



The role of genetic biases in shaping the correlations between languages and genes

Dan Dediu *

School of Philosophy, Psychology and Language Sciences, The University of Edinburgh, 14 Buccleuch Place, Edinburgh EH8 9LN, UK

ARTICLE INFO

Article history:

Received 10 February 2008

Received in revised form

20 May 2008

Accepted 21 May 2008

Available online 29 May 2008

Keywords:

Genetic bias

Language change

Computer model

ABSTRACT

It has recently been proposed [Dediu, D., Ladd, D.R., 2007. Linguistic tone is related to the population frequency of the adaptive haplogroups of two brain size genes, ASPM and Microcephalin. *Proc. Natl Acad. Sci. USA* 104(26), 10944–10949] that genetically coded linguistic biases can influence the trajectory of language change. However, the nature of such biases and the conditions under which they can become manifest have remained vague. The present paper explores computationally two plausible types of linguistic acquisition biases in a population of agents implementing realistic genetic, linguistic and demographic processes. One type of bias represents an innate asymmetric initial state (*initial expectation bias*) while the other an innate asymmetric facility of acquisition (*rate of learning bias*). It was found that only the second type of bias produces detectable effects on language through cultural transmission across generations and that such effects are produced even by weak biases present at low frequencies in the population. This suggests that learning preference asymmetries, very small at the individual level and not very frequent at the population level, can bias the trajectory of language change through the process of cultural transmission.

© 2008 Elsevier Ltd. All rights reserved.

1. Introduction

In their recent paper, Dediu and Ladd (2007) argued for the existence of a correlation between the population frequency of the derived haplogroups of two brain growth and development related genes, *ASPM* and *Microcephalin*, and the usage of linguistic tone in the language(s) spoken by that population. To this end, they used a world sample of 49 populations for which information was collected for 983 alleles and 26 linguistic features, while controlling for geographical distance and known historical linguistic relationships. The most controversial claim of the paper concerns the nature of this correlation, which is argued to be causal, due to a putative genetic bias in the processing of tone induced by the haplogroups concerned.

The exact nature of this genetic bias is not specified, but it is argued that it involves three components (Dediu and Ladd, 2007):

- (i) from *interindividual genetic differences to differences in brain structure and function*,
- (ii) from *differences in brain structure and function to interindividual differences in language-related capacities*, and
- (iii) from these to *typological differences between languages*.

* Tel.: +44 7903387241.

E-mail address: Dan.Dediu@ed.ac.uk

Any claim involving a genetic bias manifested in typological differences between languages must make reference to these three components, which represent the flow of causation from genes to language. More generally, this same argument applies to any claim of genetic biases causing cultural differences between human populations. Components (i) and (ii), concerning inter-individual variability, are generally well established for language either independently (Lenroot et al., 2007; Wright et al., 2002; Bartley et al., 1997; Thompson et al., 2001; Scamvougeras et al., 2003) or as a conglomerate in studies involving the heritability of language (Stromswold, 2001; Bonneau et al., 2004; Bishop, 2003; Fisher et al., 2003; Felsenfeld, 2002; Plomin and Kovas, 2005). Component (iii) concerns interpopulation variability and asserts that populations with different genetic structures could develop overt linguistic differences. The claim is that individual biases can be either amplified or hidden by the cultural transmission of language in a population of biased agents, making them visible or not at the language level.

Previous studies of the cultural transmission of language by biased agents are not numerous and come mainly from the field of language evolution.¹ For example, Nettle (1999) uses computer models and is mainly concerned with explaining language change

¹ Boyd and Richerson's (1985) distinction between various types of bias in cultural evolution and their treatment of directly biased transmission is relevant, but their approach seems too general to answer the questions of interest for this paper.

and the threshold problem, and includes the impact of functional biases, suggesting that they are effective in influencing the trajectory of language change. However, the study is limited to uniform populations with respect to the strength of these biases. Smith (2004) shows that the evolution of vocabulary is influenced by the “innate” biases of simulated agents (in favor of, neutral or against homonymy) and the population structure with respect to the relative frequencies of these different biases.

A new and productive framework for treating language evolution and language change is represented by the Bayesian approach (Press, 2003), where agents are considered to be Bayesian learners (Griffiths and Kalish, 2007; Kirby et al., 2007; Smith and Kirby, 2008; Hawkey, 2008), having a prior distribution over the possible languages, $P(h)$, and updating this distribution to reflect the observed linguistic data, d , to result in their posterior distribution,

$$P(h|d) = \frac{P_{PA}(d|h)P(h)}{P_{PA}(d)}$$

where h represents a hypothesis (language), $P_{PA}(d|h)$ is the probability of producing the observed linguistic data, d , under the hypothesis h , and $P_{PA}(d) = \sum_h P_{PA}(d|h)P(h)$. In this framework, the prior $P(h)$ is equated to the *learning bias* and the agent selects a single “winning” hypothesis from the posterior $P(h|d)$ to represent its linguistic knowledge. Griffiths and Kalish (2007) propose two such *learning algorithms*, namely the *sampling learner* which chooses as the “winner” a random hypothesis with a probability proportional to its posterior probability, and the *maximum a posteriori* or *MAP learner*, which chooses as the “winner” the hypothesis with the maximum posterior probability. They prove that, if certain assumptions are met, including identical agents and generations composed of a single agent, iterated learning with sampling agents is equivalent to a Gibbs sampler and always converges to the prior, while for MAP agents, the system is equivalent to an expectation-maximization algorithm, and the behavior is more complex but still largely influenced by the prior.

Kirby et al. (2007) focus on the MAP learner and show that there is a continuum of learning algorithms by proposing that learners choose the “winning” hypothesis with probability $(P_{PA}(d|h)P(h))^r$: when $r = 1$ the learner samples from the posterior (sampler), while for $r \rightarrow \infty$ the learner picks the hypothesis with the maximum posterior probability (MAP). As opposed to the sampler, $r = 1$, which invariably converges to the prior $P(h)$, the learners with $r > 1$ pick languages proportional to $P(h)^r$, deviating from the prior. Therefore, they conclude that small learning biases can be amplified by the process of cultural transmission and made manifest as universals. Smith and Kirby (2008) analyze the evolutionary stability of sampling and maximizing (MAP) against invasion by the opposing strategy and conclude that maximizing is always preferred over sampling. Moreover, assuming a fitness cost to strong priors, they show that evolution favors weak biases.

