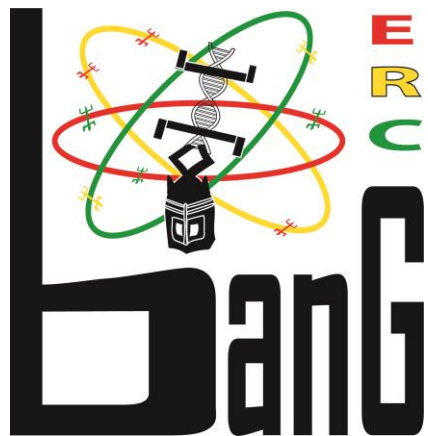




Optimizing Cognacy for Automatic Internal Classification of Languages: The Critical Role of Segmentation



Promise Dodzi Kpoglu
promisedodzi@gmail.com / promise-dodzi.kpoglu@cnrs.fr

LATTICE, Paris.
17th September, 2024.



Outline

- Introduction
 - General introduction
 - Objectives
- The Dogon languages
- Methodology
- Results
- Discussion
- Conclusion

Introduction

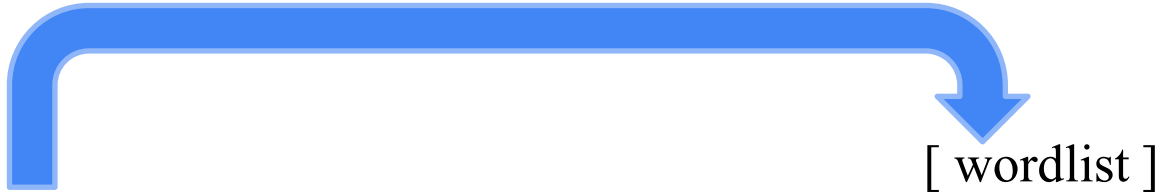
general idea

- Traditional historical linguistics approaches are anchored on the comparative method (Campbell, 2013).
- Discrete (non mutually exclusive) steps often involved in the comparative method (Weiss, 2015)



Introduction

general idea



Language	Family	Time-depth	Proto-form	Cognates
French	Romance	~1000	cane [kane]	chien [ʃiɛ̃]
Italian	Romance	~1000	cane [kane]	cane [kane]
Spanish	Romance	~1000	cane [kane]	can [kan]

Hypothesis

CM

p.c (Fabian Zuk)

Introduction

general idea

- Since the beginning of the 20th century, attempts have been made to automate some of the steps involved in the CM:
 - Cognate identification: words that are related, especially in form – see Jäger 2019
 - Establishment of sound correspondences –see List et al. 2022, Kim et al 2023
 - Proto-form reconstruction- See Meloni et al 2021 for instance
 - Time-depth estimation – See Gray & Atkinson 2003, 2006
 - Phylogenetic relationships – See Rama et al. 2018
- This paper concerns automatic methods for cognate detection.

Introduction

cognacy

- Cognates can be of different forms depending on phonetic and semantic properties.
- Phonetics:
 - Lexemes can manifest high phonetic transparency i.e. “strong cognacy”, or opaque phonetic properties i.e. “weak cognacy” (Meelen & Hill 2022: 52)
- Semantics:
 - Lexemes can manifest semantic equivalence i.e. “synonymous cognates” or a weak semantic correlation i.e. “non-synonymous” (Koch & Hercus 2013:34)
- In this paper, I consider:
 - Strict cognacy = phonetically strong, semantically synonymous e.g. Ger. *herz* vs. Eng. *heart*
 - Partial cognacy = less than the full phonetic and semantic form is available eg. Ger. *walfisch* vs. Eng. *whale*
- Not considered are synchronic forms involved in ‘dialexification’ (cf. François & Kalyan 2023):
Albanian. *gardh* ‘yard’ vs Romani. *kher* “family”

PIE root.

**g^herd^h*- ‘enclose’

Introduction

automatic cognate detection

- A critical component of automatic cognate detection involves word similarity calculation which is then followed by cognate alignment (Rama et al. 2018: 4).
- Word similarity calculation:
 - Semantic similarity: similarity based on traditional methods (comparative concepts) –Forkel et al (2018) , or on corpus properties - (Kondrak, 2001)
 - Phonetic similarity: similarity based on particular metric which can be feature-based – Kondrak (2000) or class-based (List, 2012)
- In this paper, automatic cognate detection involves comparative concepts and class-based distances.

Introduction

Automatic cognate detection

- Attempts to improve results of automatic cognate detection have not only been concerned with similarity metric innovation, but also:
 - Model-experimentations:
 - Jager et al (2014) use a SVM to automate the process, while Konojia et al (2021a) use a feed-forward neural network
 - Feature-enrichment:
 - Konojia et al (2021b) for instance introduce a feature-enriched dataset that is feed into various machine learning models

Introduction

Automatic cognate detection

- A critical factor that can influence the quality of automatic cognate detection involves the quality of data (List, 2017).
- In other words, tasks which are generally referred to as “low-level” - see Mikheev 2022: 550, are critical to the results obtainable.
- Nevertheless, to our knowledge, there is no dedicated study seeking to understand to the relationship between data quality and improvement of automatic cognate detection.

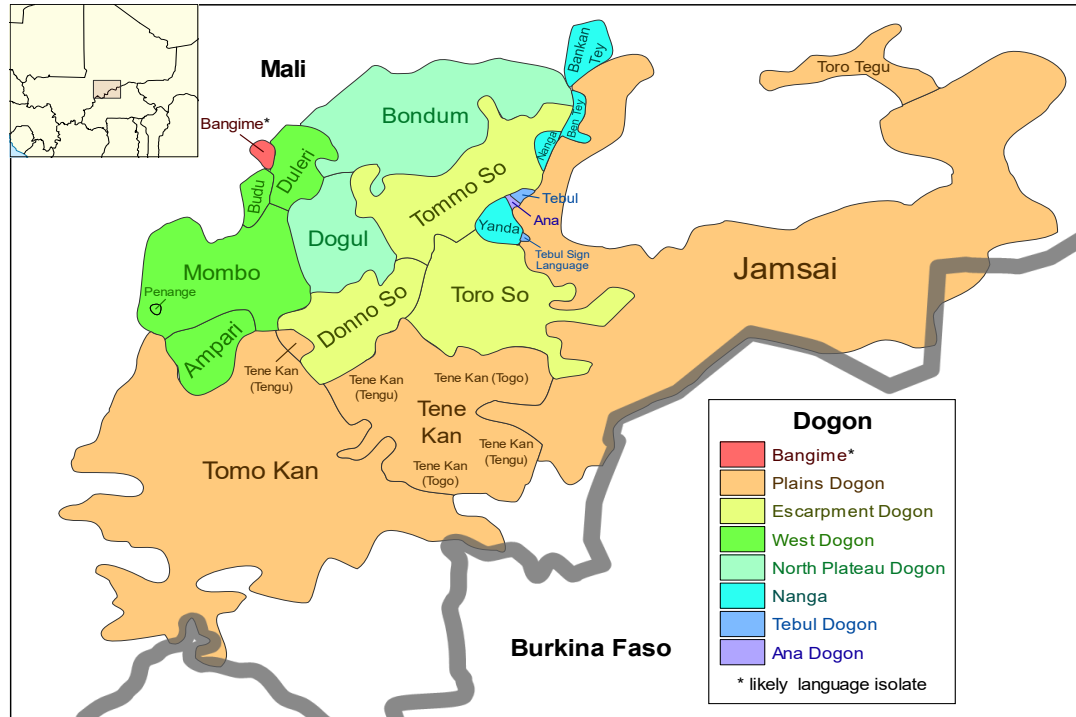
Introduction

Objectives

- The objectives of this paper are two-fold:
 - Understand how data ‘quality’ influences results obtained by automatic cognate detection techniques
 - Examine the consequences of change in data ‘quality’ on tasks further upstream in the automated historical comparative method workflow

The Dogon languages

Relevant typological features



The Dogon languages

Relevant typological features

- The vowel inventory usually consists of seven vowel qualities, short and long. Nasalized counterparts, long and short are also available.
- The consonantal inventory usually includes nasalized sonorant [wⁿ], [yⁿ] and [rⁿ] – (Heath 2015:7).
- There is often ATR harmony in the languages.
- Dogon languages are tonal languages, with the particularity of syntactic categories exhibiting a complex tonal controller/non-controller dichotomy (McPherson 2014:60)
 - Controllers trigger tonal overlay; non-controllers do not.

