Dating and Relationships



(or...)

Computational historical linguistics and long-range reconstruction

Eric Smith, Santa Fe Institute

Part of the Evolution of Human Languages project at SFI <<u>http://ehl.santafe.edu/</u>> in collaboration with the Tower of Babel project <<u>http://starling.rinet.ru/main.html</u>> Thanks to Murray Gell-Mann, George Starostin, Ilia Peiros, in memory of Sergei Starostin

Work done jointly with Tanmoy Bhattacharya, Jon Wilkins, Dan Hruschka, William Croft, Ian Maddieson, Logan Sutton, and Mark Pagel

Outline

- Goals of historical linguistics
- The classical comparative method
- Attempts at deep reconstruction
- New observations change the landscape
- Our attempts at quantitative reconstruction

Goals of historical linguistics

- To understand how languages change, and how they have changed historically
- To identify relations among languages due to common ancestry or cultural contact
- To reconstruct the languages of past speakers
- To contribute to an understanding of human populations and migrations
- To understand what is possible in language as a window on cognitive constraints

The interaction of history with process

"Proto-Turkic"

Yenisei'

Yakut Xalaj

Chuvash

Xakas Dolga "Chulym"

Old

Chulym

Tofa Tuv

 $\overline{\leq}$

- History-dependent phenomena combine *lawful dynamics* with *historical accident*
 - Accidents make branching processes -help us infer *diachronic* relations from *synchronic* variability
 - Diachronic relations assign the correct weights to processes
 "Western which act probabilistically

14 other close languages

The classical comparative method of historical linguistics: to interpret innovations

- A hypothesis of relationship among a set of languages.
- Cognate identification to correctly group elements (word roots or other morphemes) that have a common origin as evidenced by sound structure and meaning.
- Regular sound correspondences that describe relations among phonological segments in different languages across the lexicon; these may depend on phonetic context.
- Reconstruction of the ancestral sounds, morphemes and lexical items of a protolanguage that serves as a model of the ancestor of the proposed group.
- Establishing the innovations that created the descendant languages from the protolanguage.
- Classification of the languages within the family, using shared innovations to identify a structure of subfamilies.
- Construction of an etymological dictionary that traces semantic shifts and borrowings.

Language structure: many kinds of innovation can Phonology: motivate hypotheses of relatedness

- the sound system used by a language ("phones" or un-analyzed segment)
- the sound sets (phonemes) recognized as carrying distinctions in meaning
- Lexicon:
 - the map from root meanings to words (strings of phones or phonemes)
 - the overlap structure of meanings (so, what's in a word)
- Morphology:
 - "word shape": modifications to indicate case, tense, aspect, etc.
- Syntax:
 - word categories and composition rules to determine phrase structure, etc.
- Typology:
 - major structure categories of languages: word order, implicational groups Helpful glossary at <<u>http://www.sil.org/linguistics/GlossaryOfLinguisticTerms/</u>>

Rare innovations versus clusters of common innovations

- Rare innovations: single features with ~0 probability to occur by chance
 - Imply common descent or borrowing, even w/o mathematics
 - Only seen once: hard to assign probabilities from frequencies
 - Common in morpho-syntactic features
 - Useless for dating; do not support induction
- Common variations:
 - Examples: sound shift and meaning shift in core lexicon
 - Individually uninformative, but can assign probabilities from data
 - Require math to handle, but do support induction, and can be informative about dates if change processes are regular

| Morp | hology, syntax | k, typology | • |
|---|---|---|---|
| regular func | tions; differing | g represen | tations |
| Word orde | (Russian borrowin Oglum shkola my-son at-sch (S O | ng) da oe:renip tu ool is studyir V) | urar ng-3 |
| nom nomnar nomnaram nomnaramnung | a book books my books of my books | Vowel ha nom inek nomnar inekter | rmony a book a cow books cows |
| ajtyr ajtyrar men ajtyryp tur men ajtyrbas tur men | to ask I will ask I am asking I am not asking | Turkic family: ex "inflectional" an "agglutinating" | xamples of nd language TYPES |

Word lists as the key to lexical (= phonological / semantic) reconstruction

http://starling.rinet.ru/cgi-bin/response.cgi?root=config&morpho=0&basename=\data\alt\turcet&first=1

Turkic etymology :

