# DISTANCE METHODS FOR PHYLOGENETIC INFERENCE

Thesis Proposal | Philip Georgis

# DISTANCE-BASED METHODS

Jäger (2018): Global-scale phylogenetic linguistic inference from lexical resources

*"Aggregating over all families suggests that distance-based inference produces the best fit with the expert gold standard. [...] Combining both types of characters in a partitioned model always leads to better results than the two character types individually."*

# RESEARCH QUESTION

What types of distance-based methods are useful (or: improve) phylogenetic inference?

- Phonetic
  - Phonemes (IPA) or sound classes?
  - Which phonetic features to use?
    - e.g. PHOIBLE features, proposed features by Tresoldi (2016), other?
  - Weighting features or segment types? (e.g. weight consonants more than vowels, or weight certain phonetic features more than others)
  - Penalties for deletions?

# RESEARCH QUESTION

What types of distance-based methods are useful (or: improve) phylogenetic inference?

- Lexical (= cognate distance?)
  - Character-based binary lexical distances
  - PMI score *(method used by Jäger)*
  - Surprisal score → new method
    - Lower surprisal : total phoneme entropy ratio → closer cognates

- Morphosyntactic
  - How to represent features? (binary or continuous?)
  - Where to take features from? (e.g. WALS, grammars, UD corpora)
  - Which features to include?
    - → need to research which MST features are most associated with genetic relatedness

# EVALUATION

Validation with gold cognate sets to measure performance in cognate detection step

- F1 score or MCC for classification accuracy
- Purity and completeness measures for cluster identification

| Gloss | Czech | Slovak | Polish | Slovenian | Croatian | Serbian | Macedonian | Bulgarian | Russian | Ukrainian | Belarusian |
|---|---|---|---|---|---|---|---|---|---|---|---|
| many | mnoɦo | mnɔɦɔ | | mno:gɔ | mnogo | mnogo | mn̪ɔgu | mn̪ɔgo | mn̪ogo | | mn̪ɔɣa |
| many | | vɛʌa | vjɛlɛ | | | | | | | | |
| many | | | | | puno | puno | | | | | |
| many | | | duʐɔ | | | | | | | | |
| many | | | | | | | | | | baɦatɔ | |
| meat | maso | mæsɔ | mjɛ̃sɔ | mɛso: | mesɔ | mesɔ | mɛsɔ | mesɔ | mʲasɔ | mjasɔ | mʲasa |
| moon | mn̪ɛsi:ts | mɛsʲats | mjɛсɔ̃ts | me:sət̠s | mjese:t̠s | mese:t̠s | mɛsɛtʃin̪a | | | misʲatsʲ | mʲɛsʲatⱮ |
| moon | | | | lu:na | | | | ɫun̪a | ɫʲun̪a | | |
| mountain | ɦora | ɦora | gura | gɔ:ra | gora | gora | | | gora | ɦɔra | ɣara |
| mountain | | | | | pɫanina | pɫanina | pɫan̪in̪a | pɫan̪in̪a | | | |
| mouth | u:sta | u:sta | usta | u:s̪ta | u:s̪ta | u:s̪ta | us̪ta | us̪ta | | | |
| mouth | | | | | | | | | rot̪ | t̠ɔt̠ | rɔt̪ |
| name | jmɛːno | mɛnɔ | imjɛ̃ | ime: | ime | ime | imɛ | ime | imʲa | imja | jimʲa |
| neck | krk | krk | kark | | | | | | | | |
| neck | | | ʂɨja | | | | ʃija | ʃija | ʂeja | ʃɪja | ʂija |
| neck | | | | ʋra:t̪ | ʋra:t̪ | ʋra:t̪ | vrat̪ | vrat̪ | | | |
| new | novi: | nɔvi: | nɔvi | nɔʋ | noʋ | noʋ | n̪ɔv | n̪ɔv | n̪ovij | n̪ɔʋɪj | n̪ɔvi |
| night | nots | nɔts | nɔts̠ | no:tʃ | no:t̠ɕ | no:t̠ɕ | n̪ɔc | n̪ɔʃt̪ | n̪ot̠ɕ | nʲit̠ʃ | n̪ɔt̠ʂ |

# EVALUATION

Comparison of generated phylogenetic trees with gold-standard trees from experts
- e.g. Glottolog trees

## TreeDist

'TreeDist' is an R package that implements a suite of metrics that quantify the topological distance between pairs of unweighted phylogenetic trees. It also includes a simple 'Shiny' application to allow the visualization of distance-based tree spaces.

'TreeDist' primarily employs metrics in the category of 'generalized Robinson–Foulds distances': they are based on comparing splits (bipartitions) between trees, and thus reflect the relationship data within trees, with no reference to branch lengths.

## Distance metrics for ranked evolutionary trees

Jaehee Kim, Noah A. Rosenberg, and Julia A. Palacios
+ See all authors and affiliations

# EVALUATION

How to evaluate the **distance measures?**

- Correlation with mutual intelligibility measures
  - Available to some extent for:
    - Slavic (SFB C4 experimental results; Lindsay; Golubović & Gooskens, 2015)
    - Germanic (Gooskens, 2017)
    - Turkic (Lindsay)
    - …possibly others?

- Other ideas?

# FORMAL ASPECTS

First step: Get clearer idea research question and contact supervisor

Master seminar

- Meet regularly with supervisor to refine topic
- Write thesis proposal
- Give ~30 minute presentation of proposal
- Supervisor then signs thesis registration form

Thesis

- Continue meeting with supervisor
- Write up findings
- Submission: minimum 3 months, maximum 6 months from date of registration

# FORMAL ASPECTS

Second advisor

- **Bernd Möbius**
  - → if the focus would be more on phonetic aspects generally

- **Tania Avgustinova**
  - → if the focus would be more on correlation with inter-comprehension, morphosyntax, and/or Slavic as a key example family