TommoSo

gámmá=ge ‘the cat’ (N Def) = Def is a non-controller

gàmmà gém ‘black cat’ (N Adj) = Adj is a controller

The Dogon languages

Relevant typological features

- The favored syllabic structures are Cv, CvCv, CvNcv, and CvCvCv.
- Morphologically, Dogon languages have agglutinative features. Thus, various “affixes” can attach to stems.
 - Eg. adjectives in YandaDom can have an inchoative and factitive with one or more derivational suffixes involved in the process (Heath 2017:238/239).

‘hot’	<i>inchoative</i>	<i>factitive</i>
òjú	ój-jé	ój-jé-mé
- Verbs generally, have allomorphic forms (usually instantiated with an ATR harmony paradigm) with each form of the verb exponenting various TAM categories.
 - For instance while the bare stem in TebulUre typically expresses perfectivity, the A/X stem expresses imperfectivity (Heath 2023: 21/22).

TebulUre

‘abandon’

dògó	bare stem
dògé-∅	3sg simple perfective E/I stem
dógà-m-dò-∅	3sg imperfective A/X stem

Methodology

- Data preprocessing
- Word segmentation
- Cognate detection
- Cognate evaluation

Methodology

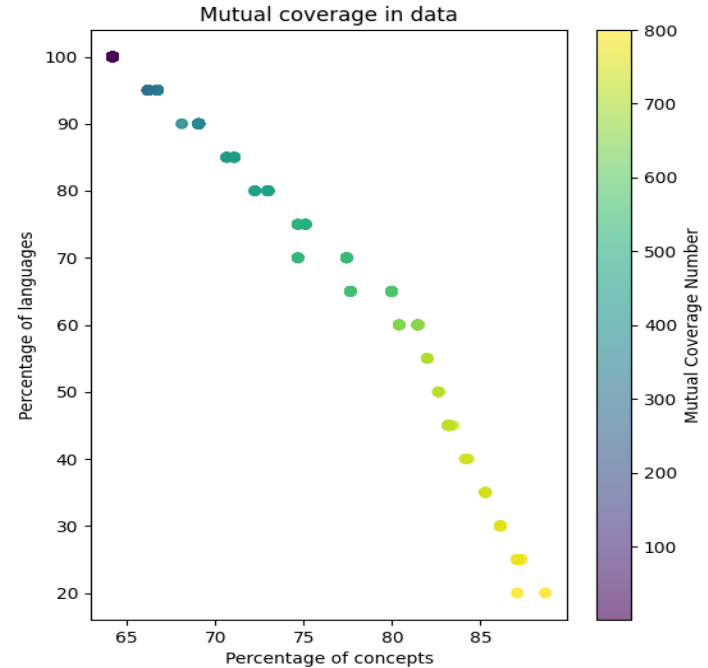
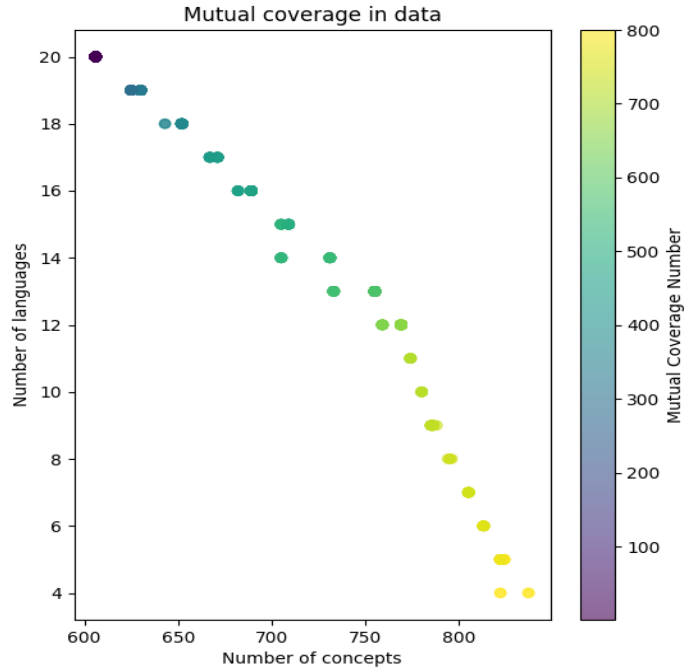
Data preprocessing

- Data is from Dogon and Bangime Linguistics project (Moran et al. 2016), curated in CLDF (Forkel et al. 2018), and available as `heathdogon` on GitHub.
- Data cleaning and conversion includes the following steps:
 - Formatting data so each form is defined with unique ID, and first row specifies language, form and concept, just as prescribed by List et al (2018)
 - Running an orthography profile to harmonize graphemes for computation, by using the segments package (Forkel et al 2019)
 - Identifying concept coverage rate in various languages and choosing languages that meet the required threshold (i.e. 288)

Methodology

Data preprocessing

- Cleaned data had 20 languages with a mutual coverage number of 288.



Methodology

Word segmentation

- Based on the grammars available for the Dogon languages, different parsing rules were applied to the data via scripts written in python.
- Four datasets.
 - No parse dataset: the raw preprocessed data
 - Phonetically parsed dataset: long vowels and tones parsed
 - Morphologically parsed dataset: phonetic and identified suffixes parsed
 - Morpho-phonotactically parsed dataset: morphological and phonotactic parsing

No parsing	Phonetic parsing	Morphological parsing	Morpho-phonotactic parsing
b ε r j i m	b è r j î m	b è r j î + m	b è r + j î + m
goatkid.ANM	goatkid.ANIM	goatkid-ANIM	“goat-kid-ANIM

Methodology

Cognate detection

- The LexStat algorithm is used for cognate detection (List, 2012).
 - Converts sequences into sound classes
 - Calculates language specific scoring schemes
 - Calculates pairwise distances
 - Undertakes Sequence clustering
 - A distance matrix can be obtained; as well as alignments, for visualization
- LexStat is available in the Lingpy library (in python).
- Full and Partial cognate detections are carried out on all four datasets – LexStat can be used for both (List et al. 2016)

Methodology

Cognate detection

Full cognate detection

ID	DOCULECT	CONCEPT	TOKENS	NOTE	COGID
769	DogulDomKundialang	(child) be born	n à l + j é		130 ³
770	Najamba	(child) be born	n à l + i + j é		130 ³
771	TommoSoTongoTongo	(child) be born	n à l + i + j é		130 ³
772	YomoSo	(child) be born	n à r + é é		133

Two things to note:

- full=1 row, 1 cogid
 - partial=1 row, multiple cogids
- b. counting cognate rows = favoring full

Partial cognate detection

ID	DOCULECT	CONCEPT	TOKENS	NOTE	COGIDS
769	DogulDomKundialang	(child) be born	n à l + j é		347 ⁴ 348 ³
770	Najamba	(child) be born	n à l + i + j é		347 ⁴ 349 ² 348 ³
771	TommoSoTongoTongo	(child) be born	n à l + i + j é		347 ⁴ 349 ² 348 ³
772	YomoSo	(child) be born	n à r + é é		347 ⁴ 350

Methodology

Cognate detection

- Cognacy results are evaluated in two ways:
 - Computing cognacy scores: statistics on cognacy
 - Assessing cognacy goodness: clustering and comparison to qualitative ‘ground-truth’
- Computing cognacy scores:
 - Number of unique cognate pairs (cogids and concepts)
 - Number of unique concepts involved in cognate pairs
 - Number of rows involved in cognate pairs
 - Proportion of data with cognacy
 - average number of cognate items per concept
- Cognacy scores favor full cognate detection as they are cogid count based.