New query Total of 2017 records 101 page Pages: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 Forward: 1 20 50 100

| Proto- Turkic | Altaic etymology | Meaning | Russian meaning | Old Turkic | Karakhanid | Turkish | Tatar | Middle Turkic | Uzbek | Uighur | Sary-Yughur | Azerbaidzhan | Tur |
|------------------|-----------------------|---|--|--|--|--------------|--|--|-------|----------------|-------------|--------------|--------------|
| *Ăb ⊞ | Altaic etymology ⊞ | hunt, chase | охота | ab (Orkh.), av (OUygh.) | av (MK) | av | aw | aw (Pav. C.) | v | aw, dial. σ | | σν | āv |
| *ab- ± | Altaic etymology ⊞ | to crowd, come together | собираться, встречаться | av- (OUygh.) | av- (MK, KB) | | | | | | | | |
| *abuč ⊞ | Altaic etymology ⊞ | handful | пригоршня | | avut (MK), avut-ča, avuč-ča (KB), avuč (Tefs.) | avuč | uč | avuč (MA, Sangl., Бор. Бад.) | xowuč | oč | oš | ovuč | ovui jan⊰ |
| *Abuč-ka 🗄 | Altaic etymology ⊞ | 1 husband, old man 2 foster- mother 3 elder sister 4 uncle | 1 муж, старик 2 кормилица 3 старшая сестра 4 дядя | avičya, abučya 1, abučga 2 (OUygh.) | avičya 1 (MK, KB) | abuš 3 dial. | abušqa, awucqa 1 dial. (Sib.) | abušqa, avušqa 4 (Abush., Sangl.) | | | | | |

Fundamental object is the history of a given word's sound and meaning

- n.b., word forms are attested; meanings are indirectly inferred, and often ambiguous
- easy to trace a form; but inadequate to infer history from forms alone



Representing sound and meaning "innovation" in the comparative method

- Suppose that some stable meaning categories can be identified
- Identify primary words for each meaning
- Try to exclude "borrowed" terms; suppose that what is left has been transmitted through vertical descent
- Identify systematic sound relations and try to infer historical sound changes
- Associate semantic innovations with inlanguage substitutions within meaning categories



Preserved meanings suggest sound maps



Sound maps help identify meaning shifts



But sound changes can depend on context

Split of Old English /k/

| Stage I | katt | keaff | kinn | | | | | | | | |
|-------------------|-------|---------------------|------------------------|--|--|--|--|--|--|--|--|
| Stage II | katt | t∫eaff | t∫inn | | | | | | | | |
| Stage III | katt | <mark>t∫</mark> aff | t∫inn | | | | | | | | |
| Split of Latin/s/ | | | | | | | | | | | |
| Stage I | ka:ra | flo:s | flo: <mark>s</mark> es | | | | | | | | |
| Stage II | ka:ra | flo:s | flo:zes | | | | | | | | |
| Stage III | ka:ra | flo:s | flo:res | | | | | | | | |

From R.L.Trask, "Language change"

- Even if true, not uniformly available, and unrealistic to rely on at larger time depths
- Eventually context discovery requires more information than the language pattern yields

And the nature of meaning shift is not mathematically understood

- Phonological and semantic constraints interact with polysemy and synonymy to structure sound and meaning change
- Semantic categories, split, join, and move in some "space" which we do not know



Attempts at long-range reconstruction: lexicostatistics and glottochronology

- Assign any sound map without penalty, but require regularity
- Exclude borrowed items in either language from consideration
- Identify fraction preserved cognates; convert to separation time (penalty)
- Attempt to fit separation times to an ultrametric structure (tree)

$$P_{\text{Preserve}}(\text{word}) = e^{-t/\tau}$$

$$\Delta t_{\rm sep} = -\tau \log (\text{frac. preserved})$$



Example: the Nostratic Hypothesis



http://starling.rinet.ru/maps/maps.php?lan=en

- Coined by Holger Pedersen (1903)
- Modern form of the hypothesis by Vladislav Ilitch-Svitych and Aharon Dolgopolsky (1960s -- present)
- Estimated 12,000-15,000 BCE



http://en.wikipedia.org/wiki/Nostratic_languages

Sub-families within Nostratic



New quantitative observations: finer process models for language change

- Rates of change in semantics, phonology, morphology
- "Point processes" associated with branching

Typology I: frequency of use and rates of change

Frequency of word-use predicts rates of lexical evolution throughout Indo-European history

Mark Pagel^{1,2}, Quentin D. Atkinson¹ & Andrew Meade¹

NATURE Vol 449 11 October 2007



Quantifying the evolutionary dynamics of language

 ${\sf Erez} \ {\sf Lieberman}^{1,2,3}\star, {\sf Jean-Baptiste} \ {\sf Michel}^{1,4\star}, {\sf Joe} \ {\sf Jackson}^1, {\sf Tina} \ {\sf Tang}^1 \ \& \ {\sf Martin} \ {\sf A}. \ {\sf Nowak}^1$

NATURE Vol 449 11 October 2007



Linguistic punctuated equilibrium: an interaction of phylogeny with underlying process?



