

# 1 **Inferring linguistic transmission between** 2 **generations at the scale of individuals**

3 Valentin Thouzeau<sup>†\*</sup>, Antonin Affholder<sup>†</sup>, Philippe Menecier<sup>†</sup>, Paul Verdu<sup>†</sup>, Frédéric Austerlitz<sup>†</sup>

4 <sup>†</sup> *CNRS, MNHN, Université de Paris, UMR 7206 Eco-Anthropologie, Paris 75016, France*

5 <sup>\*</sup> *Laboratoire de Neurosciences Cognitives, Département d'études cognitives, ENS, PSL, Research*  
6 *University, Paris, France*

## 7 **Abstract**

8 Historical linguistics strongly benefited from recent methodological advances inspired by  
9 phylogenetics. Nevertheless, no available method uses contemporaneous within-population  
10 linguistic diversity to reconstruct the history of human populations. Here, we developed an  
11 approach inspired from population genetics to perform historical linguistic inferences from  
12 linguistic data sampled at the individual scale, within a population. We built four within-population  
13 demographic models of linguistic transmission over generations, each differing by the number of  
14 teachers involved during the language acquisition and the relative roles of the teachers. We then  
15 compared the simulated data obtained with these models with real contemporaneous linguistic data  
16 sampled from Tajik speakers from Central Asia, an area known for its large within-population  
17 linguistic diversity, using approximate Bayesian computation methods. Under this statistical  
18 framework, we were able to select the models that best explained the data, and infer the best-fitting  
19 parameters under the selected models. This demonstrates the feasibility of using contemporaneous  
20 within-population linguistic diversity to infer historical features of human cultural evolution.

## 21 **1. Introduction**

22 Several recent studies use linguistic data under a computational framework aiming at  
23 reconstructing various aspects of the cultural history of human populations (Atkinson, 2011;  
24 Bouckaert et al., 2012; Gray and Atkinson, 2002; Pagel et al., 2013; Thouzeau et al, 2017). They  
25 rely on data mainly consisting in a set of presence or absence of linguistic items, within a given set  
26 of contemporaneous languages, which can be found, for example, in databases such as the World  
27 Atlas of Language Structures WALS (Dryer and Haspelmath, 2013), or the Global Database of  
28 Cultural, Linguistic and Environmental Diversity D-PLACE (Kirby et al., 2016). Thus, most studies  
29 consider languages at a macro-evolutionary scale, i.e. they deal only with differences among  
30 languages, neglecting the variability within each language. For instance, Gray and Atkinson (2002)  
31 used a set of Swadesh lists obtained for 87 languages to investigate the origin of the Indo-European  
32 linguistic family. Atkinson (2011) considered the number of phonemes used in 504 languages  
33 worldwide to test the hypothesis of a serial founder effect due to the Out-Of-Africa expansion.  
34 Reesink et al. (2009) used the linguistic diversity of the ancient Sahul continent (present-day  
35 Australia, New Guinea, and surrounding islands) among 121 languages to infer the history of the  
36 structural characteristics of these languages.

37 These approaches rely implicitly on several assumptions. They require primarily a clear  
38 separation between several differentiated languages. Nevertheless, this notion of distinct languages  
39 is often irrelevant at local scale, in particular in contexts of dialectal continuum or linguistic  
40 contacts (Heeringa and Nerbonne, 2001; Livingstone and Fyfe, 1999). Furthermore, most of these  
41 studies do not take into account within-population linguistic diversity, since traditional linguistics  
42 often considers languages as unique and coherent systems (Pateman, 1983). This assumption  
43 implies the loss of a large amount of information, knowing that the demographic phenomena  
44 occurring at population level – different population sizes, bottlenecks, expansions – are expected to  
45 play a major role in language evolution (Vogt, 2009). The inclusion of contemporaneous within-  
46 population linguistic diversity in the reconstruction of the demographic history of human

47 populations at a local scale is thus expected to open a completely new dimension in the field of  
48 historical linguistic inferences.

49 In this context, Croft (1996) argued for the replacement of the ‘essentialist’ theory of language  
50 changes by a ‘population’ approach of these changes, and later proposed a detailed review of the  
51 “evolutionary linguistic” field and underlying paradigms (Croft, 2008). Nevertheless, very few  
52 previous studies dealt with the contemporaneous within-population linguistic diversity in a  
53 historical reconstruction perspective. Rodriguez-Larralde and Barraı (2000) used surnames of  
54 telephone users in Austria as linguistic contemporaneous information, showing that Austrian towns  
55 are subdivided into five main clusters with uniform levels of endogamy. Darlu et al. (2012)  
56 reviewed the analyses of paternally-inherited family names distributions, as an analogy to Y  
57 chromosomes, for historical inferences. Verdu et al. (2017) contrasted the proportion of African  
58 words in free speech among Cape Verdean Kriolu speakers with their proportion of African genetic  
59 admixture, showing that Cape Verdean genetic and linguistic admixture processes followed parallel  
60 histories, with possible co-transmission of genetic and linguistic variation. None of these studies  
61 developed an inferential approach that would allow researchers to distinguish among different  
62 historical mechanistic models, and inferring their constitutive parameters.

63 In order to perform such historical linguistic inferences from observed linguistic data, we need to  
64 assume one or several possible models of linguistic transmission between generations, and a  
65 possible set of historical scenarios that produced the observed data. Nevertheless, there is no  
66 consensus framework that allows handling within-population linguistic diversity data, in order to  
67 infer historical scenarios and evolutionary mechanisms. It requires to first build an explicit  
68 mechanism of linguistic evolution, and then to study the range of historical scenarios that could  
69 have produced the observed linguistic data. Nevertheless, the validity of the historical conclusions  
70 will depend on the validity of the assumed mechanism. It is, therefore, crucial to first determine the  
71 most relevant mechanism of linguistic evolution of a given set of linguistic objects, in order to  
72 produce, ultimately, valid inferences.

73 We evaluated here a series of models of linguistic evolution between generations at the individual  
74 scale. We did not study the history of higher-order objects such as “the languages”, but the history  
75 of the linguistic diversity carried by individuals within a population, among which communication

76 events may occur over time. We aimed at understanding how the evolution of linguistic diversity  
77 among generations was affected by demographic parameters such as population sizes (the number  
78 of individuals of a given speech community), and thus to assess whether it was possible to infer the  
79 best demographic scenario and its corresponding parameters from a set of linguistic data.

80 Approximate Bayesian Computation methods (ABC, Beaumont et al., 2002; Tavaré et al., 1997)  
81 provide a particularly well-adapted framework to tackle this problem. In this paper, we used the  
82 recently developed Approximate Bayesian Computation via Random Forest (ABCRF) algorithm to  
83 assess, among a set of possible competing scenarios, the scenario that best explained the observed  
84 data, and to estimate the posterior parameters of this scenario (Breiman, 1999; Pudlo et al., 2016,  
85 Raynal et al., 2017).

86 For this purpose, we implemented an individual-based simulation program, which simulates the  
87 evolution of word variation among generations, under different modes of linguistic transmission.  
88 These simulated data allowed us to perform the ABCRF procedure on a real dataset from Central  
89 Asia. This dataset consisted of 30 individuals interviewed for 185 words across 10 villages in  
90 Tajikistan. These villages are known to use the same language, but with some variation across  
91 speakers (Menecier et al., 2016). We aimed at inferring the most probable models of linguistic  
92 transmission mechanisms between linguistic generations, under a demographic scenario of  
93 population-size expansion or contraction.

