death rate in joint model), we chose the one that gave us the smaller
502 area on the efficiency distribution. For the double-stall aversion model the fitting procedure
503 is the identical, and only requires to fix γ .

504 *Simplified Likelihood analysis.* We performed a (simplified) likelihood ratio analysis in order
505 to test the better quantitative performance of the combined model. The full likelihood of the

506 models analyzed here is complex, but we have defined “partial” likelihoods for the joint and
 507 the double stall aversion model only taking into account the marginal probabilities shown
 508 in Fig. 4 and 4 - Supplement 1. Hence the test evaluates for both models the goodness of
 509 the predicted correlation between the efficiency and firing rate of the lost origins and the
 510 ones of their neighbors. The likelihood ratio test quantifies how much the prediction of
 511 a certain model is better than a reference (“null”) model. We chose the the double stall
 512 aversion model as reference (equivalent to setting $\beta = 0$ in the joint model). Specifically,
 513 one evaluates

$$L_r = 2 \log \left(\frac{L_{joint}(\gamma, \beta)}{L_{DS}(\gamma, \beta = 0)} \right) = 2 (l_{joint}(\gamma, \beta) - l_{DS}(\gamma, \beta = 0)), \quad (8)$$

514 where L_X are the likelihoods of the two models and l_X are the log-likelihoods. Assuming that
 515 L_r is χ -squared distributed (this is generally the case for large samples), we could compute
 516 a P-value associated to this test.

517 **Null birth-death model**

518 We defined a null birth-death model where origin birth-death events in sister species
 519 are uncorrelated, in order to analyze the divergence of shared origins and compare it with
 520 the prediction of the evolutionary model. This model implements birth and death events
 521 uniformly, regardless of origin position and firing rate, fixing the number of events for each
 522 branch of the simulated phylogenetic tree. These values are taken from the inference reported
 523 in ref. [22] (shown in Fig. 3A of that study and in Fig. 1 - Supplement 1). The simulation of
 524 this model starts with 220 origins (the number of origins inferred for to LA2, the species at
 525 the root of the tree). Subsequently, following the structure of the *Lachancea* phylogenetic
 526 tree, the simulation proceeds as follows: (i) at each branching point the genome is copied
 527 into two daughters, (ii) for each daughter the prescribed number of random death and birth
 528 events (in this order) is generated on random origins (iii) the simulation stops when it reaches
 529 the leaves of the *Lachancea* tree.

530 **ACKNOWLEDGMENTS**

531 We are very grateful to Ludovico Calabrese, Simone Pompei and Orso Maria Romano for
 532 the useful discussions. We also thank the editors and reviewers of this manuscript for their

533 constructive and helpful feedback. This work was supported by the Italian Association for
534 Cancer Research, AIRC-IG (REF: 23258).

- 535 [1] A. C. Leonard and M. Méchali, *Cold Spring Harb. Perspect. Biol.* **5**, a010116 (2013).
- 536 [2] M. W. Musiałek and D. Rybaczek, *Cell cycle (Georgetown, Tex.)* **14**, 2251 (2015).
- 537 [3] O. Ganier, P. Prorok, I. Akerman, and M. Méchali, *Current opinion in cell biology* **58**, 134
538 (2019).
- 539 [4] D. M. Gilbert, *Science* **294**, 96 (2001).
- 540 [5] M. Méchali, K. Yoshida, P. Coulombe, and P. Pasero, *Curr. Opin. Genet. Dev.* **23**, 124 (2013).
- 541 [6] S. C. Di Rienzi, K. C. Lindstrom, T. Mann, W. S. Noble, M. K. Raghuraman, and B. J.
542 Brewer, *Genome research* **22**, 1940 (2012).
- 543 [7] J. Bechhoefer and N. Rhind, *Trends in genetics : TIG* **28**, 374 (2012).
- 544 [8] N. Rhind, S. C.-H. Yang, and J. Bechhoefer, *Chromosome research : an international journal*
545 *on the molecular, supramolecular and evolutionary aspects of chromosome biology* **18**, 35
546 (2010).
- 547 [9] M. Hawkins, R. Retkute, C. A. Müller, N. Saner, T. U. Tanaka, A. P. de Moura, and C. A.
548 Nieduszynski, *Cell Rep.* **5**, 1132 (2013).
- 549 [10] A. Baker, B. Audit, S. C.-H. Yang, J. Bechhoefer, and A. Arneodo, *Phys. Rev. Lett.* **108**,
550 268101 (2012).
- 551 [11] N. Agier, O. M. Romano, F. Touzain, M. Cosentino Lagomarsino, and G. Fischer, *Genome*
552 *Biol. Evol.* **5**, 370 (2013).
- 553 [12] C. A. Müller, M. A. Boemo, P. Spingardi, B. M. Kessler, S. Kriaucionis, J. T. Simpson, and
554 C. A. Nieduszynski, *Nature methods* **16**, 429 (2019).
- 555 [13] M. Hennion, J.-M. Arbona, L. Lacroix, C. Cruaud, B. Theulot, B. L. Tallec, F. Proux, X. Wu,
556 E. Novikova, S. Engelen, A. Lemainque, B. Audit, and O. Hyrien, *Genome biology* **21**, 125
557 (2020).
- 558 [14] R. Retkute, C. A. Nieduszynski, and A. de Moura, *Phys. Rev. E* **86**, 031916 (2012).
- 559 [15] A. P. S. de Moura, R. Retkute, M. Hawkins, and C. A. Nieduszynski, *Nucleic acids research*
560 **38**, 5623 (2010).