However, there are a number of issues with these studies, including the assumption that the prior distribution over the languages, $P(h)$ captures all the aspects of the vague concept of a *learning bias* or that the human language learning process can be approximated by a Bayesian formalism (for a critique see, for example, Hawkey, 2008). But, from the point of view of the present study, two assumptions are very relevant and open to argument: that the agents are identical and that the linguistic community is degenerated to a single teacher and a single learner.²

² Griffiths and Kalish (2007, Section 7), discuss the case of an infinite population of identical agents, where at each time step a learner sees data produced by a single random teacher, and show that such a model converges to a

The present study is a computational investigation of the effects of genetic biases on language change in structured populations of agents. It implements realistic³ demographic, genetic and linguistic processes and two types of genetic biases on language acquisition. It is well known that first language learning represents one of the proposed mechanisms of language change (Campbell, 2004) together with second language acquisition by adult learners (Ostler, 2005) and adult usage (Croft, 2000), and there is still a debate concerning their relative roles. The biases explored in this paper refer strictly to first language acquisition, in the vein of the previous studies cited above, while the effects of biases manifest in adult second language learners and adult usage will be investigated in a future study.

2. The model

The model world is composed of $m \times n$ regions, arranged in a square grid. Time is discretized in simulation years. Each region, R_{ij} , $1 \leq i \leq m$, $1 \leq j \leq n$, can support a population, P_{ij} , of a given (constant) optimal size, S_{ij}^{opt} . The current population size at time t , S_{ij}^t is attracted towards S_{ij}^{opt} , in the sense that when $S_{ij}^t > S_{ij}^{opt}$, both mortality and emigration increase, while when $S_{ij}^t < S_{ij}^{opt}$, mortality decreases and the region becomes a preferential target for immigration. The unit of the simulation is represented by an agent. Each agent has a limited lifespan, age_{max} , and the probability that an agent will die at each time step t is determined by its age, its fitness and the population pressure $S_{ij}^t - S_{ij}^{opt}$. There is a critical period up to $age_{critical}$ during which language acquisition takes place. The agent becomes sexually active at $age_{puberty}$ and can mate with another mature agent of the opposing sex. All demographic processes (mating and migrations) depend on space, in that both involve only immediately neighboring (Moore neighborhoods) regions.

For this study, $m = n = 10$, $S_{ij}^{opt} = 50$, $age_{critical} = 5$, $age_{puberty} = 15$, $age_{max} = 70$, and there are no fitness differences between agents (neutral evolution). Test runs for different world dimensions (m and n) and optimal population sizes (S_{ij}^{opt}) have shown that the behavior of the model is qualitatively the same irrespective of these parameters, but for bigger worlds and optimal sizes the random fluctuations are less important. The model is also robust with respect to the agent-related parameters $age_{critical}$ and $age_{puberty}$.⁴

Each agent has a genome composed of two independent genes, G_1 and G_2 , each with two alleles, one of which is denoted * and is of “special interest”, akin to the derived haplogroups of *ASPM* and *Microcephalin*. All four alleles are selectively neutral but they can influence the linguistic development of the agents by coding specific linguistic biases. Concerning the linguistic aspect of the model, there are two features, denoted F_1 and F_2 , each with two possible values, one of them denoted *. An agent represents its linguistic world through two probabilities, p_1 and p_2 , where p_i is the probability that F_i has value *. In general, these two linguistic features are not necessarily independent and the joint probability

(footnote continued)

state dependent on the prior. However, in this case, there is no social structure and the learning process is still essentially single teacher-single learner.

³ The qualification “realistic” must be understood by comparison with previous models and the hypothesis of interest.

⁴ 280 test runs: 54 runs, 27 for $m = n = 5$ and 27 for $m = n = 15$; 54 runs, 27 for $S_{ij}^{opt} = 25$ and 27 for $S_{ij}^{opt} = 100$; 81 runs, 27 for $age_{critical} = 1$, 27 for $age_{critical} = 3$ and 27 for $age_{critical} = 10$; 81 runs, 27 for $age_{puberty} = 1$, 27 for $age_{puberty} = 5$ and 27 for $age_{puberty} = 25$. Larger world and population sizes, and critical and puberty ages have higher computational costs while smaller values tend to be more erratic.

$p_{1,2}$ of both F_1 and F_2 having value $*$ can represent functional or cognitive relationships between them. Utterances are produced conforming to these probabilities, containing F_1^* with probability p_1 , F_2^* with probability p_2 , and the joint distribution of F_1^* and F_2^* being governed by $p_{1,2}$.

During first language acquisition, an agent samples utterances produced by the agent's own mother (n_{mother} utterances), the other linguistically mature (i.e., past $age_{critical}$) members of the agent's population ($n_{stranger}$ utterances each) and a proportion ($f_{foreigners}$) of the linguistically mature members of the neighboring populations ($n_{foreigner}$ utterances each). For the present study $n_{mother} = 100$, $n_{stranger} = 50$, $f_{foreigners} = 0.05$ and $n_{foreigner} = 10$, so that the most influential role is played by the agent's mother, followed by the agent's own speech community and lastly by the neighboring speech communities. Test runs⁵ have shown that the model is robust with respect to these parameters.

The agent computes the probabilities p_1 , p_2 and $p_{1,2}$ based on the frequencies of heard utterances containing F_1^* , F_2^* and (F_1^* and F_2^*), respectively, denoted f_1 , f_2 and $f_{1,2}$. More specifically, let us focus on F_1^* but the same applies to the other two cases, as well: the difference between the observed frequency, f_1^t , and the agent's internal probability, p_1^t , at time t , $\Delta_1^t = |p_1^t - f_1^t|$, is used to update the agent's internal probability at time $t + 1$, p_1^{t+1} , by making

$$p_1^{t+1} = \begin{cases} p_1^t + \Delta_1^t \cdot r_1^+ & \text{if } p_1^t \leq f_1^t \\ p_1^t - \Delta_1^t \cdot r_1^- & \text{otherwise} \end{cases}$$

where $0 \leq r_1^+, r_1^- \leq 1$ are the *learning rates* adjusting the weight of the evidence in favor of or against F_1^* .

The three models for the genetic biases are:

- **M₀** (*No genetic biases*): no influence from the genome on the computation of the linguistic probabilities p_1 , p_2 and $p_{1,2}$. More specifically, $r_1^+ = r_1^- = r_2^+ = r_2^- = r_{1,2}^+ = r_{1,2}^- = 1.0$ and $p_1 = p_2 = p_{1,2} = 0.5$ initially. For example, if we consider F_1^* to represent tone and G_1^* to represent the derived haplogroup of *ASPM*, this model describes the case where a carrier of *ASPM-D* is not different in any relevant respect from a non-carrier in learning about the tonality of its language;
- **M₁** (*Genes bias the initial expectation*): represents the case where genes bias language acquisition by coding for *different initial starting points*. G_1 influences F_1 , in the sense that if an agent has the G_1^* allele, then initially its $p_1 = 1.0$ (very strongly "predisposed" to expect a language of type F_1^*) and, if not, its $p_1 = 0.0$ (very strongly "predisposed" against such a language). The other parameters are as for **M₀**, namely $r_1^+ = r_1^- = r_2^+ = r_2^- = r_{1,2}^+ = r_{1,2}^- = 1.0$ and $p_2 = p_{1,2} = 0.5$ initially. The language learning process subsequently adjusts these expectancies conforming to the actual language spoken around the agent. In our example, a carrier of *ASPM-D* is born expecting its language to be tonal;
- **M₂** (*Genes bias the rate of learning*): represents the case where genes bias language acquisition by coding for *preferential rates of learning*. Initially all agents have a neutral expectancy irrespective of their genomes ($p_1 = 0.5$), but the rate of adjustment of p_1 given the linguistic evidence is asymmetric. If an agent has the G_1^* allele, then it is more ready to accept that