94 We proposed four transmission models. The “*Clonal* model” assumed that each individual learns  
95 his/her linguistic words from only one teacher. The “*Sexual model 1*” assumed that each individual  
96 learns his/her words from two teachers (one “male” and one “female”), with specific words  
97 transmitted only by males and others transmitted only by females. The “*Sexual model 2*” assumed  
98 that each individual learned his/her words from two parents (one “male” and one “female”), without  
99 specific words belonging to males or females. The transmitting parent was drawn independently for  
100 each word, so that the set of words of an individual was a recombination between the sets of his/her  
101 two parents. Finally, the “*Social* model” assumed that each individual learns his/her words from the  
102 entire population. We aimed, then, at inferring, with ABC, the best-fitting parameters under the  
103 winning scenario: linguistic mutation rates, and population sizes. We demonstrated thus the

104 feasibility of using contemporaneous within-population linguistic diversity to infer historical  
105 features in human cultural evolution.

## 106 **2. Materials**

107 We sampled cognate variation for 30 individuals from 10 Tajik villages in Central Asia (Figure 1)  
108 assuming that the individuals belonged to a single linguistic population. In contrast with our  
109 previous study, where we considered for each cognate only its most frequent variant in each locality  
110 (Thouzeau et al., 2017), we kept here the linguistic variants recorded for each individual separately.  
111 Thus, for each individual, we recorded the words used for 185 meanings from an adapted Swadesh  
112 list. Individuals were asked to state the most frequently used word for the associated meaning. We  
113 considered as “cognate” a group of words with the same etymological origin and the same meaning,  
114 such words being more likely to be related by a common ancestry (Atkinson et al., 2005). The  
115 phonetic transcriptions and classification of lexical data gathered in the field into cognates was  
116 performed by Philippe Menecier similarly as in previous work (Menecier et al., 2016; Thouzeau  
117 et al., 2017).

## 118 **3. Principle of ABCRF method**

119 ABC methods were first introduced by Tavaré et. al. (1997) and Beaumont et al (2002), aiming to  
120 encompass the limitation of Markov chains Monte Carlo (MCMC) methods. For simple models,  
121 analytical formulas may be derived to compute the likelihood of the data under a given model.  
122 However, for complex models or for large datasets, computing the likelihood may be highly  
123 difficult and/or highly time consuming. ABC methods allow circumventing these problems by  
124 approaching the likelihood instead of exactly computing its value. It is a particularly well-suited  
125 statistical framework to develop within-population linguistic historical inferences tools, allowing to  
126 specify complex and explicit processes of linguistic interactions between a large set of agents.

127 ABC consists first in defining a set of models that could fit the observed data. Each model is  
128 characterized by several parameters, such as, but not limited to, effective population sizes and time  
129 of change in population sizes. A prior distribution for each parameter is chosen by the user, and  
130 corresponds to the range of values that are realistic for this parameter. A large number of  
131 simulations are then performed in order to generate data sets under the different models. For each  
132 simulation, each parameter of the model is drawn at random in its prior distribution. Sets of  
133 summary statistics are then computed on each simulated data set, each corresponding, therefore, to  
134 a given set of model parameters. In the ABCRF method (Pudlo et al. 2016), a random forest (RF)  
135 procedure is then applied to choose the best-fitting model. In short, the aim of the RF method is to  
136 produce a set of decision trees from the simulated data sets. Each tree is built by performing a  
137 supervised categorization of the whole set of simulations, according to the models which produced  
138 those simulations, and each one using a different subset of their summary statistics. Those subsets  
139 of summary statistics are selected randomly for each tree in order to improve classifications,  
140 because using the full set of summary statistics can lead to overfitting. The curse of dimensionality  
141 is also reduced by the RF procedure (Pudlo et al. 2016).

142 Then, the full set of summary statistics are computed on the real data, and this “observed” set of  
143 summary statistics is independently evaluated by each decision tree. Each one votes for a model,  
144 and the final decision is the majority of votes from the forest. Then, an error rate is computed to  
145 assess the confidence of this final decision. At this step, several models are usually rejected by the  
146 random forest.

147 For each selected model, another random forest is then constructed to estimate the parameters of  
148 the models. The principles of ABCRF regression are analogous to the principles of ABCRF  
149 classification (Raynal et al., 2017), but in this case, the trees use the summary statistics in order to  
150 predict the value of scalar transformed parameters. The forest is then built on the simulated  
151 summary statistics, in order to estimate the mean, median, and quantiles of the distributions of the  
152 real parameters.

## 153 **4. Models**

#### 154 **4.1. Production of utterances**

155 We considered a linguistic population as a group of individuals that may potentially interact  
156 through linguistic communication. The mechanisms of linguistic communication and transmission  
157 may follow different modalities, which correspond to different models of linguistic evolution. In all  
158 cases, we considered that the unit of linguistic communication is the *utterance*, a production of  
159 words associated with a meaning (Croft, 1996).

160 We developed a general model of word transmission, which we applied in particular to the case of  
161 cognates, which correspond to words with different etymological origins that express the same  
162 meaning. For example, the Spanish word “Flor” and French word “Fleur” are two words with the  
163 same meaning (“Flower” in English) and the same etymological origin, and classified as the same  
164 cognate. The Spanish word “Mariposa” and French word “Papillon” are two words with the same  
165 meaning “Butterfly”, but with a different etymological origin. They are thus considered as different  
166 cognates. We considered here that cognates can vary among individuals within a population. This  
167 differs from the assumptions made in previous studies (Bouckaert et al., 2012; Gray et al., 2009;  
168 Thouzeau et al., 2017) where cognates are sampled at the language scale and for which individuals  
169 are considered as users rather than producers of the language.

#### 170 **4.2. Four models of acquisition of a new language**

171 We developed a new Open Source C++ simulation software *PopLingSim 2 (PLS2)*, available at  
172 <https://github.com/ValentinThouzeau/PLS2>. This software implements an individual-based  
173 forward-in-time simulation model with discrete linguistic generations, in which we assumed that  
174 populations were composed of only two types of individuals: “learners” and “teachers”. The  
175 linguistic generation time corresponded to the time required for an individual between learning a  
176 language from teachers and teaching this language to learners at the following generation. We did  
177 not specify the linguistic generation time in our models, allowing it to be completely decorrelated  
178 from the reproductive generation time, and possibly much smaller. We assumed a neutral model in  
179 the sense that, even if the number of teachers per learner varied across models, the learners selected  
180 their teachers at random in the previous generation, with equal probabilities. We assumed that the

181 rules of utterance productions of a teacher depended only on the utterances that he/she heard when  
182 he/she was a learner. We assumed that each learner chose only one word for each meaning during  
183 the learning phase. Two learners could choose the same word. After the whole learning phase, all  
184 teachers were discarded and all learners became teachers. Then, at the following generation, new  
185 learners appeared. The proportions of males and females were exactly 50%/50% at each generation.  
186 The models of linguistic acquisition differed by the number of teachers involved for each learner in  
187 his/her language acquisition process, and the relative roles of these teachers.

#### 188 **4.2.1. Clonal Model**

189 In the first model, named the “*Clonal*” model, each learner had only one teacher, which was  
190 drawn at random in the teacher population. The learner copied “in a clonal way” every word that the  
191 teacher produced. This would correspond, in genetics, to a clonal reproduction model, as observed  
192 e.g. for bacteria or for mitochondrial DNA and non-recombining regions of the Y chromosome in  
193 humans and other mammals, for instance.

#### 194 **4.2.2. Sexual Model 1**

195 In the second model, named the “*Sexual 1*” model, two different teachers (one “male” and one  
196 “female”) were drawn at random within the population for each learner respectively. The learner  
197 then copied directly the first half of the words produced by teacher 1, and the second half of the  
198 words produced by teacher 2. Thus, half of the words were always transmitted by one teacher, and  
199 the other half by the other teacher, the two different sets being always the same for all generations.

