



# Verba volant, scripta volant? Don't worry! There are computational solutions for prootword reconstruction

Liviu P. Dinu\*♣  
Ana Sabina Uban\*♣  
Alina- Maria Cristea\*♣  
Bogdan Iordache\*♣  
Teodor Marchitan\*♣  
Simona Georgescu\*♣  
Laurentiu Zoicas\*♣

\*University of Bucharest, ♣ Faculty of Mathematics and Computer Science, ♣ Faculty of Foreign Languages and Literature  
{ldinu, auban}@fmi.unibuc.ro, alinaciobanu20@gmail.com, iordache.bogdan1998@gmail.com, teodormarchitan@gmail.com {simona.g

## Abstract

We introduce a new database of cognate words and etymons for the five main Romance languages (Romanian, Italian, Spanish, Portuguese, French), the most comprehensive one to date with over 19,000 entries. We propose a strong benchmark for the automatic reconstruction of prootwords for Romance languages by applying a series of machine learning models and features on these data. The best results reach 90% accuracy in predicting the prootword of a given cognate set, surpassing existing state-of-the-art results for this task and showing that computational methods can be very useful in assisting linguists with prootword reconstruction.

## 1 Introduction and Related Work

Prootword reconstruction, consisting of recreating the words in a proto-language from their descendants in daughter languages, is central to the study of language evolution. As the foundation of historical linguistics [?, ?] and the basis for linguistic phylogeny [?, ?, ?, ?], prootword reconstruction offers

important pieces of information concerning the geographical and chronological dimensions of ancient communities [?, ?], at the same time, allowing an insight into the cognitive and cultural world of our ancestors. The traditional process of reconstructing ancient languages consists of the "comparative grammar-reconstruction" method [?, ?], and the etymological data thus obtained can be used as a source on human prehistory, corroborating the archaeological inventory [?], and providing the basis for 'linguistic paleontology' [?]. The reconstruction of a word automatically implies a reconstruction of the surrounding realities, both natural and socio-cultural. For example, the presence in different Indo-European languages of obviously related words for 'beech' or 'salmon' allowed the reconstruction of words from Proto-Indo-European and thus information about the elements of nature present in the immediate vicinity of the Indo-Europeans could be extracted. In the absence of any clear documentary or archaeological data, these lexical clues allowed the geographical identification of the Indo-European homeland, also facilitating the chronology of successive waves of separation of Indo-European languages from the common trunk.

In the case of Romance languages, although the

mother tongue - Latin - is attested, its presence in written texts is not an exhaustive source for linguistic, social, and historical analysis of the community that spoke it. It is now generally accepted that the spoken language represented a different diatactic, diaphasic, and diamesic variety from written language, used by the few educated people who decided to express themselves in writing [?]. The Latin language that we reconstruct from words inherited in Romance languages is thus the only concrete and reliable living variety of the language from which Romance languages originate, whether we call it oral/vulgar Latin or Proto-Romance. We will opt here for the name "Proto-Romance" when we refer to the language from which the Romance languages originate, as this corresponds to the concept of protolanguage and protoword [?].

Furthermore, there are still numerous clearly cognate words present in several Romance languages, whose etymons are not attested in Latin (nor in any other language from which it might have been borrowed). For example, in the case of It. *trovare* 'find', Fr. *trouver*, Cat. *trobar*, etymologists have hotly debated over the decades whether one should reconstruct the protoform \**tropare* or \**turbare* [?]. A series of cognates attested in all Romance geographical areas, like Rom. Inca 'moreover', It. *anche*, Old Fr. *anc*, Cat. *anc* etc., has triggered over 15 etymological hypotheses over the last century, still without a generally accepted solution.

