

This table contains information pertaining to individual inscriptions as a whole.

1.1 Links

Text\_ID The ID of the inscription within this corpus. This number can

be used to link this table to other tables (\$\frac{1}{2}\$ sentences.csv and

links.csv).

Reference A single bibliographical reference for this inscription, typically

its reference in the principal printed corpus in general use. For instance, the corpus aims to provide the CIL reference for Latin

inscriptions and Crawford for Sabellic.

Name The informal name of the inscription. For instance, values of

this field include Fibula Praenestina and Iguvine Tables. This

field applies to only a select few texts.

1.2. Language information

Language The main language in which the inscription is written. Values

in the corpus are Latin, Oscan, Umbrian, Old Sabellic, Messapic

and Venetic.

Language\_Family The higher-level classification of the main language of the

inscription. For instance, Latin is classified as *Indo-European::Italic::Latino-Sabellic*. This classification is maximalist and includes speculative levels (such as placing Messapic under "Balkan"), to allow for maximal search

possibilities.

Language\_Variety This is generally equivalent to "Language", but may specify a

variety on a lower taxonomic level. For instance, in line with recent research, this corpus classes Faliscan as a variety of

Latin.

Script The alphabet or script in which the inscription is written. This

is relevant mostly for Sabellic, which was written in various

scripts.

1.3. Chronological information

Date\_after The earliest possible date of this inscription.

Date\_before The latest possible date of this inscription.

#### 1.4. Provenance

Provenance The location (town or city) in which this inscription was found.

In some cases inscriptions are known to have been written in a different location to where they were found. In those cases, provenance provides the location where the text was written.

GeoID The ID allocated to this location by Trismegistos.

Latitude The latitude of this location.

Longitude The longitude of this location.

#### 1.5. Additional information

Finite\_verb [Automatically generated field]

This field returns the value "TRUE" if the inscription in question contains at least one finite verb form, "FALSE" otherwise. The purpose of this field is to search specifically for inscriptions that constitute full sentences (as opposed to

inscriptions containing only names, for instance).

Analysable\_token [Automatically generated field]

Similar to the above, but broader, this field returns "TRUE" if the inscription contains at least one token with a linguistic analysis. This allows the user to sift out inscriptions which are wholly unintelligible or, for instance, contain only single

letters.

Text\_length [Automatically generated field]

The number of individual tokens this text links to. This is essentially a count of the number of words (and clitics) that

the inscription contains.



### 2. Sentences

This table contains information pertaining to individual syntactic units (sentences) within a given inscription. Frequently, an inscription contains only a single syntactic unit.

#### 2.1. Links

Text\_ID The ID of the inscription to which this sentence pertains

(allowing this table to be linked to 1 texts.csv).

Sentence\_ID A fixed and unique number referring to this specific sentence

(allowing this table to be linked to  $\downarrow$  tokens.csv).

Sentence\_position A number which keeps track of the order of sentences within

an inscription. This counter starts from 1 in every new

inscription and is therefore not a unique ID.

#### 2.2. Sentence

Sentence This field contains the full text of the sentence in question.

Since the focus of this corpus is linguistic, not epigraphic, the resolution here is low. The text provided is a trimmed-down version of the text provided by published corpora (or papers in

which more recent inscriptions appear).

This field is meant only for convenient reference and should on absolutely no level be considered as an alternative to

consulting an epigraphic apparatus.

Section In the rare cases where a text is long enough that a subdivision

into larger units than sentences, this field keeps track of them. Currently this applies only to the Iguvine Tables (where this

field keeps track of tables running from I to VIIb).

## 3. Tokens

This table contains information on individual tokens within a syntactic unit, usually words or enclitics.

3.1. Links

Text\_ID The ID of the inscription to which this token pertains (allowing

this table to be linked to 1 texts.csv).

Sentence\_ID The ID of the sentence to which this token pertains (allowing

this table to be linked to **1** sentences.csv).

Sentence\_position A number which keeps track of the order of sentences within

an inscription. This counter starts from 1 in every new

inscription and is therefore not a unique ID.

Token\_ID A fixed and unique number referring to this specific token

Token\_position A number which keeps track of the order of tokens within a

sentence. This counter starts from 1 in every new sentence and

is therefore not a unique ID.