- 561 [16] Q. Zhang, F. Bassetti, M. Gherardi, and M. Cosentino Lagomarsino, *Nucleic acids research*
562 **45**, 8190 (2017).
- 563 [17] E. V. Koonin, *PLOS Computational Biology* (2011).
- 564 [18] T. J. Newman, M. A. Mamun, C. A. Nieduszynski, and J. J. Blow, *Nucleic acids research*
565 **41**, 9705 (2013).
- 566 [19] A. Letessier, G. A. Millot, S. Koundrioukoff, A.-M. Lachagès, N. Vogt, R. S. Hansen, B. Malfoy,
567 O. Brison, and M. Debatisse, *Nature* **470**, 120 (2011).
- 568 [20] R. S. Cha and N. Kleckner, *Science* **297**, 602 (2002).
- 569 [21] C. A. Müller and C. A. Nieduszynski, *Genome research* **22**, 1953 (2012).
- 570 [22] N. Agier, S. Delmas, Q. Zhang, A. Fleiss, Y. Jaszczyszyn, E. van Dijk, C. Thermes, M. Weigt,
571 M. Cosentino-Lagomarsino, and G. Fischer, *Nature communications* **9**, 2199 (2018).
- 572 [23] J.-X. Yue, J. Li, L. Aigrain, J. Hallin, K. Persson, K. Oliver, A. Bergström, P. Coupland,
573 J. Warringer, M. Cosentino Lagomarsino, G. Fischer, R. Dubin, and G. Liti, *Nature Genetics*
574 (2017).
- 575 [24] J.-M. Arbona, A. Goldar, O. Hyrien, A. Arneodo, and B. Audit, *eLife* (2018).
- 576 [25] D. Boos and P. Ferreira, *Genes (Basel)* (2019).
- 577 [26] C. Marchal, J. Sima, and D. M. Gilbert, *Nature Reviews Molecular Cell Biology* (2019).
- 578 [27] E. V. Koonin, *BMC Biology* **14** (2016).
- 579 [28] J. Karschau, J. J. Blow, and A. P. S. de Moura, *Physical review letters* **108**, 058101 (2012).
- 580 [29] S. P. Das, T. Borrman, V. W. T. Liu, S. C.-H. Yang, J. Bechhoefer, and N. Rhind, *Genome*
581 *Res.* **25**, 1886 (2015).
- 582 [30] D. T. Gillespie, *Journal of Computational Physics* **22**, 403 (1976).
- 583 [31] S. Matmati, S. Lambert, V. Géli, and S. Coulon, *Cell Reports* (2020).

Figure supplements for Droghetti *et al.*

“An evolutionary model identifies the main evolutionary biases for the evolution of genome-replication profiles.”

