

1 Integration of abundances and chromatin state
2 data of Long INterspersed Elements reveals
3 dynamics transitions during evolution in
4 mammalian genomes

5 Silvia Vitali^{1,3,*,+}, Enrico Giampieri^{1,+}, Steven Criscione², Claudia
6 Sala¹, Italo do Valle^{1,4}, Nicola Neretti², and Gastone Castellani¹

7 ¹*University of Bologna, Department of Physics and Astronomy,*
8 *40127 Bologna, Italy.*

9 ²*Brown University, Department of Molecular Biology, Cell Biology*
10 *and Biochemistry, Providence, RI 02906 USA*

11 ³*Current affiliation: Basque Center for Applied Mathematics,*
12 *48009 Bilbao, Spain*

13 ⁴*Current affiliation: Northeastern University, 177 Huntington Ave.*
14 *Boston, MA 02115, USA*

15 **svitali@bcamath.org*

16 *+ these authors contributed equally to this work*

17 May 6, 2020

18 **Abstract**

19 Genome ecology and evolutionary biology have being increasingly in-
20 vestigated by interdisciplinary approaches, complementing experimental
21 techniques with advanced modeling and statistical methods. Both disci-
22 plines, with distinct perspectives, have been successful in giving theoret-
23 ical insights of the processes that happen inside and shape the genomes.
24 Distinguishing between evolutionary and ecological origin of genomes pat-
25 terns is not easy, and often the two approaches dedicate to well separated
26 topics. Here, we integrate data of Long-INterspersed Elements (LINEs)
27 abundances in 46 mammalian genomes with the insertions chromatin con-
28 figuration, and their estimated age of amplification, to study the evolution
29 of LINEs ecosystem inside and together with the genome landscape. We
30 describe LINEs amplification dynamics by a birth-death process with as-
31 sumption of competitive neutrality. Then, a competition mechanism for
32 the internal promoter is introduced, spontaneously breaking the neutral
33 assumption. We show that LINEs abundances, as well as the inherent
34 model rates, cluster according to the host taxonomic order. The tempo-
35 ral variation of these rates combined with the average abundances and
36 chromatin state of LINEs copies highlights host-elements interaction and
37 taxa-specific element appearance, such as Lx, associated to the radiation

38 of the murine subfamily, and LIMA/LPB sub-families, related to primates
39 evolution.

40 **Keywords**— Long Interspersed Elements, evolution, transposons ecology, mammals,
41 Approximate Bayesian Computation method

42 Introduction

43 In the last decades an increasingly large number of genomes has been annotated [28]
44 thanks to the availability of high-throughput sequencing methods and the develop-
45 ment of algorithms and bioinformatics tools to process such data [27]. Additionally,
46 extensive studies for chromatin state characterization in different cell lines of several
47 organisms have been published [12, 11, 46], as well as results concerning genomes
48 structure [32] and genomes content statistics [15] and evolution [14]. Therefore, a
49 huge amount of data is available for many organisms nowadays.

50 Transposable elements (TEs) constitute a large portion of most species' genomes,
51 covering roughly 45% of the human genome [7]. Also known as *selfish genes*, they are
52 protein coding DNA sequences that can move from one genomics location to another
53 and/or increase their copy number within the host-genome. Horizontal transfer of TEs
54 between different organisms is also common for several sub-types of TEs.

55 Host organism and TEs interaction has been often described as a host-parasite
56 relationship. The impact over the host organism fitness of TEs appearance and ampli-
57 fication has been extensively investigated to study how TEs copies reach fixation in a
58 population and how TEs abundances are shaped [5, 41, 10, 1, 34, 35, 33]. The loss of
59 fitness due to TEs could be responsible for the appearance of mechanisms of regulation
60 and self-regulation, to limit the number and the impact over the genome functionality
61 of new TE insertions [3, 18, 42]. The appearance of these regulation mechanisms is
62 especially significant for TE sub-types which are unable to transfer from one host to
63 another.

64 The co-evolution of TEs with the genomes and the interaction between these two
65 parties make TEs a fundamental player [25, 31], heavily involved in both evolutionary
66 biology and genome ecology.

67 Interdisciplinary approaches should provide advanced modeling techniques to genome
68 ecology and evolutionary biology to extract most information as possible from the
69 available data. Advanced Bayesian modeling approaches have been already applied
70 successfully for dating the appearance of TEs in genomes, overcoming some limita-
71 tions of standard techniques such as consensus sequence divergence or phylogenetic
72 analysis [14].

73 TE sequences are classified into two classes on the basis of the transposition mech-
74 anism: DNA-transposons and retro-transposons. Their transposition mechanism is
75 respectively DNA- or RNA-mediated. Each class is further partitioned into subclasses,
76 families, subfamilies and elements, on the basis of: their structural organization, the
77 proteins they encode, the sharing of specific insertions, deletions or substitutions. A
78 TE family is composed by the ensemble of copies that share > 80% of their DNA
79 sequence [43]. From an ecological perspective a single copy could be considered as an
80 individual, and all the individuals belonging to the same family or subfamily as con-
81 stituting a species, meaning that they occupy the same ecological niche in the genome
82 [43]. Thus, all the elements belonging to the same family or subfamily can be located
83 at the same trophic level. In the present work we reserve the word *species* to the ele-
84 ment level for modeling purposes. Therefore, different elements belonging to the same
85 family or sub-family are treated as different species, each one with its own number of
86 copies. The copy number of an element constitutes the abundance, or the number of
87 individuals, of that species. Then, according to the previous classification, all the TE

88 species belonging to the same family or subfamily can be located at the same trophic
89 level.

90 A large diversity of TE sequences exists in the genomes in terms of biodiversity,
91 biomass, and abundances. The Relative Species Abundance (RSA) of TEs varies
92 significantly in different organisms and great variability in the amount of TE species
93 and abundances can be observed also at the same trophic level. There is currently
94 no agreement about the relative importance of the possible sources of this variability:
95 such as the host specific selection pressure, both at the genome and host population
96 levels, and the stochasticity of the forces acting on the individual copies [43, 34, 24].

97 The interdependence between the TE community and the host genome, together
98 with the replication mechanisms of the elements, suggests a strong parallelism between
99 TEs dynamics in the genome and species community dynamics in the ecosystem [43,
100 37]. Both the niche and the neutral theory have features convenient to describe TEs
101 ecosystem. Niche theory is based on the partitioning of resources between compet-
102 ing species [6]. In the neutral theory, the stochastic mechanisms as demographic
103 stochasticity, migration, and speciation are the most important forces shaping the
104 community [17]. However, TEs ecosystem contains some peculiarities that differenti-
105 ate it from standard ecosystems [43]. TEs create and continuously reshape their own
106 environment, because the copies that lose any transposition ability generate a large
107 part of the genomics landscape in which new copies may insert without deleterious
108 effect on the cell functionalities. Furthermore, the natural selection acts on two levels,
109 the genome level and the host one. In this regard, we distinguish *transposon ecol-*
110 *ogy*, describing the interaction of TEs with the genome and cellular environment only,
111 from *genome ecology*, which includes the interaction with the external environment
112 mediated by the host fitness [24].

113 The selection at the level of the host could eventually induce TEs to evolve traits
114 that constitute a selective disadvantage at the individual level, as for example a lower
115 transposition rate [16, 33]. Phenomena related to the molecular nature of TEs may
116 also occur, for example mutations, insertions and sequence rearrangements, which may
117 lead to functional variations of the elements.

118 Here we focus on the study of a particular family of *non-long terminal repeats*:
119 the Long Interspersed Elements (LINEs). We restrict the model to LINE family only,
120 discarding the possible interaction with different TEs entities. Future studies may
121 include information about inter family interactions, for example Short INterspersed
122 Elements (SINE) -LINE parasitism [29].

123 We model LINEs copy number distribution inside a cohort of mammalian reference
124 genomes under the hypothesis of competitive neutrality [23]. Competitive neutrality
125 represents the absence of competitive differences among different LINE species. Thus,
126 all the copies of all elements in the community could be characterized by the same
127 transposition activity, sequence divergence, and death rate [43]. The variability of the
128 *abiotic* component (in this context genes, repetitive sequences, and intracellular com-
129 ponents) of the ecosystem is further reduced by including only mammalian genomes
130 in the study.

131 LINEs are the most abundant family of TEs in mammals, in terms of biomass.
132 They belong to the retro-elements class, which means that their replication is RNA
133 mediated (figure 1). They are also incapable of horizontal transfer [1]. Full-length ele-
134 ments contain a promoter region (5'UTR), two protein coding regions (ORF1, ORF2)
135 and a poly-A tail (3'UTR). The internal promoter directs transcription initiation, and
136 permits autonomous transposition. When the transcribed RNA reaches the cytoplasm,
137 the protein encoding regions ORF1 and ORF2 are translated to an RNA-binding pro-
138 tein and a protein with endonuclease and reverse-transcriptase activities, respectively.
139 Both proteins show a strong cis-preference; consequently, they preferentially associate
140 with the RNA transcript that encoded them to produce what is called a ribonucleo-
141 protein (RNP) particle. After coming back into the nucleus, the proteins on RNA can

142 open a nick in DNA and produce a DNA copy of the template through a process termed
143 target-primed reverse transcription (TPRT). The new insertions often result in low fi-
144 delity copies of the parent LINE [8]. Some transposition events are incomplete such
145 that the inserted copy is incapable of autonomous retro-transposition; for example, L1
146 insertions are often 5'-truncated (e.g. Figure 6B of [9]). Furthermore, a transcribed
147 incomplete copy can hijack the retro-transposition machinery of autonomous copies
148 to duplicate into a new location: a process called trans-complementation. The phe-
149 nomenon of trans-complementation has been observed, for example, in LINE-1 retro-
150 elements, although it should happen at a much smaller rate than retro-transposition
151 in *cis* [Wei2001].

152 Despite occasional re-activation of inactive elements has been observed in certain
153 diseases, LINE community in mammalian genomes is mainly composed by the col-
154 lection of defective and/or silenced copies of inactive elements that reached fixation
155 in the genome. The stratification of such elements in the genome, in the absence of
156 horizontal transfer, can be used to infer the changes in time that may have occurred
157 in the dynamics of LINES.

158 Some peculiarities of LINES replication mechanism and their evolutionary history
159 inside mammalian genomes support our hypothesis that LINES community in mam-
160 mals could be successfully described by a birth-death process under competitive neu-
161 trality hypothesis. LINES evolved often on a single lineage, in particular in primates
162 [22], with a subsequent appearance of active elements, making competition between
163 different elements negligible. Coexistence of multiple L1 lineages is documented for
164 ancient LINES [38] and currently in mouse [26], where L1 frequently recruited novel
165 5'UTR sequences [40], suggesting that simultaneous activity of non-homologous pro-
166 moters does not introduce a competition between the elements. Finally, the genome
167 environment is unique to each of the TE copies and full-length L1 copies may differ
168 randomly in their level of transposition activity [4, 36]. Therefore, the stochasticity
169 at the individual level could have a significant impact on the structure of the entire
170 community, supporting the neutral approach to describe the community dynamics.

171 Here we model the distribution of LINES copy number through a Master Equation
172 approach [2, 21], under the hypothesis of competitive neutrality and an alternative
173 hypothesis of competition for the promoter region. The competition is occasional
174 between two species and induces a reduction of the birth rates of the competing species,
175 breaking spontaneously the hypothesis of competitive neutrality. The two models
176 are nested. They are tested by fitting the RSA of LINES communities through a
177 hierarchical Approximate Bayesian Computation (ABC) method. A sliding window
178 analysis is applied to study the evolution of the RSA in time by the same ABC method
179 together with the chromatin state characterization of the individual copies.

180 Results and methods are summarized in the respective sections. The parameters
181 expectation associated to several regimes of the model and the results are extensively
182 analysed in the discussion section.

183 Results

184 The RSA of LINES in 46 mammalian genomes have been fit by a negative binomial
185 and by a mixture of negative binomial through hierarchical ABC method (see, sec-
186 tion methods for details). From the same prior distribution, we obtained different
187 posteriors of the model parameters for the RSA of 42 datasets, with the exception of
188 Wallaby, Tasmania devil, Opossum, and Platypus, for which the rate of success of ABC
189 method was extremely small in comparison to the others 42. Posteriors distribution
190 and model comparison with the LINES RSA are shown in Supplementary Materials
191 (supplementary figures 6-15). The goodness-of-fit of the two models, assessed through
192 the likelihood-ratio test, shows that the models are comparable for the majority of the

193 data set. In figure 2 the expected value of the parameters are shown, color labeled by
194 the taxonomic order of the corresponding host genome.

195 LINEs community in human, chimp, rhesus macaque, mouse and rat genomes
196 have been sorted by their relative time of appearance and amplification according to
197 published genome-wide defragmentation results [14]. For each of these collections we
198 performed a sliding window analysis (see section methods for details) of RSA distribu-
199 tion patterns, average percentage of insertions in open chromatin regions, and average
200 abundance. Here we show the results for the windows length $N = 15$. Windows
201 lengths of $N = 30, 40$ produce consistent results, with smoother temporal trends.

202 The sliding windows of RSA distribution patterns have been tested with the same
203 method applied previously to the entire timeline. In figure 3 the likelihood-ratio test
204 of the two models and the mixture coefficient α of the competition model are shown.
205 The results for the other parameters are reported in the supplementary materials.

206 In figure 4 the sliding window analysis of the average abundance and the average
207 percentage of insertions in open chromatin regions [46, 12] is shown. The average
208 percentage of LINE copies belonging to open chromatin regions displays a decreasing
209 temporal trend for all the organisms for which such information was available (human
210 and mouse). The average abundance is decreasing for the primates cohort and display
211 a fast transition to larger abundances in the rodents cohort.