Although etymologists' interest in reconstructing the protolanguages has risen over the years, they still encounter numerous gaps when using exclusively the classical, manual methods [?, ?]. As the task of pro-toward reconstruction plays an important role in historical linguistics, studies have gone beyond the comparative method in an attempt to automate the process [?, ?, ?, ?]. However, the task has been recognized as difficult and challenging. Computational pro-toward reconstruction is a fairly new direction of study, and consequently even state-of-the-art approaches have limitations. Complete automation of the reconstruction process is still a desideratum. [?] proposed two systems (Jakarta and Prague) that, combined, cover the steps of the comparative method for protolanguage reconstruction, and several

other approaches to reconstruct pro-towards computationally had been attempted previously [?, ?, ?]. The work of computational biologists such as Alexandre Bouchard-Cote, Russell Gray, Robert McMahon, and Mark Pagel, and co-authors took the protoword reconstruction one step further by applying methods from computational biology to the problem of the reconstruction of language history, often in collaboration with linguists [?, ?, ?, ?]. In recent years, researchers have introduced new methods for pro-toward reconstruction, based on modern computational techniques (for example, CRF, transformers, RNN, deep learning) [?, ?, ?, ?, ?, ?, ?, ?]. The computational methods are limited today by 1) the available data (sparse, inconsistent) and 2) by the insufficiency of linguistic knowledge embedded in the systems.

The latest computational results on Romance protoword reconstruction, in particular, are reported on the database of [?], which contains 8,799 cognates set in Latin, Italian, Spanish, Portuguese, French, and Romanian (not all full cognates set). This is a revision of the dataset of [?] (used with very good results in [?]) with the addition of cognates scraped from Wiktionary.

Starting with these remarks, our main contributions are:

1. We introduce a comprehensive Romance database for protoword reconstruction by processing RoBoCoP [?], the largest Romance cognate-borrowing database obtained from electronic dictionaries with etymological information of Romanian, Italian, Spanish, Portuguese, and French.
2. We propose a strong benchmark for automatic protoword reconstruction, by applying a set of machine learning models (using various feature sets and architectures) on any cognate set of Romance languages.

The rest of the paper is organized as follows: In Section 2 we present the database that we have created and offer details about the processing steps involved; in Section 3 we introduce our approach for the automatic protoword reconstruction, along with methodological details; the results of our proposed experi-

ments are fleshed out in Section 4; and a comprehensive error analysis is described in Section 5. The last section is dedicated to final remarks.

## 2 Data

A major inconvenience in Historical Linguistics in general, and in computational approaches of protoword reconstruction in particular is the scarcity of available data. Nonetheless, in the last few years, several initiatives have been undertaken in this direction. [?] developed a database of Latin protowords, further expanded by [?] with Wiktionary data. Recently, this dataset was extensively used for several studies [?, ?, ?]. In 2023, Dinu et al. (2023) published the most comprehensive database of Romance related words, named RoBoCop. It contains cognates and etymons in five Romance languages: Italian, Spanish, Portuguese, Romanian, and French. It has already been used with good results on prominent historical linguistic tasks such as cognate identification [?], cognate- borrowings discrimination [?], and determining the borrowing direction [?].

### 2.1 The ProtoRom Database

Starting with the RoBoCoP database [?], in order to obtain cognate sets with common etymons in the five Romance languages, we filtered out the words with Latin etymology. We then created maximal tuples of words in the Romance languages with the same etymon ( $\langle w_{L_i}, e \rangle$ ), where  $L_i$  are all the languages among the five where the etymon  $e$  engendered a word, and  $w_{L_i}$  are the corresponding words in each of the languages discussed. In cases where multiple words in  $L_i$  derive from the same etymon  $e$ , we created multiple tuples ( $\langle w_{L_i}, e \rangle$ ) with all possible combinations of cognate words  $\langle w_{L_i} \rangle$  and the same etymon  $e$ . For an example of such a case see Table 1.

