# Detecting Regular Sound Changes

# in Linguistics

# as Events of Concerted Evolution

**Daniel J. Hruschka, Simon Branford, Eric D. Smith, Jon Wilkins, Andrew Meade,**

**Mark Pagel, and Tanmoy Bhattacharya**

# Supplemental Information

**Table S1.  List of estimated regular changes with support in > 90 trees.**

| Vowels | | | | Vowels | | |
|--------|-----|-----|---|--------|-----|-----|
| from | to | branch | | from | to | branch |
| ɨ | ə | CHV | | ä | i | HAK |
| ɨ | i | UIG | | ä | ɛ | TAT |
| ɨ | i | UZB | | ä | e | GAGX,TRK |
| a | ɔ | UZB | | o | ǫ | UZB |
| a | ū | CHV | | o | ö | UZB |
| a | ā | SAL,TRM | | o | u | BAS,TAT (c) |
| a | ā | KHAL | | ö | ụ | UZB |
| e | ä | BAS,TAT | | ö | əʷ | CHV |
| e | ä | GAGX,KHAL,AZB,TRK | | ö | ü | BAS,TAT (c) |
| e | ä | UIG | | ü | əʷ | CHV |
| e | ä | HAK | | ü | u | UZB |
| i | ɘ | HAK | | u | əʷ | CHV |
| i | ɘ | CHV | | ā | e | UIG (a) |
| i | ɨ | SJG (b) | | ā | ä | BAS,TAT (a) |
| i | e | BAS,TAT | | ā | ō | KRG (a) |

a. rare starting sound

b. language without sound correspondence proposal in Starostin et al.

c. phoneme swap (e.g. o to u and u to o).

| Consonants | | | | Consonants | | |
|---|---|---|---|---|---|---|
| from | to | branch | | from | to | branch |
| d | d | TOF,TUV,DOLG,JAK (a) | | ž | j | BAS,TAT |
| ɣ | v | CHV | | ž | j | SAL,KRMX,GAGX,KHAL,TRM,AZB,UIG,UZB,TRK |
| ɣ | w | QUM,KLPX,BLKX,BAS,NOGX,KAZ,KRG,TAT | | ž | j | SJG (b) |
| č | d́ | ALT | | ž | j | NOGX |
| č | ś | CHV | | ž | j | QUM |
| č | š | KLPX,NOGX,KAZ | | ž | z | HAK |
| č | h | DOLG,JAK | | ž | ʒ | KRG |
| č | š | TOF,TUV,DOLG,JAK,ALT,SHR,HAK | | q | χ | SAL (b) |
| č | s | BAS | | q | G | SAL,TRM |
| b | p | SHR,HAK | | q | G | AZB |
| b | p | CHV | | q | k | TRK |
| d | č | CHV | | q | x | DOLG,JAK |
| š | ź | CHV | | q | x | HAK |
| š | s | KLPX,NOGX,KAZ | | q | x | CHV |
| š | s | DOLG,JAK | | s | h | BAS |
| š | s | SJG (b) | | ŋ | m | CHV |
| š | s | HAK | | ŋ | n | GAGX,AZB,TRK |
| h | s | JAK | | x | k | DOLG (b) |
| ž | č | TOF,TUV,DOLG,JAK,ALT,SHR,HAK | | z | ž | DOLG,JAK |
| ž | h | DOLG | | z | r | CHV |
| ž | h | JAK | | z | ð | BAS |
| ž | ś | CHV | | | | |

**Table S2  Sound correspondences compared to historical linguistic proposals**

Each pair of rows of the table records (upper row) linguists' proposals as to regular sound changes in twenty-three Turkic languages and (lower row) corresponding predictions from the model of regular sound changes. Items coloured red correspond to regular sound changes, uncoloured items refer to ancestral phonemes being retained. In each case the upper row records linguists' proposals and the lower row reports the phonemes to which the model assigns its highest probabilities. The left columns identify the proto-phonemes favoured by linguistic analysis and the corresponding phoneme favoured by the model of regular change. The rightmost column records the count of the number of proposed changes under the two approaches. A blue border around a cell denotes a change that can be attributed to one or more of the regular sound changes inferred on the phylogeny of Figure 2 (main text)