212 In figure 5 the estimated time range of activity of LINE elements inside the human
213 genome[14] belonging to the windows range most significant for competition (largest
214 values of the mixture coefficient in figure 3) is shown. Elements are color labeled by
215 their abundance in the human genome. The presence of significant similarity between
216 5'UTRs pairs (see section for details) is highlighted for the following high and low copy
217 number pairs (or group): L1M2-L1M2c and L1MA9-L1M3a-L1M3b-L1M3c-L1M3d.

218 Hierarchical clustering of the LINEs abundances inside the 46 genomes, for the
219 elements shown in figure 5, is reported in the supplementary material (supplementary
220 figure 19). Elements are color labeled according to their abundance. Rare or abun-
221 dant classification of the elements across all the mammals included in this study is
222 consistent, except for the White Rhinoceros, which shows the opposite trend. The
223 corresponding clustering for the host species is also mostly in agreement with the
224 taxonomic classification.

225 The study of the correlation between the negative binomial parameters x, Υ ob-
226 tained with the sliding window analysis of the RSA patterns is shown in figure 6.

227 The negative binomial parameters x, Υ of both cohorts can be clearly separated
228 in two clusters, divided by the time of appearance of specific elements (L1MA/LPB
229 for primates and Lx for rodents). One cluster contains all the samples before the
230 appearance, the other all the sample after. Thus, a transition in time of the correlation
231 of the parameters can be identified for both the group of primates and rodents. The
232 same transition can be observed for the neutral model description as well as for the
233 competition model (mixture of two negative binomials), by looking at the component
234 describing the elements with large copy number.

235 The panel in figure 6 containing the parameters of the negative binomial component
236 describing the less abundant LINEs do not display significant correlation. The other
237 parameters describing the sliding windows RSA are shown in supplementary figures
238 17-18.

239 By using chromatin state assignments in human [12] and mouse [46] genomes, and
240 the coordinates of the respective LINEs insertions from RepBase, we assigned to each
241 LINE copy a chromatin configuration, distinguishing between open and closed states
242 (*euchromatin* and *heterochromatin*). In figure 7 the correlation between the number of
243 insertions in euchromatin (ECN) and the number of insertions in heterochromatin re-
244 gions (HCN) and the corresponding 2D principal component analysis (PCA) is shown.
245 Two clusters are clearly visible for mouse data set. The elements can be classified in
246 the two clusters according to the same time threshold identified in figure 6 for mouse.

247 Information about the time of appearance is included in the PCA, to highlight the
248 separation between the two clusters. The logarithms of ECN and HCN in figure 7 are
249 strongly correlated. To estimate the correlation, a linear regression between the loga-
250 rithm of the counts has been performed. This correlation corresponds to a power-law
251 relationship between the raw counts:

$$N_{Eu} = 2^{c_0 \pm \epsilon} (N_{Het})^c. \quad (1)$$

252 We obtained for human: $c = 1.18$, $c_0 = -4.58$ and $\epsilon = 0.035$, which correspond to
253 the standard error in the estimate. The correlation coefficient is $r = 0.96$ with p-value
254 $p \sim 10^{-55}$. For the most ancient group of LINE in mouse we obtained $c = 1.12$,
255 $c_0 = -4.84$, $\epsilon = 0.043$, $r = 0.96$, $p \sim 10^{-33}$, for the most recent we obtained $c = 1.03$,
256 $c_0 = -5.84$, $\epsilon = 0.084$, $r = 0.90$, $p \sim 10^{-13}$.

257 The beginning of the transition in figure 7 for the mouse data set is contemporary
258 to the transition of the parameters in figure 4 and figure 6 for the rodents cohort.
259 Regarding the primates cohort, in 7 we highlighted the same group of elements defined
260 by the transition in 6 for sake of clarity, but a sharp transition as the one observed for
261 the rodents cohort is absent.

262 Discussion

263 We modeled the way LINES populated different mammalian genomes as a birth-death
264 process of two interacting sub-types: full-length (autonomous or active) copies and
265 incomplete (non-autonomous or inactive) copies. The biological processes that char-
266 acterize LINES retro-transposition activity (replication, mutation, disappearance, ex-
267 tinction, etc.) have been described in terms of transition rates by a Master Equation
268 (ME) approach (equation 2). The stochastic processes described by the ME, with the
269 corresponding rates, are depicted in figure 1.

270 We analyzed and tested the two variants of the model proposed, with and with-
271 out competition, to describe LINES communities over 46 mammalian genomes. We
272 focused on the two realistic dynamics regimes characterized by a specific asymptotic
273 stationary solution to describe LINES RSAs: a negative binomial distribution in case
274 of competitive neutrality, and a mixture of two negative binomial distributions in case
275 of direct competition between elements.

276 The negative binomial distribution depends on two parameters: the "probability
277 of success", x , and the "number of failures", Υ . If trans-complementation does not
278 produce a relevant contribution and the system is out of equilibrium, we expect to
279 observe an RSA following a negative binomial distribution with $\Upsilon \sim 1$. Instead, if the
280 trans-complementation process is relevant we expect $\Upsilon \sim \frac{(b_I + d_A)}{b_{AI}} \gg 1$, due to the
281 experimental observation that trans-complementation events are rarer if compared to
282 retro-transposition in cis.

283 We considered that simultaneous activation of elements sharing the same promoter
284 may introduce a disadvantage for the species that compete for the molecular machinery.
285 According to our model a competitor could appear, with a certain probability, every
286 time a new active copy is created in the system. Thereafter, the extinction of one of
287 the competitors restores the neutrality.

288 Both the variants of the model reproduce the fundamental features of LINES RSA
289 patterns inside 42 of 46 mammalian genomes. The four exceptions are Wallaby, Tas-
290 mania devil, Opossum, and Platypus, which are characterized by a smaller number of
291 resident LINE elements. From an evolutionary perspective they also correspond to a
292 well defined subgroup of host genomes (marsupials and monotremes). Instead, the 42
293 genomes successfully described by the model belong to the Eutheria clade.

294 The expectation of the parameters of the competition model (mixture of two neg-
295 ative binomials) permits to separate the host species at the level of taxonomic order

296 (figure 2 and supplementary figure 5). Instead, the neutral model produces a less pro-
297 nounced separation between taxa (figure 2). The couples of parameters Υ_1 , Υ_2 and
298 x_1 , x_2 seem the best representation to discriminate the host organisms in different
299 taxonomic orders. Within our description, such couples of parameters are related by
300 the average value of the disadvantage due to competition ($\sim \frac{(b_{A,1})}{(b_{A,2})}$). The value of
301 the failure parameter remain always closed to one ($\Upsilon \sim 1$), suggesting that a pure
302 accumulation process is more convincing ($b_{AI} \approx 0, d_I \approx 0$). Hence, despite trans-
303 complementation may take place, our results suggest it is not a very relevant process
304 in shaping the RSA of the community. The shape of RSA, together with the rareness
305 of trans-complementation events, supports the idea that equilibrium between host and
306 LINE population does not hold in general ($b_A \ll d_A$), but a competition between
307 the host and LINE species takes place. In fact, in the case of equilibrium for the
308 active sub-type dynamics ($b_A \approx d_A$), numerical simulations display heavier tails in the
309 generated RSA (supplementary figure X). Such excess of abundant elements could be
310 compensated by a higher rate of excision from the genome (d_I) (supplementary figure
311 X). However, this parameter configuration results less convincing, at least in primates,
312 because more ancient LINE elements are more abundant on average than the recent
313 one (figure 4).

314 For the majority of the organisms under study, the likelihood-ratio test for the two
315 model examined was of order one, suggesting no clear preference between the assump-
316 tions of competitive neutrality and competition. The ABC method applied already
317 penalizes the model with the higher number of parameters because the phase-space is
318 larger, hence the two model results equivalently acceptable from a purely statistical
319 perspective. The law of parsimony supports the choice of the simplest theory when
320 two alternatives possess equivalent power to describe the data. However, the intro-
321 duction of competition between LINE species permits to better distinguish different
322 taxonomic orders based on our model description, suggesting that the introduction of
323 the second component in the PDF enhance the ability to extract useful information
324 from the data and a better characterization of the biological phenomena (figure 2).

325 In fact, hierarchical clustering of the LINEs abundances in different organisms
326 (supplementary figure 16) displays the same tendency to aggregate host organisms
327 belonging to the same taxa. Horizontal transfer is very uncommon for the LINE family
328 and LINEs mutually evolve with the genome with a turnover of active species. Hence,
329 closely related host species could possess more similar LINEs abundances because they
330 have been inherited more recently from a common ancestor. For this reason, it cannot
331 be excluded that the better performance of the mixture model is due to a better
332 characterization of the statistical fluctuations in the RSA that have been inherited
333 (supplementary figure X).

334 The evolution of the LINEs community has been investigated by sorting LINE
335 elements according to their time of appearance and amplification in the host genomes
336 [14] and by applying a sliding window analysis. The resulting time-dependent host-
337 specific RSA have been tested for the neutral and competition models, by applying
338 the same ABC method and prior distributions employed previously for the entire
339 timeline. Where the mixing coefficient of the mixture model is higher, ABC model
340 selection score supports the presence of competition (figure 3). The sliding window
341 analysis suggests that, at specific times during the evolution in the mammalian genome,
342 multiple concurrently active LINE subfamilies might have been in direct competition
343 for the promoter region (figures 3 and supplementary figures X). Despite the time-
344 dependent RSA is host-specific, some patterns are recurrent between the majority
345 of the host organisms under study supporting the hypothesis of inheritance of the
346 abundances (supplementary figures X). The hypothesis that competition could have
347 been shaped by the LINE 5'UTR structure is supported by the similarity (measures of
348 distance between pairwise aligned sequences) of the 5'UTR sequences in concurrently
349 active LINEs (figure 5).

350 The transition in the correlation of the parameters x , Υ in figure 6 is thus associated
351 to a transition to lower average abundance for primates and higher average abundance
352 for the two rodents. These transitions are in agreement with the trend in time of
353 the abundance (figure 4) and further supported by the transition in the chromatin
354 landscape of the LINE copies for mouse (7) and by the concurrent amplification of
355 host-specific LINE elements.

356 The average percentage of LINE copies inserted in euchromatin regions in the
357 sliding window displays a decreasing trend with time ordered age in human and mouse
358 (figure 4). However, in humans, it also shows a clear peak within the neutral time
359 interval. The average percentage of copies in euchromatin regions is bigger for the
360 windows with higher average copy number respect to the one with low copy number
361 in human and in ancient LINE species in mouse. The presence of a higher fraction of
362 rare species within the non-neutral time interval results then in agreement with the
363 lower average percentage of insertions in euchromatin observed.

364 The super linear correlation observed between the number of copies belonging
365 to open and closed chromatin regions (figure 7) lead to the interesting result that a
366 higher copy number, i.e., the sum of euchromatin and heterochromatin contributions,
367 is associated to a higher percentage of insertions in euchromatin states. A reduction
368 of the average percentage of insertions in euchromatin regions thus corresponds to
369 a reduction in the average abundance. Moreover, the presence of the same type of
370 correlation in human and mouse genomes, shared by all the most ancient elements,
371 indicates the existence of a common pattern in the chromatin landscape of LINEs in
372 mammals.

373 In agreement with such prediction, in primates, the decreasing trend of the inser-
374 tion percentage in open chromatin regions is accompanied by a decreasing trend of
375 LINE species abundance in time (figure 4). On the contrary, referred to mouse, the
376 average abundance at a certain point drastically rears up. The time point at which the
377 abundance rears up corresponds, in figure 7, to a transition to a smaller value of the
378 coefficient c_0 (see section Results), while for ancient LINE species the correlation trend
379 in mouse genome is very close to the one observed in human. Given the same number
380 of insertions in open chromatin regions, a lower value for c_0 corresponds to a larger
381 abundance, and, consequently, to a lower percentage of insertions in euchromatin.

382 The beginning of the transition in figure 7 is defined as the time of appearance
383 of the most ancient element belonging to the cluster of recent elements. In mouse
384 data set, such transition is contemporary to the transition to higher average LINEs
385 abundance shown in figure 4 and figure 6, and corresponds to the appearance of the
386 LINE family Lx. Where the amplification of the LINE family Lx is associated with
387 the murine subfamily radiation ~ 12 Myr ago according to [30, 13]. In fact, the other
388 elements characterizing this group are mainly murine specific.

389 In figure 7, referred to human, there is not a sharp transition between two different
390 chromatin state distributions as observed for mouse. This is reasonable if we look at
391 the problem from the perspective of the host organism fitness. A transition to a
392 higher average copy number (figure 4) surely have a bad impact on the host fitness, if
393 it is not compensated by some host defence mechanisms, because the probability of
394 deleterious insertions in the genome increases. The combined sharp transitions in the
395 chromatin landscape and abundances, further associated to the evolutive transition of
396 the host genome (murine subfamily differentiation), observed in mouse support this
397 idea, and are perhaps the result of the competition between the host and LINEs. In
398 human instead, a transition to a lower average copy number is observed (figure 6
399 and 4). This could be the reason why chromatin states distribution in human is not
400 affected significantly, since further changes were not necessary to preserve the host
401 fitness. Indeed, the most ancient LINE species involved in the transition depicted in
402 figure 6 are related to the evolutive differentiation of Primates, associated with the
403 amplification of LIMA/LPB subfamilies $\sim 70 - 100$ Myr ago [22].

