

Simonetta Montemagni

Istituto di Linguistica Computazionale “Antonio Zampolli” – CNR – Pisa (Italy)

**Report of the joint research activity carried out
from June 13th to July 4th 2010**

**at the University of Groningen - Faculty of Arts – Center for Language and Cognition
Groningen (CLCG) directed by Prof. John Nerbonne**

Title of Research Program

***Modelli computazionali della variazione dialettale e fattori linguistici sottostanti*
*Computational Models of Dialectal Variation and Underlying Linguistic Features***

1. Introduction

The transition from the description of the geographic distribution of individual linguistic features to a more abstract level of description intended to make generalizations on diatopic variation is now made possible by the use of dialectometric techniques that have proved particularly promising in the study of language variation in different languages and dialects, also typologically very distant: dialectometric studies have a tradition of over thirty years, since the pioneering studies of Seguy (1971) to the more recent analyses by Goebel (1984, 2005) and Nerbonne (Nerbonne et al. 1999, 2001). The greatest strength of this approach is to disregard the individual data that have contributed to the observed patterns of linguistic variation and the possibility of an “aggregate” analysis of increasingly large amounts of data, such as the entire corpus of a linguistic atlas. However, abstracting from the individual data is in danger of losing “the connection to the linguistic characterization” (Nerbonne, in press), aspect – this one - that makes dialectometric analyses not particularly interesting to the eye of the community of linguists and dialectologists. Michele Loporcaro (2009) effectively describes this view: “dialectometry measures the structural distances without passing through a rationalization of linguistic structure.”

One response to such criticism has been advanced recently by Wieling and Nerbonne (2009, 2010), who used a technique of co-clustering (called “bipartite spectral graph partitioning”) to identify dialects on the basis of aggregate large corpus of dialect and simultaneously identifying the underlying linguistic basis. In particular, through this technique it is possible to understand which factors underly the identified patterns of dialectal variation, the role played by each of them and the weight associated with them. In this way, the gap between models of linguistic variation based on quantitative analyses and more traditional analyses based on specific linguistic features is significantly reduced. Achieved results for Dutch dialects turned out to be particularly promising.

This report illustrates the application and specialization of the technique of “hierarchical bipartite spectral graph partitioning” (Wieling and Nerbonne, 2010) with respect to the dialectal corpus of the *Atlante Lessicale Toscano* (‘Lexical Atlas of Tuscany’, henceforth ALT) and discusses achieved results. The analysis focuses on the level of phonetic

variation: this is the level of analysis for which an aggregate analysis of the ALT dialectal corpus has provided divergent results compared to the analyses by Giannelli (1976, 2000) and Pellegrini (1977), as documented in Montemagni (2007, 2008). Phonetic variation in Tuscany thus provides a particularly challenging case study to test the potential of this new analysis technique to study models of linguistic variation.

2. Construction of the experimental data set

2.1. The data source

This study on Tuscan phonetic variation is based on the corpus of dialectal data of the *Atlante Lessicale Toscano* (ALT, Giacomelli et al., 2000). ALT is a specially designed linguistic atlas in which dialectal data have both a diatopic and diastratic characterization. The adjectives qualifying this linguistic atlas in its name are “lexical” and “Tuscan”. ALT is lexical in the sense that its main focus is on lexical variation but this does not exclude that it contains valuable information for what concerns e.g. phonetic or morphological variation. ALT is Tuscan in the sense that it is a regional atlas focusing on dialectal variation within Tuscany, a region where both Tuscan and non-Tuscan dialects are spoken; the latter is the case of dialects in the north, namely Lunigiana and small areas of the Apennines (so-called Romagna Toscana), which rather belong to the group of Gallo-Italian dialects.

ALT interviews were carried out in 224 localities of Tuscany, with 2,193 informants selected with respect to a number of parameters ranging from age and socio-economic status to education and culture. The interviews were conducted by a group of trained fieldworkers who employed a questionnaire of 745 target items, designed to elicit variation mainly in vocabulary, semantics and phonetics. A dialectal corpus with these features lends itself to investigations concerning geographic or horizontal (diatopic) variation as well as social or vertical (diastratic) variation: in this study we will focus on the diatopic dimension of linguistic variation. ALT is now available as an on-line resource, ALT-Web (<http://serverdbt.ilc.cnr.it/altweb/>; for more details see Montemagni et al. 2006).

ALT data were collected between 1974 and 1986, resulting in millions of responses (tokens) from the 2,193 speakers who were each asked 745 questions, corresponding to more than 84,000 different attested dialectal items (types). During the collection phase, the results of interviews carried out by the group of trained fieldworkers were revised by the head of the project, Gabriella Giacomelli, in order to guarantee comparability of collected data and reduce as much as possible potentially misleading effects deriving from fieldworker's collection techniques or transcription peculiarities.

In ALT, all dialectal items were phonetically transcribed. In order to ensure a proper treatment of these data, an articulated encoding schema was devised in ALT-Web in which all dialectal items are assigned different levels of representation: a first level rendering the original phonetic transcription as recorded by fieldworkers; other levels containing representations encoded in standard Italian orthography. In this multi-level representation scheme, dialectal data are encoded in layers of progressively decreasing detail going from phonetic transcription to different levels of orthographic representations eventually abstracting away from details of the speakers' phonetic realisation.

For the specific concerns of this study, we will focus on the representation levels of a) phonetic transcription (henceforth, PT), and b) normalised representation (henceforth, NR)

where the latter is the representation level meant to abstract away from within-Tuscany vital phonetic variation.

At the NR level a wide range of phonetic variants is assigned the same normalised form: e.g. words such as [skja'ttʃata], [skja'ttʃaθa], [skja'ttʃada], [skja'ttʃaða], [stja'ttʃata], [stja'ttʃaθa], [stja'ttʃada], [stja'ttʃaða], [stʃa'ssɛda] etc. (denoting a traditional type of bread, flat and crispy, seasoned on top with salt and oil) are all assigned the same normalised form, SCHIACCIATA. Note that at this level neutralisation is only concerned with phonetic variants resulting from productive phonetic processes: this is the case, for instance, of variants involving spirantization or voicing of plosives like /t/, as in [skja'ttʃaθa] and [skja'ttʃada]. On the contrary, there are word forms like [kaʎʎo] and [gaʎʎo] (meaning 'rennet') which are assigned distinct NRs, CAGLIO and GAGLIO respectively: this follows from the fact that the [k] vs [g] alternation in word-initial context represents a no longer productive phonetic process in Tuscany. It should also be noted that the NR level does not deal with morphological variation (neither inflectional nor derivational). This entails that words such as [skja'ttʃata] (singular) and [skja'ttʃate] (plural) as well as [skjattʃa'tina] (diminutive) are all assigned different NFs. Currently, NR is the most abstract representation level in ALT-Web.

2.2. Dialectal data selection

For this study of phonetic variation, phonetic transcription was taken as the starting point. The alignment of the different representation levels was exploited to automatically extract all attested phonetic variants of the same normalised word form (henceforth, NF). In practice, the various phonetic realisations of the same lexical unit were identified by selecting all phonetically transcribed dialectal items sharing the same NF, as exemplified in Table 1 for the normalised form SCHIACCIATA.

Location	Phonetic variant	NF
15 Vergemoli	[sca'ttʃata]	SCHIACCIATA
16 Pieve Fosciana	[sca'ttʃada]	
18 San Pellegrino in Alpe	[sca'ttʃata], [stʲa'ttʃata]	
19 Brandeglio	[sca'ttʃaθa], [stʲa'ttʃaθa]	
22 Prunetta	[stja'ttʃaθa]	
23 Orsigna	[skja'ttʃaθa], [sca'ttʃaθa], [stja'ttʃata]	
24 Spedaletto	[stja'ttʃaθa]	
25 Castello di Sambuca	[sca'ttʃada]	
28 Barberino di Mugello	[skja'ttʃata], [stja'ttʃata]	
...

Table 1 – Excerpt from the experimental data set used for this study

Since the ALT-Web normalised representation level does not abstract away from either morphological variation or no longer productive phonetic processes, we can be quite sure that phonetic distances calculated against phonetic variants of the same NF testify vital phonetic processes only, without influence from any other linguistic description level (e.g. morphology).

In particular, the whole set of 34,912 normalised forms attested in the ALT dialectal corpus was taken into account. For 20,671 normalised forms (59.20%) attested variation (if any) occurs within a single locality; on the other hand, there are 4,688 normalised forms

(13.42%) showing no phonetic variation at all, in spite of their being attested in different locations (with geographical coverage ranging from 2 to 206). Since both cases are of no value in assessing diatopic phonetic variation, they have been removed from the data set which served as the basis of this study. There remained 9,553 normalised forms having at least two different phonetic realisations and being attested in at least two different locations. The graph in Figure 1 shows the geographical coverage and the phonetic variability range for the selected 9,553 normalised forms.