The table's column headers (left to right) are:

Model Ancestor · Linguists' proposal · Yakut · Tuvan · Tofalar · Khakass · Shor · Altai · Kirghiz · Uighur · Uzbek · Kazak · Karakalpak · Nogay · Bashkir · Tatar · Kumyk · Balkar · Karaim · Turkmen · Azerbaijanian · Gagauz · Turkish · Khalaj · Chuvash · Phonological innovations

Summary rows:

| | Phonological innovations |
|---|---|
| Total linguists' phonological innovations | 526 |
| Total model's phonological innovations | 479 |

| | Model Ancestor | Linguists' proposal | Yakut | Tuvan | Tofalar | Khakass | Shor | Altai | Kirghiz | Uighur | Uzbek | Kazak | Karakalpak | Nogay | Bashkir | Tatar | Kumyk | Balkar | Karaim | Turkmen | Azerbaijanian | Gagauz | Turkish | Khalaj | Chuvash | Phonological innovations |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | a | a | a | a | a | a | a | a | e | a | ɔ | a | a | a | a | a | a | a | a | a | a | a | a | o | 3 |
| 2 | a | | a | a | a | a | a | a | a | e | ɔ | a | a | a | a | a | a | a | a | a | a | a | a | ā | o | 4 |
| 1 | | u | u | u | u | u | u | u | u | | u | u | u | u | o | o | u | u | u | u | u | u | u | u | əw | 3 |
| 2 | u | | u | u | u | u | u | u | u | | u | u | u | u | o | o | u | u | u | u | u | u | u | u | əw | 3 |
| 1 | | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | i | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | i | ə ɘ | 5 |
| 2 | ɨ | | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | i | i | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | ɨ | | ɘ ə | 4 |
| 1 | | i | i | i | i | ə | i | i | i | i | | i | i | i | i | e | e | i | i | i | i | i | i | i | ɘ ə | 5 |
| 2 | i | | i | i | i | ə | i | i | i | i | | i | i | i | e | e | i | i | i | i | i | i | i | | ɘ a | 5 |
| 1 | | e | e | e | e | i | e | e | e | e | ä | e | a | e | e | e | i | i | e | e | e | e | ä | e | e | ä | a | 8 |
| 2 | e | | e | e | e | i | e | e | e | ä | i | e | i | e | e | ä | ɛ | e | e | e | e | ä | i | e | e | ä | e | 9 |
| 1 | | o | o | o | o | o | o | o | o | o | | ŭ | o | o | o | u | u | o | o | o | o | o | o | o | vɨ u | 5 |
| 2 | o | | o | o | o | o | o | o | o | o | | ŭ | o | o | o | u | u | o | o | o | o | o | o | o | o u | 4 |
| 1 | | ü | ü | ü | ü | ü | ü | ü | ü | ü | u | ü | ü | ü | ö | ö | ü | ü | ü | ü | ü | ü | ü | i | əw | 5 |
| 2 | ü | | ü | ü | ü | ü | ü | ü | ü | ü | ü | ü | ü | ü | ö | ö | ü | ü | ü | ü | ü | ü | ü | i | əw | 5 |
| 1 | | ö | ö | ö | ö | ö | ö | ö | ö | ö | ü | ö | ö | ö | ü | ü | ö | ö | ö | ö | ö | ö | ö | e | vi u ü | 7 |
| 2 | ö | | ö | ö | ö | ö | ö | ö | ö | ö | ü | ö | ö | ö | ü | ü | ö | ö | ö | ö | ö | ö | ö | ü | əw e ü | 7 |
| 1 | | ā | ā | a | a | a | a | a | a | a | e | a | ɔ | a | a | a | a | a | a | a | a | ā | a | a | a | āa o | 23 |
| 2 | ā | | a | a | a | a | a | a | a | e | a | ɔ | ā | a | a | a | a | a | a | a | a | a | a | ā | a | 23 |
| 1 | | ū | ū | u | u | u | u | u | u | u | | u | u | u | o | o | u | u | u | ū | u | u | u | ūu | əw | 21 |
| 2 | ū | | u | u | u | u | u | u | u | | u | u | ū | ū | u | u | u | u | u | u | u | u | u | | o | 22 |
| 1 | | ī | ī | i | i | ə | i | i | i | i | | i | i | i | e | e | i | i | i | i | i | i | ī | ī | ə ə | 23 |
| 2 | ī | | i | i | i | ə | i | i | i | i | | i | i | ī | e | ī | i | i | i | i | i | i | i | i | ə | 22 |

