

Research



Cite this article: Bryden J, Wright SP, Jansen VAA. 2018 How humans transmit language: horizontal transmission matches word frequencies among peers on Twitter. *J. R. Soc. Interface* **15**: 20170738. <http://dx.doi.org/10.1098/rsif.2017.0738>

Received: 5 October 2017
Accepted: 15 January 2018

Subject Category:
Life Sciences—Mathematics interface

Subject Areas:
evolution

Keywords:
language transmission, linguistic evolution, evolution of language, Moran process, horizontal transmission, word heritability

Authors for correspondence:
John Bryden
john.bryden@rhul.ac.uk
Vincent A. A. Jansen
vincent.jansen@rhul.ac.uk

Electronic supplementary material is available online at <https://dx.doi.org/10.6084/m9.figshare.c.3986328>.

How humans transmit language: horizontal transmission matches word frequencies among peers on Twitter

John Bryden, Shaun P. Wright and Vincent A. A. Jansen

School of Biological Sciences, Royal Holloway, University of London, Egham TW20 0EX, UK

JB, 0000-0002-5301-5367; SPW, 0000-0003-0215-7530; VAAJ, 0000-0002-6518-2090

Language transmission, the passing on of language features such as words between people, is the process of inheritance that underlies linguistic evolution. To understand how language transmission works, we need a mechanistic understanding based on empirical evidence of lasting change of language usage. Here, we analysed 200 million online conversations to investigate transmission between individuals. We find that the frequency of word usage is inherited over conversations, rather than only the binary presence or absence of a word in a person's lexicon. We propose a mechanism for transmission whereby for each word someone encounters there is a chance they will use it more often. Using this mechanism, we measure that, for one word in around every hundred a person encounters, they will use that word more frequently. As more commonly used words are encountered more often, this means that it is the frequencies of words which are copied. Beyond this, our measurements indicate that this per-encounter mechanism is neutral and applies without any further distinction as to whether a word encountered in a conversation is commonly used or not. An important consequence of this is that frequencies of many words can be used in concert to observe and measure language transmission, and our results confirm this. These results indicate that our mechanism for transmission can be used to study language patterns and evolution within populations.

1. Introduction

Language use is constantly in flux and language evolution can happen at many spatial and temporal scales. Historical evidence shows how population groups experience wholesale changes in word usage and language syntax across many generations [1–5]. A broad theoretical background has been developed which explains how these large-scale and dynamic language patterns can be generated by language change at the individual level [2–12]. These studies assume that language elements are repeatedly transmitted between individuals in a population, and then use mathematical models or computer simulations to show that a macroscopic language pattern is generated from iterations of this individual behaviour. This makes it plausible that macroscopic changes follow from an accumulation of individual transmission events. However, these are ‘plausibility arguments’ [9] and most theoretical efforts to explain language evolution suffer from not having been confronted with data, and are often unverifiable [13]. The origins and mechanism of the evolution of language—arguably the most distinctive form of human behaviour—remain a mystery.

Darwin noted the similarity between biological and linguistic evolution [14]. This similarity inspired Labov [15,16] in explaining linguistic change. Although the similarity in homology of descent between the two processes is similar, in biological evolution the mechanism of descent is the transmission of genetic material. The mechanism of linguistic change is much harder to pinpoint. Of course children acquire their first language from parents or caretakers, but in a later phase children's language use diverges from that of their original, and adults change their language use, indicating transmission of language elements between

speakers [15,16]. It has been posited that words transmit like alleles [17], but evidence for this hypothesis has so far been scarce.

At an individual level, we adopt elements of our language throughout our lives. As children we acquire the majority of our language from our parents, but as we grow older we increasingly pick up language from our peers [1,15,16]. This form of cultural transmission between peers is called horizontal transmission [18]. While language acquisition early in life (known as vertical transmission) can be easily observed, the effect of horizontal transmission later on is more subtle and more difficult to detect. It has been known for several decades that word-usage patterns, as well as other linguistic variables, are imitated between interlocutors [19–23]. This imitation can be transient or reflective. This is due to people mirroring language while conversing or talking about similar conversation topics. To look for lasting changes we need to look for iterated transmission where people adopt words and use them in other conversations, which has been observed under laboratory conditions [11]. How language elements transmit in a lasting way between peers in natural situations is hard to measure, in part because there is a weak effect per conversation.