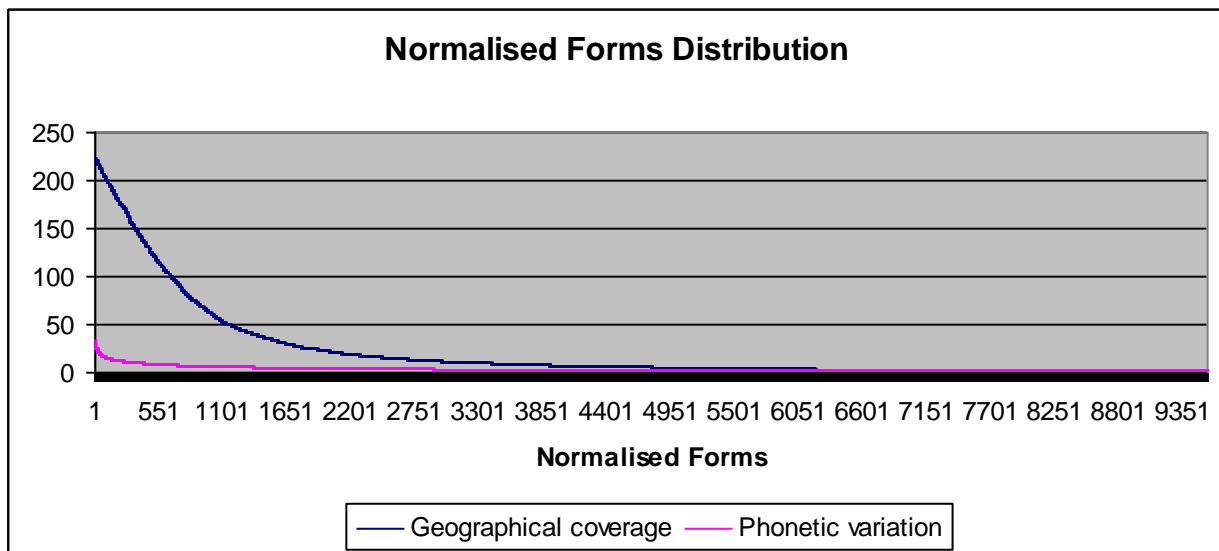


Figure 1 - Geographical coverage and phonetic variability range for the selected 9,553 normalised forms

Geographical coverage ranges between 2 and 223, and phonetic variability between 2 and 34: it should be noted, however, that within this set only 3,397 normalised forms (i.e. 35.55%) are attested in at least 10 different locations and only 1,920 show a phonetic variability range greater than 4 (corresponding to only 20.09% of NFs).

For the specific concern of this study, the following constraints have been enforced for the definition of the experimental data set: we focussed on dialectal items with a geographical coverage of at least 100 locations and showing at least 5 phonetic variants, corresponding to 523 normalised forms (5.47% of the whole sample).

The selected data set included adjectives, nouns and verbs. Due to the fact that in ALT verbal answers can be represented by different inflected forms (typically, the infinitival form, but also third person singular of the present indicative, or past participle) which are not always explicitly marked, we removed them from the experimental dataset in order to prevent potential noise deriving from verbal morphology. In this way, the set of selected normalised forms was reduced to 444 (4.64% of the whole set of diatopically varying NFs), including adjectives and nouns in the form of both single words and multi-word expressions. Note that selected multi-word expressions were represented by “frozen” word combinations, thus not showing variability due to the insertion/deletion of constituents.

In order to test the representativity of the selected sample of 444 NFs with respect to the whole set of normalised forms having at least two phonetic variants attested in at least two locations (used in Montemagni, 2008), we measured the correlation between overall phonetic distances and phonetic distances focussing on the selected sample which turned

out to be very high, with $r=0.923$. We can thus conclude that the selected sample can be usefully exploited to reliably study the patterns of phonetic variation in Tuscany.

2.3. Using atlas data as a corpus

From what has been said so far, it should be clear that here we are using atlas data in quite a peculiar way. Although this study is based on atlas data, it uses them as a corpus. This is to say that the dialectometric analysis of Tuscan phonetic variation here is not based on a predefined set of questionnaire items which were specifically designed to investigate the geographic distribution of phonetic features. Rather, it took the whole set of ALT attested lexical items, which were elicited from informants with quite different (mainly, lexico-semantic) purposes, and used it for studying phonetic variation.

By using atlas data as a corpus, the problem of inherently subjective feature selection is significantly reduced, thus providing a more “realistic” linguistic signal (Szmrecsanyi, to appear, p.3). On the other hand, by using atlas data as a corpus one of the main advantages usually ascribed to atlas-based studies, namely the areal coverage of dialectal items, can no longer be taken for granted. As we have seen, the areal coverage of attested NFs ranges from 1 to 223 locations: to overcome this potential problem, we enforced a minimal areal coverage threshold, corresponding to 100 locations (see section 2.2).

2.4. Dialectal data preparation

Having defined the extra-linguistic constraints which guided the definition of the experimental data set, all phonetic variants of the selected normalised forms were extracted. Extracted phonetic variants were enriched with information about the informants who attested them and were converted to IPA representation.

2.4.1. Extracting informants information

In previous studies based on ALT data (Montemagni 2007, 2008), phonetic variants of the same NF were used in a purely “categorical” way. This appears as a simplification, since the coordinates of each ALT item are not restricted to the location in which it was attested but also include the informants who testified it. This entails that for each attested phonetic variant we also know the number of informants who attested it, together with their socio-cultural profile.

For the specific concerns of this study, two different versions of the selected data set were generated, the first one containing frequency information associated with each phonetic variant token, and the second one also providing for each attested phonetic variant informants’ features (age, education, profession). The format of the two versions of the data set is documented in Appendix 6.1.

2.4.2. ALT-CDI to IPA conversion

The phonetic alphabet used in the ALT project was a geographically specialized version of the *Carta dei Dialetti Italiani* (CDI) transcription system (Grassi et al., 1997), henceforth referred to as CDI-ALT. This choice was in line with the Italian tradition of dialectological studies, which preferred the CDI transcription system with respect to the International

Phonetic Alphabet (IPA). Nowadays, this choice needs to be revised to make the ALT corpus usable by the wider international community of dialectologists and linguists. For this reason, the whole ALT corpus of phonetically transcribed data was converted into IPA.

Appendix 6.2 provides the correspondence table between the CDI-ALT and IPA notations. In most part of the cases, a 1:1 correspondence can be found:

- it can be the case that a CDI-ALT phonetic segment combined with one diacritic corresponds to an individual IPA phonetic segment (e.g. [e̞] > [ɛ]);
- the reverse can also occur, whenever an individual CDI-ALT phonetic segment is converted into a IPA phonetic segment combined with a diacritic (e.g. [ɾ] > [ɾ̞]).

Interestingly enough, there are three different cases (highlighted in grey in the table) in which two different CDI-ALT segments are assigned the same IPA representation: this is the case of the weakened realization of palatal affricates, e.g. [t͡ɕ] and [t͡ɕ̞], whose representation coincides with the representation of the [ʃ] and [ʒ] CDI-ALT segments, i.e. the voiceless and voiced postalveolar fricatives [ʃ] and [ʒ].

In ALT, phonetically transcribed data were represented through a hybrid encoding schema including both compositional and atomic representations which, depending on the task, were automatically converted into each other (Montemagni and Paoli 1989-90, pp. 36-43).

Compositional representations (see column 6 in the Table in Appendix 6.2) encode each phonetic symbol with a basic sign which may be further specified through one or more diacritics (conveying information, for instance, about stress or nasality of vowels). This representation type was particularly convenient for inputting and editing ALT data since all different phonetic symbols (about 110) are encoded by means of a restricted number of codes (36 basic signs and 9 diacritics) belonging to the first 128 ASCII codes and which can be directly accessed through the computer keyboard. To be more concrete, the compositional representation of a word like [stja'ttʃaθa], or [stja'ččáta] in terms of the CDI-ALT notation, is <sti4aCCa8t5a> where letters represent basic signs and numbers diacritics: in the case at hand, '5' marks the spirantization of the voiceless dental occlusive, '8' indicates the stress and '4' represents a semivowel sound. This type of representation is particularly convenient for both sorting and retrieval tasks: in fact, if basic signs only are considered, it is possible to generalise over phonetic variants. Consider as an example the compositional representation of the word forms [stja'ttʃaθa]/[stja'ččáta] and [stja'ttʃata]/[stja'ččáta]: <sti4aCCa8t5a> and <sti4aCCa8ta> respectively. In both cases, the sequence of basic signs is the same, i.e. <stiaCCata>; this entails that a query starting from this sequence of bases will retrieve both of them.