Methodology

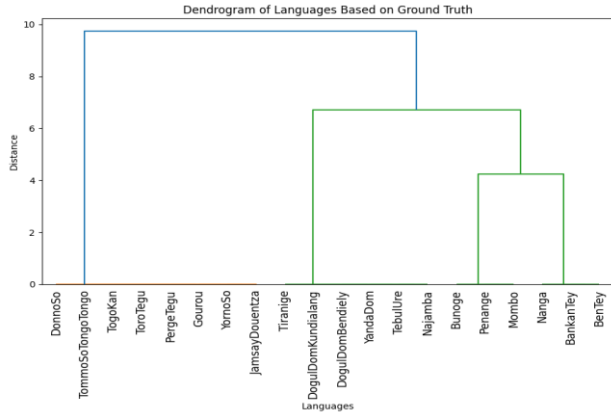
Cognate evaluation

- Assessing cognacy goodness is done via a three-step process:
 - Computing distance matrix
 - Clustering distance matrix
 - Comparing resulting clusters to a ‘field-linguist ground truth’
- Computing distance matrix:
 - aggregated computed distances via LexStat for all languages
- Clustering distance matrix
 - processing distance matrix into a condensed form
 - Hierarchical clustering of distance matrix using average-linkage clustering
- Comparing resulting clusters to ground truth
 - Clusters obtained are compared to a qualitatively generated organization of Dogon languages (Heath 2012)

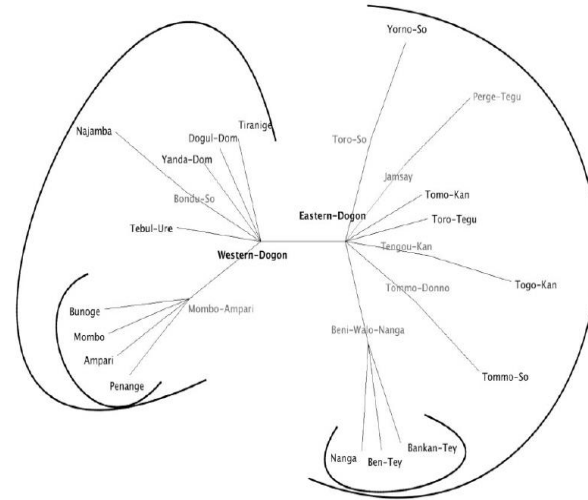
Methodology

Cognate evaluation

Generated cluster vs. field linguist ground truth



VS.



Moran & Prokić (2013:12)

● Computed scores:

- Adjusted Rand score: Measures the similarity between clusterings, adjusting for random chance
- Normalized mutual information score: Quantifies shared information between the true and predicted clusters
- Fowlkes-Mallows score: Balances precision and recall of cluster assignment
- Homogeneity score: Checks whether each cluster contains only samples from one class
- Completeness: Ensures that all samples from a class are in the same cluster
- V-Measure: Combines homogeneity and completeness into a single score

Results

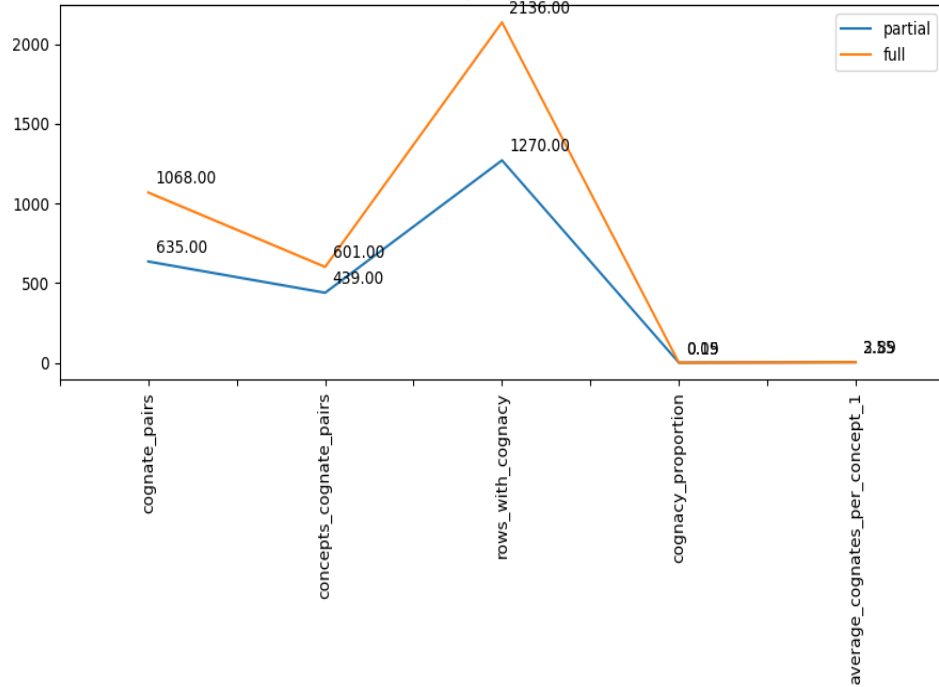
- For purposes of recall:
 - Parsing types: no parsing, phonetic parsing, morphological parsing, morphophonological parsing
 - Cognacy type: full cognacy vs. partial cognacy
 - Cognacy performace: cognacy scores vs. cognacy goodness
- First : results of parsing type vs. cognacy type
- Then : results of cognacy type vs. cognacy performace

Results

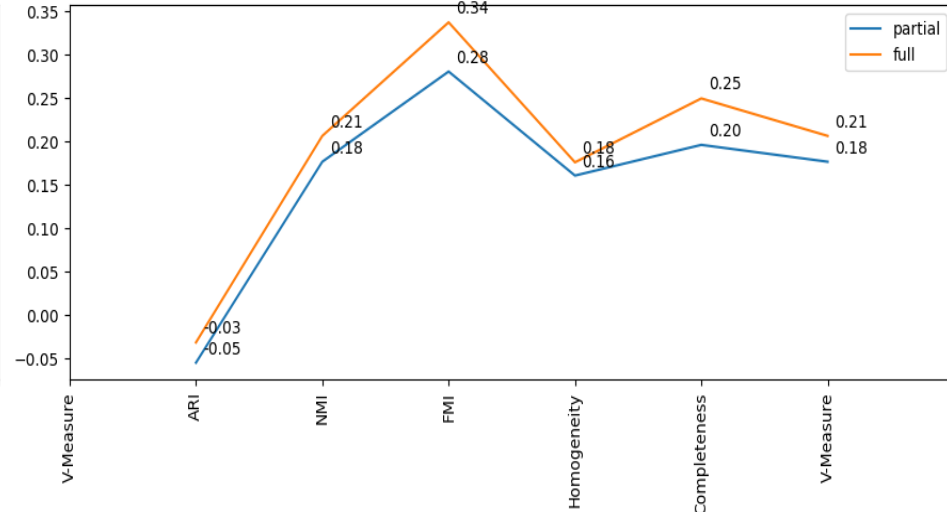
No parsing

- Full cognacy performs best on unparsed data.

Cognacy stats



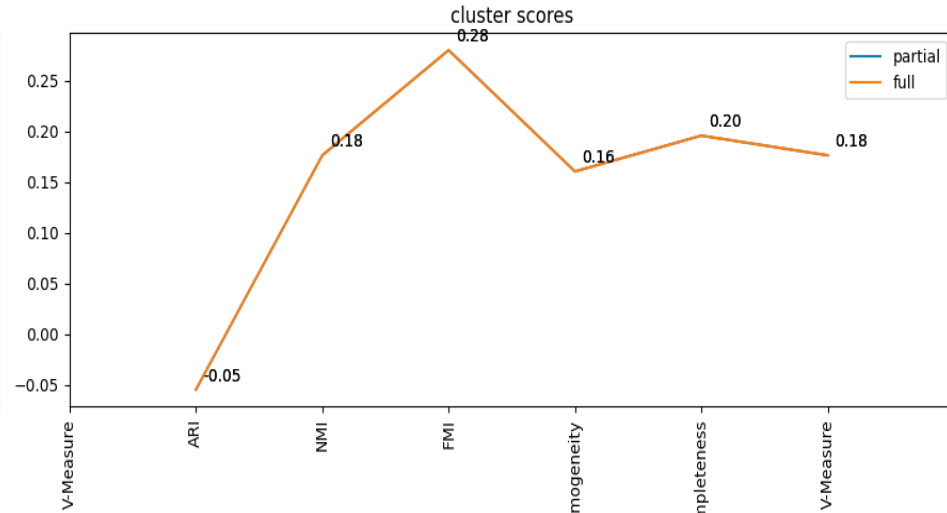
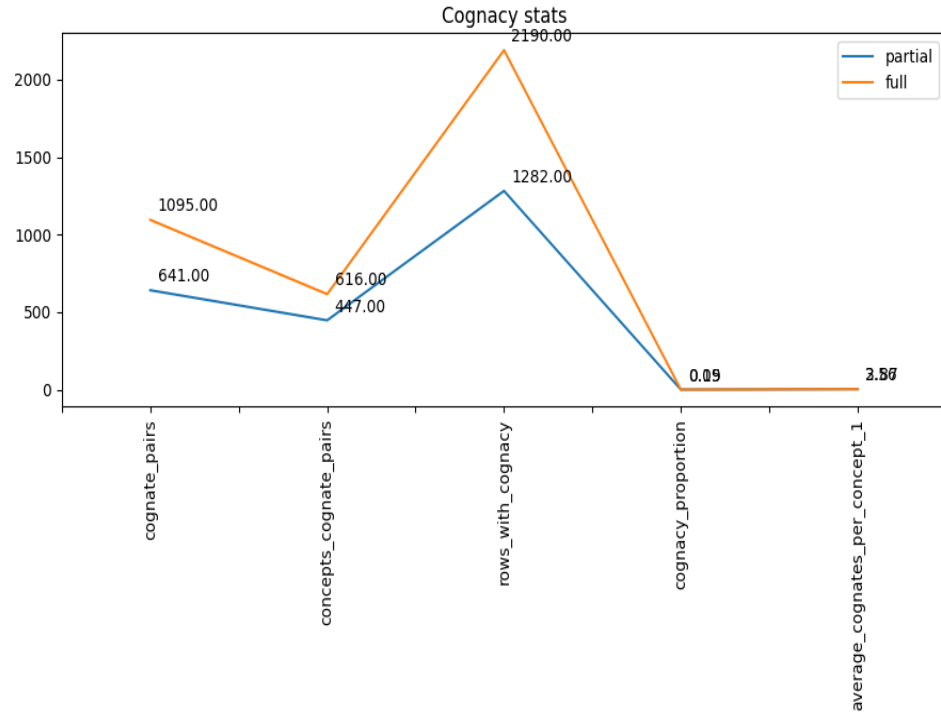
cluster scores



Results

Phonetic parsing

- Full cognacy performs best on scores, but both cognacy types have equal goodness.