#### 200 **4.2.3. Sexual Model 2**

201 In the third model, named the “*Sexual 2*” model, two different teachers (one “male” and one  
202 “female”) were attributed to each learner at random. For each word, the learner copied at random  
203 either the word from teacher 1 or teacher 2, with equal probabilities ( $\frac{1}{2}$ ,  $\frac{1}{2}$ ). Thus, no particular  
204 word had a teacher-specific transmission; every word was transmitted from one of the two teachers  
205 chosen at random. This is analogous, in genetics, to a sexual reproduction model with free  
206 recombination.

#### 207 **4.2.4. Social Model**

208 In the fourth model, named the “*Social*” model, for each meaning each learner copied a word  
209 drawn at random from all the words produced by all the teachers in the population. Thus, each



210 learner learned his/her set of words randomly from the entire speech community, or rather, from all  
211 possible utterance variants of teachers for a given meaning at a given generation.

### 212 **4.3. Mutation model**

213 For each model, we assumed that errors could occur during the transmission of each word,  
214 leading to the creation of a completely new word. We denoted such errors “linguistic mutations”.  
215 The mean mutation rate per linguistic generation  $\mu_L$  was drawn in a log-uniform prior distribution,  
216 between  $10^{-6}$  and  $10^{-1}$  mutations per word per generation. For each word, its mutation rate was  
217 subsequently drawn in a beta distribution with mean  $\mu_L$  and shape parameter  $\beta = 2$ , allowing us to  
218 simulate a set of words each with different rates of change over time.

### 219 **4.4. Historical scenario**

220 We focused here on a single linguistic population, defined as the number of individuals that  
221 contributed significantly to the currently observed linguistic diversity, where the utterances of a  
222 sample of individuals have been obtained using a linguistic questionnaire in the final generation.  
223 This linguistic population evolved under a historical scenario (Figure 3), in which it was first of  
224 constant size  $N_0$  individuals involved in linguistic communications during  $t_0 = 5 \times N_0$  generations. As  
225 we visually checked, this time was sufficient to reach an equilibrium between the production of  
226 linguistic diversity through mutation, and the reduction of this diversity through random sampling  
227 (i.e. linguistic drift). Then, this population underwent an instantaneous change of population size,  
228 reaching a new size  $N_1$ , and the population remained at this size during  $t_1$  generations. This model  
229 allowed simulating a range of histories, depending on the relative values of the parameters  $N_0$  and  
230  $N_1$  and on the value of  $t_1$ .  $N_0$  and  $N_1$  were drawn in a uniform prior distribution bounded between  
231 100 and 1000 individuals, this upper bound being set to limit the large computation-time  
232 requirements for completing forward-in-time simulations. These prior distributions reflected the  
233 uncertainty in the number of individuals that contributed significantly to the linguistic diversity  
234 observed in the sampled population. The size of this ancestral population was indeed completely  
235 unknown. Indeed, even if some information could be obtained on the census size of the current  
236 population, it likely does not reflect the ancestral linguistic census sizes. Time  $t_1$  was drawn in a

237 uniform prior distribution between 0 and 1000 generations. The median, the minimum, the  
238 maximum, and the 5% quantiles of the priors are summarized in Table 1.

## 239 5. Analyses

### 240 5.1. Simulations

241 For each model, we performed 10 000 simulations using our newly-developed software  
242 *PopLingSim 2 (PLS2)*. We parallelized the simulations using 250 cores of the cluster station  
243 *Genotoul*, amounting to approximately 90 000 CPU hours. Most of this computation time was spent  
244 during the phase to reach equilibrium between mutation and drift at  $t_0 = 5 \times N_0$  generations.

245 During the process of sampling words from our simulations, we simulated missing values by  
246 transforming cognates drawn at random into missing values; the total number of simulated missing  
247 values was set to the number of missing values in the real data set, to avoid the bias they may  
248 induce in the following ABC procedures.

### 249 5.2. Summary statistics

250 We constructed a new set of population linguistic summary statistics. We computed  $p_{i,j}$ , the  
251 proportion of individuals using the word  $i$  of the meaning  $j$ , and then computed the linguistic  
252 diversity  $D_j = 1 - \sum_i p_{i,j}^2$ , analogous to genetic diversity (Nei, 1987). We also computed chi-square  
253 values, over 200 pairs of randomly sampled meanings:  $\chi^2_i = \sum_j (O_{ij} - E_{ij})^2 / E_{ij}$ . The observed value  $O_{ij}$   
254 corresponds to the number of individuals for which a pair of words  $j$  is observed for meaning  $i$ . The  
255 expected value  $E_{ij}$  corresponds to the number of individuals for which this pair of words  $j$  would be  
256 observed for meaning  $i$  if the words were randomly distributed among individuals. We computed  
257 correlation coefficients values, over 200 pairs of meanings randomly sampled, as  $r = (p_{i,i'jj'} - p_{i,j} p_{i',j'})$   
258  $/ [ p_{i,j} (1 - p_{i,j}) p_{i',j'} (1 - p_{i',j'}) ]^{1/2}$ , with  $p_{i,i'jj'}$  the proportion of pairs of individuals using the word  $i$  of  
259 the meaning  $j$  and the word  $i'$  of the meaning  $j'$ . We then computed the frequency spectrum of the  
260 number  $i$  of words per meaning,  $F_i$ .

261 Then, we computed across all words:

- 262 - The mean linguistic diversity,  $\bar{D}$ ;
- 263 - The range of linguistic diversity,  $R(D)$ ;
- 264 - The variance of linguistic diversity,  $V(D)$ ;
- 265 - The mean number of words per meaning,  $\bar{N}$ ;
- 266 - The variance of the number of words per meaning,  $V(N)$ ;
- 267 - The mean number of different words between two individuals,  $\bar{X}$ ;
- 268 - The range of the number of different words between two individuals,  $R(X)$ ;
- 269 - The variance of the number of different words between two individuals,  $V(X)$ ;
- 270 - The mean of the chi-square values,  $\bar{\chi}_2$ ;
- 271 - The range of the chi-square values,  $R(\chi_2)$ ;
- 272 - The variance of the chi-square values,  $V(\chi_2)$ ;
- 273 - The mean of the correlation coefficients values,  $\bar{r}$ ;
- 274 - The range of the correlation coefficients values,  $R(r)$ ;
- 275 - The variance of the correlation coefficients values,  $V(r)$ ;
- 276 - The minimum of the frequency spectrum,  $\min(F)$
- 277 - The maximum of the frequency spectrum,  $\max(F)$
- 278 - The mean of the frequency spectrum,  $\bar{F}$
- 279 - The mode of the frequency spectrum,  $\text{mode}(F)$
- 280 - The range of the frequency spectrum,  $R(F)$
- 281 - The 25th quartile of the frequency spectrum,  $F_{25}$
- 282 - The median of the frequency spectrum,  $F_{50}$
- 283 - The 75th quartile of the frequency spectrum,  $F_{75}$
- 284 - The three axis of the linear discriminant analysis of the previous statistics,  $LD_1, LD_2, LD_3$ .

### 285 **5.3. Power analysis on simulated data**

286 We performed a power analysis of the model selection procedure, to evaluate the impact of the  
287 number of sampled individuals, the number of sampled words, and the number of simulations on  
288 the prior error rate, i.e. the number of cases in which a wrong model was selected among the four  
289 possible models by the ABCRF model-choice procedure. This was done for a total of 61 situations

290 in which we varied the number of sampled individuals between 2 and 100, the number of sampled  
291 words per individuals between 2 and 300, the number of simulations between 1,000 and 10,000.  
292 This maximal value of 10,000 simulations was due to the high computational cost of forward-in-  
293 time simulations. In each case, we computed the prior error rate through cross-validation using the  
294 function *abcrf* of the R package *abcrf*. This procedure considers, in turn, each simulation under each  
295 competing model as a pseudo-observed data, and performs the ABCRF model-choice using all other  
296 simulations in the reference table for training the random forest.

