

# Macro-Comparative Linguistics In The XXIst Century: State Of The Art And Perspectives

**George Starostin**

Center for Comparative Linguistics, RSUH

Handout presented at the International Conference  
"Comparative-Historical Linguistics Of the XXIst Century", March 20-22, 2013

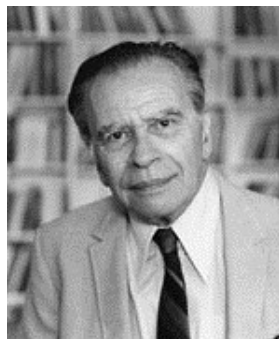
## General purposes of the talk:

- concentrate more on issues of method than on particular «long-range» hypotheses;
- suggest possible ways of integrating (reintegrating?) macro-comparative studies into «mainstream» historical linguistics.

## Macro-comparative studies today:

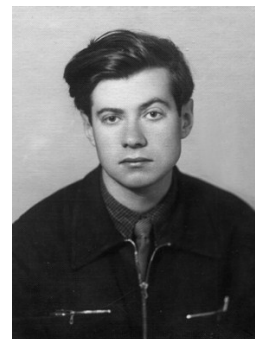
- pushed towards the fringes of historical linguistics;
- no fundamental unity between «macro-comparativists» as such, even those working on the same hypotheses («Dolgopolsky's Nostratic» vs. «Bomhard's Nostratic», «Militarev's Afro-Asiatic» vs. «Ehret's Afro-Asiatic», etc.);
- in very dire need of «rehabilitation» due to progress in adjacent disciplines preoccupied with human prehistory (archaeology, genetics, etc.);
- may be complemented with, but not replaced by intrinsically ahistoric statistical studies;
- need to operate more openly on a general probabilistic basis.

## Two main types of macro-comparative research:



«Greenberg-style»:

- based on «mass comparison»
- emphasis on classification
- «naively probabilistic»



«Illich-Svitych style»:

- based on the comparative method
- emphasis on reconstruction
- oriented at «proof» of relationship

## Common for both:

**no intrinsic difference between «short-range» and «long-range»**

## Principal issues in macro-comparative linguistics

### 1. QUANTITY vs. QUALITY

«*The more, the better*»:

- huge etymological dictionaries (*EDAL, Nostratic Dictionary, Sino-Caucasian*, etc.)
- advantage: easier to demonstrate systematic regularities
- flaw: high risk of accumulating «noise» from chance resemblances

«*Less is more*»:

- small, concise sets of comparanda (Sino-Austronesian, Dene-Yeniseian, etc.)
- flaw: difficult to confirm the presence of regular correspondences
- advantage: easy to evaluate the evidence *in toto* by outside specialists

*Recommended solution*: Size does not matter, as long as there are clearly stated and unbiased criteria (phonetics; semantics; lexicostatistics; distribution in daughter branches, etc.) for ranging the evidence.

### 2. GRAMMAR vs. LEXICON

«*Paradigmatic morphology*» argument:

- generally supported in various «mainstream» Western models
- advantage: highly resistant to borrowing, highly unlikely to arise by chance
- advantage: easily presentable, of significant interest to synchronic linguistics and typology
- flaw: notoriously unstable whenever languages undergo phonetic or typological restructuring
- flaw: very small amounts of «hard data»; easy to slip from systemic comparison into free-flowing speculation

«*Basic lexicon*» argument:

- generally supported in the «Moscow school» (cf. *Global Lexicostatistical Database*)
- flaw: may be borrowed in large quantities under specific sociolinguistic conditions
- flaw: allegedly «discredited» through an association with crude, obsolete glottochronological models
- advantage: due to fewer interdependencies, more resistant even in situations of rapid phonetic/typological change
- advantage: as a rule, supplies enough data to make realistic classifications even on very high levels

*Recommended solution*: Morphology need not matter; may not be used as grounds for a universal and easily formalizable classificatory standard; almost (but not necessarily) useless for macro-comparative purposes. Some grammatical morphemes may be used on par with the «quasi-grammatical» elements of the Swadesh wordlist (personal and demonstrative pronouns, negation markers, etc.).