3.2. Token

Token The token itself.

Token\_clean [Automatically generated field]

This field is intended for ease of searching. It reproduces the token stripped of all special or non-alphabetic characters, and

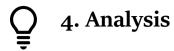
in transliteration where a non-Latin symbol is used.

#### 4.5 Syntax

Syntactic information is provided on the level of Token, not Analysis, because it concerns the way in which tokens relate to each other. It is therefore tricky to incorporate multiple possible syntactic analyses into the design of a corpus, and this has been avoided.

Relation Under construction

Head Under construction



This table provides morphological information. This information is provided on a level subordinate to tokens, because in an obscure or fragmentary text it is possible for a given token to have multiple plausible analyses.

4.1. Links

Text\_ID The ID of the inscription to which this analysis pertains

(allowing this table to be linked to 1 texts.csv).

Sentence\_ID The ID of the sentence to which this analysis pertains (allowing

this table to be linked to 1 sentences.csv).

Sentence\_position A number which keeps track of the order of sentences within

an inscription. This counter starts from 1 in every new

inscription and is therefore not a unique ID.

Token\_ID The ID of the token to which this analysis pertains (allowing

this table to be linked to 1 tokens.csv).

Token\_position A number which keeps track of the order of tokens within a

sentence. This counter starts from 1 in every new sentence and

is therefore not a unique ID.

Analysis\_ID A fixed and unique number referring to this analysis.

4.2. Lemmatisation

Lemma The lemma to which the current token belongs.

Note that the lemma field does not disambiguate homonyms. Thus, searching *dico* will return forms of *dicare* and *dicere* without fear or favour. This avoids the need to resort to arbitrary and confusing disambiguations such as *dico\_2*. If one wishes to find only forms of *dicere*, use this field in combination with "morphological type" to specify the

conjugation.

The lemma is the dictionary form in Latin. For other Italic languages it is undeclined root (so *deyk* in Sabellic), to avoid

having to reconstruct unattested forms.

Lemma\_simplex The simplex form of a compound lemma. For instance, the

"Lemma" entry for the token perfecerit is perficio, while the

"Lemma\_simplex" is simply facio.

Lemma\_frequency [Automatically generated field]

The number of times this particular value occurs in the "Lemma" field. As noted above, be aware that "Lemma" does not discriminate between homonyms.

Language

The language this token is analysed as belonging to. This field may differ from the "main language" of the inscription as a whole (provided in 1 texts.csv), as inscriptions often contain tokens in multiple languages. In addition, the same token may have different possible analyses in different languages.

#### 4.3. Morphology

Since the aim of this corpus is to facilitate comparison between different ancient languages, it strives for a maximally descriptive (and therefore linguistically accurate) terminology in its morphological analyses. For ease of reference, the field POS\_code converts the conventions employed here into standard POS tags.

Where a dimension of morphological description is not applicable for a given token, it is annotated simply as "N.A.".

Morphological\_type

This field keeps track of inflectional declensions and conjugations, usually based on the final or characteristic sound of the stem.

The terms used are descriptive, and based on the final or characteristic sound of the stem. The Latin "first declension" thus appears as "A-stem". This terminology is less arbitrary than the traditional terms and facilitates intercomparison between the languages involved.

Part\_of\_speech

The part of speech of the token analysed (noun, verb, etc.).

Stem

This field keeps track of basic stem categories which host morphological inflection. These include the present and perfect stems of Latin verbs, but also, for instance, the degrees of comparison of adjectives.

Tense

Tense includes any category denoting temporal distinctions. For the Latin present stem this would include simple (e.g. "amo"), future ("amabo") and past ("amabam"). Note that the perfect indicative is also analysed as the "simple" (unmarked) tense form, where as "past" is reserved for the pluperfect.

Mood

The mood of the token analysed.

In line with some traditional grammatical analysis, it has been found convenient to group nominal forms of verbs (participle, gerundive, infinitive, supine) in this category as well: although these are not modal categories, they are in complementary

distribution with modal exponents.

Diathesis The grammatical voice or diathesis of the token analysed.

Person The grammatical person of the token analysed.

Number The grammatical number of the token analysed.

Gender The nominal class of the token analysed. In the Italic languages

this is usually (but not always) equivalent to the traditional category of "gender", although note Etruscan, for instance, where an animate is frequently opposed to an inanimate class.