Atomic representations (see column 2 in the Table in Appendix 6.2), on the other hand, show a 1:1 correspondence between CDI-ALT phonetic symbols and computer codes; they were typically used for on-screen display and printing. So, to keep with the <sti4aCCa8t5a> example, the combination of each base together with its diacritics is encoded through a symbol which uniquely identifies it (e.g. t5>t̞).

The CDI-ALT to IPA conversion started from the compositional representation described above and was performed on the basis of 158 ordered conversion rules encoded as PERL regular expressions.

2.5. Reference data preparation

In the proposed analysis method, the phonetic variants recorded for each dialectal variety are compared with those attested in a reference variety (Wieling and Nerbonne, 2009), be it a standard language or a proto-language from which investigated dialects originate. The only prerequisite is that the reference language data should be available in the same transcription system as the dialectal material under study.

For the specific concerns of this study, two different reference languages have been selected, namely standard Italian and Latin, where the former is taken to originate from specific varieties of Tuscan dialects whereas the latter can be seen as a kind of proto-language from which Tuscan dialects originate. Different sets of experiments were performed by using respectively standard Italian and Latin as a reference.

For what concerns Italian, the standard Italian phonetic realization of selected normalised forms was manually encoded. Note that due to the historical relatedness between standard Italian and Tuscan dialects all phonetic variants attested in the reference language are also attested as phonetic variants attested in some ALT location.

For what concerns Latin, we started from a subset of the 444 selected normalised forms (see section 2.2). In this case, the areal coverage constraint was increased to 150 locations, resulting in a 340 normalised forms from which verbs and multi-word expressions (e.g. *luna piena* ‘full moon’, *al sole* ‘in the sun’) have been removed (for a total of 40 NFs). For the remaining 300 cases, we looked for the Latin etymology, if any.

To this specific end, we used the Italian etymological dictionary by Manlio Cortelazzo and Paolo Zolli, “Dizionario Etimologico della Lingua Italiana”, 4 volumes, Bologna, Zanichelli, 1979.

All selected normalised forms were looked up in the reference dictionary and were classified as follows:

1. Latin etymology;
2. diminutive/augmentative forms of Latin words (provided that the used suffix is Latin as well);
3. complex derivative of a Latin word (e.g. *brinata* ‘hoarfrost’ from Lat. *pruinam* ‘brina’ Eng. ‘frost’ ; *castagnaccio* ‘cake made out of chestnut flower’ from Lat. *castāneam* ‘castagna’ Eng. ‘chestnut’);
4. complex etymology (e.g. *albicocca* ‘apricot’);
5. uncertain or unknown etymology (e.g. *afa* ‘sultriness’ or *bischero* ‘fool’);
6. non Latin etymology (e.g. *grullo* ‘silly’, *grattugia* ‘grater’);
7. onomatopoeic words.

Only cases 1 and 2 above were selected as a basis of this case study based on Latin as a proto-language.

For the IPA encoding of Latin words, the grapheme to IPA conversion rules were based on Allen (2004).

3. Method: adaptations and customizations

A general description of the bipartite spectral graph partitioning method is provided in Wieling and Nerbonne (2009, 2010). The method can be seen as articulated into the following steps:

1. obtaining, for each investigated location, the attested realization(s) of a given phonetic segment in a reference variety. The reference variety can be either a standard language or a proto-language;
2. bipartite spectral graph partitioning of the data matrix Locations x Phonetic_features resulting from step 1);
3. for each identified cluster of linguistic varieties, identify the most relevant features characterising it with respect to other clusters of varieties.

In this section, we briefly summarise the peculiarities of the ALT dialectal data set and illustrate the customizations of the method which have been performed to deal with them.

3.1. The ALT dialectal data set

The data set which has been selected for this study has been described in detail in section 2.2. Among its main features, it is worth pointing out here that the selected sample includes nouns and adjectives, both single words and multi-word expressions, testifying both productive phonetic processes and phonotactic processes. The focus is on the phonetic representation, including diacritics for a total of 109 different phonetic symbols.

3.2. Adaptations and customizations of the method

3.2.1. Enriching phonetic segment pairs with context information (Step 1)

In the adopted clustering method, each dialectal variety is described in terms of the attested phonetic realizations of a given phonetic segment with respect to a reference variety. Attested phonetic realizations are represented in terms of segment pairs where the phonetic segment attested in a given location L_x is associated with its realization in the reference variety R (either standard Italian or Latin):

[phonetic_realization_in_R]:[phonetic_realization_in_L_x].

For each selected NF (see section 2.2 above), phonetic segment pairs are obtained by aligning the phonetic realization in the reference variety R against the phonetic realizations recorded in the investigated varieties using the Levenshtein algorithm: alignments were induced by enforcing the syllabicity constraint on the basis of the PMI-based Levenshtein distance measure (Wieling et al., 2009).

Due to the fact that in the ALT dataset the same segment pair could originate from different phonetic processes, we decided to enrich the representation of segment pairs with contextual information, as exemplified below:

[L_ctx|phonetic_realization_in_Ref|R_ctx]:[L_ctx|phonetic_realization_in_L_x|R_ctx]

Consider as an example the segment pairs involving a consonant C and its geminated counterpart, i.e. [C]-[C:]. As it can be observed in the following examples, the same segment pair could result from different phonetic processes:

1. consonantal lengthening in intervocalic position

a. [536] 225 Italiano

[6] 101 Pieve Santo Stefano

a _ b a tʃ i o

a _ b a tʃː i o

b. [829] 225 Italiano

[628] 198 Piancastagnaio

a b e t e

a bː e t e

2. palatalization + consonantal lengthening

a. [1013] 225 Italiano

[847] 198 Piancastagnaio

a l b e r o

a j bː e r u

b. [1141] 225 Italiano

[872] 198 Piancastagnaio

f a r i n a _ d o l tʃ e

f a r i n a _ d o j tʃː e

3. phonotactic lengthening (in word initial position)

a. [567] 225 Italiano

[87] 107 Rosignano Marittimo

a _ k a z o

a _ kː a z o

b. [536] 225 Italiano

[24] 107 Rosignano Marittimo

a _ b a tʃ i o

a _ bː a ʃ i o

For example, the same pair [tʃ]~[tʃ:] appears both in 1.a and 2.b, as the result of different phonetic processes, namely consonantal lengthening occurring in intervocalic position and palatalization of preconsonantal /t/ followed by lengthening of the following consonant. Note that the involved phonetic phenomena show quite a different areal distribution.

For the time being, the representation of context includes:

- vowel (V);
- consonant (C);
- both (i.e. matching vowel and consonant, encoded as B);
- indel (-);
- word boundary (_);
- unknown (?).

In principle finer-grained distinctions can be resorted to in the representation of context information, with the danger of increasing the data sparseness problem.

3.2.2. Constraints on extracted segment pairs (Step 1)

In the previous section, we have seen that extracted phonetic segment pairs are formalised as follows:

[L_ctx|phonetic_realization_in_Ref|R_ctx]:[L_ctx|phonetic_realization_in_L_x|R_ctx]

Extracted segment pairs represent the basis of this study: i.e. each dialectal variety spoken in a given location is described in terms of the set of phonetic realizations of the underlying phonetic segment in the reference variety. This entails that extracted segment pairs should testify productive phonetic processes only. For this reason, from the set of extracted segment pairs we pruned out the segment pairs attested for single words only. This is the case, for instance, of the segment pair **V|n|C:V|r|C**, originating from the comparison of phonetic variants of the word *fanfarone* 'boaster' which included among its phonetic variants also *farfarone*: here, we cannot exclude that the attested variation is lexically driven, i.e. it originates from an assimilation process.

The parameters which could be used to define the set of extracted segment pairs thus include:

- the minimum number of words from which the same phonetic segment pairs could be extracted (at least 2 on the basis of what it was said above);
- the number of locations with respect to which the same phonetic segment pair has been attested.

3.2.3. Treatment of multiple responses (Step 1)

In a dialectal corpus of atlas data, in principle the same questionnaire item can be assigned multiple responses, attested either by the same informant or by different informants belonging to the same community.

In the case of ALT data, for each attested response type to the same questionnaire item within the same location we also know the number of informants who attested it. This is to say that also the frequency of occurrence of a given response in a given location is available in the ALT data set.

In previous dialectometric studies of Tuscan dialectal variation (Montemagni 2007, 2008), the treatment of multiple responses was carried out along the lines suggested by Nerbonne and Kleiweg (2003), where the distance was computed “between sets of strings where the sets represent alternative lexicalizations. The basic idea is that we average the distances between the individual strings where we consistently choose pairs in a way that minimizes the distance measure”.