We curated the obtained data, with the help of linguists. In the process, we discarded sets that contained irrelevant or erroneous information, e.g.: erroneous lexical forms (e.g. Lat. *videre* 'see' - It. *vedere* - Fr. *voir* - Ro. *videa* (correct: *video*); included a

verb form in any mood other than the infinitive (e.g. Lat. *videre* - Sp. *veas* (subjunctive) / *viendo* (gerundive) / etc.); retained the reflexive form of a verb (e.g. Lat. *ponere* 'put' - It. *porre* - Sp. *ponerse* (poner + reflexive pronoun *se*), etc.); or contained words derived on Romance ground (e.g. Lat. *dens* 'tooth' - It. *dente* - Ro. *dintos* (= *dinte* + suff.- *os*), etc.).

We were able to apply manual corrections for all these errors for the smaller subset of entries in the database that have a cognate in each of the five languages. For the rest of the full database ProtoRom, we applied a semi-automatic correction by lemmatizing the cognate words, using the default lemmatizers implemented in the spaCy library for each of the Romance languages. In all experiments described in the rest of the paper, we use the lemmas of the cognates instead of the original forms found in the dictionary.

In addition to the correct series thus retained, we integrated the database created by Reinheimer- Rippeanu [?], a high quality collection of cognate series manually selected from the etymological dictionaries of each Romance language, some of which still not digitized (which probably explains why certain cognate sets from this collection were not among ones in the RoBoCoP database). We thus obtained a new database of cognate sets.

The proposed database contains a total 39,973 full or partial cognate sets along with their etymons. For the experiments in this paper, we focus on the 19,222 entries with at least 2 cognates. We choose this subset in order to ensure the robustness of our experiments, focusing on Latin etymons that engendered at least two cognates in two different languages, and we ignore the entries with only one cognate for a given etymon. Going further, this restricted dataset will be referred to as ProtoRom<sup>3</sup>. A cognate set is composed of a tuple of words in different languages with a common etymon, where the tuple can be either a full set of 5 cognates or a partial set of 2 to 4 cognates, where the cognate in one or more of the languages is missing (the Latin etymon did not produce an attested word in these languages according to our sources).

There are 1,245 full cognate sets in the database, the rest being partial cognate sets. To facilitate

distinguishing between the two settings, we name the first one ProtoRom-all5, and the second one ProtoRom. When we leave out one of the languages, we can obtain more full sets of 4-tuples (sets with at least 4 cognates) as follows: 1,480 if we leave out Italian, 2,493 if we leave out French, 1,489 when we leave out Portuguese, 1,504 when we leave out Spanish, and 1,946 by leaving out Romanian. The statistics detailing the number of partial cognate sets in all combinations are shown in Table 2.

ProtoRom is the largest database of cognate sets for Romance languages so far, significantly exceeding the widely used database for this task [?], containing 8,799 cognate sets of Romanian, French, Italian, Spanish, Portuguese words and the corresponding Latin form (which, in turn, is an extension of Ciobanu and Dinu (2018)’s original database).

### 3 Methodology and Experiments

#### 3.1 Experimental Setting

For our experimental trials, we consider two settings: In the first one, we limit our dataset to only the full cognate sets (i.e. 5-tuples of cognates from each of the five languages, that originate from the same Latin etymon), while in the second one we consider all cognate sets (with at least two cognates from different languages, per etymon, as previously mentioned). The second setting uses the full breadth of our proposed dataset (ProtoRom-all5), whereas the first one is a strict subset (ProtoRom).

**Data splitting.** In order to train and validate our models, we split our datasets into 80% : 10% : 10% train-dev-test subsets. Because of the nature of the cognate sets, generating a language-level stratified split is a non-trivial task. Since a Latin etymon can produce more than one reflex in a given language, we end up with  $\prod_i \max(1, n_{L_i})$  cognate sets for a given etymon, where  $n_{L_i}$  is the number of reflexes generated by that etymon in language  $L_i$ .