#### 297 **5.4. Model selection on real data**

298 Before model selection, we performed a goodness-of-fit test to check whether the simulations  
299 were able to produce data close to the real data using the function *gfit* from the R package *abc*  
300 (Csilléry et al., 2012). We performed model selection using the R package *abcrf* with the RF  
301 algorithm and the function *abcrf* (Pudlo et al., 2016). We graphically checked if a forest of 500 trees  
302 allowed a convergence of the error rate. We computed the variables importance, indicating which  
303 variables have the most predictive power. We also performed a cross-validation analysis using an  
304 out-of-bag approach implemented in the function *abcrf* of the package *abcrf*, evaluating how the  
305 algorithm *a priori* distinguished between the four models.

306 For the selected model, we then selected the 100 simulations which were closest to the real data,  
307 based on the Euclidean distance of the statistics that were standardized for a mean of 0 and a  
308 variance of 1. We then tested whether the random forest algorithm was able, in this region of  
309 simulated data close to the real data, to correctly select the true model.

#### 310 **5.5. Parameters estimation on real data**

311 We used the RF algorithm with the function *regAbcrf* of the package *abcrf* (Raynal et al., 2017)  
312 to estimate the expectation, median, variance and 5% quantiles of the parameters  $N_1$ ,  $N_0$ ,  $t_1$ ,  $\mu_L$  and  
313 of the composite-parameters  $N_1 \times \mu_L$ ,  $N_0 \times \mu_L$  and  $t_1 \times \mu_L$ . Note that the RF algorithm does not estimate  
314 the entire posterior distribution of the parameters directly, but estimates the quantiles of this  
315 distribution instead.

## 316 **6. Results**

### 317 **6.1 Power analysis**

318 Using simulated data under the four competing linguistic transmission models, we showed that an  
319 increase in the number of words sampled beyond 185 words increased moderately the power of the  
320 analyses (Figure 4). We found also that the decrease in error followed an exponential decay profile  
321 (Figure S1). Increasing the words sampling effort by several orders of magnitude would therefore  
322 be necessary to significantly reduce model selection error. Increasing the number of sampled  
323 individuals beyond 30 individuals increased only slightly the statistical power of the analysis  
324 (Figure 5), which converged towards a limit value. An increased sampling effort on the number of  
325 individuals could also, therefore, only moderately reduce the model selection error. Finally, we  
326 showed that the model-selection prior error rate converged with 10000 simulations (Figure 6),  
327 which indicated that increasing the number of simulations could not lead to a lower error.

### 328 **6.2. Model selection**

329 Using the goodness-of-fit test, we verified that there was no significant differences between the  
330 real and simulated datasets ( $p$ -value = 0.71, with 1000 replications). We performed the RF analysis  
331 using 500 trees, and verified graphically that the error rate converged. The number of trees voting  
332 for the second model was 487 out of 500 (Table 2). The RF analysis thus rejected the *Clonal*,  
333 *Sexual 2* and *Social* models, and selected the *Sexual 1* model for the real data with a posterior  
334 probability  $p = 94.4\%$ .

335 The variable importance analysis showed that the main statistics used by the RF procedure to  
336 select the models were the first two axes of the linear discriminant analysis  $LD_1$  and  $LD_2$ , the  
337 variance of the number of different words between two individuals  $V(X)$ , and the variance of the  
338 correlation coefficients values  $V(r)$  (Figure S2).

339 The cross-validation analysis on simulated datasets (Table 3) indicated a good *a priori*  
340 differentiation between the *Clonal* model and other models, with about 76% of simulated datasets  
341 under this clonal model correctly assigned to the true model. Similarly, the *Sexual 1* model was  
342 correctly attributed for about 76% of the simulated datasets. On the other hand, the *Sexual 2* model  
343 and the *Social* model were difficult to distinguish *a priori*, as the simulated datasets from these two  
344 models were arbitrarily attributed to one or the other by the cross-validation procedure.

345 The RF algorithm assigned to the correct model 100% of the simulations produced by the *Sexual*  
346 *1* model which were closest to the real data. Compared to the global cross-validation results, this  
347 indicated that the method performed better in selecting the correct model in the region of the  
348 parameter space occupied by the real data, than in the entire space occupied by simulations.

### 349 **6.3. Parameter estimation**

350 For the selected model (*Sexual 1*), we could estimate the linguistic mutation rate ( $\mu_L$ ) on the real  
351 data: the quantiles of its posterior distribution were much narrower than the quantiles from its prior  
352 (Table 4). We estimated that this rate ranged between  $1.61 \times 10^{-4}$  and  $1.50 \times 10^{-3}$  mutations per cognate  
353 per linguistic generation at the 95% credibility level. Conversely, we could not estimate the  
354 demographic parameters ( $N_1$ ,  $N_0$ , and  $t_1$ ), for which posterior quantiles did not differ substantially  
355 from prior quantiles. However, we could estimate the composite parameters  $N_{1 \times \mu_L}$ ,  $N_{0 \times \mu_L}$  and  $t_{1 \mu_L}$ ,  
356 for which posterior quantiles were substantially narrower than those of their respective priors. There  
357 was no clear evidence of expansion or contraction, since the confidence intervals of  $N_{1 \times \mu_L}$  and  $N_{0 \times \mu_L}$   
358 overlapped.

## 359 **7. Discussion**

360 We built here individual-based models simulating the linguistic evolution of a population,  
361 under given demographic scenarios, considering four possible types of linguistic transmission

362 between generations. We used an ABCRF framework (Pudlo et al, 2016, Raynal et al, 2019) to  
363 compare the simulated data with a real dataset of 30 individuals in Central Asia typed for 185  
364 words, in order to estimate which model fitted best the data and estimate the parameters of the  
365 selected model.

366 ABC relies on approximating the likelihood of the data by that of a set of summary  
367 statistics, *a priori* informative about the historical process to be inferred. ABC was initially  
368 developed with summary statistics explicitly linked to the parameters of interests, and therefore  
369 highly informative for accurate ABC inference (Tavaré 1997). However, for most case studies, it is  
370 not known, *a priori*, which summary statistics will be informative for ABC inference (Blum et al.  
371 2013). Several complex statistical approaches have been developed, therefore, to select *a priori* sets  
372 of relevant summary statistics for ABC inference and to overcome the dimensionality curse and  
373 parameter posterior identifiability issues, which result from considering very large numbers of  
374 summary statistics, possibly correlated and unevenly informative (Csilléry et al. 2012; Blum et al.  
375 2013; Prangle 2015). Importantly, ABCRF model choice inference is unaffected by the  
376 dimensionality curse faced by most other ABC model-choice frameworks, as each decision tree is  
377 built with random subsets of summary statistics (Pudlo et al. 2016). However, the accuracy of  
378 model parameter inference in ABC, whether using RF or another approach, still relies on finding  
379 minimal subsets of highly informative summary statistics (Raynal et al. 2017), which therefore  
380 requires empirical case-by-case testing of novel sets of summary statistics.

381 A main advantage of the ABC framework is its high flexibility, which will allow researchers,  
382 in future work, to include more sophisticated models with additional parameters of interests to  
383 linguistic evolution. Moreover, ABC offers a model selection procedure that has no equivalent  
384 under an analytic framework, and it offers also the possibility to compute the credibility interval of  
385 the inferred parameters, which would require a fully stochastic approach in an analytical  
386 framework. We showed, first, that some of our models were able to produce simulated data close to  
387 the contemporaneously observed data. Therefore, our approach implements realistic individual-  
388 based linguistic transmission, consistent with the observed linguistic diversity of the sampled  
389 populations.