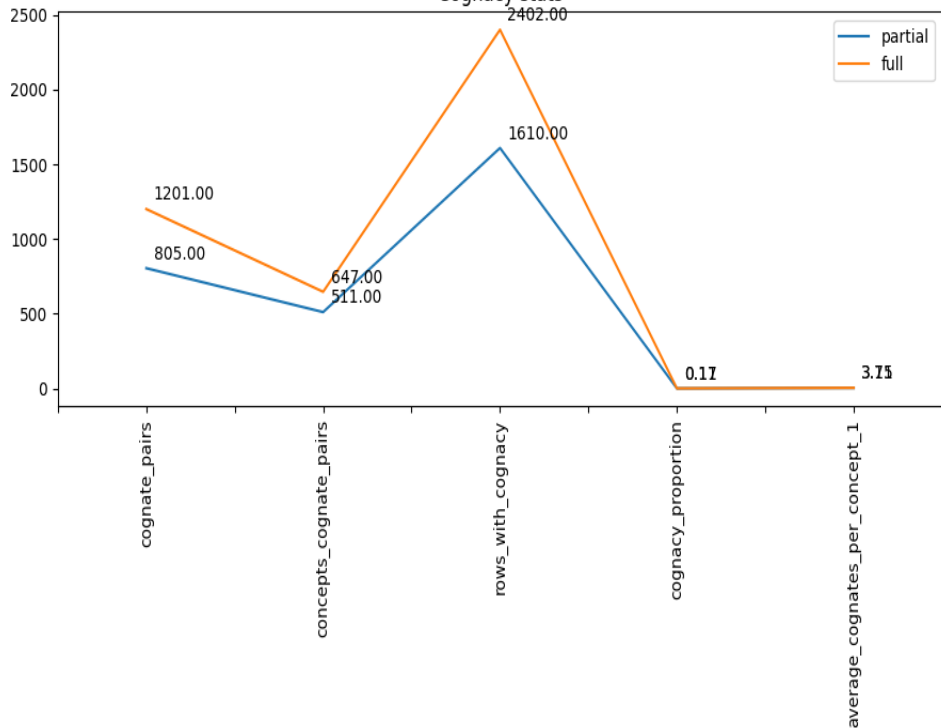


Results

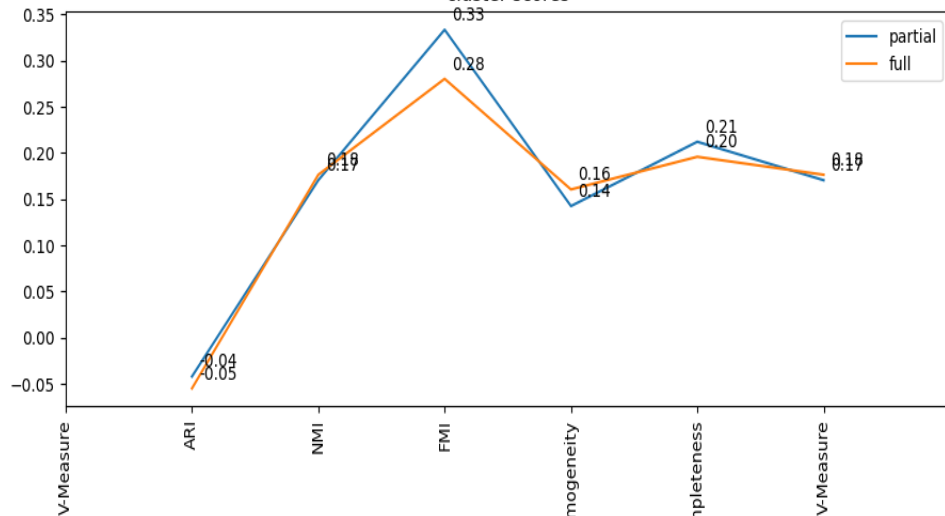
Morphological parsing

Full cognacy has better scores, but partial cognacy has better goodness.

Cognacy stats



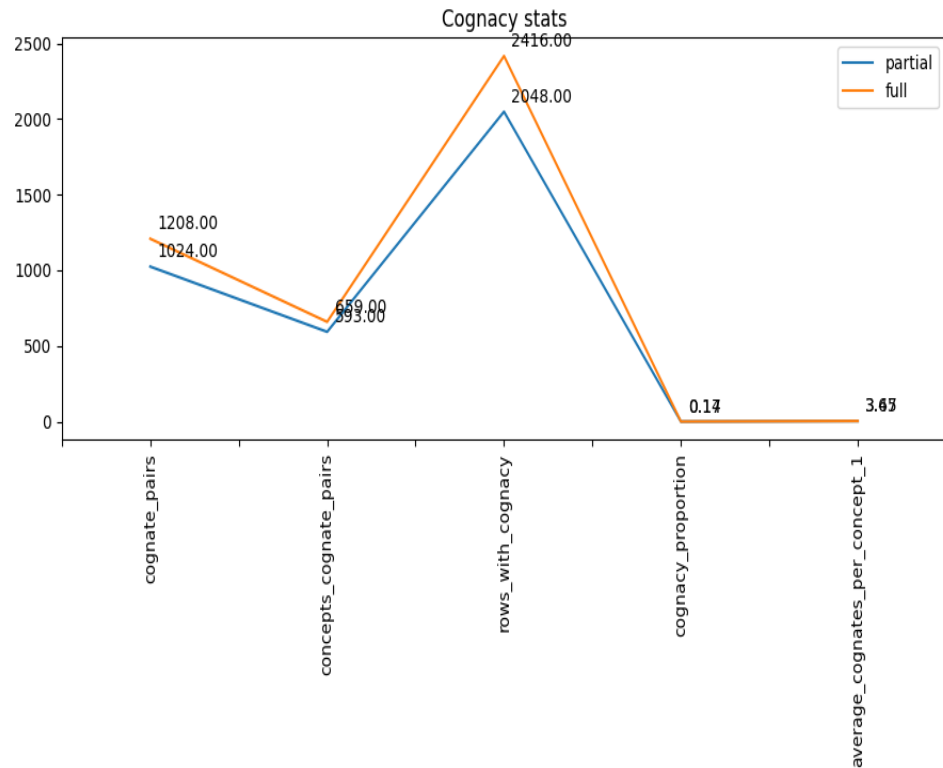
cluster scores



Results

Morpho-phonotactic

- Full cognacy performs better both in terms of scores and goodness.