A possible clue to the mechanism of language element transmission lies in the observation that speakers often demonstrate probability matching: if different variants of a word or phoneme exist in a population, learners tend to match the frequency of these variants in their language use [16]. This indicates that the process of transmission does not just involve the adding of words to a lexicon, but the frequency with which these words are used is somehow stored and internalized.

Here, we will provide evidence of horizontal language element transmission. Our method detects lasting changes in language due to conversations between online individuals. However, to eliminate transient effects that can happen within conversations, we detect transmission by looking for changes in language sent to third parties which were not involved in the original conversations. To detect this weak signal, we need to use a large corpus of online conversations. The transmission of language elements is often assumed to be analogous to the spread of genetic traits [5]. We therefore use techniques from the toolbox developed within evolutionary biology on the interface between population genetics and linguistics [18,24]. We study horizontal language transmission by investigating the change in the use of words following exposure to the language of other people. This assumes that, beyond simply having a lexicon, we have some internal language representation which influences which words we choose and how often we use them [24]. We cannot directly observe this representation, but we can infer it from word-usage frequencies in a person's outgoing communication [16,25–27]. We will show here how it is possible to identify a change in the representation over time and then show, using advanced statistical methods, that this change happens due to conversations with another individual.

We will use a simple model for the internal representation of language which incorporates transmission of language between individuals. Because our aim is to study how word frequencies change, this highly simplified internal representation does not place any specific importance on grammar, syntax or word order. We simply treat communication as a multiset or a 'bag of words' [28]: how often a person uses a word is reflected by the number of copies of the word in their bag. Word instances received from conversation partners can occasionally replace

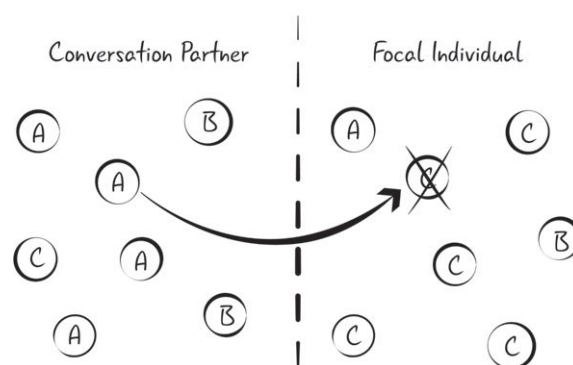


Figure 1. An osmosis-like process for horizontal language transmission used in our model. The two halves of the diagram show the internal language representations of two individuals as bags of words. The figure shows how an individual in our framework copies and stores a word from their conversation partner; an instance of word A is incorporated, replacing an instance of word C. The number of instances of a particular word defines how likely someone is to use the word in a given situation. In our model of this process, each bag contains s words; user i sends words to user j at a rate r_{ij} and the recipient replaces a randomly chosen word in their bag with a received word with incorporation rate α . Since the likelihood of a word being replaced depends on its frequency in the bag, word frequencies change similarly to osmosis in that over time the frequencies of words in both halves will tend to equilibrate.

other words in the bag, changing the internal representation and allowing the frequency of stored words to change in response to conversation (figure 1). This model forms a Moran process and can be analysed using well-understood techniques [29]. Our analysis of the model (see the electronic supplementary material) shows how the word frequencies used will equilibrate over time towards the frequencies received from conversation partners in a way that is very similar to osmosis (figure 1). The model predicts that an individual's word-usage patterns change through conversations with others and that this change will manifest itself in the word frequencies that the individual then uses to other people. Although in this model language changes in response to all language received, the effect of a conversation with a particular conversation partner will leave its mark, even if this conversation is only a relatively small part of all their conversations.

2. Results

We first show that word frequencies used by an individual change in response to the language used by a conversation partner, as predicted by our model. We studied a dataset of conversations formed from a sample of 200 million messages sent publicly between users of the Twitter web site [30] (see Methods). To eliminate any transient imitation that others have found in online communication [22,23], we excluded any mutually directed messages between a pair being studied in our analysis. Motivated by the result from our model that the difference between users is important, we looked at the influence that the difference between a focal user and their partner's early usage of a word has on any later change of the focal user's usage of the word. As this is mathematically related to the heritability of genetic traits [31], we dub this *word heritability*. Over the 1000 words tested (see Methods), we found that mean word heritability was significantly greater for pairs of users that had sent each other messages than for

