

TLG Technical Note 001: Greek Word Definition

Authored: Nick Nicholas, TLG
Maintained by tlg@uci.edu
Created: October 1999
Last Revised: 2002-11-23

The TLG maintains a word index of the Greek words occurring in its texts. Versions of the word index appear on the CD ROMs published by the TLG, and on the online TLG Search Engine. The following documents what constitutes a distinct word in the compilation of the TLG word index.

1. Definitions

Different versions of the TLG word index are referred to in the following:

- **CD ROM #D** is the CD ROM created at the Packard Humanities Institute for the TLG in 1992.
- **CD ROM #E** is the CD ROM created at the TLG in 1999. The software to create the CD ROM was developed afresh, and the resulting CD ROM differs in several details from its predecessor.
- The **Web index** is the TLG online version of the word index.

2. Word delimiters

A new word is assumed to occur whenever a blank, a dash (— ; Beta code _), a punctuation sign, or a new line occurs in the source text.

Punctuation is defined as one of the following characters:

- . (period: .)
- , (comma: ,)
- ; (Greek question mark: ;)
- : (Greek raised dot: ·)
- %4 (exclamation point: !)
- %10 (dicolon: :)
- %100 (Roman semicolon: ;)

However, %1 (Roman question mark: ?), which may indicate a doubtful letter, is not considered to be punctuation.

Nor are quotation marks, which may on occasion be used as editorial brackets, or as continuing block quotation markers.

Comma is only considered punctuation when not followed immediately by a letter; otherwise, it is considered a hypodistole.

On the Web Index, a new word is also assumed to occur whenever an apostrophe occurs in the source text; the current TLG orthographic norms do not admit word-internal apostrophe. Thus, δ'ο(/s δ'ός is considered two words, δ' and ο(/s. This also applies where the apostrophe has supplanted the coronis: γ'οὐν γ'οὐν (normally γοῦν) =N γοῦν is still analysed as the two words γ' and οὐν.

3. Normalization

3.1. Punctuation

Punctuation is normally stripped out of the word. (See below for exceptions.)

3.2. Hyphenation and Non-Text

Hyphenated words are joined; words in non-text brackets ({}) are indexed separately from words outside non-text brackets, and do not interfere with hyphenation. For example,

A) /N- {STR.} ἄν- στρ.
Q̄RWPOS θρωπος

is indexed as A) /NQ̄RWPOS ἄνθρωπος, and the marginal note STR στρ is indexed separately.

There is one exception: {27 }27 (formerly occurring in work 1595.107), indicating apograph textual emendation, are treated as brackets, and are thus ignored in extracting words. For example, ARIS{27T}27WN ἀριστ*ων is indexed as ARISTWN, not ARISWN, T. (The escape code has been reassigned as of August 2000, so this exception is no longer relevant.)

All other non-text brackets by definition contain text extraneous to the main text (e.g. marginalia, stage directions, titles, editorial emendations); so their content is never concatenated with main text words, even if no space intervenes.

3.3. Roman Script

Words in other languages are ignored in the word index; this currently concerns only words tagged as being in Roman script, since words from other languages transliterated into Greek are not yet tagged distinctly. Thus, in a text like *)APOMA/SAR &[1addidi]1 \$E)/FH Ἀπομάσαρ (*addidi*) ἔφη, only the ‘Greek’ words A)POMA/SAR ‘Abu-Masar’ and E)/FH are indexed.

3.4. Beta escapes

Characters other than elision markers, the Greek alphabet, breathings, accents, iota subscripts, and commas—i.e. almost all beta escapes—are also ignored. For example, "6\$10A)\$%5N?<1[H/]R>1#13 «ἀϗ[ή]ρ* is indexed as A)NH/R ἀνήρ.

Characters within a word in non-text or Roman font are also ignored (but see *Partial Words* below.)

3.5. Diacritics

Accents are regularized: grave accents become acute, only the first accent in a word is retained, and words are converted to lower case. Thus,

- A)NH\R ἀνήρ is indexed as A)NH/R ἀνήρ,
- A) /NQ̄RWPO/S ἄνθρωπος as A) /NQ̄RWPOS ἄνθρωπος,
- *STUMFALI/A Στυμφαλία as STUMFALI/A στυμφαλία, and
- *E*S*T*I*A ΕΣΤΙΑ as ESTIA εστια.