Results

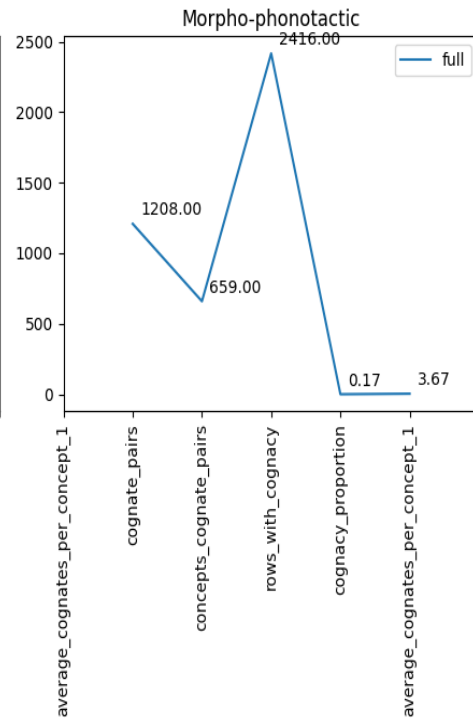
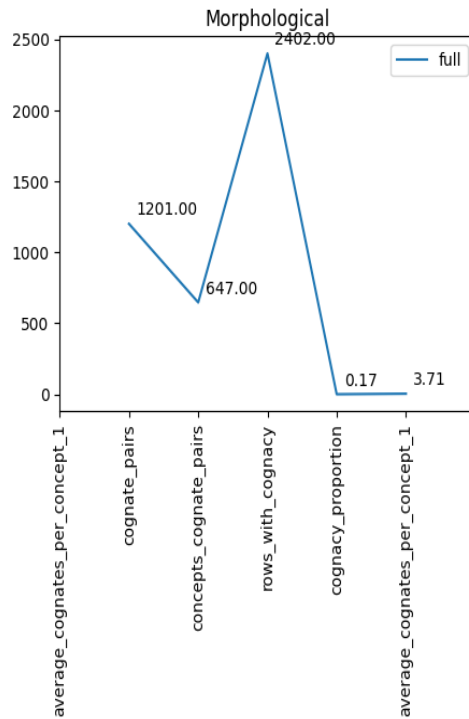
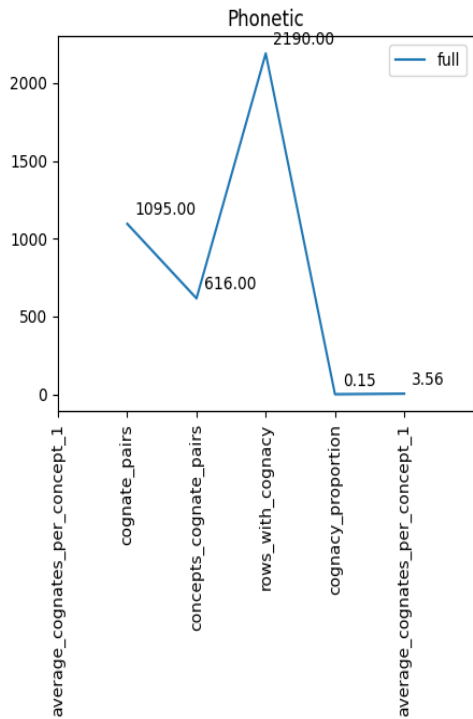
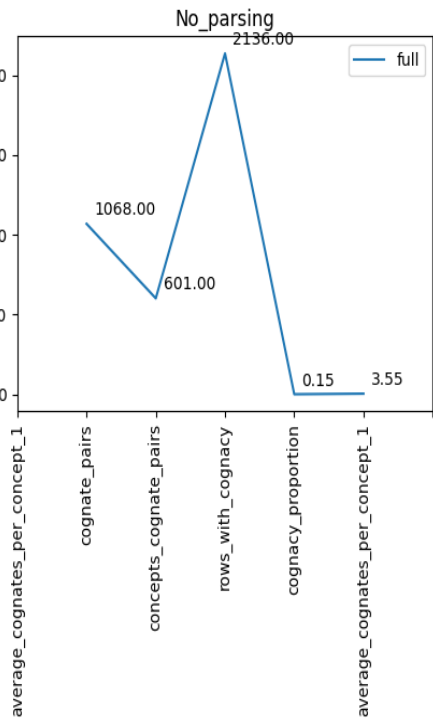
Interim summary 1

- Full cognacy produces better results when data is not parsed.
- Similar performance for both full and partial cognacy when data is parsed phonetically.
- Partial cognacy seems to perform better when data is morphologically parsed.
- Partial cognacy performs better once phonotactic information is introduced into morphological parsing.
- It is fascinating to note the correlation between linguistic level and cognacy type:
 - Phonetic = full cognacy
 - Morphological = partial cognacy
- Once phonotactic information is entered into morphological parsing, the trend reverses.

Results

Full cognacy, cognacy statistics

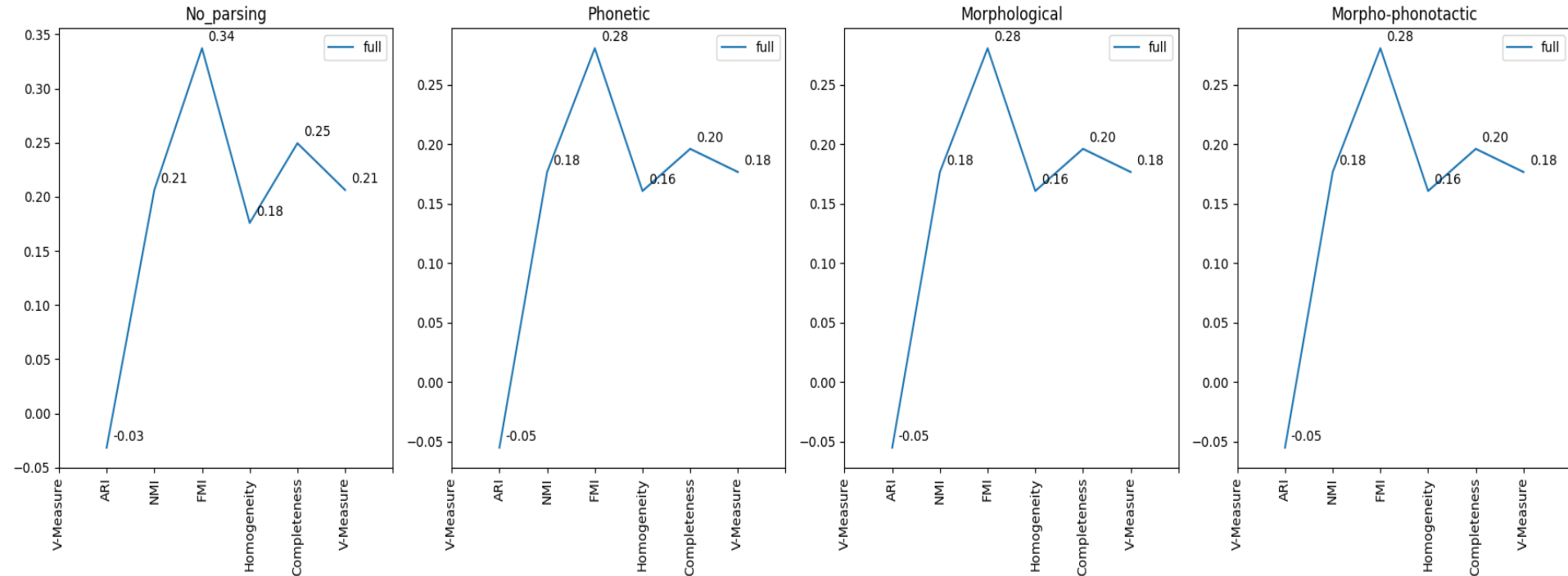
- Cognacy scores improve with increased parsing.



Results

Full cognacy, cluster goodness

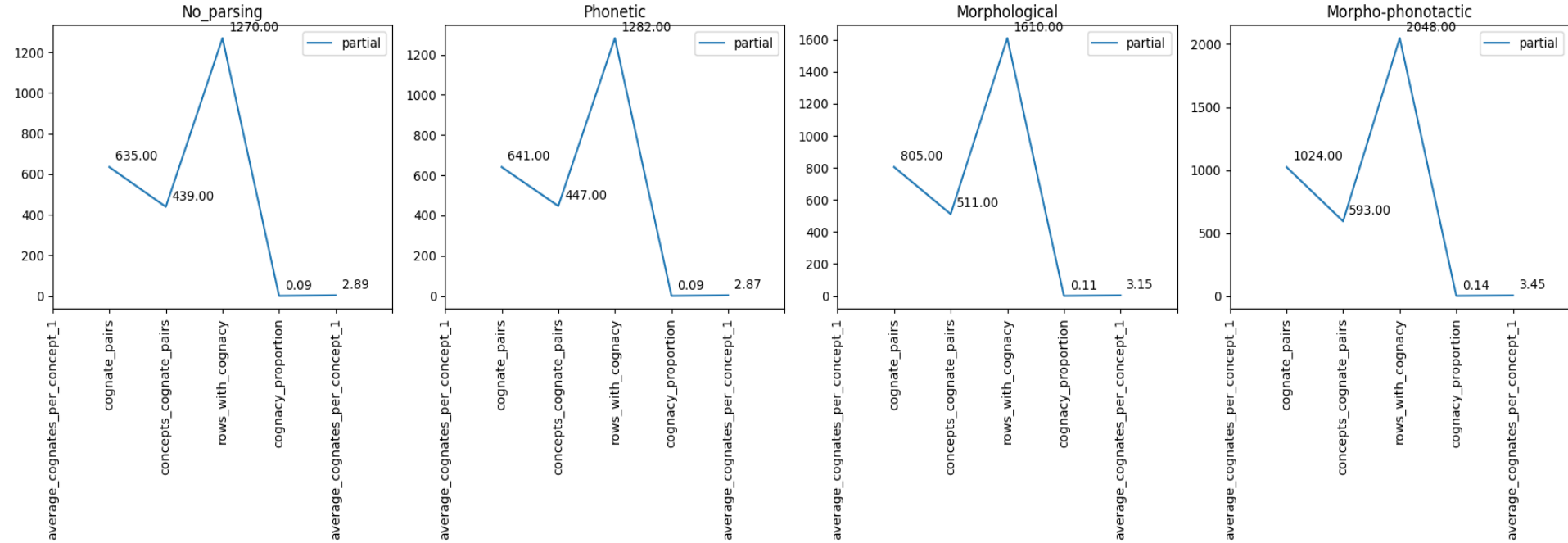
- Cognacy goodness decreases and stagnates with increased parsing.