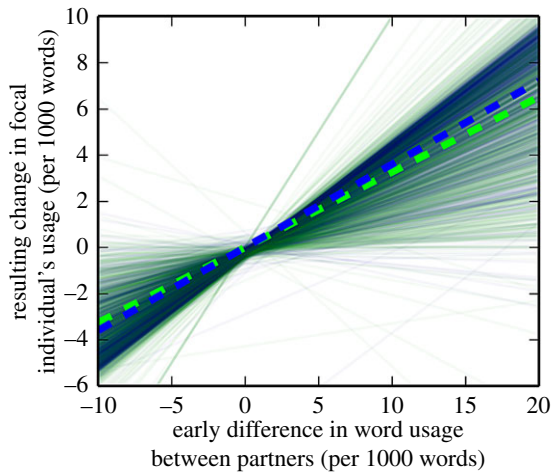


Figure 2. Word heritability between conversing partners is greater than that for non-conversing partners. For each test word, we plot regressions (see Methods) for data from conversing partners (blue solid lines) and non-conversing partners (green solid lines). The regression lines were superimposed by translucently plotting lines for each regression, interleaving between the two datasets. We found relatively high levels of word heritability in non-conversing partners due to word usage changing at population levels. A Mann–Witney U -test indicated that the slopes for conversing partners tend to be steeper than those for non-conversing partners ($p_{\text{MW}} < 9.5 \times 10^{-10}$). The two dashed lines (same colours) are slopes regressed over data collected for all of the words; the difference between these values was $W = 0.0340$, which is a measurement of word heritability due to Twitter conversations. We tested that $W > 0$ using a bootstrap ($p_B < 0.001$, see Methods).

control pairs that had not (figure 2). This indicates that an individual changes their word usage towards that used by their conversation partner.

Within our model, when a focal individual encounters word instances used by another individual, a proportion of these incoming word instances will be incorporated, replacing word instances within the focal individual's internal representation. We dub the proportion of word instances incorporated as the *incorporation rate* (α), and have developed a method to measure this rate. To do this, we implemented the model as a stochastic process. Focusing on an individual's usage of a word, we maintain a probability distribution of the word's frequency in the bag of words. We update this distribution with input received by the user according to the incorporation rate α , and then optimize α to maximize the likelihood of our observed frequencies of that word produced by the user (see the electronic supplementary material §4 for precise details). We tested 1000 different words (see Methods) and found the most likely value of α for each word.

It is important to find out if the incorporation rate of a word is dependent in any way on the frequency of usage of a word [32]. If the relationship is neutral, then studies of language change can make measurements over many words in concert. Given the heavy tailed distributions of word usage characterized by Zipf's law, one might expect that instances of more commonly used words are more likely to be incorporated than those less commonly used. Interestingly, we found that the rate of a word instance being taken up in our model is independent of word frequency across a wide range of word frequencies (figure 3). This indicates that we are as likely to adopt an instance of a frequent word as much as we are to adopt an instance of an infrequent (and therefore conversation specific) word. This suggests that we have found a perspective

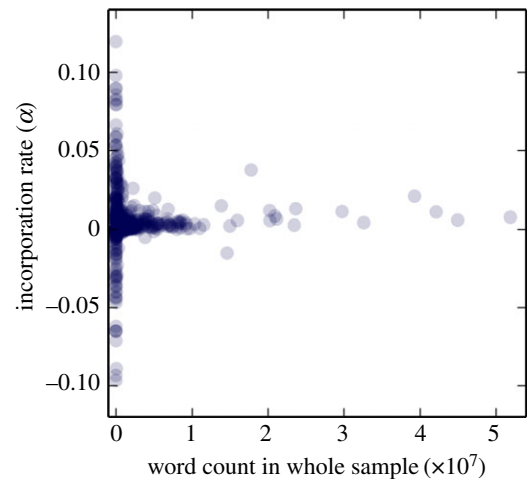


Figure 3. The rates with which words are incorporated is independent of usage frequency. Each circle is a word's incorporation rate (circles have translucency of 30%). Linear regression finds no correlation between a word's usage count (in our whole sample) and the incorporation rate (two-tailed Pearson correlation coefficient: $r^2 = 0.00040$, $p = 0.54$). The mean value of the word incorporation rate α is 0.0043, which we found to be significantly greater than zero ($p = 0.0083$, bootstrapping with 10 000 resamples of 100 values, and calculating the proportion of resamples with mean greater than zero). The high variance for very low frequencies is due to sampling effects. (Online version in colour.)

whereby word transmission is a neutral process; a view consistent with some models that generate the heavy tailed distributions of word frequencies predicted by Zipf's law [29].