404 Conclusion

405 The present research is a first attempt to answer some of the fundamental questions
406 concerning LINEs dynamics and co-evolution with the genome by combining different
407 data types through an interdisciplinary approach. Data of LINEs abundance in 46
408 mammalian genomes [39] have been integrated with data about the relative time of
409 appearance and amplification of the elements inside a sub-group of host genomes [14],
410 and with data about the current chromatin state of the LINEs insertions inside human
411 and mouse genomes [46, 12].

412 The mechanism of competition proposed between LINE species is independent
413 of the chromatin state distribution of the copies, but acts at the level of the LINE
414 species affecting their abundances. Instead, the chromatin state of LINE copies and
415 average abundance should reflect the interaction of LINE species with the host, by
416 mechanisms of silencing (for example methylation) and self-regulation (for example
417 selection of elements with lower birth rates or with specific genomics region preference
418 of the new insertions).

419 The analysis shows that LINEs abundances, as well as the inherent average birth-
420 death rate obtained through the model, cluster according to the host organism taxo-
421 nomic order. Model selection and promoter similarity analysis support the idea that
422 ancient sub-groups of LINEs could have been in direct competition within mammalian
423 genomes. Furthermore, the variation in time of the model parameters, combined with
424 the average abundances and chromatin state of LINEs copies, displays evidences of
425 host-elements interaction and features highlighting taxa-specific element appearance,
426 such as Lx, associated to the radiation of the murine subfamily, and LIMA/LPB sub-
427 families, associated to primates evolution. The sliding window analysis shows that
428 the decreasing trend of euchromatin percentage of insertions is shared by primates
429 and rodents cohort. Chromatin information highlights also a super linear correlation
430 between the insertions in open chromatin regions and insertions in closed chromatin re-
431 gions. This type of correlation indicates that a increasingly small percentage of copies
432 inserted in open chromatin regions will be associated to an increasingly small abun-
433 dance of the element (sum of the number of element insertions in open chromatin and
434 in closed chromatin regions). In fact, Lx appearance, associated to an increase in the
435 average LINEs abundance in rodents, is also characterized by a transition to a different
436 correlation regime of the number of insertions in different chromatin configurations.

437 We believe that interdisciplinary research could positively contribute to improve
438 field specific research methodologies, and possibly complete and enrich the perspective
439 of the specialists. The possibility to map in different systems (hosts) the evolution of
440 genomics elements, with careful considerations about the nature of such elements,
441 could become a powerful tool to understand present and past dynamics of the entities
442 that contribute in shaping the genome, and to bridge evolutionary biology and genome
443 ecology investigations. Despite the relative simplicity of the data types included in the
444 analysis, the results obtained are in agreement with several independent studies [22,
445 38, 26, 30] and encourage the application of similar approaches by integrating more
446 refined information at the copy level, such as DNA sequence, chromosome location
447 and genome patterns (CG/AT content, chromatin configuration, etc.), to investigate
448 specific questions at the edge of evolutionary biology, genome ecology, and transposon
449 ecology fields.

450 Methods

451 Data sources

452 LINE abundances were calculated using RepeatMasker annotation (<http://www.repeatmasker.org>)
453 [39] for human genome build hg19 and 45 other mammalian species. LINE consensus

454 sequences were downloaded from RepBase [19, 20] (<http://www.girinst.org>). Chrono-
455 logical ordering of LINEs in human, chimp, rhesus macaque, mouse and rat genome
456 was derived from genome wide defragmentation results [14]. Chromatin structure data
457 are available for mouse [46] and human [12]. The employed chromatin state assignment
458 was conducted by using ENCODE chromatin models from the ChromHMM method
459 [11].

460 Chromatine state assignment

461 Chromatin structure data were used to assign each LINE copy to open or closed chro-
462 matin state by the knowledge of their coordinates in the reference genome. Open and
463 closed chromatin states were defined and located according to the available classifica-
464 tion for mouse [46] and human [12] for the germ line. Multiple assignments of the same
465 genome region have been treated by classifying the combination of states into open,
466 weakly open and closed chromatin, depending if the assigned configurations belong
467 mainly to one of these groups. We tested grouping weakly-open chromatin population
468 with both open chromatin and closed chromatin to assess if this choice affected our
469 results, but did not observed any significant difference. Thus, the weakly-open and
470 the unknown state have been included in the closed chromatin group, which encloses
471 most of the LINE copies.

472 Sliding window analysis

473 LINEs species are sorted by their relative time of appearance and amplification in the
474 genome according to specific host genome analysis [14] (different genomes may show
475 variations of this time series). By dividing the time series into intervals (windows),
476 each of them containing a fixed number of elements, a different realization of the LINEs
477 ecosystem is obtained. By sliding the window, a different sub-sample of elements active
478 in a distinct evolution stage of the genome can be selected, representing a picture of
479 the LINE community in a different evolutionary stage of the genome. The time series
480 of the windows depicts several variations of the LINEs community in time in terms of:
481 RSA patterns; average abundance; and average percentage of insertions in different
482 chromatin configurations. The RSA patterns are described by the same ABC method
483 applied to the entire timeline.

484 Stochastic model

485 The LINEs dynamics inside the genome is described via the following two-dimensional
486 Master Equation:

$$\begin{aligned} \frac{dP}{dt} = & (E_{n_A}^- - 1)b_A n_A P + (E_{n_A}^+ E_{n_I}^- - 1)d_A n_A P \\ & + (E_{n_I}^- - 1)b_I n_A P + (E_{n_I}^- - 1)b_{AI} n_A n_I P \\ & + (E_{n_I}^+ - 1)d_I n_I P. \end{aligned} \quad (2)$$

487 where we consider $P \equiv P(n_A, n_I, t)$ and the Van Kampen step operators [21]
488 $E_n^\pm f(n, m) = f(n \pm 1, m)$, the lower index n indicates the variable on which the op-
489 erator acts, the upper index $+$ or $-$ determines the direction of the unitary change
490 of the value of the variable n . n_A and n_I represent respectively the number of full-
491 length copies (autonomous in self replication) and the number of defectives copies
492 (non-autonomous or inactive) of a specific element in a genome. Each term in on the
493 right-hand side of equation 2 represents one of the stochastic process examined to de-
494 fine LINEs dynamics: b_A, d_A , rate of birth and death of full-length copies; b_I, d_I , rate
495 of birth and death of defective copies; b_{AI} , rate of birth of defective copies by trans-
496 complementation. The stationary distribution of the bidimensional model (equation

2) doesn't have a closed form, however, it is possible to compute the marginal distributions, corresponding to the n_A and n_I species. The distribution $P(n_I)$, to which it will be referred as P_{n_I} , corresponds to the *relative species abundance* (RSA) of the LINE species, which represents the probability to observe a species with a certain number of *individuals* (the copy number) inside a community (the genome).

When equilibrium is reached for both active and inactive copies, the RSA corresponds to a negative binomial distribution. In fact, taking n_A as a constant, the equation 2 for n_I describes a well-known ecological neutral model, successfully applied to coral reefs and rain forests [45]. If equilibrium for active species does not hold ($b_A \ll d_A$) and excision and trans-complementation processes are neglected, the stationary solution P_{n_I} is still approximated by a negative binomial distribution, obtained when the "absorbing state" is reached ($n_A = 0$). The two regimes can be distinguished because the expected value of the parameters are different for biological considerations (see section discussion).

A spontaneous breaking of the competitive neutrality assumption is introduced by assuming that contemporary activation by the same promoter region of two LINES reduces the birth rates b_A , b_I and b_{AI} in equation 2 by a factor $n_1/(n_1 + n_2)$ where n_1 and n_2 are the full-length copies of the two competing elements. A smaller birth rate will result in a lower abundance for the competing LINE species, in comparison to the elements not affected by the competition mechanism, and will induce deviations from the expected distribution by generating a bimodal behavior. We surmise that the distribution arising from this type of competition is a mixture of two negative binomials (equation 3) for a range of parameters compatible with real data (confirmed by numerical simulation):

$$P_{RSA} = \alpha \cdot P_{rare} + (1 - \alpha) \cdot P_{abund}, \quad (3)$$

where α is the mixture coefficient, related to the probability that two elements compete, and with P_{rare} and P_{abund} representing respectively the RSA of rare elements.

The negative binomial distribution depends on two parameters: the "probability of success", x , and the "number of failures", Υ . The expectation value of the copy number for the LINES RSA is defined by the parameters of the negative binomial distribution by the relation:

$$\langle n_I \rangle = \frac{x}{1 - x} \cdot \Upsilon \quad (4)$$

where, in the case of competition, the parameters of the distribution of the abundant species, which contains the majority of the LINE species, can be employed.

Hierarchical approximate Bayesian computation

A hierarchical approximate Bayesian computation (ABC) method has been implemented in Python to fit the RSAs included in this study. The method can be schematized by the following procedure: (i) a set of parameters is taken from uninformative prior distributions to build a negative binomial or a mixture of two negative binomials distribution representing the RSA distribution; (ii) a sample of abundances is generated according to such distribution; (iii) the set of parameters taken from the priors is accepted if the distance between the generated sample and the empirical LINES abundance of at least one of the 46 genomes is under a certain threshold, (iv) the procedure (i-iii) is repeated $\sim 10^6$ times to build a posterior of the parameters; (v) the posterior is used to replace the uninformative prior with a more informative one, to increase the statistics and reduce computation efforts. The procedure (i-iv) is then repeated for each of the 46 data set to build a posterior distribution of the parameters specific for each of the data set. The posterior distribution of the parameters represents the probability that a certain set of parameters describes the data according to the model. The expectation value of the parameters with respect to the posterior

545 distribution is thereafter employed to describe the data. Model selection is performed
546 according to the likelihood-ratio test information criterion, approximated to be the
547 ratio of successes (number of set of parameters accepted over the total number of set
548 tested) for the two models.

549 Numerical simulations

550 To test if the dynamical model can generate a negative binomial distribution be-
551 yond the given assumptions, we performed numerical simulations. We used Gillespie
552 algorithm for the active copies dynamics, for the inactive copies dynamics we used
553 the tau-leap algorithm for the case $b_{AI} = 0$ and a hybrid algorithm when trans-
554 complementation is considered. The hybrid algorithm was chosen instead of the Gille-
555 spie one to reduce the time of computation. It consists in the estimation of the
556 expected number of inactive copies by ordinary differential equation (ODE) numerical
557 integration at the beginning of each time interval $\bar{n}_{I,t}$. Then, tau-leap algorithm is
558 applied to generate a stochastic increment associated to the time interval $\Delta n_{I,t}$. The
559 number of inactive copies at the end of the time interval is thus determined by the
560 sum $n_{I,t+1} = \bar{n}_{I,t} + \Delta n_{I,t}$. Oracle comparison to the theoretically correct Gillespie
561 algorithm was performed to test the accuracy of the hybrid simulation method. Sim-
562 ulations of the competition mechanism in both the regimes ($b_{AI} = 0$ or $b_{AI} > 0$)
563 were performed to check if the solution was compatible with a mixture of negative
564 binomials.

565 *More details about these methods can be found in the Supplementary Materials (sup-*
566 *plementary figures 1-4).*

567 Acknowledgements

568 The results presented in this work are part of the PhD Thesis of SV [44]. This work has
569 been supported by the Italian Ministry of Education at University of Bologna (Alma
570 Mater Studiorum), Department of Physics and Astronomy (DIFA), by the Basque
571 Center for Applied Mathematics (Bilbao, Spain) and in part by the following NIH
572 grants: R56 AG050582-01 to N.N. and F31AG050365 to S.W.C.. S.W.C. was also
573 supported by the NIH Institutional Research Training Grant T32 GM007601. We also
574 acknowledge IMforFuture EU project and HARMONY EU project.

575 Author contributions statement

576 N.N. and G.C. conceived the study, S.W.C. and I.F.V. prepared the data sets, E.G.,
577 S.V. and C.S. implemented the ABC pipeline and simulations, S.V. performed and
578 developed the analysis.

579 Author competing interest statement

580 All authors reviewed the manuscript and declare no conflict of interest.

581 References

- 582 [1] G. Abrusán and H. J. Krambeck. “Competition may determine the di-
583 versity of transposable elements”. In: *Theoretical Population Biology* 70.3
584 (2006), pp. 364–375.