In this study, given the availability of token frequency information, different options have been provided to deal with multiple phonetic variants of the same NF within the same location, namely:

1. average over multiple phonetic variants tokens. In this way, the token frequency is used to determine the importance of each variant within a given location;
2. average over multiple phonetic variants types. In this way, the relative frequency of individual types is ignored and each option is weighted equally. This is the option followed in previous studies;
3. “majority vote”, i.e. only the most frequent phonetic variant is considered for a given location.

The availability of these different options can help exploring the role of frequency in the study of dialectal variation, which still represents an open issue worth being investigated. As stated in Wieling and Nerbonne (2009, p.30), “while it stands to reason that more frequently encountered variation would signal dialectal affinity more strongly, it is also the case that inverse frequency weightings have occasionally been applied (Goebel, 1984), and have been shown to function well”.

3.2.4. Single segment pairs (Step 1)

The analysis can be based either on all extracted segment pairs or on a subset of them. In the previous version of the method, the latter case was handled by specifying a given phonetic segment with the results that all segment pairs including it on either side (i.e. in either the target or reference location) were selected for the analysis. By doing in this way it would have been impossible to focus on specific linguistic phenomena, since the extracted data would have included pairs relating to different phonetic phenomena.

Consider, for example, the plosives in Tuscan dialectal variation: both voicing and devoicing of plosives are attested as productive processes in Tuscany, though with a different geographic distribution:

1. **devoicing of plosives in intervocalic position**
V|g|V#V|k|V attested wrt the Tuscan words *aghetto* and *agaiolo*;
2. **voicing of plosives in intervocalic position**
V|k|V#V|g|V: a more productive process wrt the previous one, attested wrt words such as *vicolo*, *albicocca*, *bacherozzolo*, *capocollo*, *ciuco*, *grattacacia*, *idraulico*, *oca*, *radica*, *rancico*, *ricotta*, *rustico*, *strabico*, etc..

The same applies to other phenomena such as lengthening and shortening of plosives:

3. **lengthening of plosives in intervocalic position**
V|t|V#V|t : |V: attested wrt *ditale* and *sito*;

4. shortening of plosives in intervocalic position

V|t : |V#V|t|V: a more productive process wrt the previous one, attested wrt words such as *bottiglia*, *aghetto*, *bigotto*, *bruschetta*, etc.

In order to make it possible to focus on specific phonetic phenomena, in the new version of the method when one or more phonetic segment(s) are specified, all segment pairs including them on the reference side are selected for the analysis.

3.2.5. Representativeness vs Distinctiveness (Step 3)

Wieling and Nerbonne (2009) calculate the importance of each phonetic segment pair by combining two different features, i.e. 'representativeness' and 'distinctiveness', where the former indicates the proportion of varieties in a given cluster which contain the sound correspondence and the latter indicates how prevalent a segment pair is in its own cluster as opposed to other clusters.

To be able to rank the segment pairs based on their distinctiveness and representativeness, these two values need to be combined. Different options have been experimented with by Wieling and Nerbonne (2009, 2010), namely:

- a) taking the average of both values;
- b) weigh distinctiveness twice wrt representativeness.

Different ways of combining the two values have been experimented with the ALT data set. Due to the strong similarity holding between the investigated dialectal varieties, it appears that option b) above leads to uninteresting results; the first option is better but still includes some noisy data. Experiments are being carried out to identify the best balance between the two scores wrt the linguistic peculiarities of ALT data.

4. Current directions of research

The report documents the activity carried out during the Short Term Mobility stay at the University of Groningen, which focussed on two main lines of research: a) preparation of the data set; b) adaptation and customization of the method with respect to the peculiar problems posed by the ALT data set. Currently, experiments are being carried out both with standard Italian and Latin as reference varieties and results are being compared with Tuscan dialectological literature.

Preliminary results were presented at the XI Congresso della Società di Linguistica e Filologia Italiana (SILFI) which was held in Napoli on 5-7 October 2010, through a joint contribution entitled "Patterns of language variation and underlying linguistic features: a new dialectometric approach" by Simonetta Montemagni, Martijn Wieling, Bob de Jonge and John Nerbonne. The presented poster is attached at the end of the report.

Il Fruitore
Simonetta Montemagni

Il Proponente
Vito Pirrelli

5. References

- Allen, W. S. (2004), *Vox Latina — a Guide to the Pronunciation of Classical Latin* (2nd ed.). Cambridge University Press.
- Giacomelli G., L. Agostiniani, P. Bellucci, L. Giannelli, S. Montemagni, A. Nesi, M. Paoli, E. Picchi, T. Poggi Salani (a cura di) (2000), *Atlante Lessicale Toscano*, Lexis Progetti Editoriali, Roma.
- Giannelli L. (2000), *Toscana*, Pacini Editore, Pisa (1976, prima edizione).
- Goebel H. (1984), *Dialektometrische Studien: Anhand italo-romanischer, rätoromanischer und galloromanischer Sprachmaterialien aus AIS und ALF*, Max Niemeyer, Tübingen.
- Goebel H. (2005), *La dialectométrie corrélatrice. Un nouvel outil pour l'étude de l'aménagement dialectal de l'espace par l'homme*, "Revue de linguistique romane", 69, pp. 321-367.
- Grassi C., A. Sobrero, T. Telmon (1997), *Fondamenti di Dialettologia Italiana*, Roma-Bari, Laterza.
- Loporcaro M. (2009), *Profilo linguistico dei dialetti italiani*, Roma-Bari, Laterza 2009.
- Montemagni S., M. Paoli (1989/1990). *Dalla parola al bit (e ritorno): percorsi dall'inchiesta sul campo alla banca dati dell'ALT*, in AA.VV., "Quaderni dell'Atlante Lessicale Toscano", VI/VIII, pp. 7-52.
- Montemagni S., M. Paoli, E. Picchi (2006), *ALT Web: l'Atlante Lessicale Toscano in rete*, in Bruni F., Marcato C. (a cura di), *Lessicografia dialettale: ricordando Paolo Zolli*, Atti del Convegno di Studi, Venezia, 9-11 dicembre 2004, Editrice Antenore, Roma-Padova, Tomo I, pp. 209-241.
- Montemagni S. (2007), *Patterns of phonetic variation in Tuscany: using dialectometric techniques on multi-level representations of dialectal data*, in P. Osenova et al. (a cura di), *Proceedings of the Workshop on Computational Phonology at RANLP-2007* (26 September 2007, Borovetz, Bulgaria), pp. 49-60.
- Montemagni S. (2008), *The space of Tuscan dialectal variation. A correlation study*, "International Journal of Humanities and Arts Computing" (Special Issue on "Language Variation Studies and Computational Humanities"), Edinburgh University Press, Oct 2008, Vol. 2, No. 1-2, pp. 135-152.
- Nerbonne J., W. Heeringa, P. Kleiweg (1999), *Edit Distance and Dialect Proximity*, in Sankoff D., Kruskal J. (a cura di), *Time Warps, String Edits and Macromolecules: The Theory and Practice of Sequence Comparison*, Stanford, CSLI Press, pp.v-xv.
- Nerbonne J., W. Heeringa (2001), *Computational comparison and classification of dialects*, "Dialectologia et Geolinguistica. Journal of the International Society for Dialectology and Geolinguistics", 9, pp. 69-83.
- Nerbonne J., P. Kleiweg (2003), *Lexical Distance in LAMSAS*, in J. Nerbonne, W. Kretzschmar (Eds.), "Computers and the Humanities (Special Issue on Computational Methods in Dialectometry)", 37(3), pp. 339-357.
- Nerbonne J. (in press), *Various Variation Aggregates in the LAMSAS South*, in stampa in C. Davis, M. Picone (eds.), *Language Variety in the South III*, Tuscaloosa, University of Alabama Press.
- Pellegrini G.B. (1977), *Carta dei Dialetti d'Italia*, Pisa, Pacini Editore.
- Séguy J. (1971), *La relation entre la distance spatiale et la distance lexicale*, "Revue de Linguistique Romane", 35, pp. 335-357.
- Szmrecsanyi B. (to appear). *Corpus-based dialectometry – a methodological sketch*, "Corpora", 6(1), Edinburgh University Press.