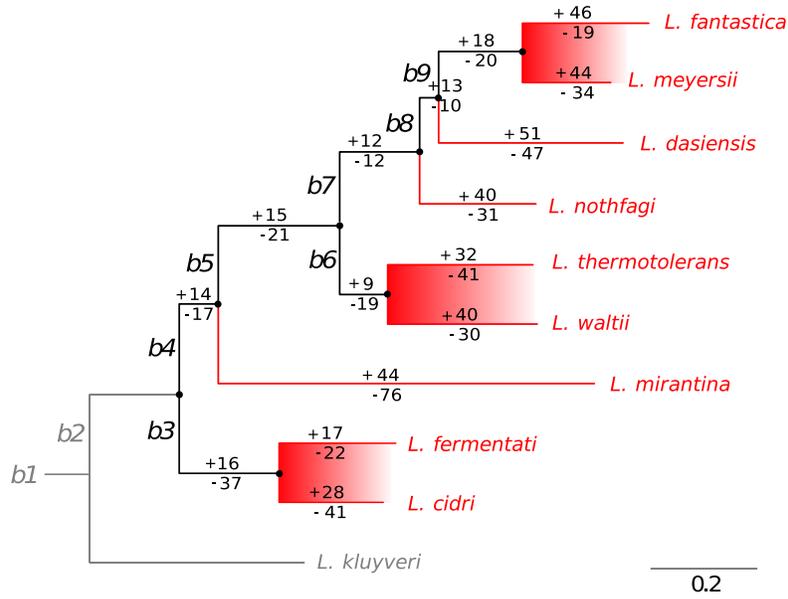


FIG. 1 - Supplement 1. **The phylogenetic tree of the ten *Lachancea* yeasts clade. taken from ref [22], Fig.3A.** *L. kluveri* was used as the outgroup species. Hence evolutionary events that occurred on both the *L. kluveri* and the b2 branches (grey lines) could not be retraced. As a consequence, our simulations of the model were not possible for the b2 and *L. kluveri* branches, and it was possible to simulate nine species instead of ten. Internal branches, labeled b3 to b9, and terminal branches are drawn in black and red, respectively. The number of origin gains (with plus sign) and losses (with minus sign) were estimated for each branch of the tree in ref. [22]. The six sister species, which belong to the three closest pairs of species, are highlighted with the red shaded areas.

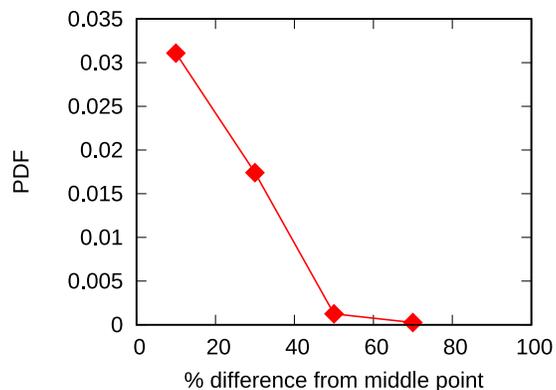


FIG. 1 - Supplement 2. **The majority of new origins are born within a 20% distance from the midpoint of the associated interval.** The plot shows the empirical distribution of the fractional distance from the midpoints of nearby origins for newborn origins of the *Lachancea* clade. More than half of all the new born origins is less than 20% far away from the midpoint of the inter-origin interval where they are born. This means that for an ideal 50 Kbp interval, more than half of the birth events would occur in positions between 20 and 30 Kbp, which is remarkably close to the midpoint position of 25Kbp. This result justifies the simplified choice of placing newborn origins at midpoints in our models.

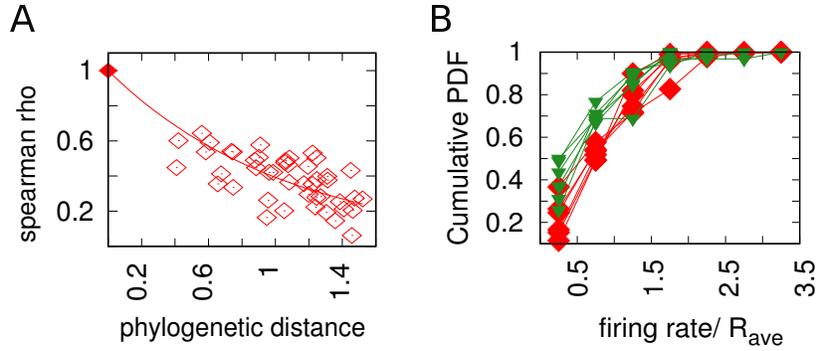


FIG. 1 - Supplement 3. **Experimental data on the evolutionary change of firing rates process.** **A:** The firing rates Spearman correlation coefficient ρ between sets of corresponding origins decreases with increasing phylogenetic distance between species. Each point in the plot represents a pair of species. The x axis reports the phylogenetic distance between the two species, while the y axis reports the Spearman correlation between the sets of normalized firing rates for corresponding origins between the two species. Empty squares represent the analysis carried out with *Lachancea* clade yeasts, while the symbol with coordinates (0,1) represents the fact that that non-distant species must have $\rho = 1$ **B:** Cumulative probability distribution of the normalized firing rates of newly-gained origins (green triangles, for the six sister species) compared to the all the extant origins (red squares, for the six sister species). This plot shows that all the functions are very similar. This results is compatible with the assumption of resampling of firing rates over evolution taken for the model (see appendix).