- 585 [2] Animesh Agarwal et al. “On the precision of quasi steady state assump-
586 tions in stochastic dynamics”. In: *The Journal of Chemical Physics* 137.4
587 (2012), p. 044105. DOI: 10.1063/1.4731754.
- 588 [3] S. Boissinot and A. V. Furano. “Adaptive evolution in LINE-1 retro-
589 transposons”. In: *Journal of Molecular Biology and Evolution* 18 (2001),
590 pp. 2186–2194.
- 591 [4] B. Brouha et al. “Hot L1s account for the bulk of retrotransposition in the
592 human population”. In: *Proceedings of the National Academy of Sciences*
593 *USA*. 100, 2003, pp. 5280–5285.
- 594 [5] B. Charlesworth and C. H. Langley. “The population genetics of *Drosophila*
595 transposable elements”. In: *Annual review of genetics* 23 (1989), pp. 251–
596 287.
- 597 [6] J. M. Chase and M. A. Leibold. *Ecological Niches: Linking Classical and*
598 *Contemporary Approaches*. Chicago: University Press, 2003.
- 599 [7] International Human Genome Sequencing Consortium. “Initial sequencing
600 and analysis of the human genome”. In: *Nature* 409.6822 (2001), pp. 860–
601 921.
- 602 [8] R. Cordaux and M.A. Batzer. “The impact of retrotransposons on human
603 genome evolution”. In: *Nature Review Genetics* 10.10 (2009), pp. 691–703.
- 604 [9] S. W. Criscione et al. “Transcriptional landscape of repetitive elements in
605 normal and cancer human cells”. In: *BMC genomics* 15 (2014), p. 583.
- 606 [10] G. Deceliere, S. Charles, and C. Biémont. “The dynamics of transposable
607 elements in structured populations”. In: *Genetics* 169.1 (2005), pp. 467–
608 474.
- 609 [11] J. Ernst and M. Kellis. “ChromHMM: automating chromatin-state dis-
610 covery and characterization”. In: *Nature Methods* 9 (2012), pp. 215–216.
- 611 [12] J. Ernst and M. Kellis. “Discovery and characterization of chromatin states
612 for systematic annotation of the human genome.” In: *Nature biotechnology*
613 28.8 (2010), pp. 817–825.
- 614 [13] A. V. Furano et al. “Amplification of the Ancient Murine Lx Family of
615 Long Interspersed Repeated DNA Occurred During the Murine Radiation”.
616 In: *Journal of Molecular Evolution* 38 (1994), pp. 18–27.
- 617 [14] J. Giordano et al. “Evolutionary history of mammalian transposons deter-
618 mined by genome-wide defragmentation”. In: *PLoS Computational Biology*
619 3.7 (2007), pp. 1321–1334.
- 620 [15] Paci Giulia et al. “Characterization of DNA methylation as a function of
621 biological complexity via dinucleotide inter-distances”. In: *Philosophical*
622 *Transactions of the Royal Society A: Mathematical, Physical and Engi-*
623 *neering Sciences* (2016). DOI: [http://doi.org/10.1098/rsta.2015.](http://doi.org/10.1098/rsta.2015.0227)
624 0227.
- 625 [16] J. S. Han and J. D. Boeke. “A highly active synthetic mammalian retro-
626 transposon”. In: *Nature* 429 (2004), pp. 314–318.
- 627 [17] S. P. Hubbell and L. B. de Água. “The unified neutral theory of biodiver-
628 sity and biogeography: reply.” In: *Ecology* 85.11 (2004), pp. 3175–3178.

- 629 [18] M. Imbeault, P. Helleboid, and D. Trono. “KRAB zinc-finger proteins
630 contribute to the evolution of gene regulatory networks.” In: *Nature* 543
631 (2017), pp. 550–554. DOI: <https://doi.org/10.1038/nature21683>.
- 632 [19] J. Jurka. “Rebase Update: A database and an electronic journal of repet-
633 itive elements”. In: *Trends in Genetics* 16.9 (2000), pp. 418–420.
- 634 [20] J. Jurka et al. “Rebase Update, a database of eukaryotic repetitive ele-
635 ments”. In: *Cytogenetic and Genome Research* 110 (2005), pp. 462–467.
- 636 [21] N. G. van Kampen. *Stochastic Processes in Physics and Chemistry*. Am-
637 sterdam: North-Holland, 1981.
- 638 [22] H. Khan, A. Smit, and S. Boissinot. “Molecular evolution and tempo of
639 amplification of human LINE-1 retrotransposons since the origin of pri-
640 mates”. In: *Genome Research* 16.1 (2006), pp. 78–87.
- 641 [23] S. Linquist et al. “Applying ecological models to communities of genetic
642 elements: The case of neutral theory”. In: *Molecular Ecology* 24.13 (2015),
643 pp. 3232–3242.
- 644 [24] S. Linquist et al. “Distinguishing ecological from evolutionary approaches
645 to transposable elements”. In: *Biological Reviews* 88.3 (2013), pp. 573–584.
646 ISSN: 14647931.
- 647 [25] Z. Lippman et al. “Role of transposable elements in heterochromatin and
648 epigenetic control”. In: *Nature* 430 (2004), pp. 471–476. DOI: <https://doi.org/10.1038/nature02651>.
- 650 [26] M. L. Mears and C. A. Hutchinson. “The evolution of modern lineages of
651 mouse L1 elements”. In: *Journal of Molecular Evolution* 52 (2001), pp. 51–
652 62.
- 653 [27] L. Milanesi et al. “Trends in modeling biomedical complex systems”. In:
654 *BMC bioinformatics* 10 (2009). DOI: [doi:10.1186/1471-2105-10-S12-](https://doi.org/10.1186/1471-2105-10-S12-I1)
655 [I1..](https://doi.org/10.1186/1471-2105-10-S12-I1)
- 656 [28] “National Center for Biotechnology Information (NCBI)”. In: *Bethesda*
657 *(MD): National Library of Medicine (US), National Center for Biotech-*
658 *nology Information* (1988). URL: [Available%20from:%20https://www.](https://www.ncbi.nlm.nih.gov/)
659 [ncbi.nlm.nih.gov/](https://www.ncbi.nlm.nih.gov/).
- 660 [29] I Ogiwara et al. “Retropositional parasitism of SINEs on LINES: identi-
661 fication of SINEs and LINES in elasmobranchs.” In: *Molecular Biology*
662 *and Evolution* 16.9 (Sept. 1999), pp. 1238–1250. ISSN: 0737-4038. DOI:
663 [10.1093/oxfordjournals.molbev.a026214](https://doi.org/10.1093/oxfordjournals.molbev.a026214). eprint: [https://academic.](https://academic.oup.com/mbe/article-pdf/16/9/1238/9593471/mbe1238.pdf)
664 [oup.com/mbe/article-pdf/16/9/1238/9593471/mbe1238.pdf](https://academic.oup.com/mbe/article-pdf/16/9/1238/9593471/mbe1238.pdf). URL:
665 <https://doi.org/10.1093/oxfordjournals.molbev.a026214>.
- 666 [30] E. Pascale, E. Valle, and A. V. Furano. “Amplification of an ancestral
667 mammalian LI family of long interspersed repeated DNA occurred just
668 before the murine radiation.” In: *Proceedings of the National Academy of*
669 *Sciences USA* 87 (1990), pp. 9481–9485.
- 670 [31] L.M. Payer and K.H. Burns. “Transposable elements in human genetic
671 disease.” In: *Nat Rev Genet* 20 (2019), pp. 760–772. DOI: [https://doi.](https://doi.org/10.1038/s41576-019-0165-8)
672 [org/10.1038/s41576-019-0165-8](https://doi.org/10.1038/s41576-019-0165-8).

- 673 [32] Michael Rosenthal et al. “Bayesian Estimation of 3D Chromosomal Structure from Single Cell Hi-C Data”. In: *bioRxiv* (2018). DOI: 10.1101/316265. eprint: <https://www.biorxiv.org/content/early/2018/05/07/316265.full.pdf>. URL: <https://www.biorxiv.org/content/early/2018/05/07/316265>.
- 674
675
676
677
- 678 [33] A. Le Rouzic, T. S. Boutin, and P. Capy. “Long-term evolution of transposable elements”. In: *Proceedings of the National Academy of Sciences USA* 104.49 (2007), pp. 19375–19380.
- 679
680
- 681 [34] A. Le Rouzic and P. Capy. “The first steps of transposable elements invasion: Parasitic strategy vs. genetic drift”. In: *Genetics* 169.2 (2005), pp. 1033–1043.
- 682
683
- 684 [35] A. Le Rouzic and P. Capy. “Population genetics models of competition between transposable element subfamilies”. In: *Genetics* 174.2 (2006), pp. 785–793.
- 685
686
- 687 [36] M. C. Seleme et al. “Extensive individual variation in L1 retrotransposition capability contributes to human genetic diversity.” In: *Proceedings of the National Academy of Sciences USA* 103 (2006), pp. 6611–6616.
- 688
689
- 690 [37] F. Serra, V. Becher, and H. Dopazo. “Neutral Theory Predicts the Relative Abundance and Diversity of Genetic Elements in a Broad Array of Eukaryotic Genomes”. In: *PLoS ONE* 8 (2013), p. 6.
- 691
692
- 693 [38] A. F. Smit et al. “Ancestral, mammalian-wide subfamilies of LINE-1 repetitive sequences.” In: *Journal of molecular biology* 246.3 (1995), pp. 401–417.
- 694
695
- 696 [39] AFA Smit, R. Hubley, and P. Green. “RepeatMasker Open-4.0”. In: <http://www.repeatmasker.org> (2013).
- 697
- 698 [40] A. Sookdeo et al. “Revisiting the evolution of mouse LINE-1 in the genomic era.” In: *Mobile DNA* 4.1 (2013).
- 699
- 700 [41] C.J. Struchiner, M.G. Kidwell, and J.M.C. Ribeiro. “Population dynamics of transposable elements: copy number regulation and species invasion requirements”. In: *Journal of Biological Systems* 13.4 (2005), pp. 455–475.
- 701
702
- 703 [42] W. Sun et al. “Pathogenic tau-induced piRNA depletion promotes neuronal death through transposable element dysregulation in neurodegenerative tauopathies.” In: *Nat Neurosci* 21 (2018), pp. 1038–1048. DOI: <https://doi.org/10.1038/s41593-018-0194-1>.
- 704
705
706
- 707 [43] S. Venner, C. Feschotte, and C. Biéumont. “Dynamics of transposable elements: towards a community ecology of the genome”. In: *Trends in Genetics* 25.7 (2009), pp. 317–323.
- 708
709
- 710 [44] S. Vitali. “Modeling of Birth-Death and Diffusion Processes in Biological Complex Environments.” In: *Ph.D. Thesis, Department of Physics and Astronomy, University of Bologna, Bologna, Italy* (2018).
- 711
712
- 713 [45] I. Volkov et al. “Neutral theory and relative species abundance in ecology.” In: *Nature* 424.6952 (2003), pp. 1035–1037.
- 714
- 715 [46] F. Yue et al. “A comparative encyclopedia of DNA elements in the mouse genome.” In: *Nature* 515.7527 (2014), pp. 355–64.
- 716

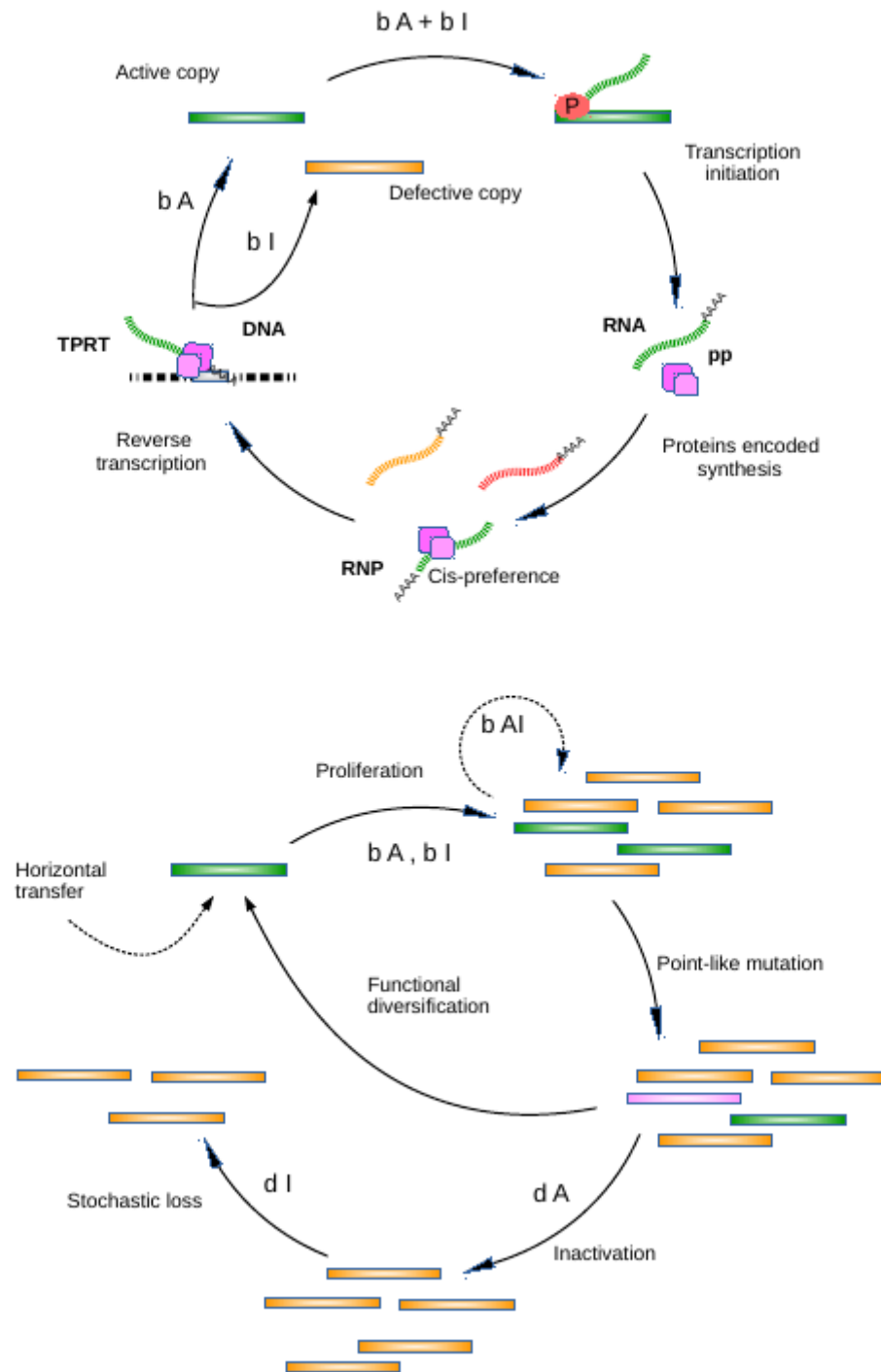


Figure 1: **LINEs transposition and activity cycle diagram.** Diagram of the birth-death process of defective (orange) and full-length copies (green) from a full-length master copy by retro-transposition and cycle of element amplification in the genome. Occasional point-like mutation (or other) generate a new element species (pink) which start the cycle again and eventually become a competitor.