Results

Partial cognacy, cognacy statistics

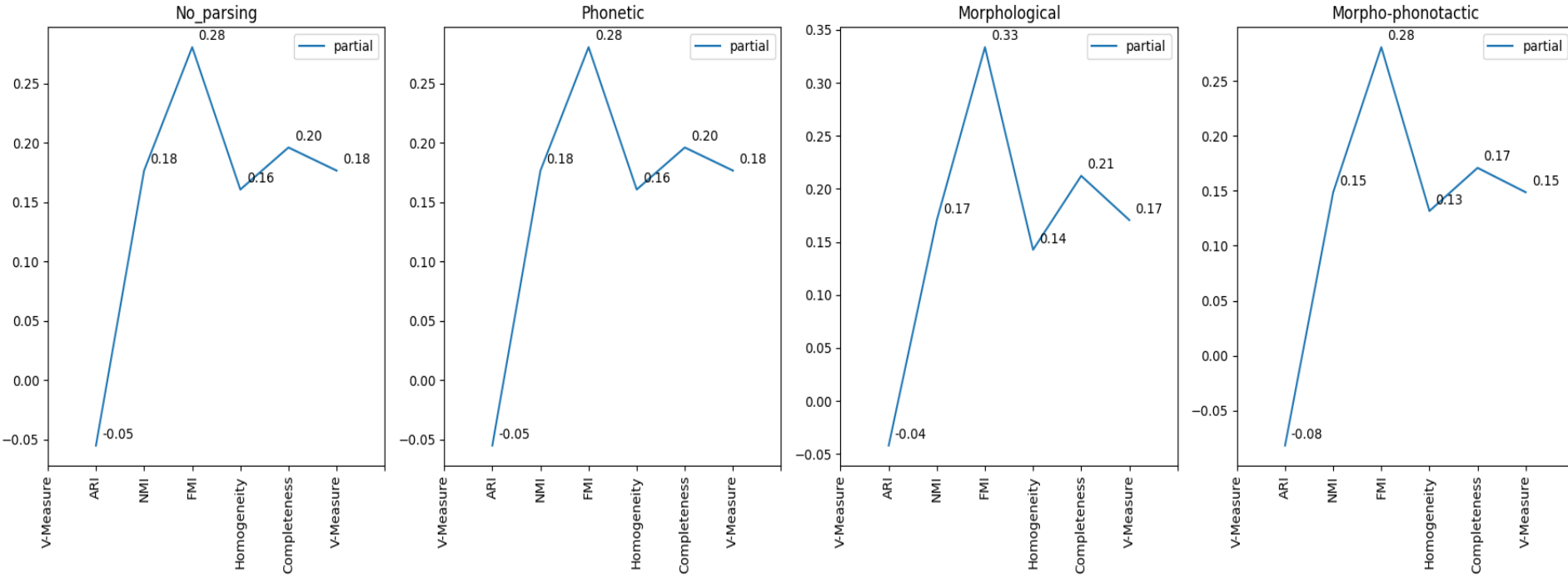
- Cognacy scores improve with increased parsing.



Results

Partial cognacy, cognacy goodness

- Goodness improves with morphological parsing but degrades with morpho-phonotactic.



Results

Interim summary 2

- For both full and partial cognacy, increased parsing increases cognacy scores.
- Cognacy goodness is not uniform:
 - Full cognacy: goodness drops with parsing, and then stagnates
 - Partial cognacy: on the average goodness increases beyond phonetic parsing, but degrades during morpho-phonotactic parsing (total of goodness measures = 0.95, 0.95, 0.98, 0.8 respectively)
- There seems to be a parsing threshold trigger for degraded performance of partial cognacy.
- “Degrader” trigger threshold seems to be at the phonotactic level

Results

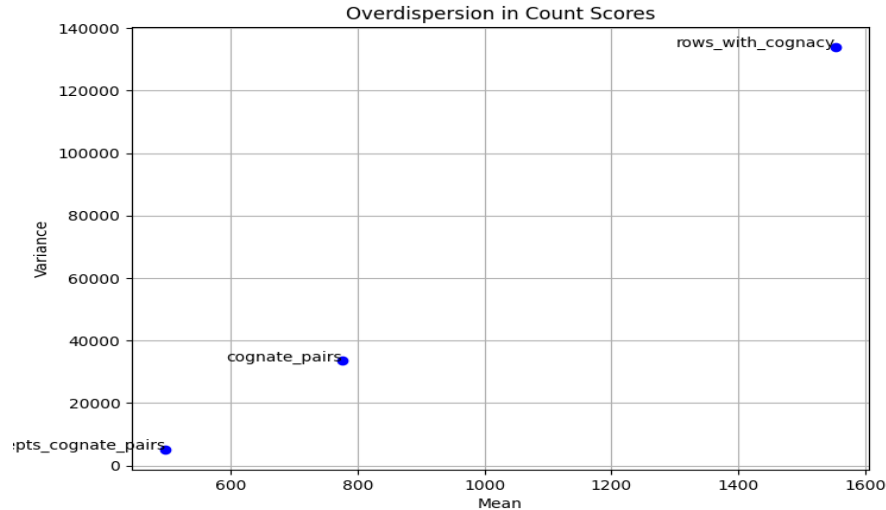
Discussion

- Given that cognacy statistics inherently favor full cognacy, cognacy goodness seems a better measure of performance.
- Results show that full cognacy performance can be characterized as unparsed vs parsed:
 - Performance degrades with parsing and stagnates no matter parsing level
- Partial cognacy on the other hand indicates an increased performance with increased parsing
 - This is however conditioned by a phonotactic threshold (noted degradation for morphonotactic parsing)
- Two facts are thus to be noted:
 - Full cognacy favors unparsed data
 - Partial cognacy favors parsed data, albeit with a phonotactic threshold