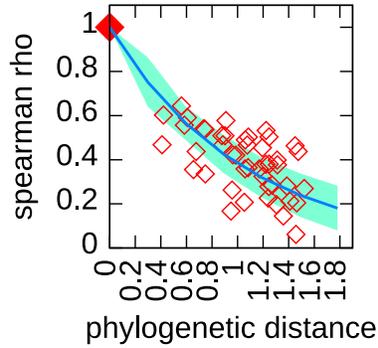


FIG. 1 - Supplement 4. **The decaying trend of the spearman correlation coefficient define a characteristic time for the firing rate resample.** For each pair of species we compute the spearman correlation coefficient between the set of normalized firing rates belonging to corresponding origins. The figure shows the results of this analysis. The red empty points refer to experimental data, each dot is a pair of species, the x coordinate is the phylogenetic distance between them while the y one is the value of the spearman correlation coefficient. The squared dot in (0,1) is a fictitious point placed to remark that the spearman coefficient between non-distant species must be 1. The blue line represent the results of a simulation (1000 runs, where we only implemented an unbiased death process) with $R = 0.92t^{-1}$ and the light blue area the standard deviation. We fixed the value of R by fitting this specific trend, and indeed the simulations that use this value of R show a remarkable agreement with the experimental trend. For the algorithm details see methods and appendix.

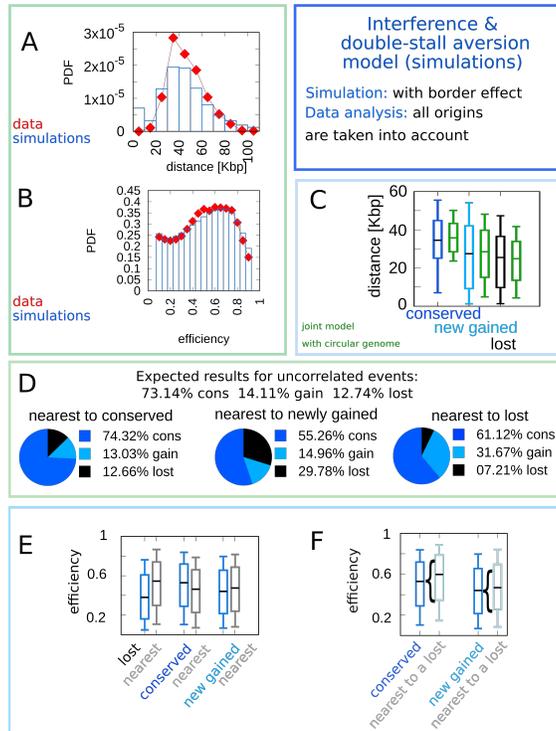


FIG. 3 - Supplement 1. **Linear chromosomes do not alter significantly the model outcomes.** We simulated eight linear chromosomes (the number of chromosomes of the majority of *Lachancea* species), with length equal to one eighth of the average genome size. We have modified the model so that the birth probability at the chromosomes ends is biased by the single stall probability (as double stalls are not possible). The plot shows the results of the simulations (100 runs) of the model. The main difference is visible in the distance distribution shown in panel A. The correlations shown in panels C-F only display minor quantitative changes. In the model, the accumulation of origins towards the chromosome ends is due to the fact that single stall events are more prone to happen than double stalls. Biologically, the region involving the last origin before telomeres is specific, and additional mechanisms such as telomerase or homologous recombination could repair stalled forks [31].