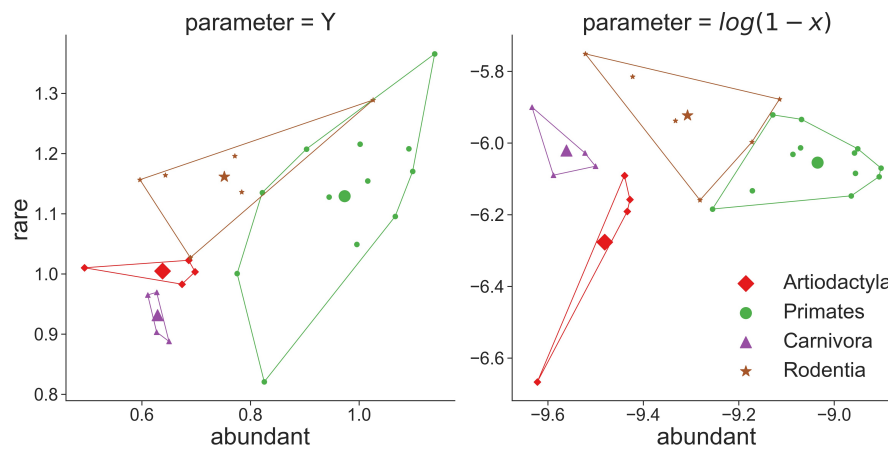


Figure 2: Competition model description clusters different mammalian orders. Set of optimized parameters (posteriors expectation value) obtained by ABC method for the competition model (mixture of two negative binomials). The couples Υ_{rare} , Υ_{abund} and $\log(1 - x_{rare})$, $\log(1 - x_{abund})$ are shown for different hosts. Convex hull of same host order parameters are highlighted by straight lines. The average of the parameters per host order are shown by the larger markers.

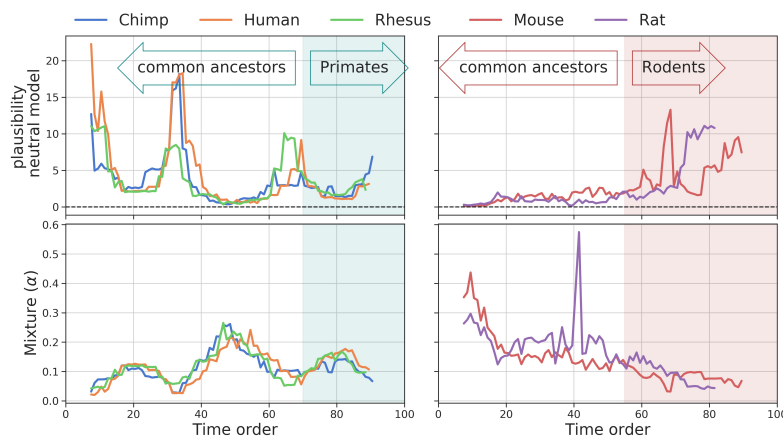


Figure 3: Comparison of the two model performance during the evolution of the LINEs ecosystem in primates and rodents. Upper panels: plausibility of the neutral model calculated in term of log-likelihood ratio test of the neutral model over competition model) by the sliding window approach. Lower panels: estimation of the mixture coefficient by the sliding window approach. Larger values for the mixture coefficient α are associated to lower plausibility of the neutral model. The portion of elements highlighted in different colors belong to different clusters in figure 6.

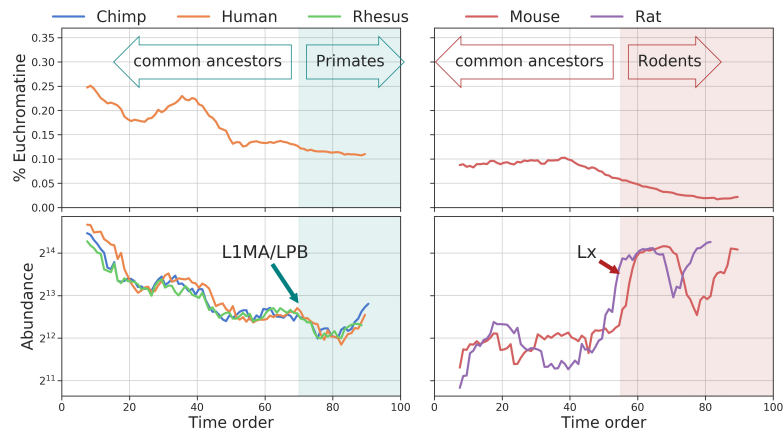


Figure 4: **Sliding window of LINE insertions percentage inside euchromatin and expected LINEs abundance in primates and rodents.** Upper panels: percentage of LINE copies inserted in euchromatin regions calculated by the sliding window approach. Lower panels: LINEs average copy number (sum of euchromatin and heterochromatin insertions) calculated by the sliding window approach. The portion of elements highlighted in different colors belong to different clusters in figure 6.

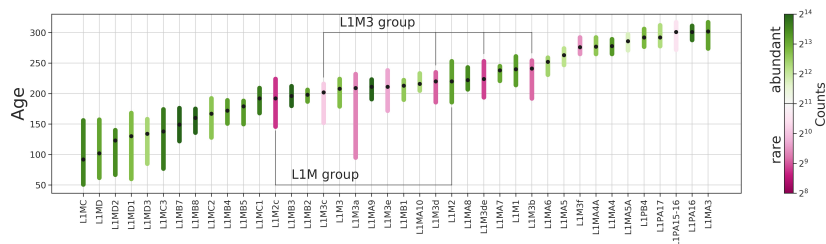


Figure 5: **5'UTR similarity between competing LINE retrotransposons in human.** The available consensus sequences of the 5'UTR of LINEs in the human genome have been aligned pairwise, with ClustalW2. The minimum distance is achieved between couples (or groups) of elements with similar ages and having high and low copy number respectively. Range of activity shown by the bar length. Abundance shown by the color legend. Significant similarity between 5'UTRs is observed for the following high and low copy numbers groups: L1M2-L1M2c, L1M3a-L1M3b-L1M3c-L1M3d.

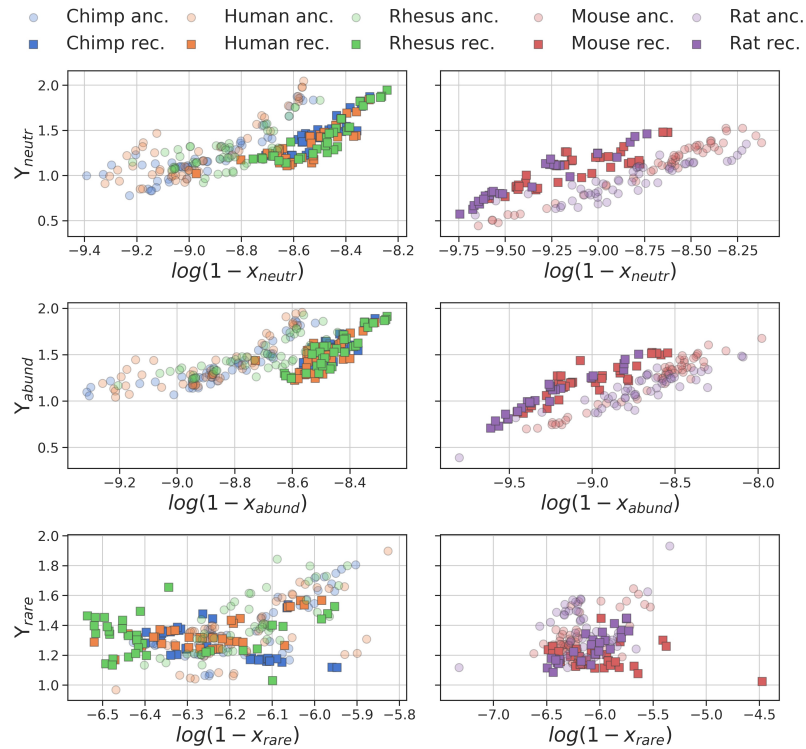


Figure 6: **Space of parameters of the two model tested by the sliding window approach in primates and rodents.** The space of parameters describing sliding window ecosystem of LINEs in human, chimpanzee, rhesus macaque (left panels) and mouse, rat (right panels) is shown. x and Y parameters are correlated by the expected value (mean) of the distribution. Upper panels refer to the neutral model, middle panels refer to the group of the mixture model with largest copy number, lower panels refer to the group of less abundant elements. Circles indicated the most ancient elements. Transition between the two cluster are associated with specific LINE species appearance.