Results

Discussion

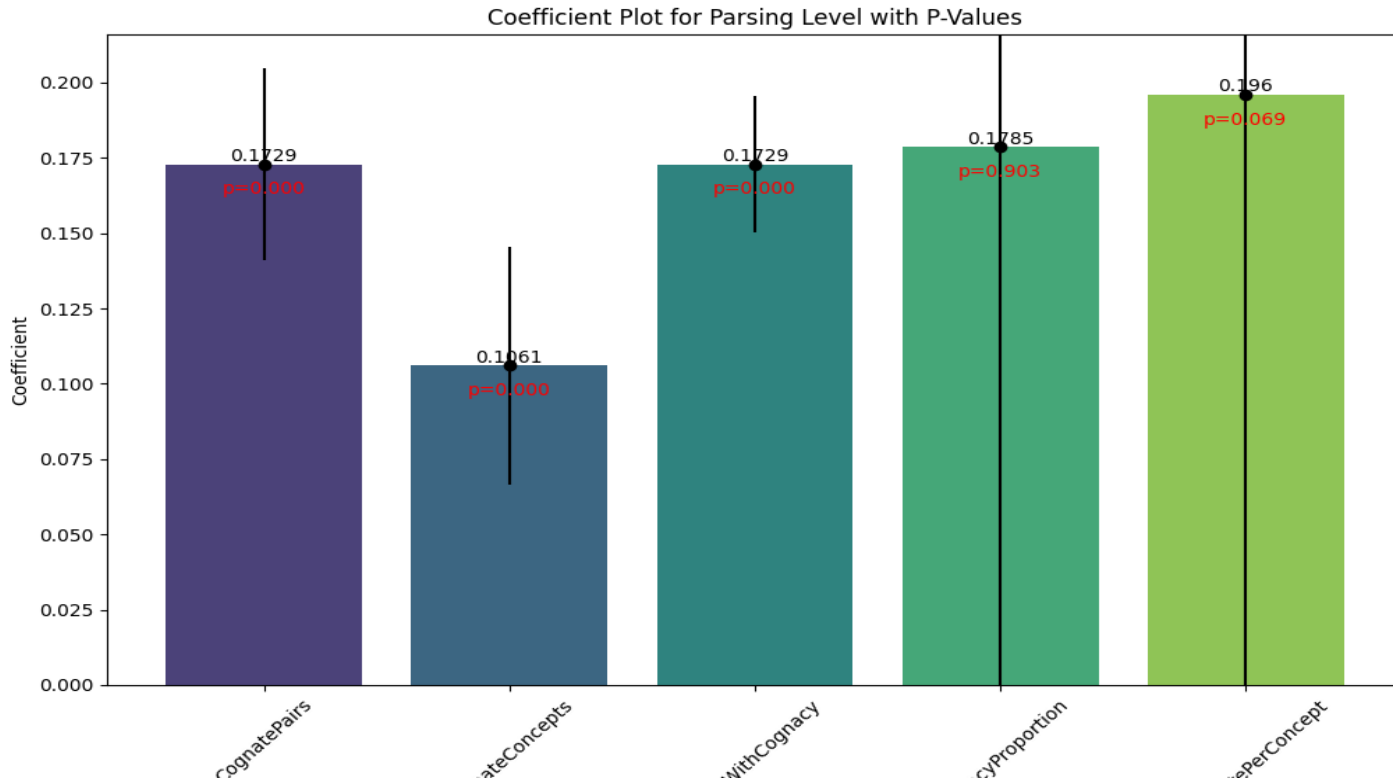
- Results can be modelled via:
 - Negative binomial regression: cognate pairs, cognate concepts, rows with cognacy.
 - Binomial logistic regression: cognacy proportion.
 - Linear regression: average cognate per concept.
 - Linear regression: cognacy goodness.



Results

Discussion

Positive coefficients confirm positive correlation between parsing and cognacy scores.



Results

Discussion

- Coefficient interpretation:
 - Parsing unit: none =>phonetic =>morphological =>morpho-phonotactic
- For variables modelled with a linear regression model, coefficient represents change in mean of variables.

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$$

- Percentage change is represented by:

$$\text{Percentage change} = \beta_i \times \frac{100}{\text{Mean of } X_i}$$

- A one-unit change in parsing level is thus associated with a 6 percent increase in average cognates per concept (mean=3.09, coefficient = 0.19).

Results

Discussion

- For variables modelled with a negative binomial regression or a binomial logistic regression, coefficients represent change in log of expected count for a one-unit change in parsing level.

$$\log(\lambda_i) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$$

- Exponentiated coefficients (IRR) indicate that for a one-unit increase in parsing level, expected cognacy score increases by a factor IRR holding all other variables constant (Buis 2010).

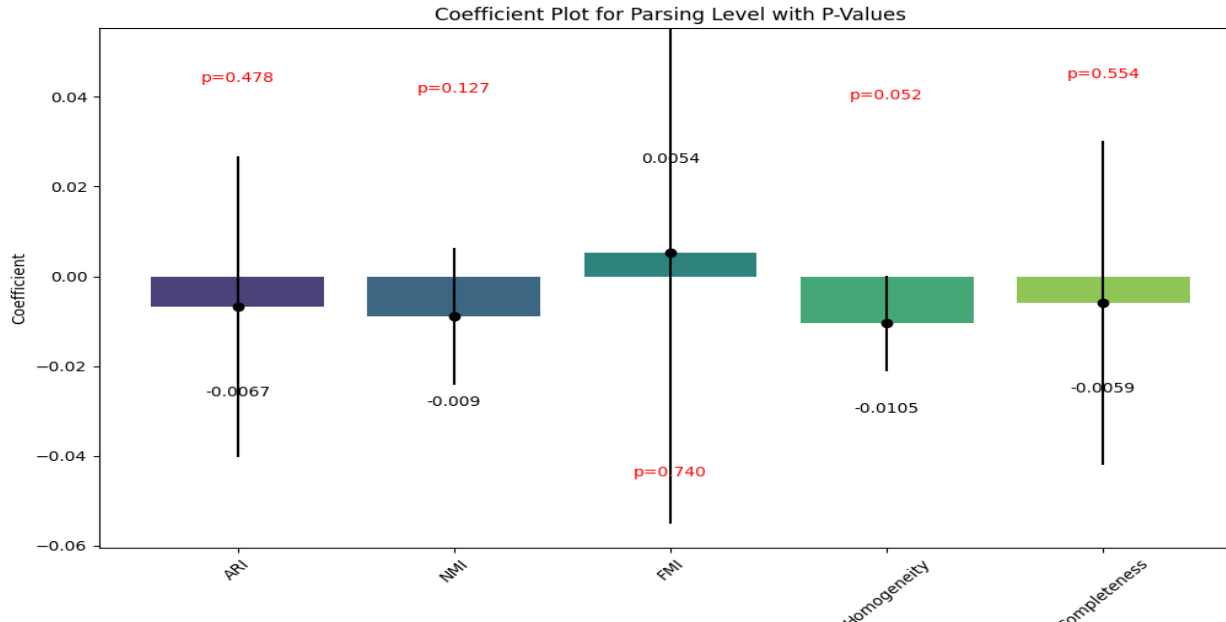
$$\text{Exp}(\beta_j) = e^{\beta_j}$$

Variable	Coefficient	Exponentiated coefficient(IRR)	Percentage Change(IRR-1)*100
Cognate pairs	0.168182	1.183152	18.3%
Cognate concepts	0.106065	1.111894	11.2%
Rows with cognacy	0.166776	1.181490	18.1%
Cognacy proportion	0.178540	1.195470	19.5%

Results

Discussion

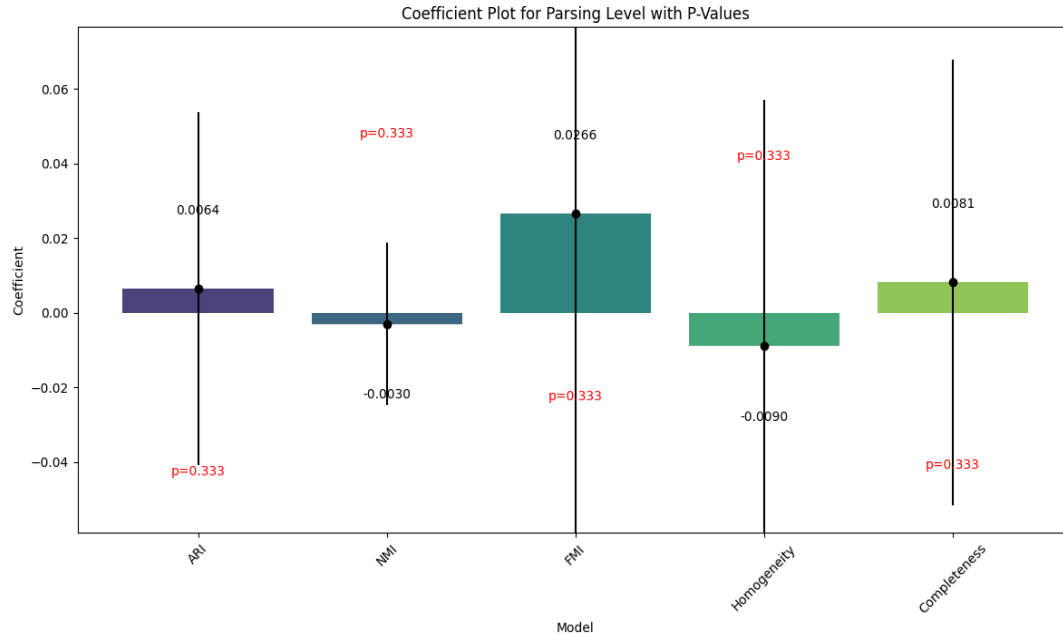
- The slightly negative coefficients confirm the overall rise-and-degrade relation between parsing and partial cognacy goodness.
- A one unit change in parsing level associates with the coefficient value changes in the means of the goodness scores.



Results

Discussion

- Modelled without morpho-phonological parsing, coefficients of partial cognacy goodness see improvement.
- A one unit change in parsing level now associates with a mostly positive change in means of goodness scores.