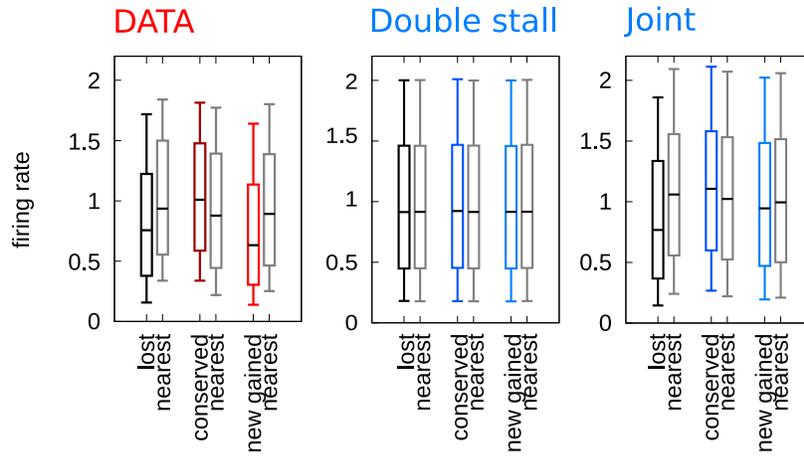


FIG. 4 - Supplement 1. **The efficiency mechanism is necessary to reproduce the correlation between firing rates and evolutionary events.** Comparison between the firing rates-events correlation for experimental data, double stall aversion model and joint model. Only the joint model can reproduce this correlation, which is observed in experimental data. The reason is that in the double-stall aversion model the evolution of firing rates is uncoupled from the origins birth-death dynamics.

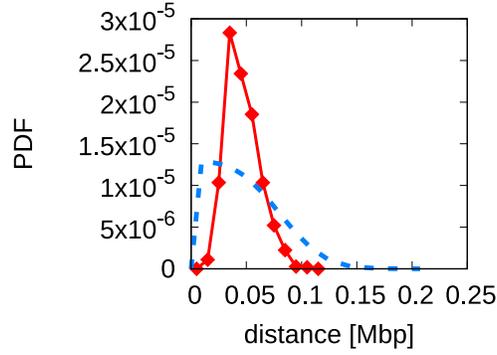


FIG. 4 - Supplement 2. **Analytical predictions for the inter-origins distance distribution falsify the scenario whereby interference alone drives replication-program evolution.**

The plot shows a comparison between the empirical inter-origin distance distribution (red line, diamonds) and the analytical prediction from the scenario of origin birth-death driven by interference alone (blue dotted line, see appendix for the calculation). The predicted distribution does not match the empirical one, thus the scenario can be rejected because it fails to reproduce a crucial feature of the data.