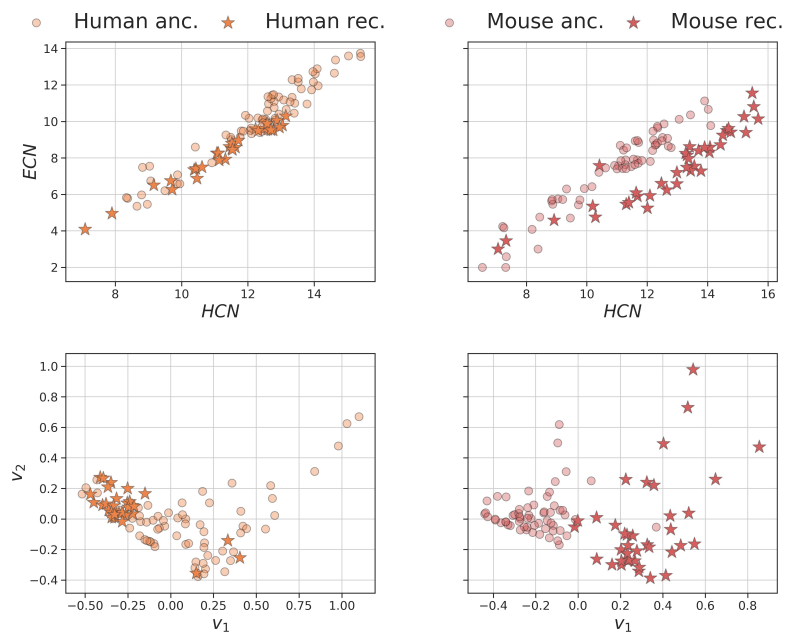
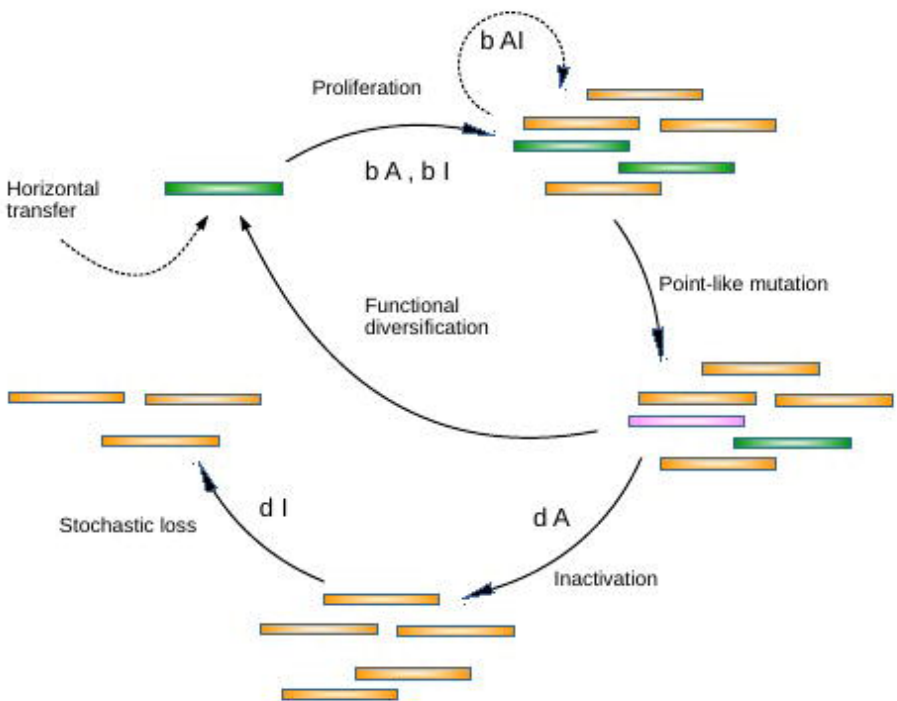
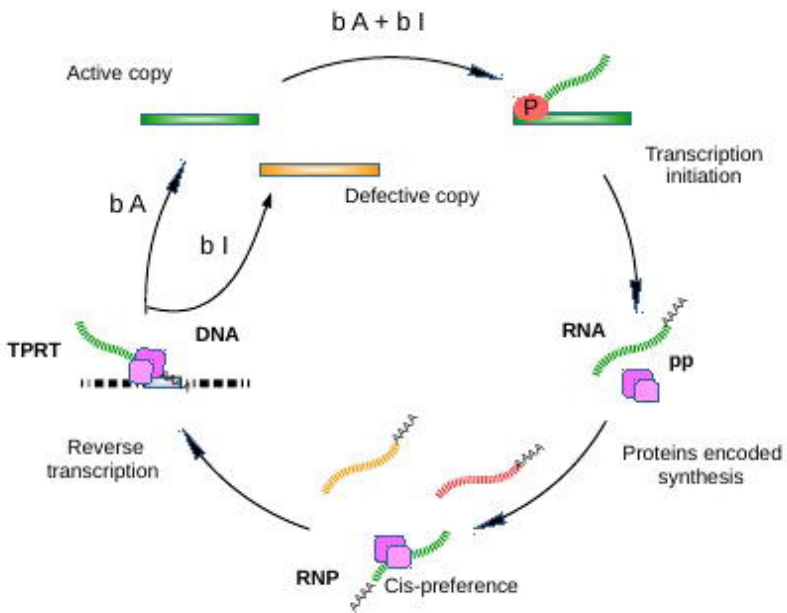
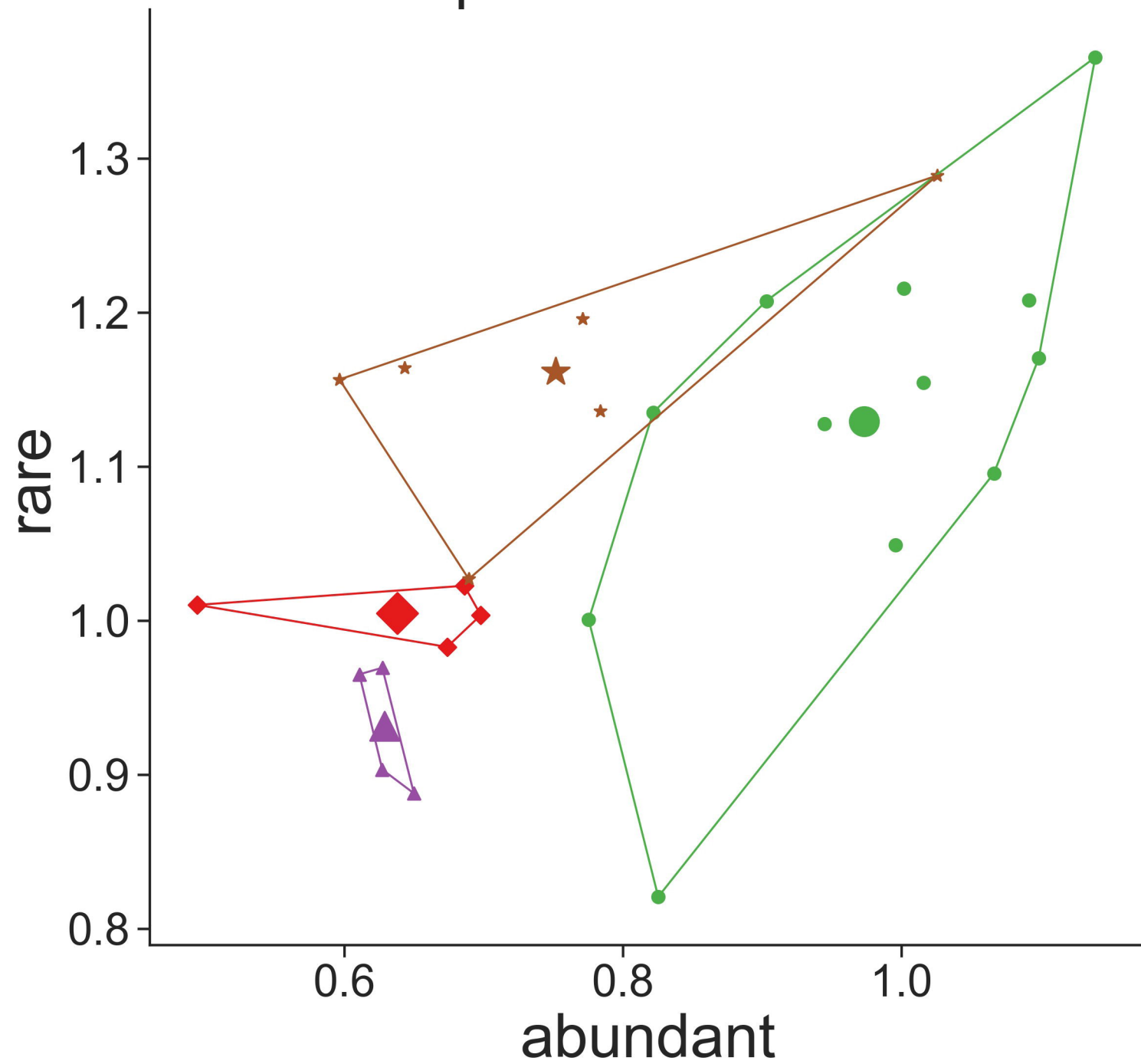
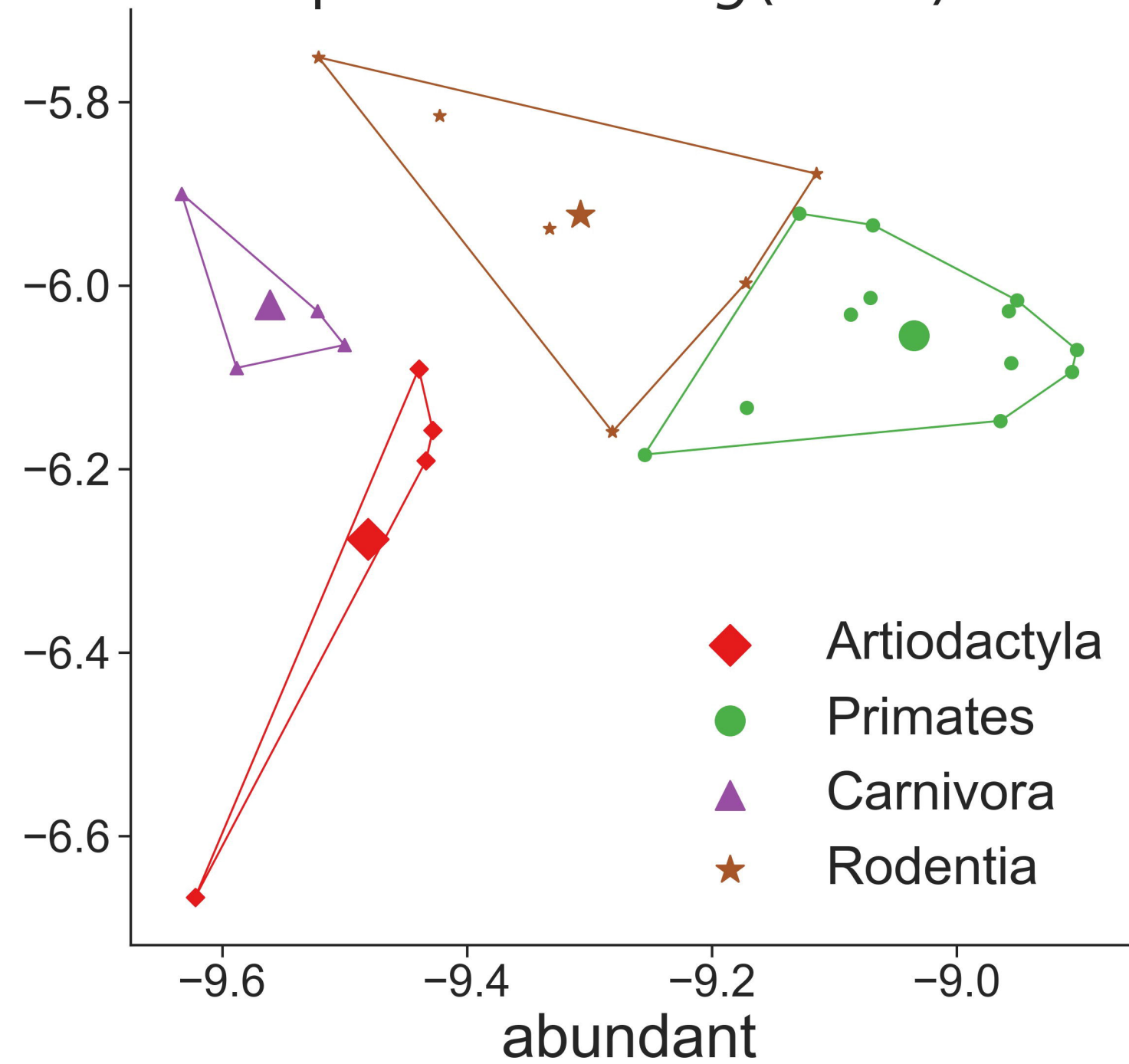