Our finding that the incorporation rate of a word is not dependent on the word's usage frequency means that we can study transmission of many words in concert. We can therefore investigate the prediction, by our model (electronic supplementary material, equation S1 in §3) and others which use a Moran process [33], that the frequencies of usage of two communicating individuals will converge exponentially over time. We did this by investigating if the Bray–Curtis similarity [34] of pairs of users increases over time according to the number of messages sent between the two users. We found a highly significant, positive correlation between the change in the proportion of word instances shared between two users and the number of messages sent between them, as well as a close quantitative fit with our model (see Methods) and the data (figure 4). We tested our transmission model against a null model ($\alpha = 0$) using the Akaike information criterion (AIC), finding essentially no support for the null model compared with the transmission model (see Methods [35]). The value of the word incorporation rate, α , found was 0.01, a similar order of magnitude to the mean incorporation rate found in figure 3. These measurements indicate that we incorporate approximately one in every 100–200 words that we experience.

3. Discussion

Our results demonstrate that humans adopt lasting changes in their language usage upon conversation. These changes are consistent with the existence of an internal representation of word frequencies, where words are incorporated in a Moran process. We found that the per-encounter rate at which words are incorporated is independent of how commonly the word is used. We also found that this per-encounter rate is

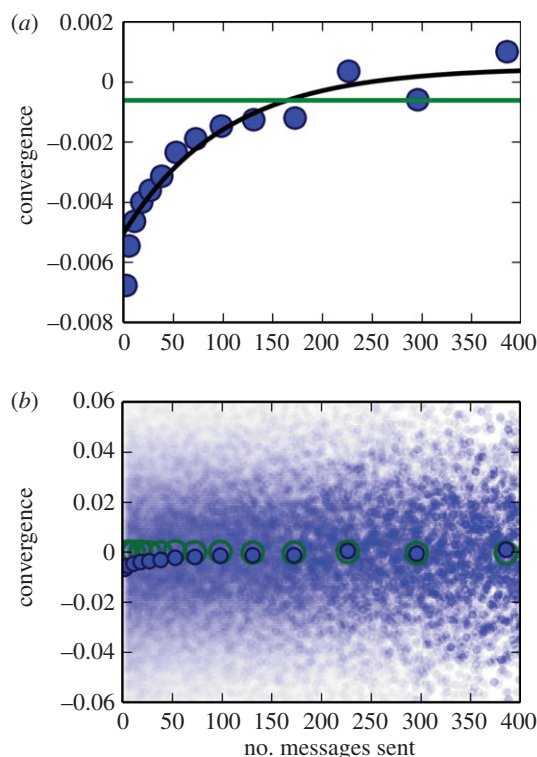


Figure 4. The more messages were sent between two users, the more their language converged. (a) Plot of the means of bins of conversation pairs (binned along the x -axis showing x , y means of each bin) and fitted models (black line is the transmission model, green horizontal line is the null model; see Methods). The fitted line of our model crosses zero at approximately 310 messages sent. (b) Illustration of the large variance in the data (unbordered translucent circles which are superimposed). The convergence of 500 conversation pairs (sampled with replacement) are plotted per bin on the x -axis (bordered blue circles). Control values are also shown (bordered green circles).

greater than zero, rejecting the null model where the per-encounter rate is equal to zero. This means that we have developed a method whereby transmission can be detected and measured on changes of individual word frequencies, or many words in concert. Put together, this means that the more two individuals converse, the more they will use similar language outside their conversations. A corollary of this is that the word usage of two isolated, or weakly connected, groups will drift apart on this time scale.

The use of large quantities of data, gleaned from online conversations, allows us to detect evidence for an underlying process of language transmission. Through identifying this process, we fill a gap in our understanding of how language is shaped and evolves [5,10,13,25,36]. We demonstrate a process which has subtle effects at the individual level (figure 4b). However, when this process is iterated many times within a population, large-scale social patterns can develop. For instance, it follows from our results that groups which interact more with one another will share similar and distinctive language patterns, which is borne out by evidence from online conversations [30]. The relatively high level of word heritability among non-conversing partners (figure 2) indicates that iterated transmission happens at a large scale in populations, which may explain increased regularization of language found among larger populations [37–39] while smaller populations are the most susceptible to language change [40,41]. Furthermore, our model and methods can be extended to infer dates which correspond to changes in

language usage of groups. This approach can be further used to study dynamical changes in population structure, and where possible link these changes to genetic changes, especially regarding whether groups have become more integrated, more isolated or have changed in size [3,42–50].