the language is of type F_1^* than of the opposite type, meaning that evidence favoring F_1^* is accepted as stronger than equivalent evidence against F_1^* . This readiness will be denoted as the value of the bias, β , varying between 0.0 (extremely strong tendency towards F_1^*) to 1.0 (no tendency towards F_1^*). More specifically, $r_1^+ = r_2^+ = r_2^- = r_{1,2}^+ = r_{1,2}^- = 1.0$ and $p_1 = p_2 = p_{1,2} = 0.5$ initially, but $r_1^- = \beta$. In our example, a carrier of *ASPM-D* is born without any special expectancy concerning the tonality of its language, but it is more inclined to accept the data in favor of tonality than against it.

There are two parameters of interest:

- the *initial frequency* of G_1^* in the population,⁶ denoted v : it can take any value between 0.0 (total absence of the G_1^* allele from the population) to 1.0 (total absence of the alternative allele, G_1 , from the population). Due to computational costs, nine equally spaced values were considered, $v \in \{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9\}$, with the extremes of 0.0 and 1.0 excluded as uninteresting due to lack of genetic variation;
- only for model **M₂**, the *value of the bias*, denoted β . As discussed above, $\beta = r_1^-$ and can take any value between 1.0 (no bias, fully equivalent to model **M₀**) to 0.0 (extreme preference for F_1^* completely discarding any evidence to the contrary). Due to computational costs, seven values were considered based on preliminary exploratory runs, suggesting a denser sampling of weaker biases: $\beta \in \{0.1, 0.5, 0.8, 0.85, 0.9, 0.95, 0.99\}$.

It must be highlighted that F_2 and G_2 are used as controls. The initial frequency of G_2^* is not a parameter and was fixed at 0.5. The two genes are considered independent, as well as the two linguistic features. The present paper assumes that the $*$ allele is dominant. Incomplete dominance of $*$ would mean that the heterozygous phenotype is intermediate, which effectively means a weaker bias, while the recessiveness of $*$ would lower the effective frequency of the biased individuals. Therefore, the dominant case was the only one investigated.

For each of the 99 cases (all three models and possible combinations of parameter values), 20 independent runs were performed. For each run, there are two types of measures of interest: those reflecting the overall correlations between genetic diversity, linguistic diversity and geography, on one hand, and the specific correlations between the frequencies of the alleles and the frequencies of linguistic features across populations, on the other. For the first type of measures, *Mantel* (1967) correlations⁷ involving the genetic distances (Nei, 1972; considering both G_1 and G_2), linguistic distances (Euclidean distances on the space of both features F_1 and F_2 ; Dediu and Ladd, 2007) and geographic distances (Euclidean distances between regions) between populations were computed: *GenGeo* (genetic and geographic distances), *LingGeo* (linguistic and genetic distances), *GenLing* (genetic and linguistic distances) and *GenLingGeo* (genetic and linguistic distances controlling for geographic distances).

GenGeo reflects the degree to which geographic distance between populations accounts for their genetic (dis)similarity. It is usually positive due to the effects of geography on the dispersal of two populations split from a single ancestor population and contact between populations, respectively. Likewise, *LingGeo* reflects the degree to which mere geographic distance accounts

⁵ 81 runs, 27 for each of the following parameter values: (i) $n_{mother} = 0$, $n_{stranger} = 100$, $f_{foreigners} = 0.05$ and $n_{foreigner} = 0$ —the agent acquires language only from its own linguistic community; (ii) $n_{mother} = 10$, $n_{stranger} = 10$, $f_{foreigners} = 0.5$ and $n_{foreigner} = 100$ —the neighboring linguistic communities have the strongest impact on language acquisition; and (iii) $n_{mother} = 50$, $n_{stranger} = 50$, $f_{foreigners} = 0.05$ and $n_{foreigner} = 50$ —equal weighting of all three types of learning models. As expected, for case (ii) the strength of the correlations is reduced and the bias is manifest only if very strong and infrequent ($\beta \leq 0.5$, $v \leq 0.5$; see below for notations).

⁶ Given that the alleles are selectively neutral and independent, the only evolutionary process affecting their frequencies is random drift in large populations, so that these frequencies tend to remain constant during the simulations.

⁷ Computed using the ZT software (Bonnet and Van de Peer, 2002).

for the differences between the languages spoken by the two populations and concerns both historical linguistic relatedness (descent with modification from a common ancestor, given that related languages tend to inhabit neighboring regions) and language contact (linguistic borrowing across language boundaries). *GenLing* reflects the degree to which the genetic and linguistic (dis)similarities between two populations tend to correlate and it is expected that most of this correlation is explained by the subtending geography (Dediu, 2007; Poloni et al., 1997). This is so because, as discussed above, geography conditions both genetic and linguistic (dis)similarities and their residual correlation is captured by *GenLingGeo*. However, if there is a relatively strong causal biasing of language by genes, it would be expected that the residual *GenLingGeo* is larger than in the purely neutral case. Therefore, the main interest in studying these Mantel correlations is that they are widely used in the literature (e.g., Jobling et al., 2004) and they might also offer a first clue to genetic biasing for language.

The second type of measures concerns specifically the hypothesis of a causal relationship between biasing alleles and linguistic diversity and are represented by Pearson correlations between frequencies across populations: F_1F_2 (F_1^* and F_2^*), G_1G_2 (G_1^* and G_2^*), F_1G_1 (F_1^* and G_1^*), F_1G_2 (F_1^* and G_2^*), F_2G_1 (F_2^* and G_1^*) and F_2G_2 (F_2^* and G_2^*). F_1F_2 reflects the typological correlations between the two linguistic features, whereby languages tend to have certain combinations of values for these features. For example, if one takes F_1 to represent the order of Object and Verb (with two possibilities: *OV*, like Turkish, and *VO*, like Gulf Arabic; Dryer, 2008b) and F_2 as the order of Adposition and Noun Phrase (with two main possibilities: *prepositions*, like English, and *postpositions*, like Japanese; Dryer, 2008a), then these two features are strongly correlated, with 427 *OV* and *postposition* languages (41.3%) and 417 *VO* and *preposition* languages (40.3%) out of a sample of 1033 languages (Dryer, 2008c; see this also for a discussion of the explanations). However, in our case, by design there is no relationship between F_1 and F_2 ($p_{1,2} = 0.5$) and, therefore, any correlation $F_1F_2 \neq 0$ reflects particular linguistic events (splits and language contact). G_1G_2 represents the correlation between the two loci of interest but, by design, these two loci are independent so that any non-null correlation between them reflects accidental events. The correlations F_iG_j , for $i, j \in \{1, 2\}$, reflect the relationship between the linguistic feature F_i and gene G_j and form the main focus of this paper. Such a correlation can be non-null for a variety of reasons, including random drift, migrations and genetic biasing. By design, there is no causal relationship between F_i and G_j except for $i = j = 1$ for models \mathbf{M}_1 and \mathbf{M}_2 : it is expected that the correlation F_1G_1 for these models to be non-null, depending on the other relevant parameters. The interest is to identify the regions in the parameter space which produce significant and large correlations between F_1 and G_1 , which, in turn, would allow the detection of genetic biasing on language.