Total linguists' phonological innovations  108
Total model's phonological innovations  108

## Table S3. Translation of dataset transcription to International Phonetic Alphabet

| Dataset | IPA | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | v | v | ī | īiː | ū | uː |
| q | q | r | r | c | ɕ | w | w |
| a | a | n | n | ε | ε | ō | oː |
| b | b | c | e | ē | eː | h | h |
| ɨ | ɨ | m | m | g | g | ẹ | ɪ |
| k | k | d | d | ä | æ | ụ | ŭ |
| u | u | ü | y | ā | æː | ọ | ʊ |
| p | p | δ | ð | ɣ | ɣ | əw | əw |
| ɔ | ɔ | j | j | ś | ɕ | d́ | ʄ |
| i | i | z | z | č | ʧ | f | f |
| o | o | ö | ø | ǯ | ʤ | ɨː | ɨː |
| G | ɢ | ə | ə | ž | ʒ | ń | ɲ |
| ā | aː | ū | yː | ŋ | ŋ | ȫ | øː |
| x | x | š | ʃ | χ | χ | ź | ʐ |
| ə | ə | s | s | l | l | ŭ | ÿ |
| | | t | t | əw | əw | ŭ | ŭ |

**Supplemental Data**

Data: Etymologically aligned and phonetically coded data for the Turkic languages is available from the Evolution of Human Languages project at: http://starling.rinet.ru/downl.php?lan=en. The Turkic languages are contained inside the altaic.exe file (filename = turcet). This study used the first 225 words.

Table S1, related to Figure 2. List of 73 regular sound changes that were detected in at least 90% of the trees in the posterior sample. Note: The model detected an average of 74.27 regular sound changes across the set of trees in the posterior sample. The extra 1.27 events per tree (74.27-73) are regular sound changes that occur infrequently in the posterior sample, and so fail to make the list in Table S1.

See Table S1

Table S2, related to Figure 3. Comparison of inferred regular sound changes to historical linguistic inferences.

See Table S2

Table S3, related to Figure 2 and Figure 3. Translation of dataset transcription to International Phonetic Alphabet

See Table S3

**Supplemental Experimental Procedures**

Concerted evolution. Concerted evolution is a term normally used to describe the process by which copies of duplicated genes, alleles or repetitive elements within a species come to be similar at many different sites of their sequences. But concerted change might be a general feature of evolving systems, including language and culture. In language evolution, concerted or regular change affects copies of the same sound (or phoneme) that appear in different words. Linguists have suggested that this regular change occurs because to fulfill their function of word discrimination, instances of the same phoneme must be both independent occurrences and also repeated realizations of the *same* functional element within the language's sound system. The resulting changes, of many instances of one phoneme in parallel to some other phoneme, yield regular sound correspondences between pairs or groups of languages.

The complexity in an instance of concerted evolution can range from a simple covariance of a few elements, such as amino acids that jointly determine a protein fold, to patterns of sound change in language that might affect groups of phonemes, to rules for conditional change

involving many segments across a syllable or a word. We model regular sound change in a language family at the level of independent sounds (the most common kind), to demonstrate the large effect that the change model has on the interpretation of patterns in data. Some sound changes in our data set cannot be incorporated in this simple form of regular shifts, either because different instances are truly independent or because they depend on multi-phoneme conditions for which our model does not test. For instance, some kinds of change depend upon 'context' defined as other features that might or might not be present in a word (see Box and Discussion in main text). If the contextual conditions are rare, our model will likely treat these as sporadic changes. If they are common, however, our model will likely treat them as a regular change but not label the context of that change. The statistical modelling of complex concerted evolution brings many challenges of representation as well as simulation, and is left to future work.