The process of transmission demonstrated here, being peer to peer in nature, forms a basis for horizontal transmission [18]. Indeed, our results reject a model that human language use can solely be explained by vertical transmission as we have shown that horizontal transmission does take place. Furthermore, the mechanism of lasting transmission we have identified can go beyond horizontal transmission and may underlie vertical transmission, whereby children acquire vocabularies from their parents, and oblique transmission, whereby children acquire vocabularies from older generations. From this perspective, we propose that vertical transmission can work in much the same way as horizontal transmission but with an inequality between parents and children whereby parents are much less likely to pick up words from their children than vice versa. With an understanding of both forms of transmission, the model and evidence that we have presented can be applied to understand how word frequencies can change across several generations of a population.

Language transmission is a cognitive process with an underlying neurological mechanism. Our evidence that word frequencies are transmitted from person to person points to insights which can inform neuroscience about the sorts of brain structures, mechanisms and memory that are necessary for language uptake and storage, and may be awaiting discovery. For example, an internal, mutable representation of word frequency suggests a reinforcement process and directs neuroscientists towards plasticity theories; a conclusion supported by various studies showing a role for plasticity and/or Hebbian learning in language therapy [51], acquisition [52] and processing [53,54].

There are no genes for words, or other specific language features, yet languages change in a way that is very reminiscent of biological evolution. These similarities to biological evolution suggest that within language evolution there is an analogous unit to the gene, even if we do not know what this unit is. Here we shed some light on the nature of this unit by showing how word frequencies can be stored and passed on. We argue that word frequencies can be passed on vertically, horizontally or obliquely. This forms a quantifiable basis for studying descent with modification of language: a requirement for language evolution.

4. Methods

4.1. Data acquisition

We used conversations between users recorded on the social networking site Twitter. Online conversations on social networks allow the observation of natural, everyday language within its social context in a way that more formal, written media does not. The informal style of this language, and its short, back-and-forth nature, makes it much closer in form and appearance to spoken language than most other forms of written language. Communication on Twitter replicates the heterogeneity in usage that is found in spoken language [12,23,30]. The ubiquity of the use of online social media for human interaction allows the gathering of these data at a large scale and in quantities that are not normally achieved for spoken language. While there are likely

to be differences, Twitter conversations are more like regular conversations than other, written forms of communication.

The data were recorded from the Twitter website during December 2009. A snowball sampling process was used to gather users as follows: for each user sampled, all their tweets that mentioned other users (using the '@' symbol) were collected directly from their profiles, meaning that we expect to have recorded a full history of their tweets at the time they were sampled. Any newly referenced users were added to a list of users from which the next user to be sampled was picked. Starting from a random user, conversational tweets (time-stamped between January 2007 and November 2009) were sampled, yielding over 200 million messages from over 189 000 users. We ignored messages that were copies of other messages (so-called retweets, which are identified by a search for tweets beginning with 'RT').

4.2. Test words

The following methods used a list of 1000 different test words (see the electronic supplementary material). These words were selected randomly from the complete collection of all text in the sample.

4.3. Word heritability analysis

Messages in a conversation were temporally split into 'early' and 'late' halves around the median time. An 'early sample' was created by randomly sampling 1000 words from the amalgamated early tweets. This was repeated with the amalgamated late tweets to create a 'late sample'.

Word heritability was measured by regressing over a series of points: each calculated on the basis of a single given word, and a randomly shuffled pair of users. For the first axis of the regression, we recorded the difference in the first user's usage of the word compared with that of the other user during the two early halves. For the second axis, we recorded the amount which the first user changed their usage of that word over time between their early and late halves. Two regressions were plotted for each word: one for conversing partners and one for non-conversing partners.

To test for significance, we carried out a bootstrap analysis by generating two resamples of 500 K points from the conversing and non-conversing datasets and regressed a line through each sample. We then measured the difference between the two slopes and recorded the proportion (reported in the main text as p_B) of the 1000 bootstrap resamples for which the slope for non-conversing individuals exceeded the slope for conversing individuals. To test that we had used enough resample points, we confirmed that similar results could be achieved with smaller resample sizes.

In all we recorded approximately 500 million data points between conversing partners. To generate controls, we randomly

generated pairs of users and checked that they had never sent one another messages in our dataset. We used 9 million pairs for our control, which was sufficient to capture its distribution for our bootstrap analysis and for the Mann–Whitney U -test.