Conclusion

- This study has sought to understand the effects of morphological segmentation on the quality of computer –assisted comparative historical work.
- It has shown that, for unsegmented-unparsed data, full cognacy is the most adequate method; segmented data is best fitted for partial cognacy detection (or vice-versa).
- It has been demonstrated that while these effects are easy to observe in cognate scores, in cognacy goodness tests, the relationship is nuanced.
- For cognacy scores:
 - The higher the linguistic level of parsing, the higher the cognates detected
 - But, higher cognate numbers does not necessarily translate into higher goodness
- For cognacy goodness:
 - Goodness stagnates with higher parsing for full cognacy
 - Goodness accelerates with higher parsing for partial cognacy, but degrades with morpho-phonotactic information

Conclusion

- The suggestion then is that, cognate detection is further optimized when data is fed with linguistically richer information, with controls set for thresholds.
- The results obtained from this study are pertinent for two reasons:
 - Optimal methods detected in this study are being transferred to workflows defined for the BANG project.
 - Results obtained will influence data preprocessing techniques to be adopted during (new) data integration.
 - Other projects employing computational techniques can be inspired by the results.
- Results nevertheless raise few questions:
 - Which specific feature triggers goodness degradation (CV, CVC, CVCC etc.)?
 - Are there any relationships between ground-truth variables and parsing variables?
For eg. will a ground-truth constructed on morphological paradigms have any consequence for morphologically parsed data as opposed to other parsings?

Thank you

Bibliography

- Atkinson, Q. D., & Gray, R. D. (2006). How old is the Indo-European language family? Illumination or more moths to the flame. In *Phylogenetic methods and the prehistory of languages* (pp. 91-109).
- Blum, F., & List, J.-M. (2023). Trimming phonetic alignments improves the inference of sound correspondence patterns from multilingual wordlists. In *Proceedings of the 5th Workshop on Computational Typology and Multilingual NLP* (pp. 52-64). Association for Computational Linguistics.
- Buis, M. L. (2010). Stata tip 87: Interpretation of interactions in nonlinear models. *The stata journal*, 10(2), 305-308.
- Campbell, L. (2013). *Historical linguistics*. Edinburgh University Press.
- Cioabanu, A. M., & Dinu, L. P. (2014, June). Automatic detection of cognates using orthographic alignment. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)* (pp. 99-105).
- Creissels, D. (2019). Morphology in Niger-Congo languages. In *Oxford Research Encyclopedia of Linguistics*.
- Coxe, S., West, S. G., & Aiken, L. S. (2009). The analysis of count data: A gentle introduction to Poisson regression and its alternatives. *Journal of personality assessment*, 91(2), 121-136.
- Dimmendaal, G. J. (2008). Language ecology and linguistic diversity on the African continent. *Language and Linguistics Compass*, 2(5), 840-858.
- François, A., & Kalyan, S. (2023). Dialexification: A tool for studying cross-linguistic patterns of semantic change. In *16th International Cognitive Linguistics Conference*.
- Forkel, R., List, J. M., Greenhill, S. J., Rzymiski, C., Bank, S., Cysouw, M., ... & Gray, R. D. (2018). Cross-Linguistic Data Formats, advancing data sharing and re-use in comparative linguistics. *Scientific Data*, 5(1), 1-10.
- Forkel, R., Moran, S., List, J. M., Greenhill, S. J., Ashby, L. C., Gorman, K., & Kaiping, G. (2019). *Segments: Unicode standard tokenization routines and orthography profile segmentation* [Software Library, Version 2.1.3].
- Gray, R. D., & Atkinson, Q. D. (2003). Language-tree divergence times support the Anatolian theory of Indo-European origin. *Nature*, 426(6965), 435-439.
- Hantgan-Sonko, A. (2019). *Linguistic support for an early Dogon diffusion*. Poster presented at the Peopling History of Africa conference, Geneva, Switzerland.
- Hantgan-Sonko, A., & List, J. M. (2022). Bangime: Secret language, language isolate, or language island? A computer-assisted case study. *Papers in Historical Phonology*, 7, 1-43.
- Heath, J. (2015). *A grammar of Togo Kan (Dogon language family, Mali)*.
- Heath, J. (2017a). *A grammar of Bunoge (Dogon, Mali)*.
- Heath, J. (2017b). *A grammar of Tebul Ure (Dogon, Mali)*.
- Heath, J. (2023). *A grammar of Tebul Ure (Dogon, Mali)*.
- Hochstetler, J. L., Durieux, J. A., & Durieux-Boon, E. I. (2004). *Sociolinguistic survey of the Dogon language area*. SIL International.
- Jäger, G., List, J. M., & Sofroniev, P. (2017, April). Using support vector machines and state-of-the-art algorithms for phonetic alignment to identify cognates in multi-lingual wordlists. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers* (pp. 1205-1216).

Bibliography

- Kanojia, D., Dabre, R., Dewangan, S., Bhattacharyya, P., Haffari, G., & Kulkarni, M. (2021). Harnessing cross-lingual features to improve cognate detection for low-resource languages. *arXiv preprint arXiv:2112.08789*.
- Kanojia, D., Sharma, P., Ghodekar, S., Bhattacharyya, P., Haffari, G., & Kulkarni, M. (2021). Cognition-aware cognate detection. *arXiv preprint arXiv:2112.08087*.
- Kim, Y. M., Chang, K., Cui, C., & Mortensen, D. (2023). Transformed protoform reconstruction. *arXiv preprint arXiv:2307.01896*.
- Koch, H., & Hercus, L. (2013). Obscure vs. transparent cognates in linguistic reconstruction. In R. Mailhammer (Ed.), *Lexical and structural etymology* (pp. 33-51). De Gruyter.
- Kondrak, G. (2000). A new algorithm for the alignment of phonetic sequences. In *1st Meeting of the North American Chapter of the Association for Computational Linguistics*.
- Kondrak, G. (2001). Identifying cognates by phonetic and semantic similarity. In *Second Meeting of the North American Chapter of the Association for Computational Linguistics*.
- List, J. M. (2012). Multiple sequence alignment in historical linguistics. In *Proceedings of ConSOLE* (Vol. 19, pp. 241-260).
- List, J. M. (2017). *Historical language comparison with LingPy and EDICTOR*. <https://doi.org/10.5281/zenodo.1042205>
- List, J. M., Forkel, R., & Hill, N. W. (2022). A new framework for fast automated phonological reconstruction using trimmed alignments and sound correspondence patterns. *arXiv preprint arXiv:2204.04619*.
- List, J. M., Lopez, P., & Baptiste, E. (2016, August). Using sequence similarity networks to identify partial cognates in multilingual wordlists. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)* (pp. 599-605).
- List, J. M., Walworth, M., Greenhill, S. J., Tresoldi, T., & Forkel, R. (2018). Sequence comparison in computational historical linguistics. *Journal of Language Evolution*, 3(2), 130-144.
- List, J.-M. (2012a). LexStat: Automatic detection of cognates in multilingual wordlists. In *Proceedings of the EACL 2012 Joint Workshop of Visualization of Linguistic Patterns and Uncovering Language History from Multilingual Resources* (pp. 117-125). Stroudsburg.
- List, J.-M. (2012b). SCA: Phonetic alignment based on sound classes. In M. Slavkovic & D. Lassiter (Eds.), *New directions in logic, language, and computation* (pp. 32-51). Springer.
- List, J.-M., Greenhill, S., & Forkel, R. (2023). *LingPy Documentation Release 2.6*.
- List, J.-M., Walworth, M., Greenhill, S. J., & Tresoldi, T. (2017). The potential of automatic word comparison for historical linguistics. *PLOS ONE*, 12(1), 1-18.
- McPherson, L. E. (2014). *Replacive grammatical tone in the Dogon languages* (Doctoral dissertation, UCLA).
- Meelen, M., Hill, N. W., & Fellner, H. (2023). What are cognates?. *Papers in Historical Linguistics*, vol. 7, (pp. 44-80).
- Mikheev, A. (2022). Text segmentation. In R. Mitkov (Ed.), *The Oxford handbook of computational linguistics* (pp. 549-564). Oxford University Press.
- Moran, S., & Prokić, J. (2013). Investigating the relatedness of the endangered Dogon languages. *Literary and Linguistic Computing*, 28(4), 676-691.
- Moran, S., Forkel, R., & Heath, J. (2016). *Dogon and Bangime linguistics*. Max Planck Institute for the Science of Human History.
- Rama, T., List, J. M., Wahle, J., & Jäger, G. (2018). Are automatic methods for cognate detection good enough for phylogenetic reconstruction in historical linguistics? *arXiv preprint arXiv:1804.05416*.
- Weiss, M. (2015). The comparative method. In *The Routledge handbook of historical linguistics* (pp. 127-145). Routledge.