We propose a random split methodology that achieves the following properties: A Latin etymon and all of its cognate sets are not allowed to be part

of more than one split; the raw number of cognate sets (i.e. entries in the dataset) follows the 80 : 10 : 10 distribution; the distribution of unique Latin etymons is also 80 : 10 : 10 ; for each of the five languages; and computing the distribution of unique reflexes in that language yields the same ratio across the splits. In other words, if we only keep the Latin etymons and their reflexes in only one language, we obtain a monolingual task with the same 80 : 10 : 10 split.

In order to perform these splits, we construct for each Latin etymon a 5-dimensional vector  $(n_{L_i})_i$  using the previous definition of  $n_{L_i}$ . In order to obtain a split of ratio  $0 < p < 1$ , we want to select such vectors that, when summed together, equal  $p \cdot (N_{L_i})_i$ , where  $N_{L_i}$  is the total number of unique reflexes from language  $L_i$ . In other words, we face a task equivalent to a five-dimensional knapsack problem, which is not feasible given the large total capacities. Considering that these vectors contain particularly small values, and are somewhat uniformly distributed, plus the large capacities that we have to fill, we are able to randomly select etymons and their associated cognate sets and add them to any of the three splits, as long as they fit. This approach yields the original split distribution with some small deviations (< 1%) .

Also note that after splitting the ProtoRom-all5 dataset, containing only the full cognate sets, we can use it as a starting point for splitting the rest of the ProtoRom dataset, thus ensuring that no training examples from one setting leaks into the validation of the other one.

**Features.** The proposed approaches can be split into two main categories: models for reconstructing the orthographical representation of the protowords using the orthographical form of modern cognates, and models that reconstruct the phonemic representation from phonetic transcriptions of modern cognates. Our extracted dataset essentially provides the necessary examples for the former, while for the latter we employ the eSpeak library to automatically generate the phonemic representations.

## 3.2 Models

We use a variety of machine learning models, including classical, neural, and transformer- based (pre-trained and trained from scratch for the task). We include methods used in previous papers on the topic and evaluate them on our larger dataset in order to provide a benchmark for the task of protoword reconstruction for Romance languages.

We experiment with a variety of models, including pre- trained large language models (LLMs) and current state- of- the- art models for protoword reconstruction with various architectures (probabilistic RNN, character- level transformer) adapted to our new database, as well as original solutions. In this way, we aim to provide a benchmark for the task of protoword reconstruction.

**CRF + reranking** We used an approach that relies on conditional random fields (CRFs), based on the method proposed by [?]. Firstly, we applied a sequence labeling method that produces the form of the Latin ancestors, for each modern language. The modern words are the sequences, and their characters are the tokens. We used character n- grams from the input words as features. We employed pairwise sequence alignment (Needleman and Wunsch, 1970) between modern words and protowords to obtain the labels for each token. Secondly, we defined several ensemble methods to take advantage of the information provided by all languages, in order to improve performance. We employed fusion methods based on the ranks in the n- best lists and the probability estimates provided by the individual classifiers for each possible production, in order to combine the outputs of the classifiers (n- best list of possible protowords) and to leverage information from all modern languages. For each word in the productions list, we multiply the rank of it with the confidence score given by the CRF model for each language; we sum up the multiplication scores for each word in the list and then rerank the productions based on these results.

**Probabilistic LSTM** We conducted experiments using a combination of recurrent neural networks with different dynamic programs and expectation- maximization techniques, as described in [?]. The overall system can be split in two stages: a) a mod-

elling stage, where we model the evolution of words by making small character- level edits to the ancestral form; for each language in the study, the distribution over newly created words is computed; b) an expectation- maximization stage, where the ancestral form is inferred; using words sampled from the posterior distribution, the expected edit count is computed and further used by the character- level recurrent neural network in order to optimize the next round of samples; the final reconstruction is the maximum likelihood word forms. This model requires a full tuple of cognates to be passed as input, so we only compute results for experiments on the ProtoRom- all5 set. Like the original authors, we only apply this model on the phonemic forms of words, since the probability distributions of edit operations used in the algorithm rely on a set of manually set features for each phoneme that are not similarly available for orthographical characters.