4.4. Convergence analysis

The convergence analysis required a method that measures language similarity between pairs of users. We used the Bray–Curtis similarity measure [34] because it takes frequency into account rather than simply binary presence/absence [55]. Words are converted to lower case and stripped of punctuation (see [55] for more information). We divided each of the two users' language into early and late time periods and sampled 1000 words (with replacement) from each time period. To measure convergence data points, we calculated the Bray–Curtis similarity between the samples from the two late time periods and subtracted the Bray–Curtis similarity between the samples from the two early time periods. For the control data points, we took early and late samples from the complete time period without division.

The transmission model fitted to the convergence data points was from the electronic supplementary material, equation (S1):

$$y = c_1 + c_2 e^{-\alpha x}.$$

The null model was with $\alpha = 0$, which was simply

$$y = c_3.$$

Fitting was done against the points sampled for display in figure 4 using a least-squares method. The values found were: $c_1 = 0.000478$, $c_2 = -0.00552$, $\alpha = 0.00982$ and $c_3 = -0.000617$. The AIC was calculated as

$$\text{AIC} = 2k - 2 \sum_i \ln[\text{pdf_norm}(y_i; \mu_i, \sigma^2)],$$

where k is the number of parameters in the model, y_i are the model predictions and μ_i are the corresponding data points, σ^2 is the variance of the data points and pdf_norm is the probability distribution function of the normal distribution. We found $\text{AIC}_{\text{transmission}} = 1\,535\,774$ and $\text{AIC}_{\text{null}} = 1\,536\,263$, which means there is essentially no support for the null model in light of the transmission model [35].

Data accessibility. Data and scripts for plotting figures have been uploaded as part of the electronic supplementary material.

Authors' contributions. All of the authors contributed equally to the work.

Competing interests. We declare we have no competing interests.

Funding. S.P.W. was supported by a Royal Holloway, University of London Reid Scholarship. J.B. was supported by the Economic and Social Research Council (grant no. ES/L000113/1).

Acknowledgement. Thanks to Yaniv Garber for artwork on Figure 1.

References

- Bloomfield L. 1933 *Language*. Chicago, IL: University of Chicago Press.
- Dunn M, Terrill A, Reesink G, Foley RA, Levinson SC. 2005 Structural phylogenetics and the reconstruction of ancient language history. *Science* **309**, 2072–2075. (doi:10.1126/science.1114615)
- Lieberman E, Michel J-B, Jackson J, Tang T, Nowak MA. 2007 Quantifying the evolutionary dynamics of language. *Nature* **449**, 713–716. (doi:10.1038/nature06137)
- Gray RD, Drummond AJ, Greenhill SJ. 2009 Language phylogenies reveal expansion pulses and pauses in Pacific settlement. *Science* **323**, 479–483. (doi:10.1126/science.1166858)
- Pagel M. 2009 Human language as a culturally transmitted replicator. *Nat. Rev. Genet.* **10**, 405–415. (doi:10.1038/nrg2560)
- Nowak MA, Komarova NL, Niyogi P. 2001 Evolution of universal grammar. *Science* **291**, 114–118. (doi:10.1126/science.291.5501.114)
- Nowak MA, Komarova NL, Niyogi P. 2002 Computational and evolutionary aspects of language. *Nature* **417**, 611–617. (doi:10.1038/nature00771)
- Steels L, Kaplan F. 2002 Aibo's first words: the social learning of language and meaning. *Evol. Commun.* **4**, 3–32. (doi:10.1075/eoc.4.1.03ste)
- Castellano C, Fortunato S, Loreto V. 2009 Statistical physics of social dynamics. *Rev. Mod. Phys.* **81**, 591–646. (doi:10.1103/RevModPhys.81.591)
- Chater N, Christiansen MH. 2010 Language acquisition meets language evolution. *Cogn. Sci.* **34**, 1131–1157. (doi:10.1111/j.1551-6709.2009.01049.x)
- Kirby S, Griffiths T, Smith K. 2014 Iterated learning and the evolution of language. *Curr. Opin. Neurobiol.* **28**, 108–114. (doi:10.1016/j.conb.2014.07.014)
- Eisenstein J, O'Connor B, Smith NA, Xing EP. 2014 Diffusion of lexical change in social media. *PLoS*