390 We also provided inferences on some features of linguistic history of Tajik-speaking  
391 individuals, selecting the most plausible mechanisms of linguistic transmission among the  
392 competing options tested, and estimating the parameters of the selected models for our sample. The  
393 high posterior probabilities of the *Sexual* model 1, compared to the other models, indicated that the  
394 mechanisms of linguistic acquisition followed, in this study-case, a process of linguistic  
395 transmission from two teachers with their own vocabulary. In other words, we inferred that these  
396 individuals did not learn their basic vocabulary from only one individual, nor from two individuals  
397 without “sex”-specific vocabulary, nor from the whole speech community. We estimated that they  
398 learn their vocabulary from two individuals with “male”-specific and “female”-specific words. This  
399 linguistic-transmission mechanism may reflect the fact that Tajik populations are cognatic (Krader,  
400 1966), i.e. they inherit social status and material goods from their two parents. This symmetric role  
401 of parents in cultural transmission across generations appears thus to be reflected linguistically, as  
402 learners appear to receive specific words from both parents. Future studies on populations with  
403 other lineage and kinship descent systems, such as patrilineal or matrilineal descent rules, will allow  
404 better understanding how social-descent rules and features may influence linguistic transmission  
405 processes in a given population.

406 Our estimates of the mean linguistic mutation rate of the words of the Swadesh list in this  
407 population ranged between  $10^{-4}$  and  $10^{-3}$  mutations per word per generation. Interestingly, the  
408 mutation rate estimated here fell in the same range as the mutation rate estimated in previous  
409 macro-evolutionary linguistic studies (Pagel et al., 2007). Considering that languages at a global  
410 scale emerge from the interactions among individuals, we may thus hypothesise that the mutation  
411 rate estimated globally emerges from the mutation rate at a local scale. Under this assumption,  
412 further studies could investigate whether macro-evolutionary linguistic processes (i.e. processes  
413 occurring at the scale of a whole language or a linguistic variety), may also emerge from micro-  
414 evolutionary linguistic processes (i.e. at the scale of the individuals within a population).

415 Population genetics effective population sizes estimated in Central Asia differ from census  
416 population sizes (Palstra et al. 2012). Similarly, the estimated linguistic population size in our  
417 model did not necessarily reflect the real size of the community, and effective linguistic populations  
418 are possibly much smaller in size than empirical groups of speakers. Our posterior estimates of the



419 number individuals that contributed significantly to the observed linguistic diversity did not differ  
420 from the priors of the simulations. It meant that our method could not directly estimate the number  
421 of individuals in the current and ancestral linguistic populations, but only synthetic parameters such  
422 as  $N_0\mu$ . In this context, a perspective might be to design specific summary statistic to improve our  
423 ability to infer the number of individuals that contributed significantly to the observed linguistic  
424 diversity. Another promising approach might be to sample individuals in the population at different  
425 moments in time, separated by at least several decades, analogously to what is done in population  
426 genetics, where it is the most efficient method for estimating recent population sizes, independently  
427 of mutation rate (Foll et al., 2014).

428         In this study, unlike most other studies focusing on within-population linguistic diversity  
429 (Baxter et al., 2009; Danescu-Niculescu-Mizil et al., 2013; Kandler et al., 2010), we only used  
430 contemporaneous linguistic diversity. This method allowed us to perform historical inferences based  
431 only on sampling campaigns conducted in existing populations. The amount of available  
432 information depended only on the sampling effort, and not on the availability of dense historical  
433 records, which are unavailable for numerous languages. It would be of great interest in future works  
434 to be able to distinguish among the *Sexual 2* model (with only two teachers) and the *Social* model  
435 (with a whole community as a teacher). As we showed in the power analysis, increasing the  
436 sampling effort (in terms of number of individuals or in term of number of words) was not sufficient  
437 to reliably distinguish between these two models, using our set of summary statistics. As for the  
438 inference of demographic parameters, developing new summary statistics and/or designing multi-  
439 generational studies might be the best solution to further distinguish among closely related  
440 linguistic transmission modes in future work.

441         Our approach could be extended in several other ways. First, the linguistic acquisition  
442 models that we proposed here did not integrate the particular constraints of communication  
443 processes. In particular, we assumed a neutral production of variants without any constraints on  
444 linguistic communication. Some evolutionary linguists would argue for an integration of the  
445 particularity of languages as communication systems, associated with a strong set of constraints  
446 (Beckner et al., 2009). Indeed, individuals maximize their probability of being understood, while  
447 minimizing their communication costs, two features that strongly affect linguistic evolutionary

448 processes (Tamariz and Kirby, 2015). These constraints are particularly strong in the case of  
449 phonological, morphological, or syntactical systems, and we may wonder whether lexical variants  
450 are also subjected to these constraints. If so, these particularities of linguistic systems may be at  
451 odds with inferences based on a model of neutral evolution, and should thus be taken into account  
452 for a more accurate model of linguistic evolution at the individual scale, for historical inferences  
453 purposes.

454 It will also be of interest to study the transmission model of other types of linguistic objects,  
455 for instance focusing on other types of words such as food lexica, or very recently acquired  
456 technological lexica. Those different types of words could be transmitted differently, and our results  
457 could be different in the case of these particular lexical elements. Other types of linguistic data  
458 could also be obtained, like phonetic productions or syntactic rules, and it could be then assessed  
459 whether these linguistic elements are transmitted or not in the same way as the words of the  
460 Swadesh list. In addition, individuals may know more variants than those they use most frequently.  
461 It may then be possible in future works to model also the evolution of word usage, in order to take  
462 into account a greater part of the lexical diversity of languages.

463 We assumed also that linguistic transmission occurs between two linguistic generations,  
464 ignoring the impact of communications between individuals of the same generation. Moreover, we  
465 did not take into account global media such as books, radio, internet, or television. We will thus  
466 consider in further investigations several alternative models of language evolution, where the  
467 acquisition of language results from a series of interactions between individuals, who would update  
468 their language during each conversation.

469 Finally, note that the formalism of our models are close to the formalism of population  
470 genetics. This should allow us to develop joint inferences coupling genetic and linguistic data for  
471 the same set of populations and individuals. However, some theoretical limits remain. We may  
472 wonder whether a speech community (a “linguistic population”) is identical to a reproductive group  
473 (a “genetic population”). It is far from obvious that human reproductive boundaries overlap  
474 language boundaries among human groups. A joint model between genetics and linguistics would  
475 thus require clarifying and articulating rigorously the concepts of population genetics with the  
476 concepts of population linguistics to propose robust joint inferences.

## 477 **8. Acknowledgements**

478 We thank the Genotoul bioinformatics platform (Toulouse, Midi-Pyrénées) for providing help,  
479 computing and storage resources. V.T. was financed by a PhD grant from the French ‘Ministère de  
480 l’Education Nationale, de l’Enseignement Supérieur et de la Recherche’. V.T. and F.A. received a  
481 travel grant from the NEFREX project funded by the European Union (People Marie Curie Actions,  
482 International Research Staff Exchange Scheme, call FP7-PEOPLE-2012-IRISES). This work was  
483 also partially funded by the Agence Nationale de la Recherche grant DemoChips (ANR-12-BSV7-  
484 0012).

## 485 **9. Bibliography**

- 486 Atkinson, Q., Nicholls, G., Welch, D., & Gray, R. (2005). From words to dates: Water into wine,  
487 mathemagic or phylogenetic inference? *Transactions of the Philological Society*, 103(2), 193–  
488 219.
- 489 Atkinson, Q.D. (2011). Phonemic Diversity Supports a Serial Founder Effect Model of Language  
490 Expansion from Africa. *Science* 332, 346–349.
- 491 Baxter, G.J., Blythe, R.A., Croft, W., and McKane, A.J. (2009). Modeling language change: An  
492 evaluation of Trudgill’s theory of the emergence of New Zealand English. *Language Variation and*  
493 *Change* 21, 257.
- 494 Beaumont, M.A., Zhang, W., and Balding, D.J. (2002). Approximate Bayesian computation in  
495 population genetics. *Genetics* 162, 2025–2035.
- 496 Beckner, C., Blythe, R., Bybee, J., Christiansen, M.H., Croft, W., Ellis, N.C., Holland, J., Ke, J.,  
497 Larsen-Freeman, D., and Schoenemann, T. (2009). Language is a complex adaptive system:  
498 Position paper. *Language Learning* 59, 1–26.