Case The grammatical case of the token analysed.

Category This is a wastebasket field for any categorial distinction not

captured by the above. It assigns alphabetic indices (A, B, C...). Currently this is applicable only to the Latin imperatives "ama" and "amato", which are labelled imperative A and imperative B in the corpus. However, the distinction between the two Etruscan genitives, for instance, would also be annotated in

this field.

POS\_code [Automatically generated field]

To facilitate intercomparability with other projects, this field translates the morphological analysis provided above into

standard 9-slot POS tags.

Tags [Automatically generated field]

This field is similar to the above, in that it provides a succinct summary tag with morphological information, but it does so

in human-readable form.

For instance, *perf-3-sg* or *M-nom-pl*.

As is apparent from the first example, for the sake of brevity and convenience some categories have default values (thus "indicative", "active" and "present" are not explicitly

signalled).

#### 4.4. Semantic information

Meaning A short English translation of the token in question.

Meaning\_category This organises a number of lemmata into very broad semantic categories. The corpus currently contains:

- PROPER (for proper names e.g. *Manius, Zeus, Roman*)
- RELATION (e.g. daughter, son-in-law, mother)

- HUMAN (e.g. slave, praetor, flute-player)
- BODY (e.g. head, hand, liver)
- ANIMAL (e.g. horse, sheep, lion)
- VICTUALS (e.g. water, wine, meal)
- NUMBER (e.g. three, twice, fourth)
- POSSESSIVE (e.g. your, his, my)
- BUILDING (e.g. house, statue, bridge)
- MATERIAL (e.g. *gold*, *silver*, *wood*)
- COLOUR (e.g. white, black, red)
- TIME (e.g. day, year, hour)
- GEOGRAPHY (e.g. sea, river, field)
- POLITICS (e.g. vote, assembly, plebs)
- LAW (e.g. fine, judge, punishment)
- ECONOMY (e.g. money, debt, sell)
- RELIGION (e.g. goddess, sacrifice, pray)

Meaning\_subcategory

This is currently used only to distinguish different categories of proper name, and contains the values PERSONAL, TOPONYM and THEONYM.

Classical\_Latin\_equivalent

This field contains a close cognate or semantic equivalent of the current lemma in orthographically standardised classical Latin. This field is more consistent and regularised than the "Translation" field above, and therefore facilitates intercomparison between different languages.

For instance, one may wish to find all occurrences of a first person pronoun in Italic languages. One can do so simply by searching "ego" in this field.

This field also standardises Old Latin forms to classical Latin. Thus the Old Latin verb *to*, *tare* "to steal", for instance, becomes *furor* in this field.

Classical\_Latin\_form

Equivalent to the above, but the classical Latin equivalent is here appropriately inflected, rather than simply giving a dictionary form.

TAM\_analysis

Information on the tense, aspect and modality of verb forms. This column contains the data used for Pitts (2020) and is limited to finite verb forms in Sabellic.

#### 4.5. Phonology and orthography

Standard\_aligned

This field contains a standardised form of the current token which has been aligned (using hyphens) with the attested form of said token in the following field. For instance, for analysis ID #1469381, these fields read:

Standard\_aligned: cu-raverunt

Form\_aligned: coeraveront

Since the alignment guarantees that both strings are the same length, they can be searched with regular expressions to gather data on sound changes (e.g. monophthongisation) or orthographic variants (e.g. the use of gemination, or the use of k, q and c).

Currently these fields contain data for Latin and Sabellic. However, the field Standard\_aligned contains somewhat different data for these languages. In the case of Latin, the "standard" form is simply the classical Latin form, equivalent to the field Classical\_Latin\_form (except aligned). For Sabellic, however, the field contains one of the older etymological forms suggested in Unterman (2000).

This field is under construction: currently the Latin alignment is automatic (using LingPy), although the Sabellic data have been aligned manually.

Form\_aligned

See above.

# Links

This table links to <u>texts.csv</u> and serves to link the corpus to extensive metadata and bibliographical information hosted by Trismegistos.

Text\_ID The ID of the text in this corpus (links to 1 texts.csv)

Trismegistos\_ID The ID of this text in the Trismegistos project.