- ONE 9, e113114. (doi:10.1371/journal.pone.0113114)
13. Hauser MD, Yang C, Berwick RC, Tattersall I, Ryan MJ, Watumull J, Chomsky N, Lewontin RC. 2014 The mystery of language evolution. *Front. Psychol.* **5**, 401. (doi:10.3389/fpsyg.2014.00401)
 14. Darwin C. 1883 *The descent of man and selection in relation to sex*. London, UK: John Murray.
 15. Labov W. 2001 *Principles of linguistic change volume 2: social factors*. New York, NY: John Wiley & Sons.
 16. Labov W. 2010 *Principles of linguistic change volume 3: cognitive and cultural factors*. New York, NY: John Wiley & Sons.
 17. Realì F, Griffiths TL. 2010 Words as alleles: connecting language evolution with Bayesian learners to models of genetic drift. *Proc. R. Soc. B* **277**, 429–436. (doi:10.1098/rspb.2009.1513)
 18. Cavalli-Sforza LL, Feldman MW. 1981 *Cultural transmission and evolution: a quantitative approach*, vol. 16. Princeton, NJ: Princeton University Press.
 19. Brennan SE. 1996 Lexical entrainment in spontaneous dialog. *Proc. ISSD* **96**, 41–44.
 20. Pickering MJ, Garrod S. 2004 Toward a mechanistic psychology of dialogue. *Behav. Brain Sci.* **27**, 169–190.
 21. Gallois C, Ogay T, Giles H. 2005 Communication accommodation theory: a look back and a look ahead. In *Theorizing about intercultural communication* (ed. WB Gudykunst), pp. 121–148. Thousand Oaks, CA: Sage.
 22. Danescu-Niculescu-Mizil C, Gamon M, Dumais S. 2011 Mark my words!: linguistic style accommodation in social media. In *Proc. of the 20th Int. Conf. on World Wide Web, WWW '11*, Hyderabad, India, 28 March–1 April 2011, pp. 745–754. New York, NY: ACM.
 23. Tamburrini N, Cinnirella M, Jansen VAA, Bryden J. 2015 Twitter users change word usage according to conversation-partner social identity. *Soc. Netw.* **40**, 84–89. (doi:10.1016/j.socnet.2014.07.004)
 24. Wang WS-Y. 1976 Language change. *Ann. N. Y. Acad. Sci.* **280**, 61–72. (doi:10.1111/j.1749-6632.1976.tb25472.x)
 25. Pagel M, Atkinson QD, Meade A. 2007 Frequency of word-use predicts rates of lexical evolution throughout Indo-European history. *Nature* **449**, 717–720. (doi:10.1038/nature06176)
 26. Michel J-B *et al.* 2011 Quantitative analysis of culture using millions of digitized books. *Science* **331**, 176–182. (doi:10.1126/science.1199644)
 27. Newberry MG, Ahern CA, Clark R, Plotkin JB. 2017 Detecting evolutionary forces in language change. *Nature* **551**, 223–226. (doi:10.1038/nature24455)
 28. Salton G, McGill MJ. 1983 *Introduction to modern information retrieval*. New York, NY: McGraw-Hill.
 29. Blythe RA. 2012 Neutral evolution: a null model for language dynamics. *Adv. Complex Syst.* **15**, 1150015. (doi:10.1142/S0219525911003414)
 30. Bryden J, Funk S, Jansen VAA. 2013 Word usage mirrors community structure in the online social network Twitter. *EPJ Data Sci.* **2**, 3. (doi:10.1140/epjds15)
 31. Falconer DS, Mackay TFC. 1995 *Introduction to quantitative genetics*, 4th edn. New York, NY: Longman.
 32. Church KW. 2000 Empirical estimates of adaptation: the chance of two Noriegas is closer to P/2 than P2. In *Proc. of the 18th Conf. on Computational linguistics, Saarbrücken, Germany, 31 July–4 August 2000*, vol. 1, pp. 180–186. Stroudsburg, PA: Association for Computational Linguistics.
 33. Blythe RA, McKane AJ. 2007 Stochastic models of evolution in genetics, ecology and linguistics. *J. Stat. Mech.* **2007**, P07018. (doi:10.1088/1742-5468/2007/07/P07018)
 34. Bray JR, Curtis JT. 1957 An ordination of the upland forest communities of southern Wisconsin. *Ecol. Monogr.* **27**, 325–349. (doi:10.2307/1942268)
 35. Burnham KP, Anderson DR. 2002 *Model selection and multimodel inference: a practical information-theoretic approach*. Berlin, Germany: Springer.
 36. Croft W. 2000 *Explaining language change: an evolutionary approach*. London, UK: Pearson Education.