- 499 Bouckaert, R., Lemey, P., Dunn, M., Greenhill, S.J., Alekseyenko, A.V., Drummond, A.J., Gray,  
500 R.D., Suchard, M.A., and Atkinson, Q.D. (2012). Mapping the Origins and Expansion of the Indo-  
501 European Language Family. *Science* 337, 957–960.
- 502 Breiman, L. (1999). Random forests. UC Berkeley TR567.
- 503 Croft, W. (1996). Linguistic Selection: An Utterance-based Evolutionary Theory of Language  
504 Change. *Nordic Journal of Linguistics* 19, 99.
- 505 Croft, W. (2008). Evolutionary Linguistics. *Annual Review of Anthropology* 37, 219–234.
- 506 Csilléry, K., François, O., and Blum, M.G.B. (2012). abc: an R package for approximate  
507 Bayesian computation (ABC): *R package: abc*. *Methods in Ecology and Evolution* 3, 475–479.
- 508 Danescu-Niculescu-Mizil, C., West, R., Jurafsky, D., Leskovec, J., and Potts, C. (2013). No  
509 country for old members: User lifecycle and linguistic change in online communities. In  
510 Proceedings of the 22nd International Conference on World Wide Web, (ACM), pp. 307–318.
- 511 Darlu, P., Bloothoof, G., Boattini, A., Brouwer, L., Brouwer, M., Brunet, G., Chareille, P.,  
512 Cheshire, J., Coates, R., Dräger, K., et al. (2012). The Family Name as Socio-Cultural Feature and  
513 Genetic Metaphor: From Concepts to Methods. *Human Biology* 84, 169–214.
- 514 Dryer, M.S., and Haspelmath, M. (2013). The World Atlas of Language Structures Online  
515 (Leipzig: Max Planck Institute for Evolutionary Anthropology).
- 516 Foll, M., Poh, Y. P., Renzette, N., Ferrer-Admetlla, A., Bank, C., Shim, H., Malaspinas, A.S.,  
517 Ewing, G., Liu, P., Wegmann, D. & Caffrey, D. R. (2014). Influenza virus drug resistance: a time-  
518 sampled population genetics perspective. *PLoS Genet*, 10(2), e1004185.
- 519 Gray, R.D., and Atkinson, Q.D. (2002). Language-tree divergence times support the Anatolian  
520 theory of Indo-European origin. *Geophysical Research Letters* 29.
- 521 Gray, R.D., Drummond, A.J., and Greenhill, S.J. (2009). Language phylogenies reveal expansion  
522 pulses and pauses in Pacific settlement. *Science* 323, 479–483.
- 523 Heeringa, W., and Nerbonne, J. (2001). Dialect areas and dialect continua. *Language Variation*  
524 and Change 13, 375–400.

- 525 Kandler, A., Unger, R., and Steele, J. (2010). Language shift, bilingualism and the future of  
526 Britain's Celtic languages. *Philosophical Transactions of the Royal Society B: Biological Sciences*  
527 *365*, 3855–3864.
- 528 Kirby, K.R., Gray, R.D., Greenhill, S.J., Jordan, F.M., Gomes-Ng, S., Bibiko, H.-J., Blasi, D.E.,  
529 Botero, C.A., Bower, C., Ember, C.R., et al. (2016). D-PLACE: A Global Database of Cultural,  
530 Linguistic and Environmental Diversity. *PLOS ONE* *11*, e0158391.
- 531 Krader, L. (1966). *Peoples of central Asia* (Indiana University [1966]).
- 532 Livingstone, D., and Fyfe, C. (1999). Modelling the evolution of linguistic diversity. *Advances in*  
533 *Artificial Life* 704–708.
- 534 Menecier, P., Nerbonne, J., Heyer, E., and Manni, F. (2016). A Central Asian Language Survey.  
535 *Language Dynamics and Change* *6*, 57–98.
- 536 Nei, M. (1987). *Molecular Evolutionary Genetics* (Columbia University Press).
- 537 Pagel, M., Atkinson, Q.D., and Meade, A. (2007). Frequency of word-use predicts rates of  
538 lexical evolution throughout Indo-European history. *Nature* *449*, 717–720.
- 539 Pagel, M., Atkinson, Q.D., S. Calude, A., and Meade, A. (2013). Ultraconserved words point to  
540 deep language ancestry across Eurasia. *Proceedings of the National Academy of Sciences* *110*,  
541 8471–8476.
- 542 Palstra, F. P., & Fraser, D. J. (2012). Effective/census population size ratio estimation: a  
543 compendium and appraisal. *Ecology and evolution*, *2*(9), 2357-2365.
- 544 Pateman, T. (1983). What is a language? *Language & Communication* *3*, 101–127.
- 545 Prangle, D. (2015). Summary statistics in approximate Bayesian computation. arXiv preprint  
546 arXiv:1512.05633.
- 547 Pudlo, P., Marin, J.-M., Estoup, A., Cornuet, J.-M., Gautier, M., and Robert, C.P. (2016).  
548 Reliable ABC model choice via random forests. *Bioinformatics* *32*, 859–866.
- 549 Raynal L., Marin J.-M., Pudlo P., Ribatet M., Robert C.P., and Estoup A. (2019). ABC random  
550 forests for Bayesian parameter inference. *Bioinformatics* *35*:1720-1728.
- 551 Reesink, G., Singer, R., and Dunn, M. (2009). Explaining the Linguistic Diversity of Sahul  
552 Using Population Models. *PLOS Biology* *7*, e1000241.

- 553      Rodriguez-Larralde, and Barraï (2000). Elements of the surname structure of Austria. *Annals of*  
554 *Human Biology* 27, 607–622.
- 555      Tamariz, M., and Kirby, S. (2015). Culture: Copying, Compression, and Conventionality.  
556 *Cognitive Science* 39, 171–183.
- 557      Tavaré, S., Balding, D.J., Griffiths, R.C., and Donnelly, P. (1997). Inferring Coalescence Times  
558 from DNA Sequence Data. *Genetics* 145, 505–518.
- 559      Thouzeau, V., Menecier, P., Verdu, P., and Austerlitz, F. (2017). Genetic and linguistic histories  
560 in Central Asia inferred using approximate Bayesian computations. *Proc. R. Soc. B* 284, 20170706.
- 561      Verdu, P., Jewett, E.M., Pemberton, T.J., Rosenberg, N.A., and Baptista, M. (2017). Parallel  
562 Trajectories of Genetic and Linguistic Admixture in a Genetically Admixed Creole Population.  
563 *Current Biology* 27, 2529-2535.e3.
- 564      Vogt, P. (2009). Modeling interactions between language evolution and demography. *Human*  
565 *Biology* 81, 237–258.
- 566      Wang, J. (2005). Estimation of effective population sizes from data on genetic markers.  
567 *Philosophical Transactions of the Royal Society of London B: Biological Sciences* 360, 1395–1409.
- 568      Wang, J. (2016). A comparison of single-sample estimators of effective population sizes from  
569 genetic marker data. *Mol. Ecol.* 25, 4692–4711.