### 3. INTERNAL SOLUTIONS vs. EXTERNAL EVIDENCE

*Precedence of internal reconstruction over external comparison*

- generally supported by «narrow» specialists in particular groups / families, skeptical of «bold» external hypotheses
- flaw: lack of proper comparative basis prone to lead to speculation (at worst, to fantasy) instead of science

*Precedence of external comparison over internal reconstruction*

- generally supported in «mass comparison» as well as Neogrammarian-based macro-comparative models (Nostratic, etc.)
- flaw: combined with the «quantity over quality» principle (q.v.), prone to lead to violation of probability expectations

*Recommended solution: No bias should be detectable.* Competing solutions should be weighed probabilistically against an accumulating bank of evidence for diachronic change (typology of phonetic change, semantic change, language contact, etc.) rather than intuition and pure theory.

### 4. FORMAL OBJECTIVITY vs. SUBJECTIVE JUDGEMENT

*Formal algorithms must play a crucial role in historical analysis of language data*

- stimulated by the respective advances in genetics and other «hard science»
- popular mostly with non-linguists («Gray Lab», Mark Pagel's group, etc.) but not exclusively so (D. Ringe; S. Wichmann's ASJP, etc.)
- advantage: seemingly objective procedures that (a) filter out subjective bias on part of the researcher; (b) significantly facilitate working with large numbers of languages
- flaw: any algorithm is only as good as the data fed into it; perfect objectivity not attainable for that reason (e. g. «sloppily» constructed wordlists; poor treatment of semantics; relying on subjective etymological cognacies, etc.)

*No formal algorithm is robust enough to consider all the important arguments*

- most historical linguists tend to be skeptical about trusting automated judgements
- possible alternative: «manual» quantification of evidence with elements of objective analysis applied to etymologies generated on a «subjective» basis
- flaw: situations that cannot be easily «cracked» by automatic algorithms often cannot be easily «cracked» manually as well; quantifying «anomalies» and «irregularities» is not an easy task in any respect
- advantage: personalized approach to different families allows adjusting for their typological peculiarities
- advantage: possibility to «manually» check data against the accumulating bank of evidence for diachronic change

*Recommended solution: Use both.* «Manual» and «automated» procedures of etymological and statistical analysis should be applied in tandem, in order to simplify each other's performance and help identify the stable and the questionable sectors of reconstruction and classification.

## **Suggested general recommendations of future research**

### **[1] Work mainly with evidence that may be quantified:**

- Swadesh-type wordlists allowing for manual as well as automatic analysis
- accurately organized etymological databases

### **[2] Build up reference corpora of typological evidence:**

- unified databases for types of phonetic change
- unified databases for types of semantic change
- unified databases for loanword typology

### **[3] Try to develop and apply universal standards and reference frames:**

- for lexicostatistics (rigidly defined standard wordlists)
- for etymology (unified databases)

## **Final conclusions**

- many of the listed principles/recommendations not properly implemented even for «short-range» families;
- this may be ignored/downplayed in relatively trivial «short-range» situations where data are abundant and daughter languages have not shifted to different typologies;
- this may not be ignored for macro-comparative purposes, where establishing reliable systems of regular phonetic correspondences becomes harder due to lack of quality data and must be compensated for with other check factors;
- not a single macro-comparative hypothesis satisfies all of the listed principles / recommendations, but some are better than others and may be ranged accordingly (e. g. Nostratic > Austric, Austric > Nilo-Saharan, etc.);
- at least some current and future efforts of the Moscow School and those projects that it coordinates («Evolution of Human Languages», «Global Lexicostatistical Database») targeted at significant methodological improvements, including use of «smart lexicostatistics» and empirically-based filters for possible semantic shifts.