- Wieling M., Nerbonne J. (2009), *Bipartite spectral graph partitioning to co-cluster varieties and sound correspondences in dialectology*, in M. Choudhury, S. Hassan, A. Mukherjee, S. Muresan (a cura di), *Proceedings of the 2009 Workshop on Graph-based Methods for Natural Language Processing*, pp. 26–34.
- Wieling M., J. Prokic, J. Nerbonne (2009), *Evaluating the pairwise alignment of pronunciations*, in L. Borin, P. Lendvai (Eds.), *Language Technology and Resources for Cultural Heritage, Social Sciences, Humanities, and Education*, pp. 26–34.
- Wieling M., Nerbonne J. (2010), *Hierarchical spectral partitioning of bipartite graphs to cluster dialects and identify distinguishing features*, in *Proceedings of the 2010 Workshop on Graph-based Methods for Natural Language Processing (TextGraphs-5)*, ACL, Uppsala, Sweden, July 16, 2010.

6. Appendixes

6.1. *ALT-RuG L04 data files: syntax for location and informants codes*

ALT-RuG L04 data files with frequency information

: 100 Caprese Michelangelo
- aJino
- aJino
- aJino
...

Location line

Es.

: 100 Caprese Michelangelo

it contains the following information types:

Numeric_location_id (100) and place_name (Caprese Michelangelo)

Phonetic variant tokens are reported one for each line starting with a dash (“-“). Frequency information can be reconstructed by counting the occurrences of the same phonetic variant token.

ALT-RuG L04 data files with informant details

: 100 Caprese Michelangelo-1-1
f;81;1898;2;9;2;1
- aJino

Location line

Es.

: 100 Caprese Michelangelo-1-1

it contains the following information types:

Numeric_location_id (100) place_name (Caprese Michelangelo) inquiry_id (1) informant_id (1)
where

- the value of inquiry_id is numeric
- the value of informant_id is alphanumeric

Note that the numbering of informants reflects age (older informants are assigned lower alphanumeric identifiers (i.e. informants labelled as 1 or A are older than 5 and C respectively). The ordering of identifiers follows the age ranking.

Informant line

Es.

f;81;1898;2;9;2;1

Fields separated by “;”

For each informant the following information types are provided:

- sex (f/m)
- age of the informant at the time of the interview
- year of birth

- the two information types are useful due to the fact that ALT interviews were carried out in a time span of 20 years
- education encoded as follows:
 - 1: illiterate or semi-literate;
 - 2: primary school (not necessarily completed);
 - 3: middle school (not necessarily completed);
 - 4: so-called “Istituto Professionale” which is a type of high school providing secondary education oriented toward more practical subjects, enabling the students to start searching for a job as soon as they have completed their studies, typically after 3 years instead of 5 (not necessarily completed);
 - 5: high school (not necessarily completed);
 - 6: university degree (not necessarily completed);
- current and past profession, encoded as follows: in the first profession field it is reported the current profession, whereas the other two contain past professions, if any. The profession codes are provided below:
 - 1: farmer, farmhand, shepherd
 - 2: craftsman
 - 3: trader
 - 4: executive or auxiliary employee
 - 5: manager, concept employee, nurse
 - 6: teacher, freelance
 - 7: unskilled worker, apprentice
 - 8: skilled worker
 - 9: non-professional status (student, housewife, retired)

6.2. ***CDI-ALT to IPA conversion: correspondence table***

	CDI-ALT	IPA-full conversion	IPA suprasegm.	IPA-simplified	CDI-ALT compositional representation	Notes
1.	ə	ə		ə	@	
2.	é	ə	Stress	ə	@8	
3.	a	a		a	a	
4.	ä	æ		æ	a3	
5.	ǣ	æ̃		æ	a37	
6.	□	æ̃	Stress	æ	a378	
7.	ǣ	æ	Stress	æ	a38	
8.	ã	ã		a	a7	
9.	ǣ̃	ã	Stress	a	a78	
10.	á	a	Stress	a	a8	
11.	b	b		b	b	
12.	ḃ	β		β	b5	
13.	č	tʃ		tʃ	C	
14.	ć	ʃ		ʃ	C5	See n.106
15.	d	d		D	d	

	CDI-ALT	IPA-full conversion	IPA suprasegm.	IPA- simplified	CDI-ALT compositional representation	Notes
16.	d'	dʲ		dʲ	D	
17.	đ	ǯ		ǯ	d5	
18.	e	e		e	e	
19.	ẹ	ẹ		e	e0	
20.	ẽ	ẽ		e	e07	
21.	ẽ̌	ẽ̌	Stress	e	e078	
22.	ẹ̌	ẹ̌	Stress	e	e08	
23.	ẹ	e		e	e1	
24.	ẹ̃	ẽ̃		e	e17	
25.	ẽ̌	ẽ̌	Stress	e	e178	
26.	ẹ̌	e	Stress	e	e18	
27.	ẹ	ɛ		ɛ	e2	
28.	ẽ̃	ẽ̃		ɛ	e27	
29.	ẽ̌	ẽ̌	Stress	ɛ	e278	
30.	ẹ̌	ɛ	Stress	ɛ	e28	
31.	□	ẹ̌		ẹ̌	e4	
32.	é	e	Stress	e	e8	
33.	f	f		f	f	
34.	ǵ	ɖʒ		ɖʒ	G	
35.	g	g		g	g	
36.	ǥ	ɣ		ɣ	g5	
37.	ǵ̌	ʒ		ʒ	G5	See n.107
38.	□	gʲ		gʲ	g6	
39.	h	h		h	h	
40.	i	i		i	i	
41.	ị	ɿ		i	i2	
42.	ĩ̃	ĩ̃		i	i27	
43.	ĩ̌	ĩ̌	Stress	i	i278	
44.	ị̌	ɿ̌	Stress	i	i28	
45.	ị̇	j		j	i4	
46.	ĩ	ĩ		i	i7	
47.	ĩ̌	ĩ̌	Stress	i	i78	
48.	í	i	Stress	i	i8	
49.	ǵ̌	gʲ		gʲ	J	
50.	č	kʲ		kʲ	j	
51.	k	k		k	k	

	CDI-ALT	IPA-full conversion	IPA suprasegm.	IPA- simplified	CDI-ALT compositional representation	Notes
52.	k'	x		x	k5	
53.	l	l		l	l	
54.	l'	ʎ		ʎ	L	
55.	ɬ	ɭ		ɭ	ɭ6	
56.	m	m		m	m	
57.	n	n		n	n	
58.	ń	ɲ		ɲ	N	
59.	ɳ	ɳ		ɳ	n1	
60.	ṇ	ɳ		ɳ	n6	
61.	o	o		o	o	
62.	ṽ	ṽ		o	o0	
63.	õ	õ		o	o07	
64.	ṽ̃	õ	Stress	o	o078	
65.	ó	ṽ	Stress	o	o08	
66.	ṽ	o		o	o1	
67.	õ	õ		o	o17	
68.	ṽ̃	õ	Stress	o	o178	
69.	ó	o	Stress	o	o18	
70.	ṽ	ɔ		ɔ	o2	
71.	õ	õ		ɔ	o27	
72.	ṽ̃	õ	Stress	ɔ	o278	
73.	ó	ɔ	Stress	ɔ	o28	
74.	ö	ø		ø	o3	
75.	õ	ø		ø	o37	
76.	ṽ̃	ø	Stress	ø	o378	
77.	õ	ø	Stress	ø	o38	
78.	ó	o	Stress	o	o8	
79.	ṽ̃̃	ɹ		ɔ	o9	
80.	ó̃̃	ɹ	Stress	ɔ	o98	
81.	p	p		p	p	
82.	ɸ	ɸ		ɸ	p5	
83.	ɖ	ɖ		ɖ	R	
84.	r	r		r	r	
85.	ɾ	ɾ		ɾ	r1	
86.	s	s		s	s	
87.	S	z		z	S	

	CDI-ALT	IPA-full conversion	IPA suprasegm.	IPA- simplified	CDI-ALT compositional representation	Notes
88.	t	t		t	t	
89.	t̃	t̃		t̃	T	
90.	ɛ	θ		θ	t5	
91.	u	u		u	u	
92.	u̇	u̇		u	u2	
93.	ü̇	y		y	u27	
94.	ũ̇	ũ	Stress	u	u278	
95.	ú̇	u̇	Stress	u	u28	
96.	ü	y		y	u3	
97.	ũ	ỹ		y	u37	
98.	ũ̃	ỹ	Stress	y	u378	
99.	ũ̇	y	Stress	y	u38	
100.	u̇	w		w	u4	
101.	ũ	ũ		ũ	u7	
102.	ũ̇	ũ	Stress	ũ	u78	
103.	ú	u		u	u8	
104.	v	v		v	v	
105.	ž	z̥		z̥	w	See n.109
106.	š	ʃ		ʃ	x	See n. 14
107.	ʃ̃	ʒ		ʒ	X	See n. 37
108.	ś	ʃ̥		ʃ̥	x5	
109.	□	z̥		z̥	X5	See n.105
110.	ɟ	ɟ		ɟ	Y	
111.	ʒ̥	dʒ		dʒ	Z	
112.	z	ts		ts	z	

6.3. Selected experimental data set

This Table is still under completion.