**Character- level transformer** The next experiments conducted in this research are based on the transformer model, proposed by [?]. Some critical changes in the architecture were made in order to be able to accept our samples format: multiple modern word sequences (one for each language) correspond to a single protoform sequence. A positional encoding is applied to each individual modern word sequence before concatenation. An additive language embedding is applied to the token embeddings alongside the positional encoding in order to make a difference between input tokens of different languages.

**Pre- trained LLM (Flan- T5)** We finally evaluate the capabilities of pretrained Large Language Models (LLMs) to solve our task. While LLMs are currently obtaining state- of- the- art performance across NLP tasks, our specific goal is unlike usual tasks included in benchmarks or in training data for LLMs, and it is strongly multilingual (including one dead language), so we suspect it might be a difficult task for an LLM. We choose to use a pretrained model and fine- tune it on our own training data in order to increase its chances to perform well. We use a "base" variant of the Flan- T5 model [?], and fine- tune the model using instructions including the prompt: "What is the etymon given the following cognates:", followed by a list of cognate and language

pairs formatted as " $< L_i > : < w_i >$ " and separated by new lines, where the list of cognate words  $w_i$  is their respective languages  $L_i$  can be arbitrarily long (from 2 to 5 cognates, in the case of our experiments). For evaluation, we attempt to generate multiple output sequences, which are used as a ranking for the etymon prediction.

One limitation of pretrained LLMs that we cannot overcome through fine-tuning is its alphabet, which contains mostly characters in the Latin graphical alphabet, which means that we can only

## 4 Results

The previously described methods have been applied on both ProtoRom and ProtoRom-all5 datasets, using the orthographical form of the cognates and Latin etymon, or alternatively the autogenerated phonemic representations (where the models were able to accommodate them). We also provide a comprehensive human evaluation of the results. Linguists from our team manually analyzed the entire list of results, and we present the most significant observations regarding the models' successes and failures. The linguists did not correct the protoforms proposed by the models, but only evaluated and commented on them in relation to current knowledge in the field of historical linguistics. The metrics used include accuracy, (normalized) edit distance, and  $Cov_i$ , with  $i \in \{1, 5, 10\}$ , which stands for an extended version of the accuracy metric, where a correct prediction is one where the model found the correct etymon within the first  $i$  etymons predicted by our method (this metric is computed for models that are able to output a ranked list of predictions - Flan-T5 and CRF-based models).

### 4.1 ProtoRom-all5 Results

Results obtained on the ProtoRom-all5 set are shown in Table 3. In terms of accuracy (or  $Cov_1$ ), the best results are obtained using the orthographical forms, with the CRF-rerank model, reaching 60.4%. From the perspective of the  $Cov_i$  metrics, it is remarkable that the CRF-rerank model obtains a  $Cov_{10}$  score above 82%.

The experiments using the phonemic forms produce weaker results, with the best accuracy reaching 55.8% in the top 1 predictions scenario. Nevertheless, the CRF approach is able to achieve an accuracy close to 80% when we consider the top 10 best ranked predictions.

The probabilistic RNN models achieve very poor performances, reaching a mean edit distance of 3.11 when trained on the phonemic representations.

### 4.2 ProtoRom Results

The best accuracy when training the orthographical models is achieved in this scenario by the Transformer model, closely surpassing 73% (Table 4). As for the  $Cov_i$  metrics, the Flan model remarkably obtains a  $Cov_{10}$  accuracy score of 85.4%, and an edit distance of 0.23.

Similarly to the previous scenario, the experiments using the phonemic forms produce weaker results, with the best accuracy reaching 66.8% via the Transformer model. These results represent a collection of baselines for protoword reconstruction using our proposed dataset configurations.