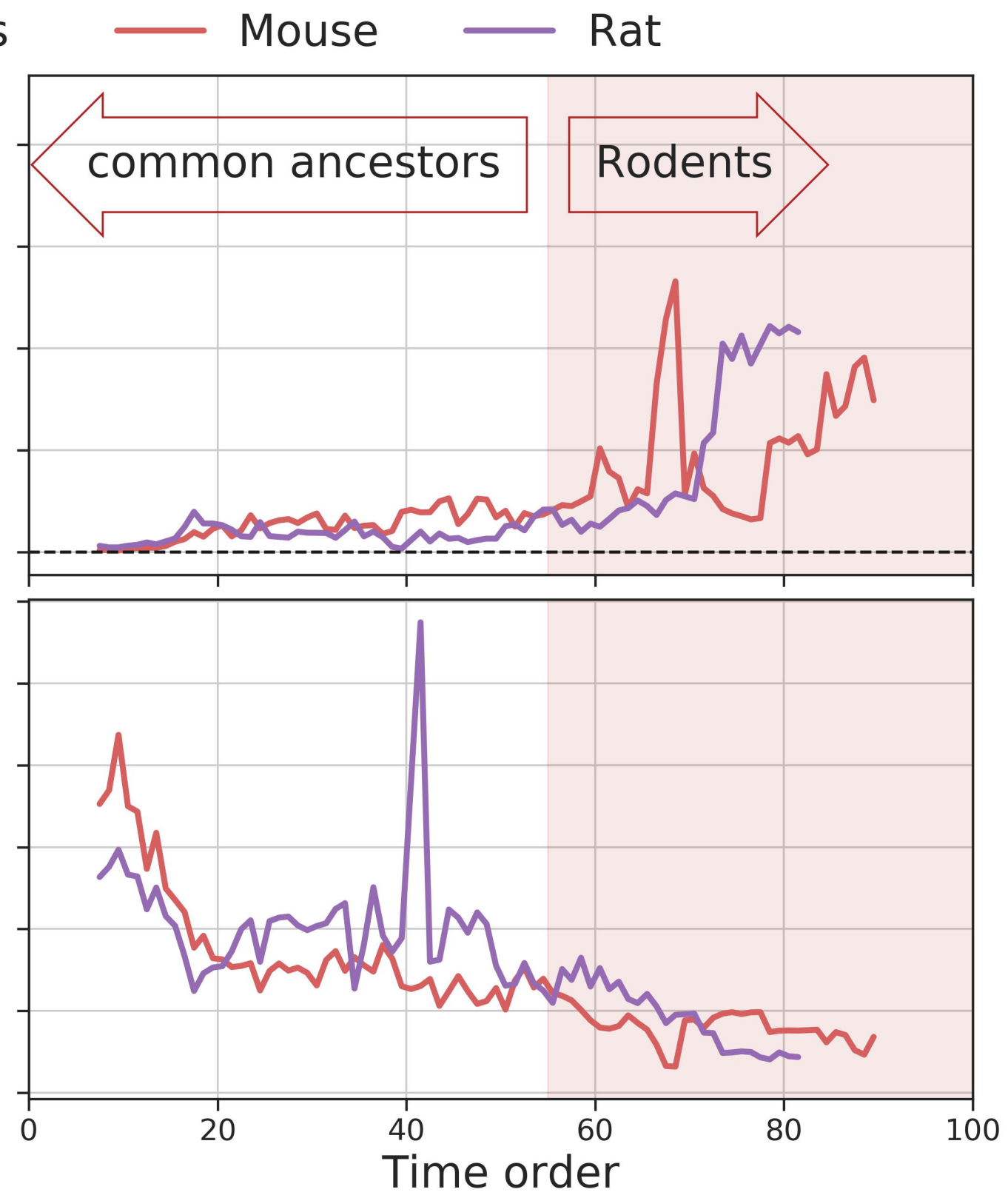
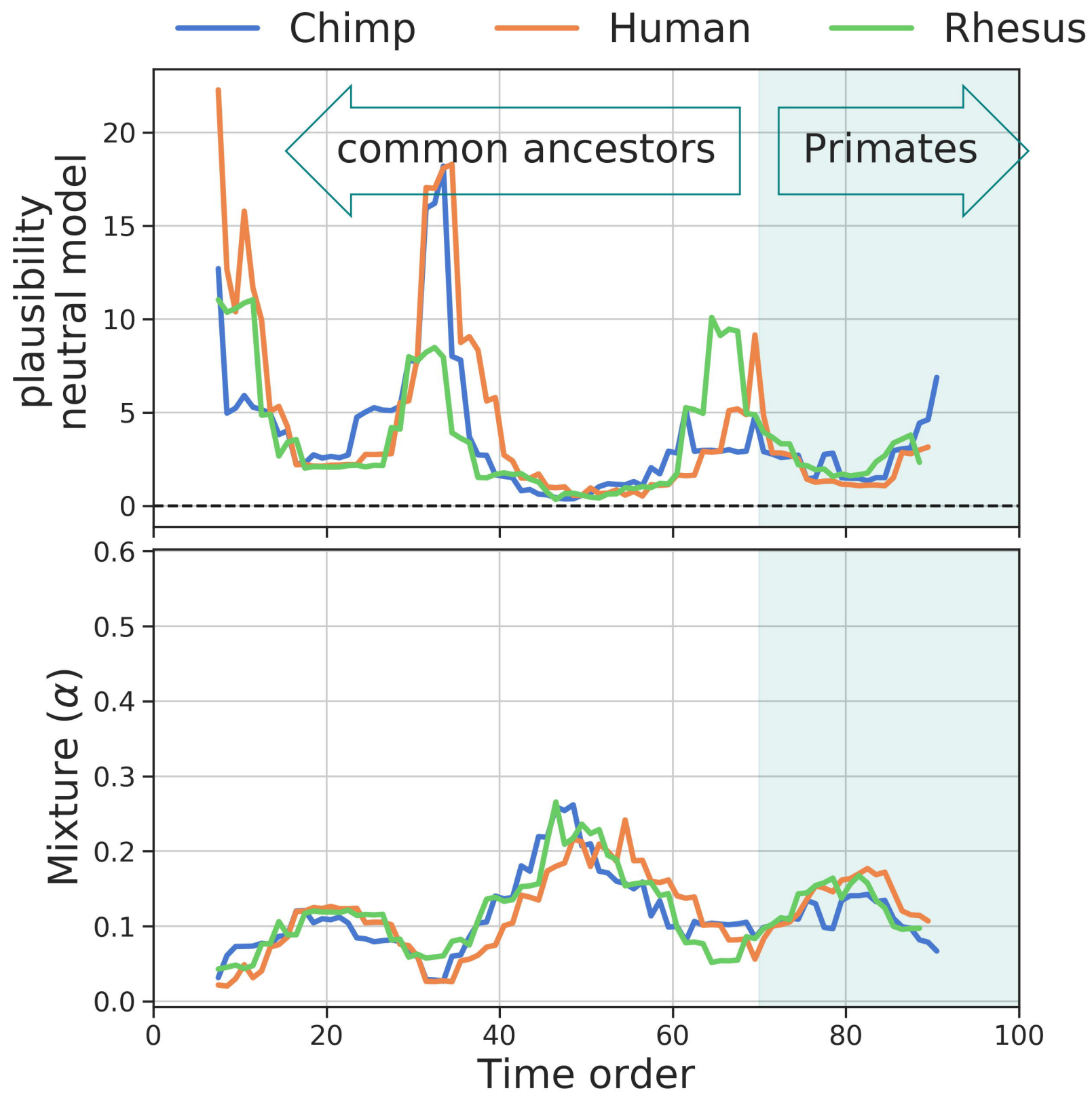
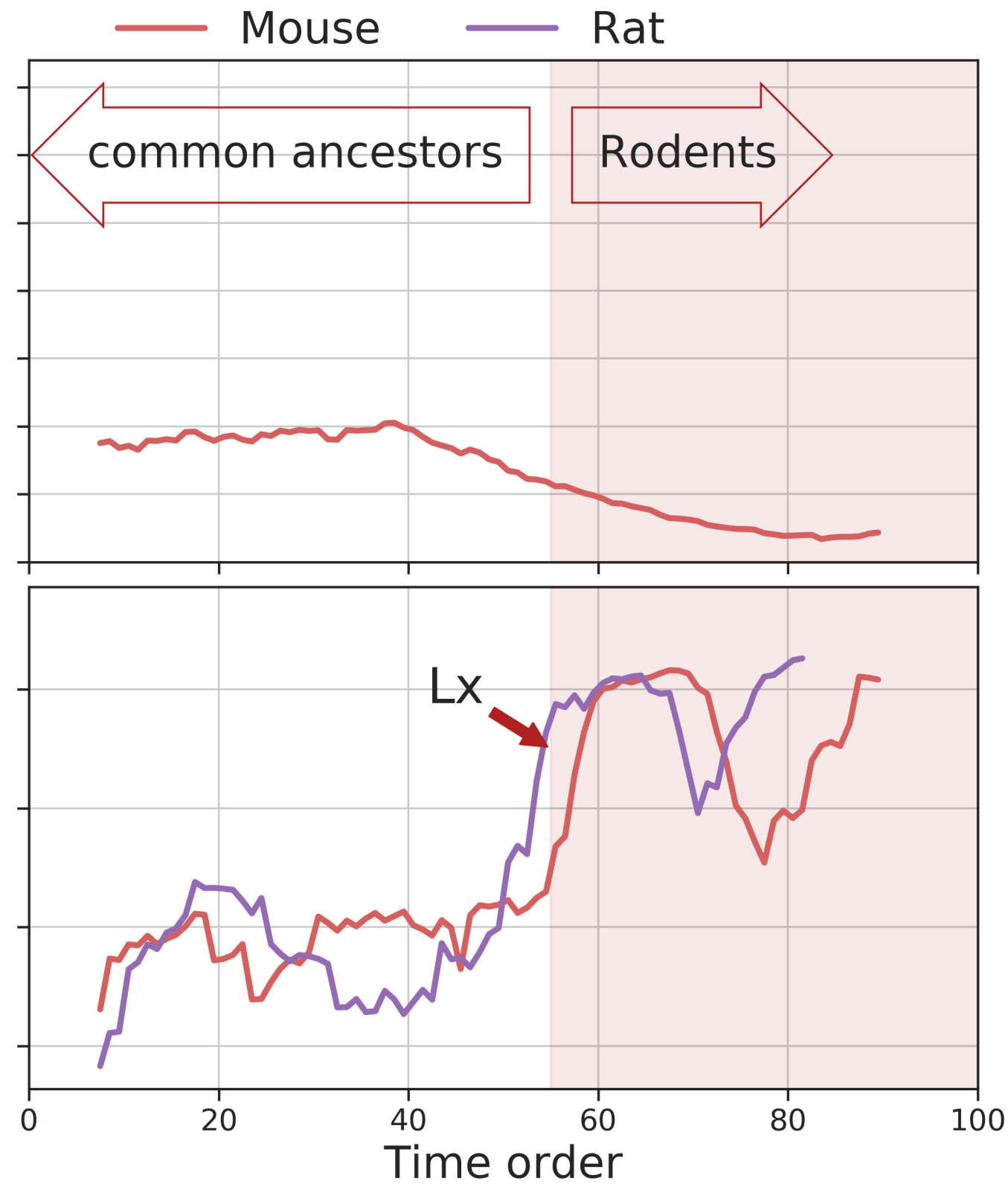
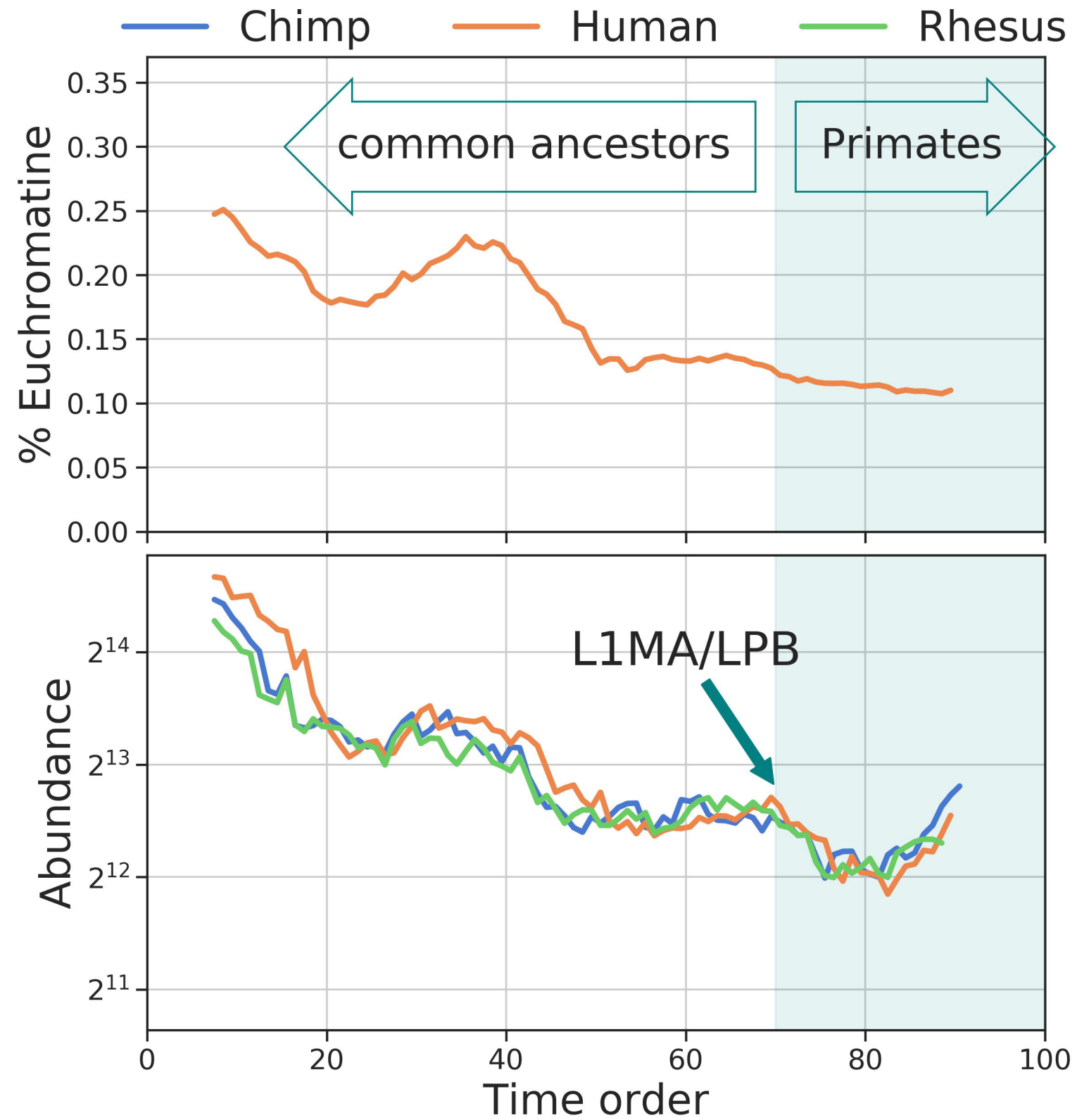