For all these correlations, both the effect size and significance were collected every 100 simulation-years for a period of 10,000 simulation-years, so that for each such correlation there resulted two time series: the effect sizes, r_t , and the p -values, p_t . The behavior of these correlations in time was measured by two related coefficients:

- $\rho(p_t) = \text{Card}\{p_t \leq \alpha\} / \text{Card}\{p_t\}$, where *Card* represents the number of elements in a set, and α is the chosen α -level (0.05 in this case). Thus, $\rho(p_t)$ represents the proportion of significant correlations across time;
- $\lambda(r_t, p_t)$ captures the idea that some correlation series tend to contain long, continuous stretches of significant correlations of

the same sign. Given the two series (r_t, p_t) , let us form a new time series y_t such that

$$y_t = \begin{cases} 0 & \text{if } p_t \geq \alpha \\ 1 & \text{if } p_t < \alpha \text{ and } r_t > 0 \\ -1 & \text{if } p_t < \alpha \text{ and } r_t < 0 \end{cases}$$

Given this new y_t series, let us define \bar{y} as the average length of a contiguous run of 1's or -1's in y_t , and then $\lambda(r_t, p_t) = \bar{y} / \text{Card}\{p_t \leq \alpha\}$.

$\rho(p_t)$ and $\lambda(r_t, p_t)$ tend to have the same behavior (Pearson's $r = 0.78$, $p < 0.01$; $\kappa_{0.01} = \kappa_{0.05} = 92.2\%$ ⁸). However, $\lambda(r_t, p_t)$ is better at identifying chaotic series, where correlations tend to be significant but alternate very rapidly between negative and positive values. On the other hand, $\rho(p_t)$ tends to assume more extreme values for correlation series which look different.

Besides $\rho(p_t)$ and $\lambda(r_t, p_t)$, the mean of the raw correlations, $\mu(r_t) = \bar{r}_t$, and their maximum absolute value, $M(r_t) = \max(|r_t|)$, are also considered. Their correlations with $\rho(p_t)$ and $\lambda(r_t, p_t)$ and with each other are good and highly significant ($0.46 \leq r \leq 0.65$, $p < 0.01$), and the concordances are also good ($76.7\% \leq \kappa_{0.01} \leq 85.5\%$; $82.2\% \leq \kappa_{0.05} \leq 88.9\%$).

3. Results

Due to the fact that ρ , λ , μ and M are not normally distributed, randomization techniques (independent samples t -test, one- and two-way ANOVA; Edgington, 1987) were used to compute the p -values with 10,000 permutations. Also, Holm's multiple hypotheses testing correction (Holm, 1979) was systematically applied and the reported p -values are adjusted; an α -level of 0.05 was used for significance decisions. All of the statistical analyses used R (R Development Core Team, 2007).

3.1. The no bias model

The *no bias* model, \mathbf{M}_0 , represents the baseline, generally accepted model for language–genes interaction, which assumes that the correlations between languages and genes are entirely due to shared demographic processes (Jobling et al., 2004; Cavalli-Sforza et al., 1994; Dediu, 2007; Poloni et al., 1997). In this case, we do not expect any correlations between particular genes and linguistic features (F_iG_j , $i, j \in \{1, 2\}$), between linguistic features (F_1F_2 , considered independent), or between genes (G_1G_2 , also independent). These hypotheses are supported by the results: the raw correlations for F_iG_j , F_1F_2 and G_1G_2 are normally distributed around 0.0 (mean, $\bar{x} = 0.0$ and standard deviation, $s = 0.14$), while *GenGeo*, *LingGeo* and *GenLing* are also normally distributed but positive and narrower ($\bar{x} \approx 0.2$, $s \approx 0.06$, $|\bar{x} - 0| \geq 2s$ for the first two, $\bar{x} = 0.07$, $s = 0.06$, $|\bar{x} - 0| \geq s$ for the third). Interestingly, *GenLingGeo* is normally distributed with $\bar{x} = 0.05$, $s = 0.06$, $|\bar{x} - 0| \approx 0.8s$, which seems to confirm the general finding in the literature that the correlations between genes and languages are mostly due to geography, as a consequence of demographic processes. Moreover, the four measures of *GenGeo* and *LingGeo* are very high, *GenLing* and *GenLingGeo* high, and the rest very low (Fig. 1, top panel⁹). Spatial proximity plays an important role in shaping both the genetic and linguistic diversities, especially in

⁸ Given two series of p -values, the concordance κ_α represents the percent of cases where the two series concord in their significance judgments for the considered α -level.

⁹ In the following, unless specified, only ρ will be reported, the other 3 measures behaving in a similar manner.

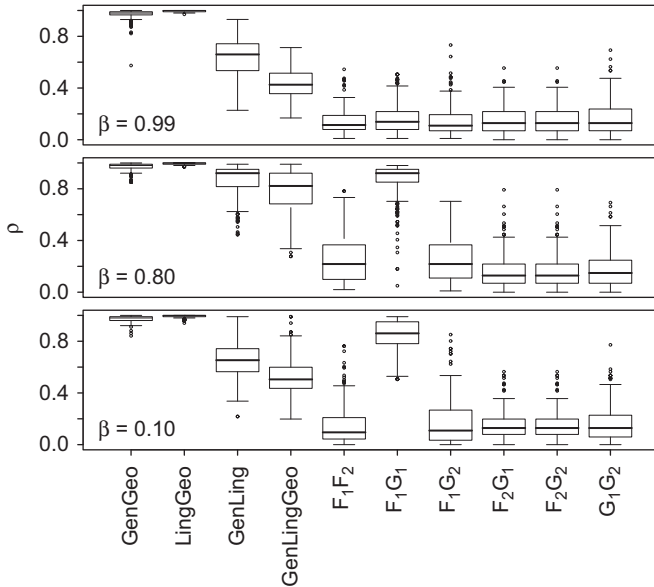


Fig. 1. Boxplots of ρ (vertical axis) for the ten measures (horizontal axis) for model \mathbf{M}_2 , when the bias, β , is very weak (0.99, top), moderate (0.80, middle) and very strong (0.10, bottom). The top panel ($\beta = 0.99$) is nearly identical to models \mathbf{M}_0 and \mathbf{M}_1 . The most important change concerns the correlation between F_1 and G_1 , which increases with increasing bias, and becomes very strong even for relatively weak biases ($\beta = 0.80$).