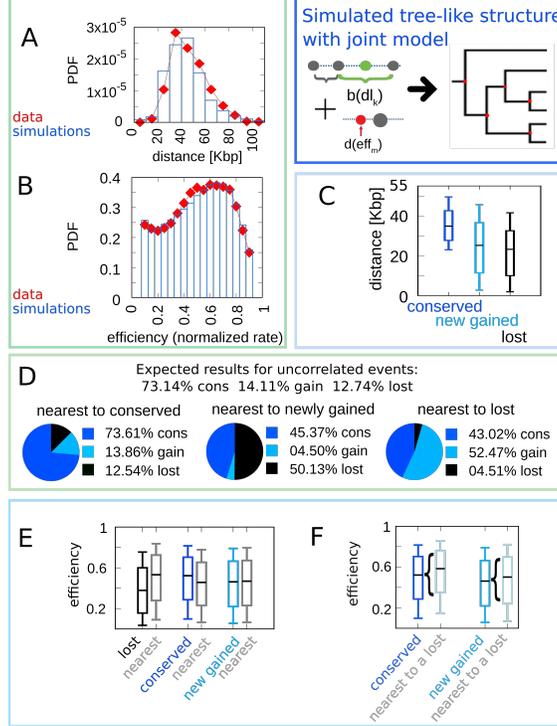


FIG. 5 - Supplement 1. **The joint efficiency/double-stall aversion model simulated on a cladogenetic structure reproduces all the results found for a single lineage.** The results refer to 100 different runs of the simulation of the joint model on the empirical tree structure, compared with empirical data. **A:** Inter-origin distance distribution in simulated species (blue bars) compared to the empirical distribution for the ten *Lachancea* species (red diamonds). **B:** Origin efficiency distribution in simulated (blue bars) *vs* empirical species (red diamonds). **C:** Box plot of the distance from the nearest origin split by evolutionary events, i.e. for conserved (dark blue), newly gained (blue) and lost origins (black), for simulated species. **D:** Fraction of origins that are nearest to conserved, newly gained and lost, for simulated species, compared to the expected result for uncorrelated events. **E:** Box plot of efficiency of lost, conserved and newly gained origins (respectively black, dark blue and blue) and their neighbors (grey), in simulated species. **F:** The efficiency of all conserved and newly gained origins compared to the ones flanking a lost origin. Box plots show the median (bar), 25-75 (box), and 10-90 (whiskers) percentiles. Panels D and F show that the model correctly reproduces the correlation between origin birth-death events over evolution and efficiency of the nearest origin. Simulation parameters (see methods): $\gamma = 2.2$, $\beta = 1.9$, overall birth and death rate $\bar{b} = 13.6\text{Mbp}^{-1}t^{-1}$, $\bar{d} = 0.61t^{-1}$ and rate of origin firing-rate reshuffling $R = 0.92t^{-1}$, where t is measured by protein-sequence divergence.

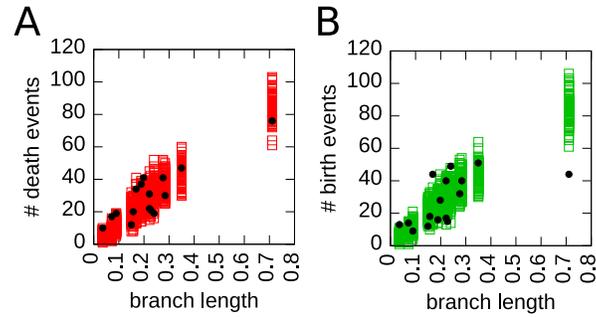


FIG. 5 - Supplement 2. Simulations and empirical data show a similar variability in number of death and birth events across branches of the tree. In each plot, a symbol corresponds to one branch of the phylogenetic tree, empty squares represent the simulations of the cladogenetic structure (100 different runs) and round black circles the experimental data. The x axis represents the branch length, while the y axis is the number of death events (panel A) or birth events (panel B) that occur in that branch. Both plots show a similar spread, supporting the idea that a fixed birth (death) rate in the simulations represents sufficiently well the fluctuations of the number of birth (death) events observed in the data.

Supplementary Files

Supplementary File 1. Results of the simplified log-likelihood tests of the joint and the double stall aversion model with the associated P-values. Positive log-Likelihood differences favor the joint model (see Methods).

Appendix

A. Estimating parameters for the evolution of origin firing rates

This section motivates the model implementation of the evolutionary dynamics of firing rates. In order to quantify the change of origin firing rates over evolutionary times, we studied how the correlation between firing rates of conserved origins behave as species diverge (Fig. 1 - Supplement 3A). To quantify the divergence, for each pair of species in the *Lachancea* clade we calculated the Spearman correlation coefficient between the sets of firing rates belonging to corresponding origins in the two species considered (normalized by the species mean firing rate). We found that the more the species are distant, the less these two sets are correlated, which means that origin initiation rates diverge during evolution and origins lose memory of their initial firing rate. The model describes the evolution of firing rates as follows. Every origin changes its firing rate by extracting a new value from the distribution of empirical normalized ones, regardless of their previous firing rate. This process is characterized by a resampling rate R , common to all the origins, which defines the probability per unit time that an origin resamples its firing rate. The slope of the correlation coefficient in empirical data defines the speed at which the origin firing rates evolve. Hence, it is possible to fit this specific slope and extract the value of R .

In order to do that, we simulated the evolutionary process with unbiased origin death and update of the firing rate. This simulation can be performed without the birth process, because the only origins that one needs to consider in computing the Spearman coefficient between two species are the conserved ones. Each simulation started from 225 origins, with firing rates randomly sampled from the empirical set of firing rates, evolved the genomes changing the firing rates with the resampling process described above and removing the origins according to the death rate estimated from the data. By performing several simulations with different values of the extracting rate R , it is possible to fit its best value. For each R tested, we ran 1000 simulations for an evolutionary time corresponding to 1.6.

After computing the Spearman correlations between snapshots at different evolutionary times, we performed an exponential fit, in order to see which value of the R parameter gave the best agreement with the experimental data, finding the best-fit value $R = 0.92$. Fig. 1 - Supplement 4 shows the trend achieved by the simulation using $R = 0.92$, and it

shows a very good agreement between experimental data and simulations.