 37. Kam CLH, Newport EL. 2005 Regularizing unpredictable variation: the roles of adult and child learners in language formation and change. *Lang. Learn. Dev.* **1**, 151–195. (doi:10.1080/15475441.2005.9684215)
 38. Lupyán G, Dale R. 2010 Language structure is partly determined by social structure. *PLoS ONE* **5**, e8559. (doi:10.1371/journal.pone.0008559)
 39. Dale R, Lupyán G. 2012 Understanding the origins of morphological diversity: the linguistic niche hypothesis. *Adv. Complex Syst.* **15**, 1150017. (doi:10.1142/S0219525911500172)
 40. Trudgill P. 2005 Linguistic and social typology: the Austronesian migrations and phoneme inventories. *Linguist. Typol.* **8**, 305–320. (doi:10.1515/lity.2004.8.3.305)
 41. Trudgill P. 2011 Social structure and phoneme inventories. *Linguist. Typol.* **15**, 155–160. (doi:10.1515/lity.2011.010)
 42. Barbujani G, Whitehead GN, Bertorelle G, Nasidze IS. 1994 Testing hypotheses on processes of genetic and linguistic change in the Caucasus. *Hum. Biol.* **66**, 843–864.
 43. Hunley K, Long JC. 2005 Gene flow across linguistic boundaries in native North American populations. *Proc. Natl Acad. Sci. USA* **102**, 1312–1317. (doi:10.1073/pnas.0409301102)
 44. Hunley K, Cabana G, Merriwether D, Long J. 2007 A formal test of linguistic and genetic coevolution in native Central and South America. *Am. J. Phys. Anthropol.* **132**, 622–631. (doi:10.1002/ajpa.20542)
 45. Hunley K, Dunn M, Lindström E, Reesink G, Terrill A, Healy ME, Koki G, Friedlaender FR, Friedlaender JS. 2008 Genetic and linguistic coevolution in Northern Island Melanesia. *PLoS Genet.* **4**, e1000239. (doi:10.1371/journal.pgen.1000239)
 46. Kutanan W, Ghirotto S, Bertorelle G, Srithawong S, Srithongdaeng K, Pontham N, Kangwanpong D. 2014 Geography has more influence than language on maternal genetic structure of various northeastern Thai ethnicities. *J. Hum. Genet.* **59**, 512. (doi:10.1038/jhg.2014.64)
 47. Longobardi G, Ghirotto S, Guardiano C, Tassi F, Benazzo A, Ceolin A, Barbujani G. 2015 Across language families: genome diversity mirrors linguistic variation within Europe. *Am. J. Phys. Anthropol.* **157**, 630–640. (doi:10.1002/ajpa.22758)
 48. Srithawong S, Srikumool M, Pittayaporn P, Ghirotto S, Chantawannakul P, Sun J, Eisenberg A, Chakraborty R, Kutanan W. 2015 Genetic and linguistic correlation of the Kra-Dai-speaking groups in Thailand. *J. Hum. Genet.* **60**, 371–380. (doi:10.1038/jhg.2015.32)
 49. Creanza N, Ruhlen M, Pemberton TJ, Rosenberg NA, Feldman MW, Ramachandran S. 2015 A comparison of worldwide phonemic and genetic variation in human populations. *Proc. Natl Acad. Sci. USA* **112**, 1265–1272. (doi:10.1073/pnas.1424033112)
 50. Karafet TM *et al.* 2016 Coevolution of genes and languages and high levels of population structure among the highland populations of Daghestan. *J. Hum. Genet.* **61**, 181. (doi:10.1038/jhg.2015.132)
 51. Sarasso S, Määttä S, Ferrarelli F, Poryazova R, Tononi G, Small SL. 2014 Plastic changes following imitation-based speech and language therapy for aphasia: a high-density sleep EEG study. *Neurorehabil. Neural Repair* **28**, 129–138. (doi:10.1177/1545968313498651)
 52. Kim KHS, Relkin NR, Lee K-M, Hirsch J. 1997 Distinct cortical areas associated with native and second languages. *Nature* **388**, 171–174. (doi:10.1038/40623)
 53. Chee MW, Hon NH, Caplan D, Lee HL, Goh J. 2002 Frequency of concrete words modulates prefrontal activation during semantic judgments. *Neuroimage* **16**, 259–268. (doi:10.1006/nimg.2002.1061)
 54. Wennekers T, Garagnani M, Pulvermueller F. 2006 Language models based on Hebbian cell assemblies. *J. Physiol. Paris* **100**, 16–30. (doi:10.1016/j.jphysparis.2006.09.007)
 55. Wright S. 2017 Tuning in to terrorist signals. PhD thesis, Royal Holloway, University of London, Egham, UK.