We believe the higher accuracy observed on the full dataset is simply due to the larger amount of available data. While ProtoRom-all5 is a subset that contains only complete cognate sets from each of the five studied languages (totaling 1,245 sets) the ProtoRom dataset includes sets of two, three, or four cognates, resulting in significantly more sets (19,222). This larger dataset allowed the models to learn more phonetic correspondences, thereby improving the reconstruction process. Even though they are not full sets of five cognates, the additional cognate sets in the full database seem to help the models learn more about their protowords. This learning process is closely similar to the human method of learning: with more examples, linguists can be more certain of particular correspondences or phonetic changes and can apply them in the reconstruction with much greater confidence.

## 5 Error analysis

This section is dedicated to a deeper dive into qualitatively quantifying the errors produced by the previously proposed models. Our objective is separating purely wrong predictions from "near misses", which may still provide value for linguists for the reasons discussed below.

The error analysis was manually conducted by the linguists from our team, who specialize in Romance languages. They did not modify the protoforms provided by the models in any way. Their only intervention was to distinguish forms that were genuinely erroneous from those whose differences from the dictionary form were either insignificant or represented a correct adjustment to the reality of Latin pronunciation. In the final quantitative analysis, forms in this category were therefore included in the list of correct predictions without any changes to their structure.

Through analyzing the errors, we have identified some patterns that typically reflect either an insufficient number of examples to support a particular phonetic change or the irregularity of the change itself. For example, the short tonic /u/ develops into Spanish /o/ in half of the cases, while it remains /u/ in the other half. In such scenarios, the model may not know which phonetic treatment the cognates underwent and might choose the wrong variant. Similarly, in cases of phonetic accidents, which are by nature irregular and unpredictable, the model cannot reconstruct the pre- accident form. Instead, it reconstructs the intermediate form between the classical word and its Romance descendants. Identifying and systematizing these errors can help improve future results by broadening the input with information related to sound changes.

Before analysing the errors, a few preliminary points should be made. Romance lexicography as a whole is graphocentric - it considers the written, classical Latin (CL) lexical variants as the basis for the Romance vocabulary, even though it goes without saying that vernacular languages, oral par excellence, developed from an oral language, in our case Proto-Romance (PR) [?]. In the latest methodology used in Romance etymology, developed within the DÉRom project [?], the etymological identification is based

strictly on the comparative grammar - reconstruction method, starting from the lexical forms that were used uninterruptedly in Romance languages. The lexemes attested in Classical Latin are only a written correlate, possibly further evidence of the existence of the form obtained by the methods of comparative historical linguistics.

In the light of these considerations, we find that some of the reconstructed variants classified as errors should actually be considered as positive results and evidence that the machine could work at the same level as a linguist applying traditional methods. By positive results instead of errors we mean cases - not a few - where the machine reconstructed exactly the phonetic form valid for oral Latin, at the expenses of the standard orthographical form as it is lemmatized in classical Latin dictionaries.

Cases where the word obtained and the one given by the dictionary did not completely match were automatically considered as errors, although sometimes it was not a mistake as such. Therefore, there are a number of protoforms which, although they appear in the list as inadvertences, are variants that should be taken into account with full attention by linguists. Some are no more wrong than the form in the dictionary, some are closer to the actual oral form than those provided by lexicographers, while some are exactly the form that historical linguists would have reconstructed using traditional methods based on the sound laws of each language (we discuss each case below). Therefore, protoforms obtained by the automatic methods proposed here are sometimes preferable to the lemmatized ones, and this is the most important thing we can expect from the machine.

Below, we provide a list of situations categorized as errors, but where the the automatic protoword reconstruction is either comparable or better than the version proposed by the dictionary, as it represents exactly the linguistic variant we should consider as intermediate between classical Latin and Romance languages.