The model assumes that a language's lexicon evolves primarily from changes to form-meaning links, from sporadic sound changes that occur independently within words, or from concerted, regular sound changes that occur across all words in a lexicon simultaneously. The model also assumes that when a language diverges into two languages, the two languages begin to change independently. These elements of the model— changes to form-meaning links, sporadic and regular sound changes, and the pattern of historical divergence—are treated as independent free parameters to be estimated. To assess how adding regular, concerted sound changes improve model fit and prediction, we fit the model in two ways: (1) permitting regular, concerted sound changes in addition to sporadic changes (henceforth, regular sound change model) and (2) permitting only sporadic changes (henceforth, sporadic sound change model). In addition, languages may lose entire words used to express a particular sense—with words of different origins used to express the original meaning—and these changes in cognacy are modeled as an independent stochastic process.

       <u>Likelihood modelling</u> We model the sporadic sound changes as a continuous-time Markov process, widely used in models of DNA or protein sequence evolution, where in place of the usual *4 × 4* or *20 × 20* matrices of nucleotide or amino acid transitions we erect a *62 × 62* sound transition rate matrix, denoted $\mathbf{Q_s}$. As in DNA sequence and protein models, $\mathbf{Q_s}$ is an instantaneous rate matrix whose rows sum to zero and whose main diagonal elements for row $i$ are equal to $-\sum_{j \neq i} q_{ij}$ .

We estimate the elements of $\mathbf{Q_s}$ from the data employing a reversible-jump Markov chain Monte Carlo (RJ-MCMC) procedure described elsewhere (*S1*) that allows the large number of potential parameters to be reduced to a potentially far smaller set of statistically distinct parameters, and without loss of statistical accuracy or prior knowledge on the part of investigators.

Regular sound changes of the general form denoting the $i^{th}$ sound changing to the $j^{th}$ ($i{\neq}j$ ) are modeled in a stochastic matrix $\mathbf{Q_r}$ that takes the form of an identity matrix with an $i^{th}$ diagonal element interchanged with the off diagonal position ($ji$). Pre-multiplication of any stochastic matrix $\mathbf{Q}$ (e.g., $P(D|Q_s,t)$) by $\mathbf{Q_r}$ is equivalent to adding all elements $q_{i1},\ q_{i2},...q_{ik}$ to the corresponding values of $q_{j1},q_{j2}...q_{jk}$, and then zeroing out the $q_{i1},q_{i2},...q_{ik}$.

We then use a different RJ-MCMC procedure to propose possible $\boldsymbol{Q_r}$ matrices in branches of the phylogenetic tree, thereby allowing regular changes to occur or not on a branch-specific basis. Candidate $\boldsymbol{Q_r}$ matrices are chosen by randomly selecting a pair of sounds ($i,j$) as above. These matrices are then proposed to land with uniform probability anywhere in the tree, allowing the model to estimate both the number of regular changes along branches and the position or timing of regular sound changes along a given branch. We do not place a prior probability distribution on the number of regular sound changes in the tree, but to be accepted in the Markov chain, new $\boldsymbol{Q_r}$ matrices must satisfy the requirements of the reversible-jump procedure.

The likelihood of the data for $n$ languages given the model ($M$) of phoneme evolution is calculated by first assigning possible reconstructions (or phoneme assignments) at each of the $n$-1 ancestral or internal nodes of the tree ($T$). Where $k$ regular sound changes have occurred along a branch, the contribution to the likelihood of this branch is computed by multiplying together the site likelihoods that are appropriate elements of the total transition matrix

$$\exp[Q_s(t_1 - t_0)]Q_{r_1} \exp[Q_s(t_2 - t_1)]Q_{r_2} ... \exp[Q_s(t_k - t_{k-1})]Q_{r_k} \exp[Q_s(t - t_k)]$$