Note that in ref. [22], a similar analysis was carried out in order to verify if the reprogramming of the origins firing rate has an impact on the differentiation of replication timing. The authors analyzed the origin firing time *differences* between conserved replication origins in all pairs of species, and found that this difference does not correlate with the phylogenetic distance between species. This finding is apparently in contrast with our results, which suggest that origin reprogramming increases with distance between species. We believe that this discrepancy is due to the higher sensitivity of the Spearman correlation and of the use of species-average normalized firing rates in this study.

B. The empirical data falsify the scenario where interference alone drives origin evolution

This section presents a theoretical analysis of the scenario where solely origin interference sets the evolutionary pressure on replication timing profiles. This analysis shows that a description that only takes into account the evolutionary pressure that acts on origin efficiency is not able to reproduce the origins spatial arrangement, a crucial feature in empirical yeast data. To carry out this analysis, we take a “maximum entropy” approach (see Banavar JR, Maritan A, Volkov I. Applications of the principle of maximum entropy: from physics to ecology. *J Phys Condens Matter*. 2010;22(6):063101. doi:10.1088/0953-8984/22/6/063101) and infer an effective “force potential” acting on inter-origin distance by looking at its (assumed equilibrium) distribution. Specifically, the effective potential acting on the origin efficiency starting from the empirical efficiency distribution, can be analytically computed from the following formula

$$H_{\text{eff}}(\text{eff}) = -\log(P(\text{eff})) \tag{S1}$$

where eff is the efficiency, $\text{eff} \in [0, 1]$, and $P(\text{eff})$ the efficiency probability density function.

The above potential, once given the relation between efficiency and distance between origins, Eq. 4, defines another potential $H_d(d)$ that act on the inter-origin distances. By taking the exponential of $H_d(d)$ one obtains the expected probability distribution predicted for the distances at equilibrium.

In order to find $H_d(d)$ one must to invert Eq. 4 and find $d(\text{eff})$. To accomplish this task, we have approximated the three-body interaction that gives the efficiency with a two-body

interaction. This assumption implies that each origin feels the interference of only one of his two neighbors, and is effective as long as three-origin interactions can be decomposed in two-origin components. Under this assumption, Eq. 4 becomes

$$d_{i,i+1} = -\frac{v}{\lambda} \log \left[\frac{\lambda_i + \lambda_{i+1}}{\lambda_i} (e_i - 1) \right] . \quad (\text{S2})$$

Note that origin efficiency, Eq. 4 also depends on the firing rates of the origin and its neighbor, hence, strictly speaking, one has that

$$H_d(d_{i,i+1}) = H_d(d_{i,i+1}, \lambda_i, \lambda_{i+1}) . \quad (\text{S3})$$

To eliminate the firing rates dependence we computed an effective potential H'_d on the distance, which averages the effect of the different firing rates. To this end, we used the mean value theorem for integrals, as follows,

$$H'_d(d) = \int d\lambda_i d\lambda_{i+1} P(\lambda_i) P(\lambda_{i+1}) H_d(d_{i,i+1}, \lambda_i, \lambda_{i+1}) = H_d(d_{i,i+1}, \langle \lambda \rangle, \langle \lambda \rangle) . \quad (\text{S4})$$

In other words we substituted all the firing rates with the average one $\langle \lambda \rangle = 1$, since the rates are normalized on the species average. With this simplification, going from $H_{\text{eff}}(\text{eff})$ to $H'_d(d)$ is straightforward, and gives

$$d(e) = -\frac{v}{\langle \lambda \rangle} \log[2(\text{eff} - 1)] , \quad (\text{S5})$$

and

$$H'_d(d) = H_{\text{eff}}(d(\text{eff})) . \quad (\text{S6})$$

From the potential H'_d , we can compute the prediction for the equilibrium probability distribution of inter-origin distances

$$P(d) = N \exp(-H'_d(d)) , \quad (\text{S7})$$

where N is a normalization factor. In order to use this calculation on the the data, we inferred the expected potential from the efficiency distribution, assuming that the interaction only depends on efficiency, and we then obtained the model prediction for the expected inter-origin distribution based on the efficiency profile. Comparison of this prediction with the empirical inter-origin distance distribution provides a test of the model. This procedure does not require to adjust any model parameter. Figure 4 - Supplement 2 shows the result of this

analysis. The predicted distribution does not match the empirical one. This means that any evolutionary model that assumes a bias based only on the efficiency (in other words, one that takes into account only the evolutionary pressure given by origin interference) cannot reproduce (at steady state) the correct spatial organization of replication origins.