- Protowords ending in -um instead of standard -us (lupum instead of lupus). The difference between the endings -us / -um did not properly exist in Proto-Romance, as the final consonant -

s/-m was no longer pronounced. Thus, if the etymological dictionaries provide the classical nominative form *lupus* as an etymon for Ro. *lup*, It. *lupo*, Fr. *loup* etc., but the computer reconstructs *lupum* - this latter variant is more correct from a grammatical point of view, since in general nouns are inherited from the accusative form (in our case ending in *-um*) and not from the nominative (ending in *-us*). Moreover, if it reconstructs *lupu*, this form is even more correct, being the real one, that reflects the pronunciation in the spoken language.

- The automatically reconstructed protoforms reflect phonetic features specific to Proto-Romance: monophthongation (*au* > *o*, e.g. CL *aucca* vs PR *oca*;  $\alpha \geq \alpha$ , e.g. *pena* vs *pena*;  $\alpha \geq \alpha$ , e.g. *hasitare* vs *esitare*); reduction of geminate consonants (*adductus* vs *adictum*); loss of the initial or intervocalic /h/ (*hasitare* vs *esitare*; *cohaerente* vs *coerente*); phonetic adaptation of Greek loanwords to the Latin pronunciation (*y* > *i*, e.g. CL *byzantinus* vs PR *bizantinus*, the aspirate consonants become oclusive, *th* > *t* (CL *citharodeu* vs PR *citaredu*), *ph* > *f* (CL *phalange* vs PR *falange*); assimilations (CL *admonere* vs PR *ammonire*); simplification of consonant clusters (CL *sculptore* vs PR *scultore*, *temptare* vs *tentare*, *unctura* vs *untura*); changes in the pronunciation of vowels (CL *guttu* vs PR *gotu*, *misculare* vs *mescolare*, *sicare* vs *sec(c)are*, *occidere* vs *ucidere*, *calcea* vs *calcia*).
- Certain reconstructed etyma retain accidental phonetic changes that must be presupposed for a particular geolinguistic area (Sp. *queso*, Pt. *queixo* imply the metathesis PR *caesu* instead of CL *caseu*, Ro *plop*, It. *pioppo*, Sp. *chopo* lead to the protoform with metathesis *plopu*, correctly identified by the machine, instead of CL *populus*), or for the global PR variety (Ro. *doamma*, It. *donna*, Sp. *doina*, lead to the syncopated protoform *donna*, reconstructed by the machine, instead of CL *domina*, registered in lexicography).
- The automatically reconstructed protoforms

may mirror morphologic changes that underlie the subsequent Romance developments: nouns of the 5th declension undergo a shift to the 1st declension (CL *canities* vs PR *canitia*, *species* vs *specia*); verbs shifting from middle-passive to the active voice (CL *renasci* vs PR *renascere*).

- The computer has reconstructed the oblique case forms representing the basis from which the Romance nouns were inherited (nominative *flos* vs oblique case *flore-* > Ro. *floare*, It. *fiore*, Fr. *fleur*, etc.; *civitas* vs *civitate* > Ro. *cetate*, Sp. *ciudad*, etc.), or the plural instead of the singular form, when the Romance lexemes descend from the former (sg. *capitium* vs pl. *capitia* > Sp. *cabeza*, Pt. *cabeça*).

The real errors in the experiments we developed stem primarily from lexicographic omissions or mistakes, as well as in the imprecise methodology employed by the Ibero-Romance dictionaries consulted, namely the lack of any distinction between inherited and borrowed Latin words [?]. This latter inaccuracy leads to a misinterpretation of the phonetic correspondences by the computer, given that only the inherited words, not the borrowed ones, underwent regular sound change. Therefore, if we put together Ro. *roata*, Sp. *rueda*, Pt. *roda*, with Ro. *rotatie*, Sp. *rotacion*, Pt. *rotacao*, the computer will not be able to correctly infer the correspondence *t/d/d* and will confuse it with *t/t/t*, also assuming the series *d/d/d*. Therefore, some reconstructions, especially in the case of words circumscribed only to Ibero-Romance languages, could not take this sound law into account (e.g., on the basis of Sp. *miedo*, Pt. *medo*, the computer could not reconstruct *metus*, but proposed *medus*, which is wrong). This kind of shortcomings will be easily overcome in the future, firstly by clearly establishing, in the ProtoRom database, the inheritance- borrowing distinction, and secondly by extending the input provided to the computer with a number of basic phonetic laws.