the linguistic case ($\rho_{GenGeo} < \rho_{LingGeo}$, $t(236.99) = -11.46$, $p = 2.49 \times 10^{-23}$), and it explains an important part of the language–genes correlation ($\rho_{GenLing} > \rho_{GenLingGeo}$, $t(410.91) = 16.90$, $p = 8.16 \times 10^{-48}$), but not all. The correlation between genetic and linguistic distances not explained by geographic distances, in the context of this model, suggests that judgments based only on partial correlations between distances must generally be taken with a grain of salt.

3.2. The initial expectation bias model

When the *initial expectation* type of bias is present (\mathbf{M}_1), the behavior of the model is overall very similar to \mathbf{M}_0 (Fig. 1, top panel), with the only exception of $M_{F_1G_1}$, which has a very interesting dependence on ν (Fig. 2). As before, spatial proximity is an important factor: $\rho_{GenGeo} < \rho_{LingGeo}$, $t(270.36) = -8.23$, $p = 9.41 \times 10^{-14}$; $\rho_{GenLing} > \rho_{GenLingGeo}$, $t(400.85) = 13.65$, $p = 5.95 \times 10^{-34}$. *GenGeo*, *GenLing* and *GenLingGeo* are not different between \mathbf{M}_0 and \mathbf{M}_1 ($t_{GenGeo}(403.97) = 0.14$, $p = 1.0$; $t_{GenLing}(387.06) = -0.48$, $p = 1.0$; $t_{GenLingGeo}(403.88) = -0.63$, $p = 1.0$), but *LingGeo* is ($t_{LingGeo}(336.05) = 3.49$, $p = 0.006$). For F_1G_1 , only M picks up the signature of this bias, showing a very strong dependency on ν (Fig. 2). A look at typical runs for different values of ν (Fig. 3) reveals that it is the initial time snapshot which is picked up by M , after which F_1G_1 drop rapidly. It seems, therefore, that this type of bias is easily swamped by linguistic change and is effective only for the first few generations. However, for these first generations it is very strong (correlations as high as 1.0 for $\nu = 0.1$), depending on the initial frequency of the biasing allele, ν . It can be concluded that the *initial expectation* bias, while impacting on the population's language, does not represent a plausible implementation of a linguistic genetic bias.

3.3. The rate of learning bias model

For the *rate of learning* type of bias (\mathbf{M}_2), the behavior depends on the strength of the bias, β (Fig. 1). When the bias is *extremely*

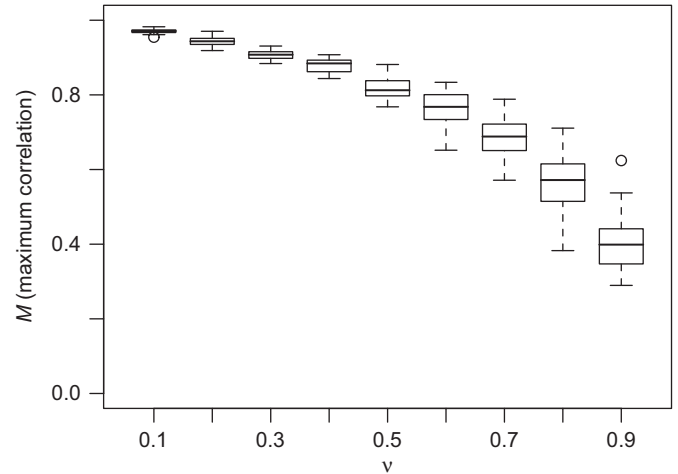


Fig. 2. The behavior of M , the biggest correlation irrespective of sign, (vertical axis) for F_1G_1 function of ν , the initial frequency of G_1^* in the population, (horizontal axis) in the case of \mathbf{M}_1 . For low ν , the maximum value of the correlations (M) is close to 1, but decreases with increasing ν .

weak, $\beta = 0.99$, the system is very similar to \mathbf{M}_0 : $t_{GenGeo}(329.58) = 1.76$, $p = 0.78$; $t_{LingGeo}(411.64) = 0.43$, $p = 1.0$; $t_{GenLing}(410.36) = 0.04$, $p = 1.0$; $t_{GenLingGeo}(411.42) = -0.22$, $p = 1.0$. When the bias is *extremely strong*, $\beta = 0.10$, its influence on the language is obvious, stable and specific (see F_1G_1 in Fig. 1): $t_{F_1G_1/F_1G_2}(343.77) = 44.48$, $p = 1.78 \times 10^{-143}$, $t_{F_1G_1/F_2G_1}(410.28) = 63.76$, $p = 9.06 \times 10^{-214}$. The distribution of the raw correlations between F_1 and G_1 depends on ν , moving from a strongly right skewed distribution with median $\bar{x} = 0.57$, $\bar{x} \approx 0.5$, $s \approx 0.16$, $|\bar{x} - 0| \geq 3s$ for $\nu = 0.1$ to another right skewed distribution with $\bar{x} = 0.37$, $\bar{x} \approx 0.4$, $s \approx 0.2$, $|\bar{x} - 0| \geq 1.5s$ for $\nu = 0.9$, going through intermediate stages of bimodality; this reflects the effects of random drift on ν and F_1 . For *intermediate* biases, the behavior of the system varies smoothly between these two extremes, with $\beta = 0.95$ more similar to $\beta = 0.99$ (and \mathbf{M}_0) and with $\beta = 0.85$ more similar to $\beta = 0.10$, suggesting that the bias starts to become manifest at $\beta \approx 0.90$. From Fig. 4, it can be seen that the effects of this type of bias depend on both the bias strength, β , and the initial frequency of the biasing allele, ν , but that even relatively weak biases ($0.85 \leq \beta \leq 0.95$) have detectable effects for certain values of ν (around 0.3), while strong biases ($0.10 \leq \beta \leq 0.50$) are detectable for any value of ν .

3.4. The behavior across models and initial frequencies

The behavior of the measures ρ , λ , μ and M relative to the initial frequency of G_1^* in the population, ν (nine levels, $\nu \in \{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9\}$) and the nine models (\mathbf{M}_0 , \mathbf{M}_1 and \mathbf{M}_2 with $\beta \in \{0.1, 0.5, 0.8, 0.85, 0.9, 0.95, 0.99\}$), was investigated using a two-way independent randomization ANOVA with 10,000 randomizations (Edgington, 1987). For the first factor¹⁰ (the initial population frequency, ν), further one-way randomization ANOVAs were conducted for each model separately with multiple comparisons corrections.