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
edera	ivy	édera	edera	hēderam	
orso	bear	órso	orso		
oca	goose	óka	oka		
occhio	eye	óḳḳjo	ok:jo	ōcūlum	
acino	grape	áčino	at̥jino		

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
acqua cotta	Lit. 'boiled water', designating a traditional Tuscan soup	ákk _χ ua kòtta	ak:wa kot:a		
albero	tree	álbero	albero	ārborem	
andito	passage	ándito	andito	ādītūm	
ape	bee	ápe	ape	āpem	
asino	ass	ásino	asino		
aspite	aspite	áspide	aspide		
aspito	aspito	áspido	aspido		
a bacio	to kiss	a bačío	a batʃio		
a caso	Random	a káSo	a kazo		
a cavalluccio	piggyback	a kávallúččo	a kaval:utʃ:o		
a solatio	in sunny	a solatío	a solatio		
abete	fir	abéte	abete	abīetem	
abeto	abeto	abéto	abeto	abīetem	
acquaio	sink	akkuájo	ak:wajo	āquārīum	
agaiolo	agaiolo	agaiòlo	agajolo		
aghetto	pond	agétto	aget:o		
agnella	lamb	ánǵella	aŋ:el:a		
al sole	sun	al sòle	al sole		
al tocco	touch	al tókko	al tok:o		
albicocca	apricot	álbikòkka	albikok:a		
allodola	lark	allòdola	al:odola		
alloro	laurel	allòro	al:oro		
altalena	swing	általèna	altalena		
arancio	orange	aránčo	arantʃo		
arcolaio	spinning wheel	árkolájo	arkolajo		
avvezzo	accustomed	avvēzzo	av:ets:o	advitiātus	
becero	yahoo	bēcero	betʃero		
borsa	bag	bòrsa	borsa		
boccia	bowl	bòčča	botʃ:a		
balzo	leap	bálzo	balʃo	bālteum	
bambola	doll	bámbola	bambola		
bazza	Bazza	bázʒa	badʒ:a		
baccano	din	bakkáno	bak:ano	bacchānal	
bacchettone	bigot	bakkettòne	bak:et:one		
bacherozzolo	bacherozzolo	bakerózzolo	bakerots:olo		
baciatura	baciatura	báčatúra	batʃatura		

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
balordo	stupid	balórdo	balordo		
bardella	pack-saddle	bardélla	bardɛl:a		
bargigli	wattles	barǵíll'i	bardʒiʌ:i		
bastone	stick	bastóne	bastone		
bastardo	bastard	bastárido	bastardo		
bazzone	Bazzone	baz̥z̥óne	badz:one		
bubbola	bells	búbbola	bub:ola	upūpulam	
bernoccolo	bump	bernókko	bernɔk:olo		
beverone	mash	beveróne	beverone		
bigotto	bigot	bigótto	bigot:o		
bighellone	loafer	bigellóne	bigel:one		
birignoccolo	birignoccolo	bírińńókko	birɪn:ɔk:olo		
bindolo	waterwheel	bíndolo	bindolo		
bordello	brothel	bordéll'o	bordɛl:o		
borraccina	stonecrop	borraččína	bor:atʃ:ina		
bottiglia	bottle	bottíll'a	bot:iʌ:a	butīculam	
brace	embers	bráče	bratʃe		
braciola	chop	bračóla	bratʃola		
braciere	brazier	bračére	bratʃere		
brinata	hoarfrost	brináta	brinata		
brindellone	Brindellone	brindellóne	brindel:one		
broccione	broccione	broččóne	broʃ:one		
bruscolo	mote	brúskolo	bruskolo		
bruschetta	bruschetta	bruskétta	brusket:a		
bischero	dawg	bískero	biskero		
biscia	snake	bišša	bɪʃ:a	bēstiam	
butto	throw	bútto	but:o		
burrone	ravine	burróne	bur:one	burrunum	
burischio	Burisch	burískjo	buriskjo		
buzzo	paunch	búz̥z̥o	budz:o		
cenci	rags	čénči	tʃentʃi		
ceppa	log	čéppa	tʃep:a		
ceppo	log	čéppo	tʃep:o	cĭppum	
cesta	basket	čésta	tʃesta		
concio	ashlar	kónčo	kontʃo	cōmptum	
cacio	cheese	káčo	katʃo	cāseum	
caglio	rennet	kál't'o	kaʌ:o	cōāgŭlum	
capo	head	kápo	kapo	căput	
cagnara	rumpus	kańńára	kaɲ:ara		

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
calabrone	hornet	kalabr _o ne	kalabrone	crabr _o ne	
caldano	brazier	kaldáno	kaldano		
calzoni	trousers	kalz _o ni	kaltsoni		
camomilla	chamomile	kámomílla	kamomil:a	camomíllam	
camposanto	cemetery	kámposáto	kamposanto		
cantuccio	corner	kantúččo	kantut _f :o		
cantonata	corner	kántonáta	kantonata		
capezzolo	nipple	kapézzolo	kapets:olo	capítium	
capone	capone	kap _o ne	kapone		
capocollo	capocollo	kápokóll _o	kapokol:o		
capomilla	capomilla	kápomílla	kapomil:a		
carbonaio	charcoal	kárbonájo	karbonajo	carbonārium	
casino	mess	kasíno	kasino	cāsīnum	
castagneto	chestnut	kástańńeto	kastan:eto		
castagnaccio	chestnut	kástańńáččo	kastan:at _f :o		
catasta	stack	katásta	katasta	catāstam	
cavaocchi	cavaocchi	káva _o kki	kavaok:i		
cavalletta	grasshopper	kavallétta	kaval:et:a		
ciccioli	greaves	číččoli	t _f it _f :oli		
cenciaio	cenciaio	čencájo	t _f ent _f aj _o		
cenciaiolo	ragpickers	čencájo _o lo	t _f ent _f aj _o lo		
cesoie	shears	čeS _o ie	t _f ezoie		
cetriolo	cucumber	četri _o lo	t _f etriolo	citrīolum	
cigli	eyelashes	číl'li	t _f il'i		
ciglia	cilia	číl'la	t _f il:a		
cigna	Cigna	čínna	t _f in:a		
chicco	bean	kíkko	kik:o		
chiocciola	snail	kij _o ččola	kjot _f :ola		
chiorba	chiorba	kij _o rba	kjorba		
chiasso	noise	kjásso	kjas:o		
chiacchierone	jay	kjákkjer _o ne	kjak:jerone		
chiorbone	chiorbone	kjorb _o ne	kjorbone		
ciocca	lock	čókka	t _f ok:a		
ciabatte	slippers	čabátte	t _f abat:e		
ciucca	ciucca	čúkka	t _f uk:a		
ciccione	fatty	číčč _o ne	t _f it _f :one		
ciucco	pacifier	čúkko	t _f uk:o		
ciuco	donkey	čúko	t _f uko		
ciuffo	tuft	čúffo	t _f uf:o		

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
cilliegia	cherry	čiliġġa	tʃiljɛdʒa	cerēseam	
cimitero	cemetery	čimitġero	tʃimitero	cimitġerium	
cinquale	Cinquale	činkuāle	tʃinkwale		
cintolino	garters	čintolíno	tʃintolino		
cintura	belt	čintúra	tʃintura		
cinturino	strap	činturíno	tʃinturino		
cipresso	cypress	čiprésso	tʃipres:o	cyparīssum	
cisoie	cisoie	čiSóġie	tʃizojɛ		
cimice	bug	čímiče	tʃimɪtʃɛ	cīmicem	
cinghia	belt	číngġia	tʃingja		
cintola	waist	číntola	tʃintola	cīnctulam	
cocomero	watermelon	kokómero	kokomero	cucūmerem	
coccinella	ladybug	kóččinġella	koʧːinɛlːa		
comare	Gossip	komáre	komare		
compagno	companion	kompánġo	kompaɲːo	companium	
compare	appear	kompáre	kompare	cōmpatrem	
coperchio	cover	kopérkġo	kopɛrkjo	copĕrculum	
covone	sheaf	kovóne	kovone		
covata	brood	kováta	kovata		
crognolo	Crognolo	króġnġolo	krɔɲːolo		
crivello	sieve	krivġello	krivelːo		
crusca	bran	krúska	kruska		
cispa	blear eyes	číska	tʃispa		
citto	citto	čítto	tʃitːo		
dolco	Pyrus	dólko	dolko		
desinare	dinner	děSináre	dezinare		
di nascosto	secretly	di naskósto	di naskosto		
di sguincio	of sideways	di Sguínčo	di zgwinʧo		
di traverso	askew	di travġerso	di traverso		
dialetto	dialect	dġalġetto	djaletːo	dġialectum	
ditale	thimble	ditále	ditale	digitāle	
denti macellari	teeth butchers	dġenti mačellári	dġenti matʃelːari		
denti occhiali	teeth glasses	dġenti okkiáli	dġenti okːjali		
esoso	hexose	eSóSo	ezɔzo	exōsum	
fosso	ditch	fósso	fɔsːo		
fango	mud	fángo	fango		
falegname	carpenter	fáleġnāme	faleɲːame		
fanfarone	braggart	fánfaróne	fanfarone		