**Revised performance scores.** Looking at the best reported predictions, we can apply the linguistic observations stated in the previous section and count which wrong predictions can be actually considered acceptable errors. Thus using these recovered

predictions, the best models’ scores would change as follows:

the orthographical Transformer accuracy for the ProtoRom dataset increases from 73.1% to 82.7% (135 out of the 575 original errors were recovered). the Flan model’s  $Cov_{10}$  accuracy on ProtoRom increases from 85.4% to 89.6% (90 out of the 311 original errors were recovered). the  $Cov_{10}$  accuracy for the orthographical CRF model trained on ProtoRom-all5 increases from 82.1% to 90.7% (11 out of the 23 original errors were recovered).

## 6 Conclusion

In this paper, we built a new dataset for automatic. protoword reconstruction, consisting of 19,222 cognate sets from five Romance languages (Romanian, Italian, Spanish, Portuguese, French). This is to date the largest database of its kind, surpassing its predecessor which totals 8,799 cognate sets.

We also proposed a series of comprehensive benchmarks ranging from deep- learning approaches, using LLMs and Transformer- based architectures, to more classical algorithms such as CRFs, some of which achieved performances of more than 85% accuracy when allowing multiple generated reconstructions.

An in- depth linguistic analysis of the erroneous reconstructions was also performed using the predictions of the best performing models. This attempt shed some light on the various categories of mistakes, out of which several could be considered acceptable. When ignoring the aforementioned acceptable errors, we were able to surpass 90% accuracies. We consider this an important distinction, since in our view similar tools should aim at assisting linguists in their scientific endeavours. Raw metrics are useful to compare computational methods, but, in order to assess their usability, a more qualitative inspection of the results should be performed. We hope through our research to incentivize further analysis.

As for future work, we are looking into an additional refinement of the current cognate sets, but also extending the database with more examples, including properly validated monolingual Latin reflexes that were excluded from our experiments for robust-

ness sake. We also intend to expand past the proposed benchmarks with more novel approaches, relying on both the proposed dataset and the additional contents of its parent database, RoBoCoP.

## Limitations

One limitation of the current work stems from the automatic generation of the phonetic representations via a third- party library (eSpeak). Although this approach was employed successfully in previous studies, the quality of the generated phonemes has a higher variance when comparing high- resourced languages to lower- resourced ones (such as Romanian, or even Latin).

Also, in this study we used the generated phonetic forms without any extra preprocessing steps, in order to have a representation of the pronunciation that is as accurate as possible. Removing phonetic markers (such as stress markers) from these representations may turn the generation task into a somewhat easier one, since currently the phonetic models are tasked with predicting the stressed sounds too.

In terms of resources, existing LLMs are mostly targeting orthographical texts, making any reasonable attempt at generating phonetic ones very difficult.

## Ethics Statement

There are no ethical issues that could result from the publication of our work. Our experiments comply with all license agreements of the data sources used. We make the contents of our package available for research purposes upon request.

## Acknowledgements

This work was supported by a mobility project of the Romanian Ministry of Research, Innovation and Digitization, CNCS - UEFISCDI, project number PN-IV- P2- 2.2- MC- 2024- 0461, within PNCDI IV. We want to thank the reviewers for their useful sug-

gestions and Diana Grigore, Cosmin Petrescu, Ioana Pintilie for their help in developing the algorithms.

## References