Languages Evolve in Punctuational Bursts

Quentin D. Atkinson,¹* Andrew Meade,¹ Chris Venditti,¹ Simon J. Greenhill,² Mark Pagel^{1,3}†

1 FEBRUARY 2008 VOL 319 SCIENCE www.sciencemag.org

Our work to quantify the comparative method

- Alignment
- Regular sound correspondence
- Minimum-bias methods for unexplained variation
- Context-dependence and information criteria
- Attempts to map the semantic space

Maximum-likelihood estimation of history and process

- Align words in daughter languages
- Propose phoneme assignments to aligned positions in the ancestor (with probabilities)
- Estimate regular correspondence of ancestor to daughter phonemes (w/ or w/o probabilities)
- Estimate random violations (with probabilities)



Alignments inferred without prior knowledge

Etymology 4 "belly"

| q | q | k | q | q | q | q | q | G | G | х | q | q | q | х | х | | х | x | q | q | q | q | q | q | q | q | q | q |
|-----|-------|------|------|------|------|------|--------------|-----|-----|----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| а | а | а | а | а | э | e | а | a | а | а | а | а | ā | i | а | | i | i | а | а | а | а | а | а | а | a | а | a |
| r | r | r | r | r | r | r | r | r | r | r | r | r | r | r | r | | r | r | r | r | r | r | r | r | r | r | r | r |
| i | i | i | i | i | i | i | i | i | i | i | | i | i | ъ | i | | i | i | i | i | i | i | i | i | i | i | i | i |
| n | n | n | n | n | n | n | n | n | n | n | n | n | n | m | n | | n | n | n | n | n | n | n | n | n | n | n | n |
| | | | | | | | | | | | i | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Εţ | ym | oloş | gy ś | 5 "t | oig, | hiş | gh" | | | | | | | | | | | | | | | | | | | | | |
| b | b | b | b | b | b | b | b | b | b | р | m | b | b | | | 1 | b | b | b | b | b | ь | m | Ъ | b | b | | b |
| е | e | ü | i | e | u | ü | e | ö | e | ö | ö | i | i | | | | e | e | i | i | i | e | i | ű | ü | i | | i |
| d | δ | j | j | j | j | j | z | j | j | z | z | j | d | | | | d | d | j | j | j | j | j | | j | j | | j |
| ü | ü | ü | e | i | u | ü | i | ü | i | ə | ü | i | i | | | i | i | i | i | i | i | e | i | | ü | i | | i |
| k | k | k | k | k | k | k | k | k | k | k | k | k | k | | | 1 | k | k | k | k | k | k | k | k | k | k | | k |
| Ff | vm | مام | TV. | 12. | í P | roo | et | nin | nlo | ,, | | | | | | | | | | | | | | | | | | |
| 1.0 | y III | uluş | 5J - | 12. | D | 1 ca | ы . , | mp | pic | r | | | | | | | | | | | | | | | | | | |
| | | m | m | 1 | | n | n | 1 | n | m | | | | m | m | 1 | | | | | | m | m | | n | 1 | m | m |
| | | e | ε | | | ä | | | i | ä | | | | ä | ē | | | | | | | ä | ä | | ä | | ä | ä |
| | | m | m | 1 | | n | n | 1 | n | m | | | | m | m | ı | | | | | | m | m | | n | 1 | m | m |
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | m |
| | | e | i | | | ä | i | | i | e | | | | ä | ē | | | | | | | ä | ä | | ä | | ä | ä |
| | | | | | | | | | | | | | | k | | | | | | | | j | j | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | - | | | | | |

Sound correspondences among the languages



Sound relations inferred from sound changes



Inferring polysemy with English as a meta-language



Sun Moon 1 Day Month

Network of polysemes in 81 diverse languages



Some concluding comments

- Historical linguistics has been a "rule-based" system, something like formal logic
- Can we re-derive such rule-based systems within principles of statistical inference?
- New high-volume data sources and new data types: a huge opportunity for computational analysis
- Languages provide another window on evolution and historical inference from molecular sequences
- BUT: Good quantitative linguistics will require collaboration, patience, and work