Inside models, *GenGeo*, *LingGeo*, *GenLing* and *GenLingGeo* generally depend on the *first factor*, ν (in the shape of a more or

¹⁰ Due to the large number of tests performed and space constraints, the p -values and test statistics were reported only for the most relevant cases. Moreover, due to the systematic application of Holm's multiple hypotheses testing correction (Holm, 1979) and the large number of tests performed, most adjusted p -values have collapsed to the extreme values (0.0 and 1.0). All results reported as significant are so for adjusted p -values less than 0.05.

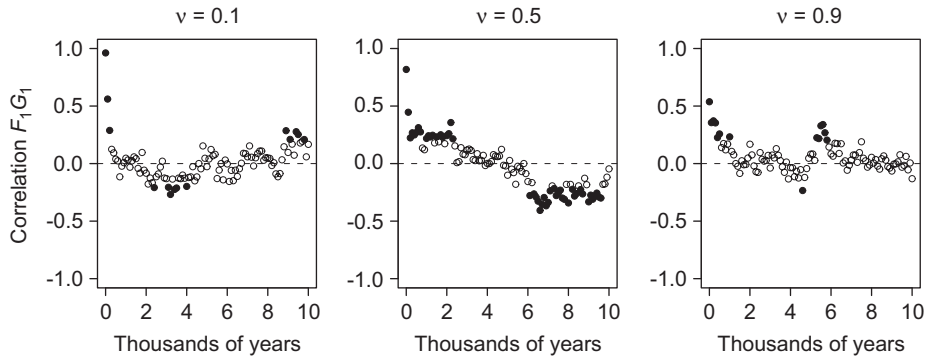


Fig. 3. The value of the correlation between F_1 and G_1 (vertical axis) function of simulation time in thousand years (horizontal axis) for model M_1 , in typical runs for $v = 0.1$ (left), $v = 0.5$ (center) and $v = 0.9$ (right). Dashed line = 0.0. Black circles (\bullet) = correlations significant at $\alpha = 0.05$, white circles (\circ) = non-significant correlations.

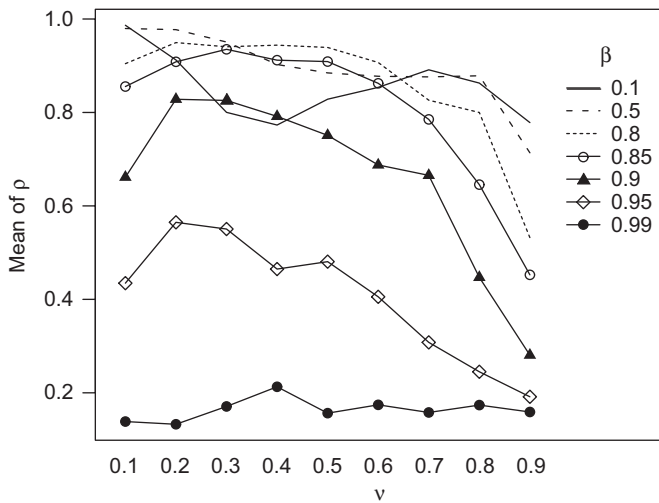


Fig. 4. The mean of ρ (vertical axes) for the correlation between F_1 and G_1 function of v (horizontal axes) and β (curves) for model M_2 . The curves for β are, from rightmost top to bottom: 0.1, 0.5, 0.8, 0.85, 0.9, 0.95 and 0.99. ρ is larger for stronger biases and there is an interaction between β and v .

less flat inverted U, “ \sim ”, with *GenGeo* and *LingGeo* very high, followed by *GenLing* and *GenLingGeo* (see also Fig. 1). Across models (*second factor*), *GenGeo* is constantly high and does not significantly differ between models, as expected, reflecting the fact that the distribution of genetic diversity is influenced only by demography, there being no causal feedback from languages to genes in the simulation. Such a feedback could be implemented as assortative mating on linguistic criteria or as linguistic group selection, but it would have added an extra level of complexity to an already complex system. However, *LingGeo*, *GenLing* and *GenLingGeo* do differ between models, mostly due to the differences between M_0 and M_1 , M_0 and M_2 ($\beta < 0.95$), M_1 and M_2 , and M_2 (small bias) and M_2 (large bias). This reflects the effects that the genetic biasing of language (F_1G_1) has on the global relationships between genetic and linguistic diversities and geographic distances. Therefore, it might be possible to devise statistical tests based on global indicators of genetic and linguistic diversity able to suggest cases where biasing effects might be at work, but this requires further study and more appropriate null models.

The other correlation of interest, F_1G_1 (see Figs. 1 and 4), depends very strongly on the model (*second factor*) with essentially all pairs of conditions being different, as expected. For M_0 it does not depend on v (*first factor*) while for M_1 , its dependency on v is picked up only by M (see Section 2.2 and

Table 1
Multiple regressions of F_1G_1 on β and v , for M_2

IV	R^2	I	β	β^2	v	v^2
ρ	0.62	0.63	2.08	-2.39	0.49	-0.79
$\sqrt{\lambda}$	0.64	0.69	1.25	-1.75	0.27	-0.50
μ	0.60	0.39	0.34	-0.49	0.93	-1.15
M	0.68	0.66	0.76	-1.06	0.33	-0.47

The regression $IV \sim I + \beta + \beta^2 + v + v^2$.

Fig. 2. For M_2 , its behavior depends on the strength of the bias, β : for an extremely weak bias ($\beta = 0.99$), it is indistinguishable from M_0 , as expected, while for stronger biases ($\beta \leq 0.95$) F_1G_1 depends on v and generally increases in strength with stronger β .

3.5. The behavior for the rate of learning model

Specifically for the *rate of learning bias* model, M_2 , the effects of the bias strength, β , and the initial population frequency of the biasing allele, v , were investigated using a two-way independent randomization ANOVA. It was found that the correlation between genetics and geography, *GenGeo*, depends only on v , not being affected by the strength of the bias, as expected. However, *LingGeo*, *GenLing* and *GenLingGeo* depend on both factors, which interact (except for ρ in the case of *LingGeo*, which depends only on β). The correlation between F_1 and G_1 depends on both factors, which also interact (see Fig. 4): it is important to note that the values of the four measures tend to increase with increasing strength of the bias (lower β).