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
faraona	guinea fowl	fáraṇa	faraona	pharaōnem	
farina dolce	cake flour	farína dólče	farina doltʃe		
favilla	spark	favílla	favil:a		
ferraio	blacksmith	ferrájo	fer:ajo		
fettina	slice	fettína	fet:ina		
fiocco	bow	fīṭkko	fjɔk:o	flōccum	
fiammifero	match	fīamífero	fjamifero	flammiŕerum	
fidanzato	boyfriend	fidanzáto	fidantsato		
filone	vein	filṇe	filone		
filare	spin	filáre	filare		
focolare	hearth	fṭkoláre	fokolare		
formaggio	cheese	formáḡḡo	formadʒ:o	formāticum	
formica	ant	formíka	formika	formīcam	
formicola	Tingly	formíkola	formikola	formīcolam	
fottio	fuck	fottío	fot:io		
fragola	strawberry	frágola	fragola	frāgulam	
frana	landslide	frána	frana	frāginam	
fregatura	swindle	frégatúra	fregatura	fricatūram	
fringuello	finch	fringuḗllo	fringwel:o		
fuliggine	soot	fulíḡḡine	fulidʒ:ine	fulīginem	
fulminante	fulminant	fulminánte	fulminante		
golpe	coup	gólpe	golpe	vŭlpem	
gota	cheek	góta	gota		
ganza	mistress	gánza	gandza	gangia	
ganzo	guy	gánzo	gandzo		
gazza	magpie	gázṛa	gadʒ:a		
gazzera	Gazzera	gázṛera	gadʒ:era		
gabinetto	cabinet	gábinétto	gabinet:o		
gallinella	hen	gállinḗlla	gal:inel:a		
ghiaia	gravel	ḡiája	gjaja	glāream	
ghiacciaia	icebox	ḡiaččája	gjaʦʃ:a		
ghiandaia	jay	ḡiandája	gjandaja	glandāriam	
ghiro	dormouse	ḡíro	giro	glīrem	
giubba	jacket	ḡúbba	dʒub:a		
ginepro	juniper	ḡinépro	dʒinepro	ienīperum	
giomella	giomella	ḡomḗlla	dʒomel:a		
girino	tadpole	ḡiríno	dʒirino	gyrīnum	
gnocchi	gnocchi	ṇṇṭkki	ɲ:ɔk:i		
gomitolo	ball	gomítolo	gomitolo	glōmitolum	

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
governo	government	govĕrno	governo	gubernum	
greppo	chasm	grĕppo	grep:o		
grappolo	cluster	gráppolo	grap:olo	clāppulum	
grasso	fat	grásso	gras:o		
grattacacia	grattacacia	grattakáča	grat:akatʃa		
grattugia	grater	grattúġa	grat:udʒa		
grembiule	apron	grembiúle	grembjule		
grillo	cricket	grillo	gril:o	grīllum	
grullo	stupid	grúllo	grul:o		
gelso	mulberry	ġĕlso	dʒɛlso	cēlsam	
gemma	gem	ġĕmma	dʒɛm:a		
guazza	dew	guázza	gwats:a	acquāceam	
guercio	one-eyed	guĕrċo	guertʃo		
idraulico	hydraulic	idráuliko	idrawliko	hydrāulicum	
imbroglio	cheat	imbróġl'o	imbroʎ:o		
imbranato	clumsy	imbranáto	imbranato		
in ghingheri	dressed up	in ġíngeri	in gingeri		
in proda	on shore	im pródá	im proda		
l' anno scorso	last year 's	l ánnu skórso	l an:o skorso		
locco	LOCKING	lókko	lok:o		
lodola	skylark	lódola	lodola	alāudam	
lampo	flash	lámpo	lampo		
lavatoio	wash	lavatòġo	lavatojo	lavatōrium	
licite	Licite	lícite	litʃite		
legnaiolo	carpenter	leńńajól'o	lej:n:ajolo		
letame	manure	letáme	letame		
luna calante	waning moon	lúna kalánte	luna kalante		
luna crescente	crescent	lúna kreššĕnte	luna kref:ente		
luna piena	full moon	lúna piĕna	luna pjɛna		
lupo	wolf	lúpo	lupo		
lucignola	wick	lučínńola	lutʃijn:ola		
lucertola	lizard	lučĕrtola	lutʃertola		
lumaca	snail	lumáka	lumaka	limācam	
livido	livid	lívido	livido	līvidum	
mento	chin	mĕnto	mento	mĕntum	
moccolo	snot	mókkolo	mok:olo	mūcculum	
mogio	dejected	móġo	mɔdʒo		
mota	mota	móta	mota	māltham	

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
macchia	stain	mákkja	mak:ja	măcŭlam	
male di capo	evil head	mál di kápo	mal di kapo		
male di testa	headache	mál di tĕsta	mal di tĕsta		
manfano	manf	mánfano	manfano		
madonnina	Madonna	madonnína	madon:ina		
maggiolino	cockchafer	maġġolíno	madʒ:olino		
magnano	locksmith	mañnáno	majn:ano		
maiale	pork	majále	majale	maiālem	
maialino	piglet	majalíno	majalino		
manciata	handful	mančáta	mantʃata		
mangiatoia	manger	manġatōja	mandʒatoja		
matassa	hank	matássa	matas:a	matāxam	
materiale	material	materjále	materjale		
mattarello	mattarello	máttarĕllo	mat:arɛl:o		
melone	melon	melóne	melone		
mestone	mestone	mestóne	mestone		
midolla	marrow	midólla	midol:a	medŭllam	
mietitura	harvest	mġetitúra	mjetitura		
migliaccio	Migliaccio	miġl'áččo	miʎ:atʃ:o	miliācium	
mirtilli	blueberries	mirtílli	mirtil:i		
mollica	crumb	mollíka	mol:ika		
moine	moine	moíne	moine		
montone	ram	montóne	montone	multōnem	
mortadella	mortadella	mortadĕlla	mortadɛl:a		
moscone	bluebottle	moskóne	moskone		
muschio	moss	múskjo	muskjo	mŭsculum	
muta	pack	múta	muta	mŭtam	
noccola	noccola	nókkola	nɔk:ola		
nottola	owl	nóttola	nɔt:ola	nŏctulam	
nottolo	owl	nóttolo	nɔt:olo		
nappone	nappone	nappóne	nap:one		
nasone	nose	nasóne	nasone		
nascondino	hide	naskondíno	naskondino		
nervoso	nervous	nervóso	nervoso	nervōsum	
nevischio	sleet	nevískjo	neviskjo	nivisculum	
noioso	boring	noióso	nojoso		
occhiali	glasses	okkiáli	ok:jali		
odori	odors	odóri	odori	odōres	
orecchio	ear	orĕkkjo	orek:jo		

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
orzaiolo	sty	orḡaiólo	ordzaioło		
pecchia	bee	pékkiḡa	pek:ja		
poggio	knoll	póġġo	pɔdʒ:o	pōdĭum	
palco	stage	pálko	palko		
palo	pole	pálo	palo	pālum	
pancia	belly	pánċa	pantʃa	pānticem	
papera	gosling	pápera	papera		
papero	gosling	pápero	papero		
padrino	godfather	padrĭno	padrino	patrĭnum	
pagliaia	Pagliaia	pal'láġa	paʎ:aja		
pagliuzza	mote	pal'lúzza	paʎ:uts:a		
paletto	pole	palétto	palet:o		
pancetta	bacon	panċétta	pantʃet:a		
pancione	paunch	panċóne	pantʃone		
paniere	basket	panniġere	pan:jere		
panzanella	panzanella	panzanélla	pantsanɛl:a		
papavero	poppy	papávero	papavero	papāverum	
pappavero	pappavero	pappávero	pap:avero		
parlata	speech	parláta	parlata		
pastone	mash	pastóne	pastone		
pastrano	overcoat	pastráno	pastrano		
pelato	peeled	peláto	pelato	pīlātus	
pelliccia	fur	pellíčċa	pel:itʃ:a		
pettata	t is	pettátata	pet:ata		
pettiroso	robin	pettirósso	pet:iros:o		
pioppo	poplar	pióppo	pjɔp:o	plōppum	
piaggia	slope	piáġġa	pjadʒ:a		
piattola	scraggy	piáttola	pjat:ola		
pignatta	pot	pińńáta	pɪn:at:a		
pipistrello	bat	pipistréllō	pipistrel:o		
pinzo	pliers	pínzo	pintso		
pollone	sucker	pollóne	pol:one		
popone	melon	popóne	popone	pepōnem	
porcino	porcine	porċĭno	portʃino	porcĭnum	
pozzanghera	puddle	pozzánġera	pots:angera	puteācula	
prese	taken	prĕse	preze		
proda	shore	próda	proda		
prezzemolo	parsley	prezzémolo	prets:emolo		
primo quarto	first quarter	prĭmo quárto	primo kwarto		