where $t$ is the total time for the branch, and $t_i$ represents the times of the $k$ regular changes. In branches where no regular sound changes occur, the likelihood is found from $P(D \mid Q_s, t)$. Then the probabilities of these possible ancestral assignments are calculated by multiplying the likelihoods of the evolutions represented by each of the $2n$-2 branches. Finally, the total likelihood is obtained by summing over all possible reconstructions at each of the $n$-1 internal or ancestral nodes in the tree. In other words,

$$L(D \mid M, T) = \sum_{\{D_{\text{int}}\}} \int \prod_{t=1}^{2n-2} \left( P(D_t \mid D_{a(t)}, Q_s, Q_{r_1}, Q_{r_2}, ..., Q_{r_m}) dQ_{r_1} dQ_{r_2} ... dQ_{r_m} \right) dQ_s \quad \text{, where } D_t \text{ and } D_{a(t)}$$

represent the possibly reconstructed phoneme sequence at the tip and ancestor of the $t^{\text{th}}$ branch, the sum runs over all possible assignments of reconstructions, the integral over all possible assignments of regular ($Q_r$)and sporadic ($Q_s$)matrices, and the product runs over the branches in the tree. Likelihoods are combined at the root of the tree weighted by their empirical phoneme frequencies. In the above notation, all priors are implicitly set to the uniform non-informative distribution. This is the standard approach to calculating the likelihood of tree topologies from discrete data such as genes, proteins or sounds, but with the addition of integrating over the collection of regular-change matrices in addition to the usual integration over the sporadic matrix.

Parameters of the models are then estimated in the usual way [e.g., (S2)] by Markov chain Monte Carlo procedures that numerically explore the posterior distribution of model parameters, given the data and the assumed prior distributions of the parameters. In practice, to improve the rate of convergence of the Markov chains, we augmented the likelihood of the sound-change model with that obtained from cognacy data for the same set of 225 words, following methods described elsewhere (S3, S4).

This model does not assume *a priori* any pattern in sound correspondence, categories of sounds (*e.g.,* vowels and consonants), or similarity between sounds. Though it would have been straightforward to incorporate known facts, this linguistic intuition is currently not stated in a

probabilistic way. This allows us to treat the raw data as the only evidence to support or deny the standard linguistic knowledge about common sound changes.

Simplifications. While the current model incorporates regular and sporadic sound changes, it also makes several simplifying assumptions. Sound changes in different positions in the same word are treated as independent, thus ignoring conditioning contexts and vowel harmony. Second, all sporadic processes are assumed to work sufficiently independently and uniformly that they can be expressed by a single probability function across the entire family. It allows us to explore the feasibility of inferring sound correspondences and sound change with no additional linguistic prior knowledge. Third, the alignment permits metathesis, but the full model explicitly ignores this and other rare processes (*S5*), replacing them with independent idiosyncratic changes in nearby sounds. The effects of non-independent evolution on various branches, *e.g.,* due to processes like contact phenomena, are ignored. Adding such within-word correlations might help resolve the discrepancies between model predictions and linguistic descriptions of sound correspondences for uvulars and velars such as $k$, $q$ and $x$ (Tables S1,S2).

## Supplemental References

S1.     M. Pagel, A. Meade, Bayesian analysis of correlated evolution of discrete characters by reversible-jump Markov chain Monte Carlo. *American Naturalist* 167, 808 (2006).
S2.     W. R. Gilks, S. Richardson, D. J. Spiegelhalter, in *Markov chain Monte Carlo in practice,* W. R. Gilks, S. Richardson, D. J. Spiegelhalter, Eds. (Chapman and Hill, London, 1996),  pp. 1-19.
S3.     M. Pagel, New approaches to lexicostatistics and glottochronology. *Time Depth in Historical Linguistics. Cambridge: McDonald Institute for Archaeological Research*, 209 (2000).
S4.     M. Pagel, A. Meade, in *Phylogenetic methods and the prehistory of languages,* P. Forster, C. Renfrew, Eds. (McDonald institute Monographs, 2006),  pp. 173-182.
S5.     J. Blevins, A. Garrett, The origins of consonant-vowel metathesis. *Language*, 508 (1998).