In order to quantitatively understand the dependency of the correlation between F_1 and G_1 on the strength of the bias, β , and the initial population frequency of the biasing allele, v , in the case of M_2 , multiple regressions of the four measures, ρ , λ , μ and M , on β and v were conducted. In all cases, the best-fitting models were quadratic in both dependent variables (λ was transformed by applying square root). The results are in Table 1 (all coefficients are significant at $p < 0.001$; R^2 are adjusted and significant at $p < 2.2 \times 10^{-16}$). The ratio of number of cases to number of dependent variables is very large (between 241.5 and 483, depending on the measure) and the skewness and kurtosis of the independent variables are within acceptable limits. Also, the examination of the residuals reveals moderate deviations from normality, nonlinearity and heteroscedasticity (Tabachnick and Fidell, 2001). The proportion of the explained variance is large (adjusted $R^2 \geq 0.60$). Focusing on ρ , the regression equation

$$\rho \approx 0.63 + 2.08\beta - 2.39\beta^2 + 0.49v - 0.79v^2$$

Table 2
Maxima for F_1G_1 function of β and ν , for M_2

Measure	Maximum	β_{max}	ν_{max}
ρ	1.16	0.44	0.31
$\sqrt{\lambda}$	0.95	0.36	0.27
μ	0.64	0.35	0.40
M	0.85	0.36	0.35

predicts that there is a unique maximum $\rho_{max} = 1.16$ for $\beta_{max} = 0.44$ and $\nu_{max} = 0.31$ (Table 2). While the actual values have large errors (for example, $\rho \leq 1.0$ by definition), the suggestion that there is a region around $\beta \approx 0.4$ and $\nu \approx 0.3$ where F_1 and G_1 correlate strongly, seems warranted. Moreover, these exact numeric values will depend on the actual model parameters (especially optimal population size and the language sampling during learning, with larger populations and stronger foreigner influence tending to mask the impact of the bias), but test runs have suggested that this behavior remains qualitatively the same. Therefore, this region in the parameter space maximizing the effects of the genetic bias seems to be optimal for the detection of *rate of learning* biases.

4. Discussion

The existence of causal correlations between inter-population genetic and linguistic diversities of the type suggested by Dediu and Ladd (2007) is potentially very important for a better understanding of the biological bases of language as well as the evolution of language and linguistic diversity (for a detailed discussion of these issues in the context of biolinguistics see Ladd et al., 2008). A very convincing support for the fact that learning biases can affect the outcome of trans-generational learning is provided, for example, by Feher et al. (2008), which reared in social and acoustic isolation song-learning male zebra finches, resulting in highly abnormal songs. Subsequently, in an iterated learning paradigm, they used these birds as models for a second generation of male birds, which were used in turn as models for the next generation of male birds, and so on. They report that the “changes in acoustic structure appeared to be directional and gradual, when observed over generations” and that “[b]y the seventh clutch, the song was indistinguishable from normal zebra finch song” (p. 424), meaning that individually very small biases recover the normal song through cultural transmission. Also, Ladd et al. (2008) discuss a specifically linguistic suggestion made by Ladefoged (1984), which compared the formant frequencies of otherwise identical 7-vowel systems of Yoruba and Italian and attributes these subtle differences to the differences in the vocal tract anatomy between the two populations, biasing their languages across generations.

However, the exact definition of what is meant by a *genetically influenced linguistic bias* is far from clear, even if intuitively this concept seem unproblematic (Hawkey, 2008). The recent Bayesian approaches to biased iterated language learning (Griffiths and Kalish, 2007; Kirby et al., 2007; Smith and Kirby, 2008), while interesting and elegant, propose a not-so-satisfying account of learning biases as representing the prior distribution over languages. In his critique, Hawkey (2008) suggests a possible classification of learning biases into *transformational biases*, affecting the outcome of learning towards the preferred variant and *biased processes*, like the *default strategy* and the *ease of learning biases*.

In this context, the present computational model suggests that, when realistic demographic, genetic and linguistic processes are

considered, the type of genetically-based linguistic bias postulated to explain the correlation between the derived haplogroups of *ASPM*, *Microcephalin* and linguistic tone (Dediu and Ladd, 2007) represents a valid mechanism shaping linguistic diversity. When no genetic bias is present, the model correctly generates the known type of correlations between genetic and linguistic diversities due to demographic processes (Jobling et al., 2004; Cavalli-Sforza et al., 1994; Dediu, 2007; Poloni et al., 1997). An *initial expectation* type of bias (akin to a “default strategy” in Hawkey’s, 2008 classification), whereby carriers are born expecting a certain linguistic state which biases the language acquisition process by changing the learner’s starting point, does not seem to be able to stably influence linguistic diversity, being easily swamped by the purely cultural transmission of language. On the other hand, a *rate of learning* type of bias (similar to a Hawkey’s, 2008 “ease of learning”), whereby carriers are born with different propensities for learning different linguistic states, can reliably link the linguistic and genetic diversities. This link is highly specific and strong for a large range of bias strengths and population frequencies of the biasing allele, which makes it possible to detect using currently available statistical methods.

These findings suggest that the hypothesis of a genetically-based linguistic bias influencing the trajectory of language change through cultural transmission in populations is supported, when a specific type of genetic bias is present (*rate of learning*). This genetic bias can be very small at the individual level and the biasing allele rare at the population level but its effects can still be amplified by cultural transmission and made manifest at the inter-population level; from a practical point of view, these results suggest that the statistical methods developed in Dediu and Ladd (2007) can be used to discover such genetic biases. However, the present model is agnostic as concerns the proximate mechanisms through which such a bias could influence language change and it is expected that various such mechanisms would be involved in different cases (sensorial, neuro-cognitive, etc.). Moreover, the model highlights the importance of cultural transmission in amplifying or swamping the effects of such biases, making any deterministic interpretations implausible. However, the present model considers only a limited set of first language learning biases, and its future extensions must also consider the effects of production and second language learning biases.

Acknowledgments

The author thanks J. Hurford, D.R. Ladd, D. Hawkey, A. Dima, M. Dowman, K. Smith, S. Kirby, M. Cysouw and an anonymous reviewer for discussions, comments and suggestions. The author was funded by a Development Trust Research Fund grant from the University of Edinburgh and an ESRC (UK) postdoctoral fellowship award.

References

- Bartley, A., Jones, D., Weinberger, D., 1997. Genetic variability of human brain size and cortical gyral patterns. *Brain* 120 (Pt. 2), 257–269.
- Bishop, D., 2003. Genetic and environmental risks for specific language impairment in children. *Int. J. Pediatric Otorhinolaryngol.* 67S1, S143–S157.
- Bonneau, D., Verny, C., Uzé, J., 2004. Genetics of specific language impairments. *Arch. Pediatr.* 11, 1213–1216.
- Bonnet, E., Van de Peer, Y., 2002. Zt: a software tool for simple and partial mantel tests. *J. Statist. Software* 7, 1–12.
- Boyd, R., Richerson, P., 1985. *Culture and the Evolutionary Process*. University of Chicago Press, Chicago, IL.
- Campbell, L., 2004. *Historical Linguistics: An Introduction*. Edinburgh University Press, Edinburgh.
- Cavalli-Sforza, L., Menozzi, P., Piazza, A., 1994. *The History and Geography of Human Genes*, Abridged paperback edition. Princeton University Press, Princeton.