	<b>Median</b>	<b>Min</b>	<b>Max</b>	<b>Variance</b>	<b>Quantile 2.5%</b>	<b>Quantile 97.5%</b>
$N_0$	550	100	1000	67645	122	978
$N_1$	550	100	1000	67645	122	978
$t_1$	500	0	1000	83490	25	975
$\mu_L$	$3.165 \times 10^{-4}$	$10^{-6}$	$10^{-1}$	$3.58 \times 10^{-4}$	$1.35 \times 10^{-6}$	$7.73 \times 10^{-2}$
$N_0 \times \mu_L$	0.150	$10^{-4}$	100	141.91	$5.25 \times 10^{-4}$	44.5
$N_1 \times \mu_L$	0.150	$10^{-4}$	100	139.05	$5.25 \times 10^{-4}$	44.5
$t_1 \times \mu_L$	0.116	0	100	129.55	$2.80 \times 10^{-4}$	42.0

**Table 1** – Summary of the prior distributions of the parameters for the four models.

<b>Clonal</b>	<b>Sexual 1</b>	<b>Sexual 2</b>	<b>Social</b>
11	487	1	1

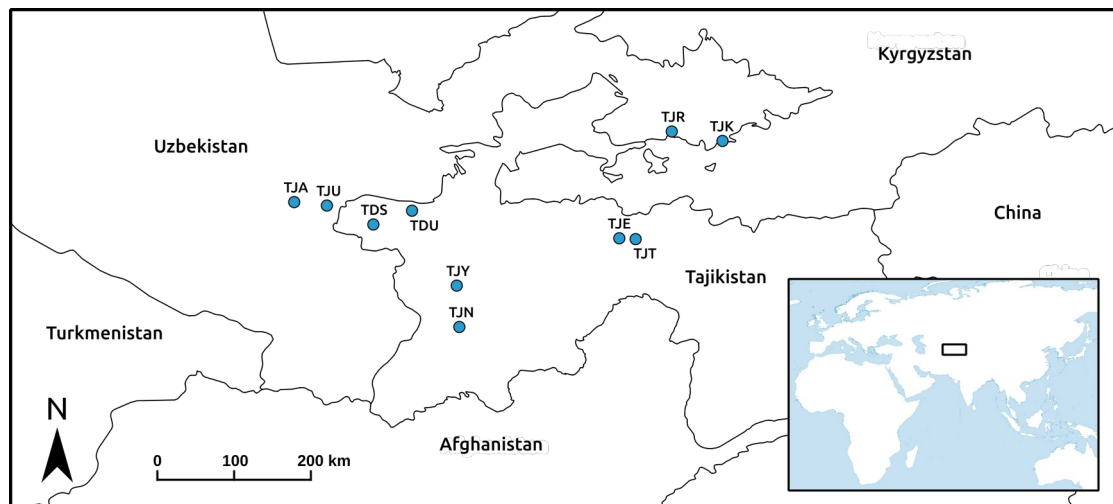
**Table 2** – Proportion of votes for the four models of linguistic evolution.

		<b>True model</b>			
		<b>Clonal</b>	<b>Sexual 1</b>	<b>Sexual 2</b>	<b>Social</b>
<b>Selected model</b>	<b>Clonal</b>	7620	1785	282	313
	<b>Sexual 1</b>	1358	7698	439	505
	<b>Sexual 2</b>	283	816	4782	4119
	<b>Social</b>	276	805	4200	4719

**Table 3** – Confusion matrices from the out-of-bag cross-validation analysis of the four models, using 10 000 pseudo-observed data.

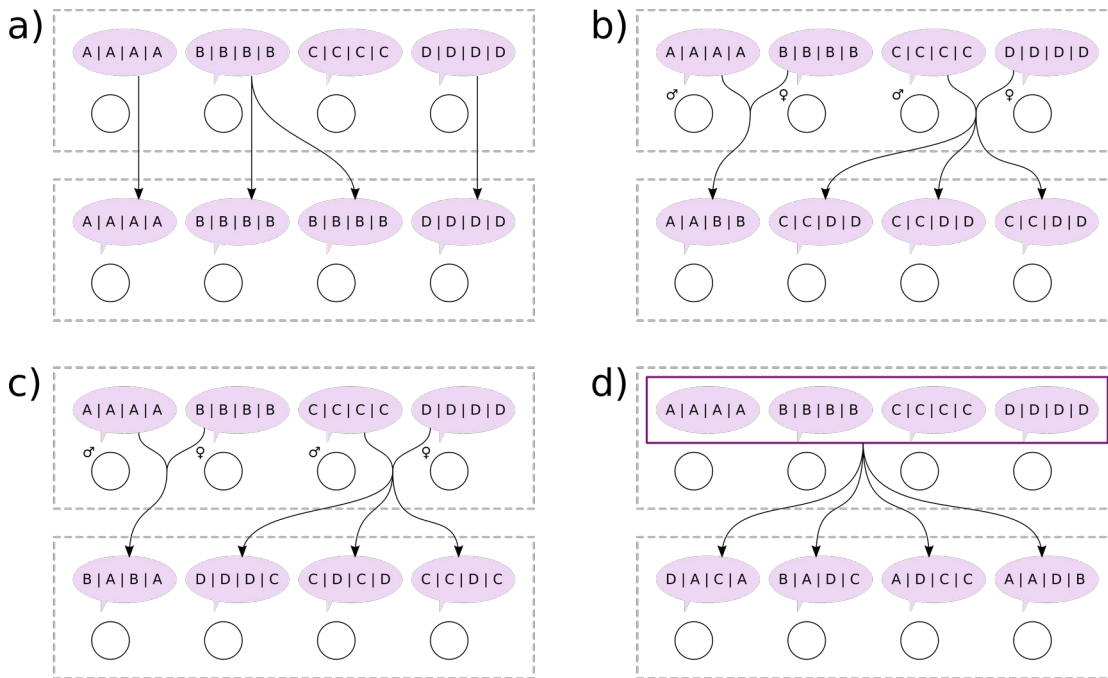
	Expectation	Median	Variance	Quantile 2.5%	Quantile 97.5%
$N_0$	489	523	70157	112	928
$N_1$	541	547	69571	105	949
$t_0$	548	541	99439	48	985
$\mu_L$	$5.42 \times 10^{-4}$	$4.14 \times 10^{-4}$	$1.8 \times 10^{-7}$	$1.61 \times 10^{-4}$	$1.50 \times 10^{-3}$
$N_0 \times \mu_L$	0.250	0.139	0.081	0.036	0.967
$N_1 \times \mu_L$	0.178	0.155	$4.35 \times 10^{-3}$	0.098	0.347
$t_1 \times \mu_L$	0.358	0.215	0.148	0.010	1.15