Figure 7: **Number of insertions in euchromatin and heterochromatin states in human and mouse.** Upper panels: Scatter plot in \log_2 scale of the number of insertions in euchromatin respect that in heterochromatin for each LINE specie. The number of insertions in euchromatin and heterochromatin states results correlated by a power law. Lower panels: PCA of the number of insertions in euchromatin, in heterochromatin, and relative time of appearance. The portion of elements highlighted in different colors belong to the two different clusters in figure 6: ancient elements (circles), recent elements (stars).

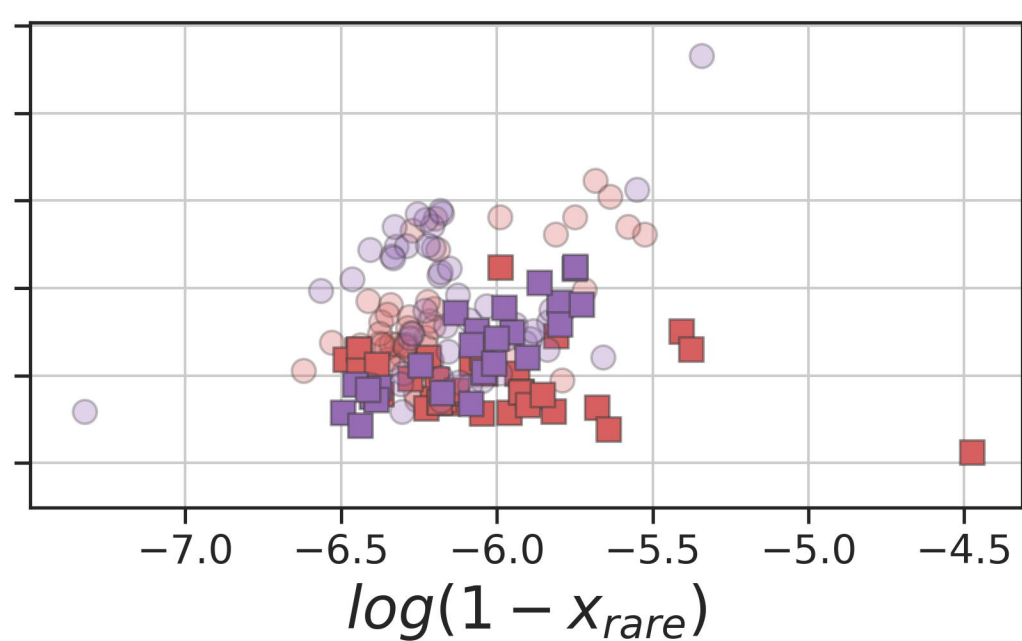
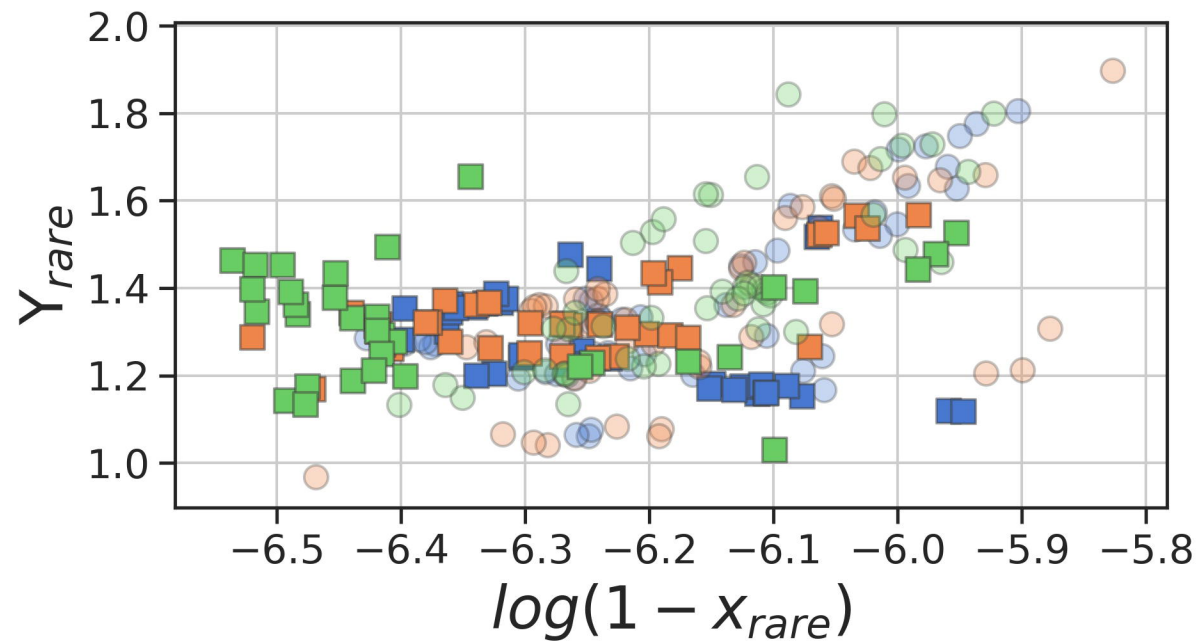
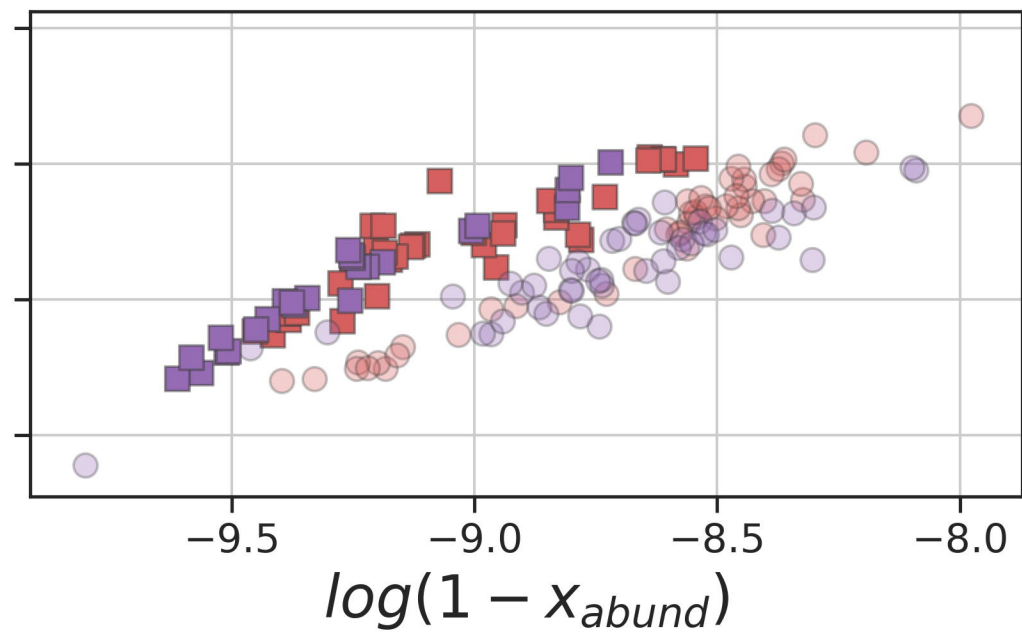
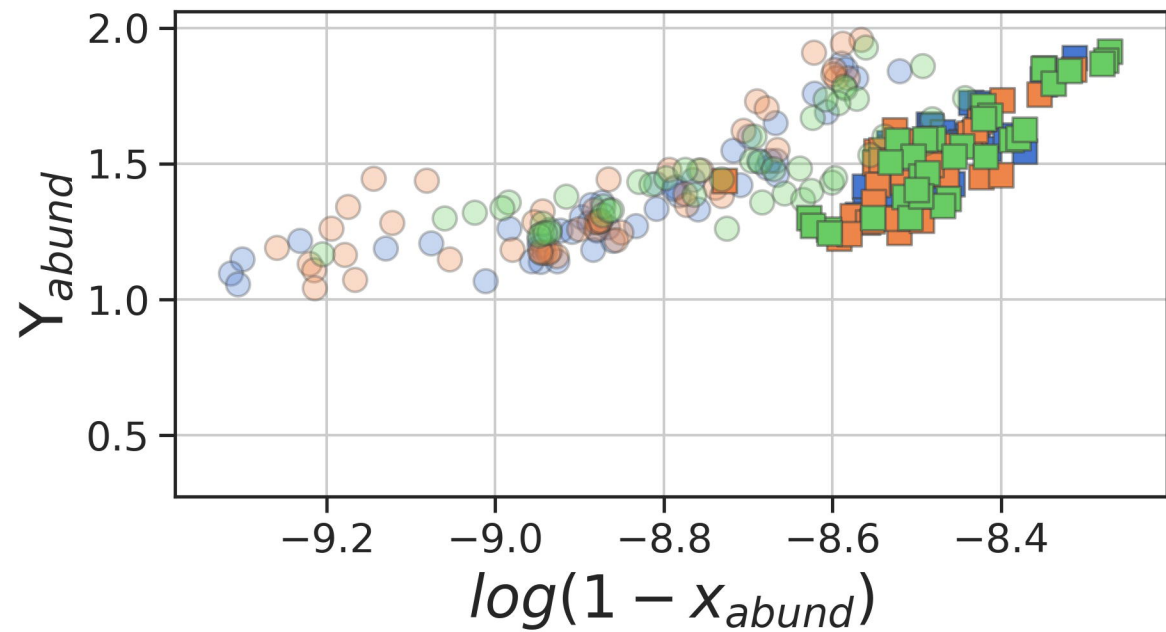
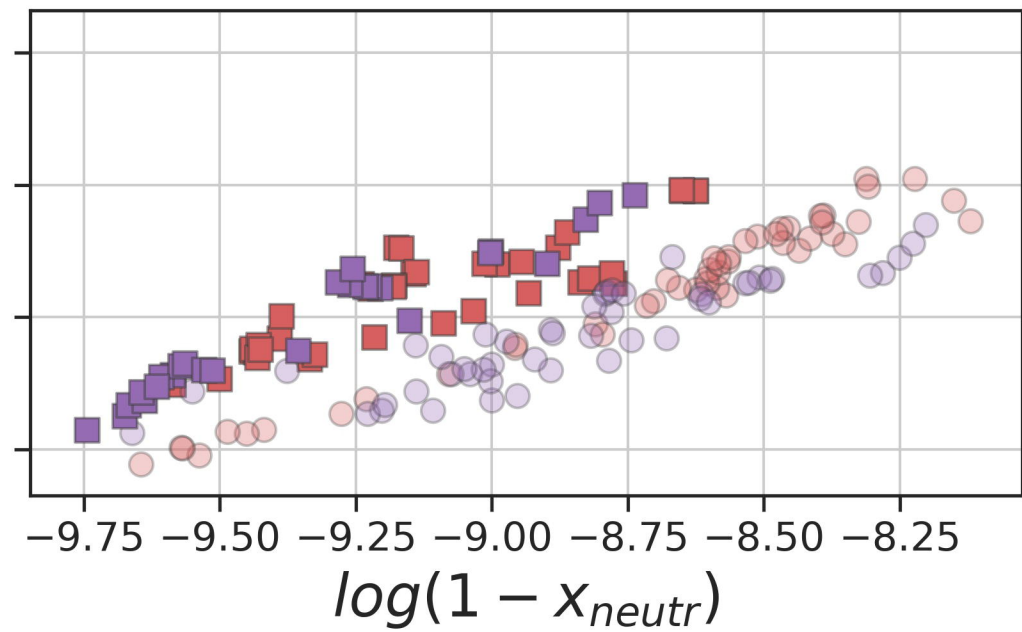
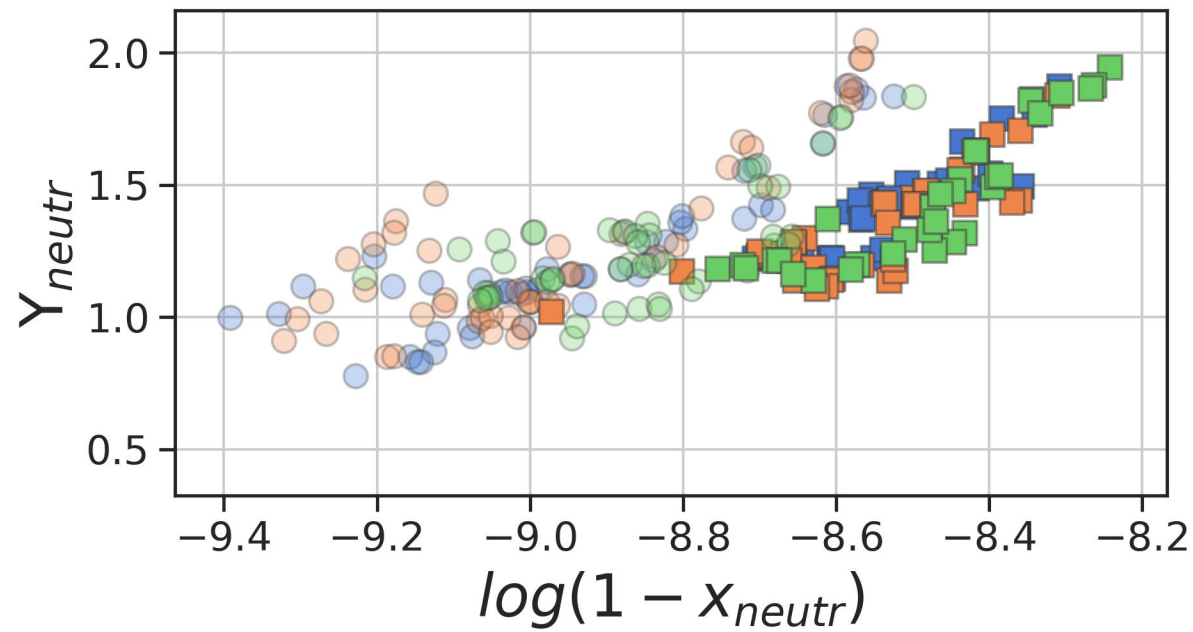


parameter = Y parameter = $\log(1 - x)$ 

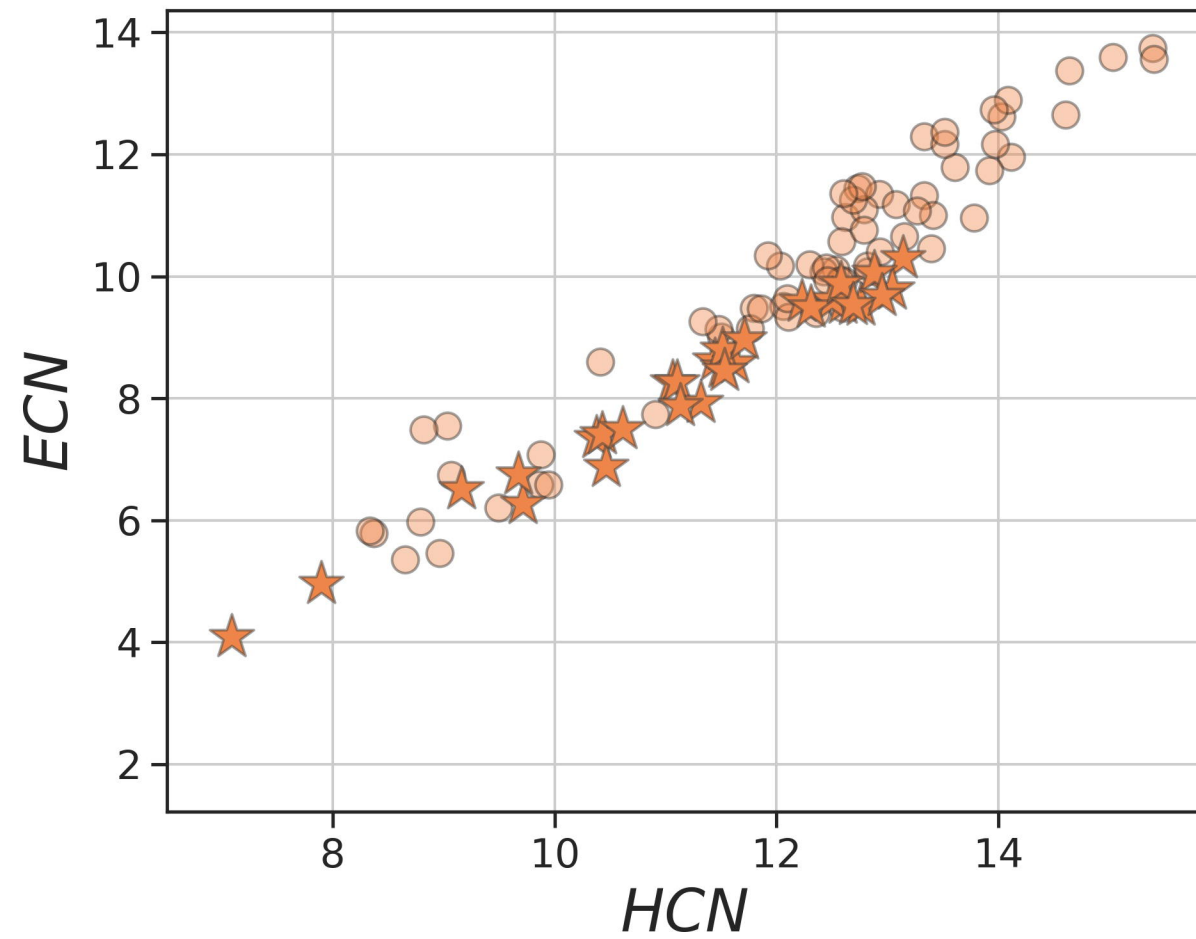




Chimp anc. Human anc. Rhesus anc. Mouse anc. Rat anc.
Chimp rec. Human rec. Rhesus rec. Mouse rec. Rat rec.



● Human anc. ★ Human rec.



● Mouse anc. ★ Mouse rec.