- Croft, W., 2000. Explaining Language Change: An Evolutionary Approach. Pearson Education Limited, Harlow, England.
- Dediu, D., 2007. Non-spurious correlations between genetic and linguistic diversities in the context of human evolution. Ph.D. Thesis, The University of Edinburgh, Linguistics and English Language.
- Dediu, D., Ladd, D.R., 2007. Linguistic tone is related to the population frequency of the adaptive haplogroups of two brain size genes, ASPM and Microcephalin. *Proc. Natl Acad. Sci. USA* 104 (26), 10944–10949.
- Dryer, M.S., 2008a. Order of adposition and noun phrase. In: Haspelmath, M., Dryer, M.S., Gil, D., Comrie, B. (Eds.), *The World Atlas of Language Structures Online*. Max Planck Digital Library, Munich (Chapter 85) (<http://wals.info/feature/85>).
- Dryer, M.S., 2008b. Order of object and verb. In: Haspelmath, M., Dryer, M.S., Gil, D., Comrie, B. (Eds.), *The World Atlas of Language Structures Online*. Max Planck Digital Library, Munich (Chapter 83) (<http://wals.info/feature/83>).
- Dryer, M.S., 2008c. Relationship between the order of object and verb and the order of adposition and noun phrase. In: Haspelmath, M., Dryer, M.S., Gil, D., Comrie, B. (Eds.), *The World Atlas of Language Structures Online*. Max Planck Digital Library, Munich (Chapter 95) (<http://wals.info/feature/95>).
- Edgington, E.S., 1987. Randomization Tests, second ed. Marcel Dekker Inc., NY.
- Feher, O., Mitra, P.P., Sasahara, K., Tchernichovski, O., 2008. Evolution of song culture in the zebra finch. In: Smith, A.D., Smith, K., Ferrer i Cancho, R. (Eds.), *The Evolution of Language*. World Scientific Publishing, Singapore, pp. 423–424.
- Felsenfeld, S., 2002. Finding susceptibility genes for developmental disorders of speech: the long and winding road. *J. Commun. Disord.* 35 (4), 329–345.
- Fisher, S.E., Lai, C.S., Monaco, A.P., 2003. Deciphering the genetic basis of speech and language disorders. *Annu. Rev. Neurosci.* 26, 57–80.
- Griffiths, T., Kalish, M., 2007. Language evolution by iterated learning with Bayesian agents. *Cognitive Sci.* 31 (3), 441–480.
- Hawkey, D.J., 2008. What impact do learning biases have on linguistic structures? In: Smith, A.D., Smith, K., Ferrer i Cancho, R. (Eds.), *The Evolution of Language*. World Scientific Publishing, Singapore, pp. 155–162.
- Holm, S., 1979. A simple sequentially rejective multiple test procedure. *Scand. J. Statist.* 6, 65–70.
- Jobling, M.A., Hurles, M., Tyler-Smith, C., 2004. *Human Evolutionary Genetics: Origins, Peoples and Disease*. Garland Science, NY.
- Kirby, S., Dowman, M., Griffiths, T.L., 2007. Innateness and culture in the evolution of language. *Proc. Natl Acad. Sci. USA* 104 (12), 5241–5245.
- Ladd, D.R., Dediu, D., Kinsella, A.R., 2008. Languages and genes: reflections on biolinguistics and the nature-nurture question. *Biolinguistics* 2 (1), 114–126.
- Ladefoged, P., 1984. 'Out of chaos comes order': physical, biological, and structural patterns in phonetics. In: Van den Broecke, M., Cohen, A. (Eds.), *Proceedings of the Tenth International Congress of Phonetic Sciences*, vol. IIB. Foris Publications, Dordrecht, Holland, pp. 83–95.
- Lenroot, R.K., Schmitt, J.E., Ordaz, S.J., Wallace, G.L., Neale, M.C., Lerch, J.P., Kendler, K.S., Evans, A.C., Giedd, J.N., 2007. Differences in genetic and environmental influences on the human cerebral cortex associated with development during childhood and adolescence. *Human Brain Mapping* 10.1002/hbm.20494.
- Mantel, N., 1967. The detection of disease clustering and a generalized regression approach. *Cancer Res.* 27 (2), 209–220.
- Nei, M., 1972. Genetic distance between populations. *Am. Naturalist* 106, 283–292.
- Nettle, D., 1999. Using social impact theory to simulate language change. *Lingua* 108, 95–117.
- Ostler, N., 2005. *Empires of the Word: A Language History of the World*. Harper Collins Publishers, London.
- Plomin, R., Kovas, Y., 2005. Generalist genes and learning disabilities. *Psychol. Bull.* 131 (4), 592–617.
- Poloni, E., Semino, O., Passarino, G., Santachiara-Benerecetti, A., Dupanloup, I., Langaney, A., Excoffier, L., 1997. Human genetic affinities for y-chromosome p49a,f/taqi haplotypes show strong correspondence with linguistics. *Am. J. Hum. Genet.* 61 (5), 1015–1035.
- Press, S.J., 2003. *Subjective and Objective Bayesian Statics*. Wiley Series in Probability and Statistics, second ed. Wiley.
- R Development Core Team, 2007. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0.
- Scamvougeras, A., Kigar, D.L., Jones, D., Weinberger, D.R., Witelson, S.F., 2003. Size of the human corpus callosum is genetically determined: an mri study in mono and dizygotic twins. *Neurosci. Lett.* 338 (2), 91–94.
- Smith, K., 2004. The evolution of vocabulary. *J. Theor. Biol.* 228, 127–142.
- Smith, K., Kirby, S., 2008. Natural selection for communication favours the cultural evolution of linguistic structure. In: Smith, A.D., Smith, K., Ferrer i Cancho, R. (Eds.), *The Evolution of Language*. World Scientific Publishing, Singapore, pp. 283–290.
- Stromswold, K., 2001. The heritability of language: a review and metaanalysis of twin, adoption, and linkage studies. *Language* 77, 647–723.
- Tabachnick, B., Fidell, L., 2001. *Using Multivariate Statistics*. Allyn & Bacon, Needham Heights, MA.
- Thompson, P., Cannon, T., Narr, K., van Erp, T., Poutanen, V., Huttunen, M., Lönqvist, J., Standertskjöld-Nordenstam, C., Kaprio, J., Khaledy, M., Dail, R., Zoumalan, C., Toga, A., 2001. Genetic influences on brain structure. *Nat. Neurosci.* 4 (12), 1253–1258.
- Wright, I., Sham, P., Murray, R., Weinberger, D., Bullmore, E., 2002. Genetic contributions to regional variability in human brain structure: methods and preliminary results. *Neuroimage* 17 (1), 256–271.