**Table 4** – Summary of the posterior distributions of the parameters, assuming a *Sexual 1* scenario.

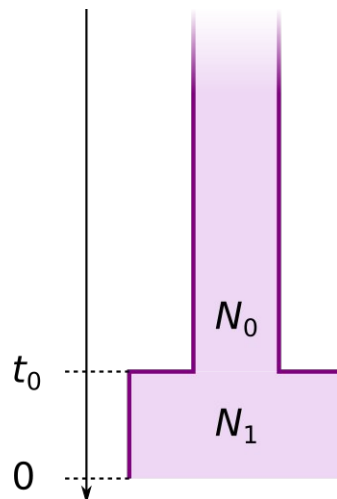


**Figure 1** – Geographical distribution of the 10 sampled units under study.

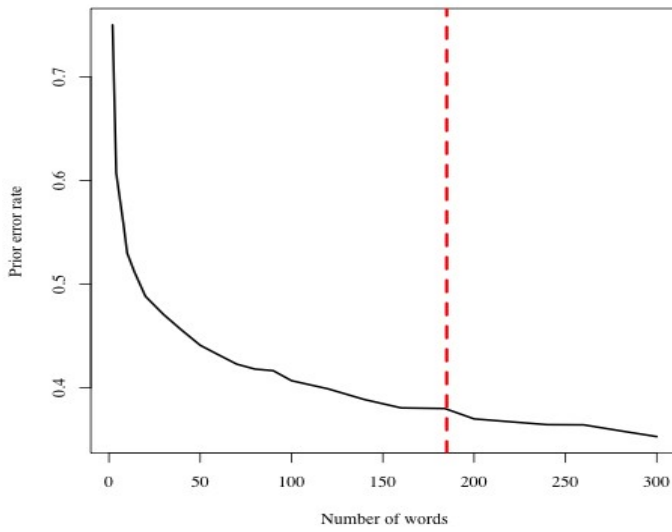




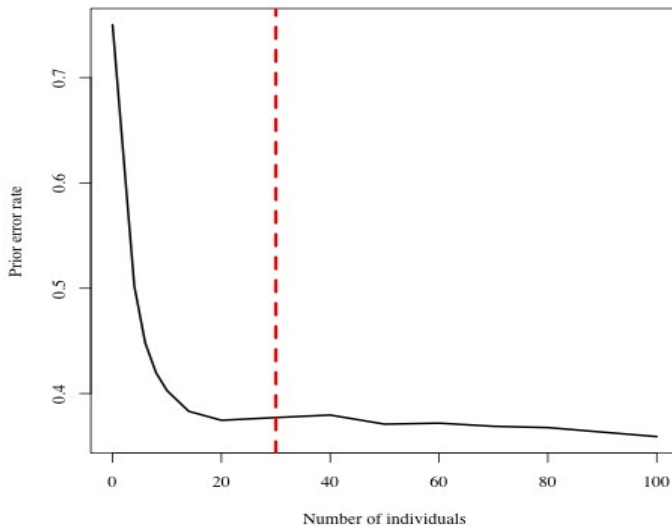
**Figure 2** – Four models of linguistic transmission between generations. Each circle represents an individual. The utterances that individuals produce depend only on the utterances that their teachers produced at the previous generation, and on the mutations induced during the transmission. Four transmission modalities were considered: (a) a “Clonal” model with only one teacher per learner, (b) a “Sexual 1” model with two teachers associated with a distinct set of vocabulary for each sex, (c) a “Sexual 2” model with two teachers without a distinct set of vocabulary for each sex, and (d) a “Social” model with the whole population as teacher for each learner.



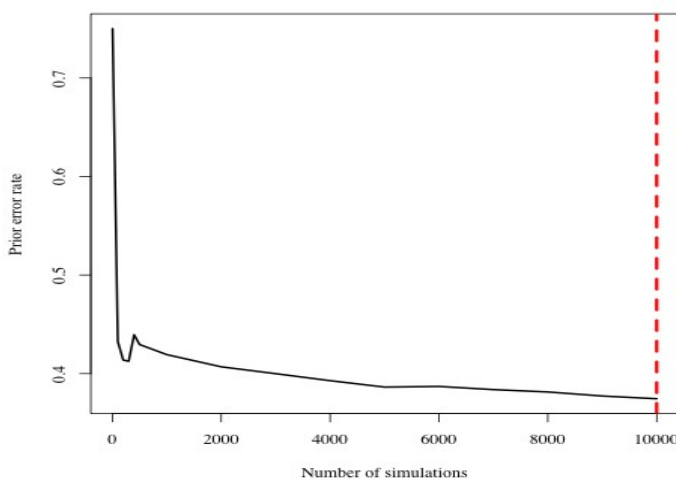
**Figure 3** – Historical scenario. If  $N_0 = N_1$ , we assumed a scenario of constant population size. If  $N_0 < N_1$ , we assumed a scenario of expansion of the population. If  $N_0 > N_1$ , we assumed a scenario of contraction of the population.



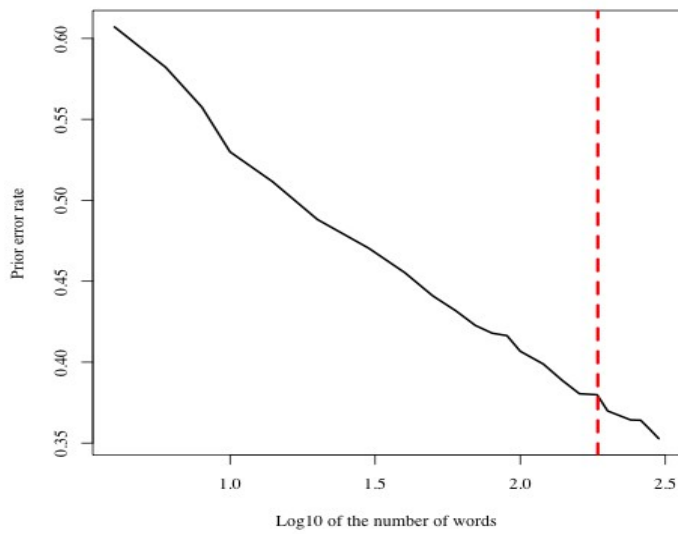
**Figure 4** – Prior error rate depending on the simulated number of sampled words, with 30 sampled individuals and 10000 simulations. The red dashed line indicates the number of words of the real sample.



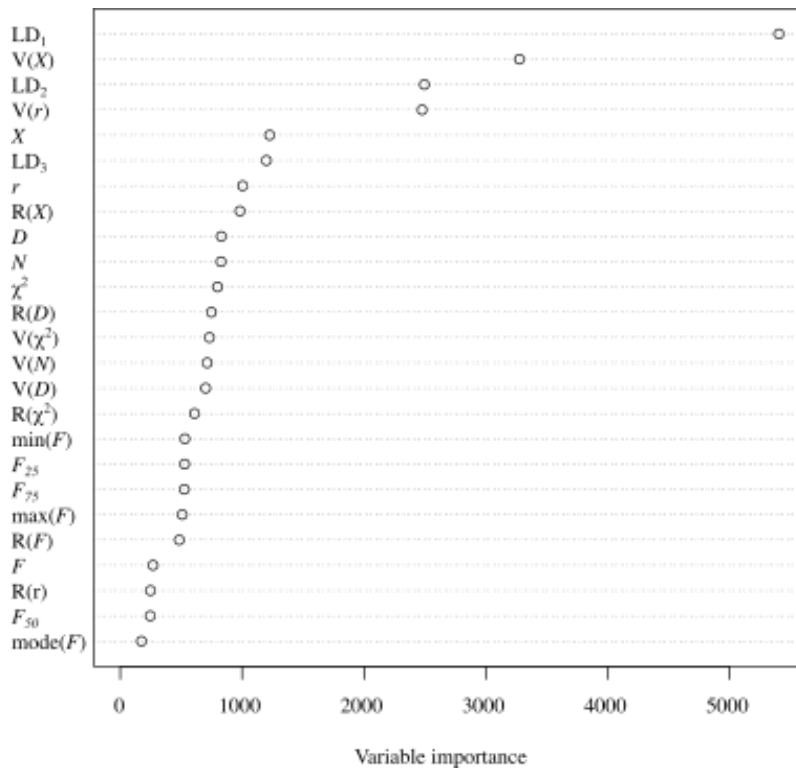
**Figure 5** – Prior error rate depending on the simulated number of sampled individuals, with 185 sampled words and 10000 simulations. The red dashed line indicates the number of individuals of the real sample.



**Figure 6** – Prior error rate depending on the number of simulations, with 30 sampled individuals and 185 sampled words. The red dashed line indicates the value used for the analyses.



**Figure S1** – Prior error rate depending on the simulated decimal logarithm of the number of sampled words, with 30 sampled individuals and 10000 simulations. The red dashed line indicates the number of words of the real sample.



**Figure S2** – Variable importance of the random forest built for the model selection.