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
prete	priest	prĕte	prete		
pendolo	pendulum	pĕndolo	pĕndolo		
pulcino	chick	pulčĭno	pultʃino	pullicĕnum	
pulenda dolce	Pulenda sweet	pulĕnda dŏlĉe	pulenda doltʃe		
pupilla	pupil	pupĭlla	pupil:a	pupĭllam	
puzzo	stink	púzzo	puts:o	pūtium	
puzzola	skunk	púzzola	puts:ola	putiolam	
rozzo	crude	rŏzzŏ	rodz:o		
rocchio	drum	rŏkkĭo	rok:jo		
radica	briar root	rádika	radika		
rancico	Rancic	ránĉiko	rantʃiko		
raspo	stalk	ráspo	raspo		
radice	root	radíĉe	raditʃe	radĭcem	
ragazza	girl	ragázza	ragats:a		
ragazzo	boy	ragázzo	ragats:o		
raganella	treefrog	raganĕlla	raganɛ:l:a		
ramaiolo	ladle	ramajŏlo	ramajolo		
raponzoli	raponzoli	rapŏnzoli	rapontsoli		
ravanelli	radishes	ravanĕlli	ravanɛ:l:i		
riccio	hedgehog	riĉĉo	ritʃ:o	erĭcium	
recinto	fence	reĉĭnto	retʃinto		
ricotta	ricotta	rikŏtta	rikot:a		
rigatino	bacon	rigatĭno	rigatino		
rimpiattino	hide and seek	rimpiattĭno	rimpjat:ino		
ronzone	Ronzone	ronzŏne	rondzone		
rustico	rustic	rústiko	rustiko	rŭsticum	
segale	rye	ségale	segale		
semola	semolina	sémola	semola	sĭmilam	
sorcio	mouse	sŏrĉo	sortʃo	sŏricem	
sodo	hard	sŏdo	sodo		
soglia	threshold	sŏll'a	sɔʎ:a	sŏleam	
salcio	willow	sálĉo	salʃo	sálicem	
salice	willow	sáliĉe	salitʃe		
sagrato	churchyard	sagráto	sagrato	sacrātum	
salamandra	salamander	salamándra	salamandra	salamāndram	
salciccia	sausage	salĉíĉĉa	salʃitʃ:a	salsĭcia	
salsiccia	sausage	salsíĉĉa	salsitʃ:a	salsĭcia	
salvastrella	Burnet	salvastrĕlla	salvastrel:a		

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
sbornia	drunkenness	Sbórni̯a	zbornja		
sbronza	drunk	Sbrón̥za	zbrondza		
scotta	sheet	skótta	skot:a		
scapolo	bachelor	skápolo	skapolo		
scaldaletto	warming pan	skaldalétto	skaldalet:o		
scaldino	Warmer	skaldino	skaldino		
scheggia	splinter	skéġġa	skedʒ:a	schĭdiam	
schacciata	crushed	skjaččáta	skjatʃ:ata		
sciocco	silly	ššókko	ʃ:ɔk:o		
sciapo	sciapo	šápo	ʃapo		
sciamannone	sciamannone	ššamannone	ʃ:aman:one		
sciamannato	sciamanno	ššamannáto	ʃ:aman:ato		
sciapito	sciapo	ššapító	ʃ:apito		
uscio	door	úššo	uʃ:o	ōstium	
scimmia	monkey	ššimmia	ʃ:im:ja		
scoiattolo	squirrel	skojaáttolo	skojat:olo	scurĭolum	
scorciatoia	shortcut	skorčatōja	skortʃatoja		
scricciolo	wren	skríččolo	skritʃ:olo		
seccatoio	squeegee	sekkatōjo	sek:atojo		
segatura	sawdust	segatúra	segatura		
sfoglia	puff pastry	sfòl'la	sfoʎ:a		
sugo	sauce	súgo	sugo		
somaro	ass	somáro	somaro		
soppressata	brawn	soppressáta	sop:res:ata		
sottana	soutane	sottána	sot:ana	subtānam	
sporta	shopping basket	spórta	spɔrta	spōrtam	
spaccone	braggart	spakkone	spak:one		
spazzatura	garbage	spazzatúra	spats:atura		
spetezza	spetezza	spetézza	spetets:a		
spigolo	corner	spígolo	spigolo	spĭculum	
spranga	bar	spránga	spranga		
sputo	spit	spúto	sputo		
sedano	celery	sédano	sɛdano		
stolto	fool	stólto	stolto	stŭltum	
stoppia	stubble	stóppia	stop:ja		
stalla	stable	stálla	stal:a		
stagnino	tinsmith	stañnino	stɔn:ino	stāgninum	
sito	site	sító	sito	sĭtum	
strolago	loon	strólagο	strolago		

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
strabico	squint	strábiko	strabiko		
stracco	tired	strákko	strak:o		
stradello	Stradello	stradêllo	stradɛl:o		
stregone	sorcerer	stregóne	stregone		
strullo	Strullato	strúllo	strul:o		
strizza	Winking	strízza	stri:ts:a		
succhiello	gimlet	sukkijéllo	suk:ʝɛl:o		
susina	plum	suSína	suzina		
topo	mouse	tópo	topo		
talpa	mole	tálpa	talpa	tálpam	
tacchino	turkey	takkíno	tak:ino		
tagliere	chopping board	tal'l'ère	taʎ:ɛre		
tarantola	tarantula	tarántola	tarantola		
tartaglione	Tartaglione	tartal'l'òne	tartaʎ:one		
tartaruga	tortoise	tartarúga	tartaruga	tartarūcam	
terriccio	soil	terriččo	ter:itʃ:o		
testone	blockhead	testòne	testone		
tincone	tincone	tinkóne	tinkone		
tirato	pulled	tiráto	tirato		
topino	Mouse	topíno	topino		
trogolo	trough	trógolo	trɔgolo		
trabiccolo	jalopy	trabikkolo	trabik:olo		
tirchio	mean	tírkjo	tirkjo		
trucioli	chippings	trúčoli	trutʃoli		
testa	head	tèsta	tɛsta	těstam	
testo	text	tèsto	tɛsto		
tuono	thunder	tùono	twono	tõnum	
uggioso	dull	uġġóso	udʒ:oso	odiōsum	
uncinetto	crochet	unčínètto	untʃinet:o		
unguanno	unguanno	unguáño	ungwan:o		
unguanno passo	unguanno step	unguáño passo	ungwan:o pas: o		
volpe	fox	vólpe	volpe	vŭlpem	
vaglio	screen	vál'lo	vaʎ:o	vállum	
vagabondo	vagabond	vagabòndo	vagabondo		
vicolo	alley	víkolo	vikolo	vīculum	
viottolo	path	vióttolo	vʝot:olo		
vitalba	clematis	vitálba	vitalba		
vitellino	calf	vitellíno	vitel:ino	vitéllinum	

Normalised form	English translation	ALT-CDI transcription	IPA transcription	Latin etymology	IPA transcription
verro	boar	věrrō	vɛrːo		
zeppa	wedge	zzěppa	tsːepːa		
zozzo	Zozzo	zzózzo	tsːotsːo		
zolla	clod	ẓẓólla	dzːolːa		
zazzera	mop	zzázzeria	tsːatsːera		
zittella	zittella	zzittělla	tsːitːɛlːa		
zolfanello	match	ẓẓolfaněllo	dzːolfanɛlːo		
ziro	Ziro	ẓẓíro	dzːiro		
zuccone	dodo	zukkōne	tsukːone		
zizzola	Zizzola	ẓẓíẓẓola	dzːidzːola		
becco	beak	běkkō	bekːo	bēccum	
caco	persimmon	káko	kako		
ciglio	edge	číl'lo	tʃiːlːo	cīlium	
grembiale	apron	grembiále	grembjale		
lumacone	snail	lumakōne	lumakone		
pagliaio	haystack	pal'lájō	paʎːajo	paleārium	
radici	roots	radíči	raditʃi	radīces	