3.6. Hypodiatolae

Since commas are potentially hypodiatolae, they are retained in the word, and are stripped out only if they do not correspond to a known instance of a hypodiatole; thus, ο(/ ,τι ὅ,τι is distinguished from

ο(/ΤΙ ὅτι. There are nine words with hypodiatolae known in our texts:

ο(/,ΤΙ, ο(/,ΤΙΡΕ, ο(/,ΤΤΙ, ΤΟ(/,ΤΕ, ο(/,ΤΕ, ο(/,ΤΟΥ, ΤΑ(/,ΤΕ, ΤΟ(/,Τ', Η(/,ΤΕ
(ὅ,τι, ὅ,τιπερ, ὅ,ττι, τό,τε, ὅ,τε, ὅ,του, τά,τε, τό,τ', ἥ,τε)

Hypodiatolae are new to CD ROM #E; in CD ROM #D, they were stripped out of the word.

3.7. Coronides

All crasis markers in words are now retained; this is a change from CD ROM #D, in which crasis markers were dropped if they occurred after the first syllable of the word. Thus, ΚΑΛΟΚΑ) /ΓΑΘΟΣ καλοκάγαθος is represented in the CD ROM #D index as ΚΑΛΟΚΑ/ΓΑΘΟΣ καλοκάγαθος, but ΚΑ)ΓΩ/ κᾶγώ remains ΚΑ)ΓΩ/; both forms retain their crases on CD ROM #E.

On CD ROM #E the coronis is not normalized: the rough coronis is treated as distinct from the smooth coronis, although these are mere notational variants. Thus, χω(ς χῶς is listed as a distinct word from χω)ς χῶς. The two coronides are conflated on the online version of the index. However, a rough breathing is considered a coronis only when following an aspirated stop; otherwise, it is an internal breathing, and left as is (e.g. *)ΑΒΡΑΑ(/Μ ᾿Αβραᾶμ)

3.8. Internal breathing marks

The internal breathing marks on double rho are removed only if they are predictable; thus, ρ)ρ(is converted to ρρ, but ρ(ρ(is left as is. For example, Α) /Ρ)Ρ(ΨΤΟΣ ᾿ῤῥῥωστος is indexed as Α) /ΡΡΨΤΟΣ ᾿ρρωστος.

This is another change from CD ROM #D, which also converted ρ(ρ(, ρ(ρ) and ρ)ρ) to ρρ.

3.9. Character Stacking

On the web index, instances of characters superposed above other characters, as tagged with the Beta escapes <10 ... >10<11 ... >11 are deemed to be textual variants, and two versions of the word are recorded in the word index: that with the above reading, and that with the below reading. Thus, the string Α)<10QH=>10<11MU/>11ΝΑΙ ᾰ^{μύ}θη^νναι is parsed as two words in the same spot, Α)QH=ΝΑΙ and Α)MU/ΝΑΙ.

3.10. Dittography

On the web index, if a bracket is known to indicate dittography or editorial or scribal deletion for a work, the indication is respected, and the marked portion of the word is excluded from the indexed word. If however the deletion bracket incorporates a word in its entirety, the word is indexed, as a variant reading. E.g. given the information that [3] 3 { } in a given work are editorial deletions, Ε) /GW[3GW] 3GE ἔγω{γω}γε is indexed as Ε) /GWGE ἔγωγε -- i.e. Ε) /GWGE ἔγωγε. However, in a text like Ε)GW\ [3DE\] 3 ΕΙ)=ΠΟΝ ἔγῳ{δε} ἔλπον, all three words are indexed.

4. Partial words

When a partial word is discovered in a text, the partial word is included in the current word index only if it contains two or more Greek letters (excluding diacritics) in sequence. This is at variance with CD ROM #D, which required three Greek letterals (including diacritics) in sequence. Thus, CD ROM #E includes word fragments like ΒΑ! βα. and excludes Α) /! ᾰ.; CD ROM #D would exclude ΒΑ! and include Α) /!.

On CD ROM #D, only the missing letter code (!) was regarded as a partial word boundary. The repertoire has been significantly expanded for CD ROM #E. A partial word boundary is discovered from the occurrence of:

- a missing letter code (!);
- a transition from Roman to Greek script within a word (though not vice versa):
 - `&lIustinian$10\S Iustinianὀς`;
- a final left bracket, followed by no Greek letterals between it and the word delimiter, *if* the preceding word is unaccented, oxytone (not barytone), or ends in a letter which is not legal for a Greek word (vowel, *s*, *n*, '):
 - `ABG[αβγ]`,
 - `$13ABG[$10_ αβγ]—`;
- an initial right bracket preceded by a space (or equivalent `^`, though not `@---` which acts as indentation), or the beginning of the line, and followed by text, disregarding any intervening beta escapes other than space or equivalents:
 - `]ABG]αβγ`,
 - `^16]1ABG)αβγ`,
 - `] %ABG]†αβγ`;
- a left bracket followed by a right bracket without any intervening characters or beta-escapes:
 - `ABG[] αβγ[]`,
 - `[2]2ABG ()αβγ`,
 - `AB[]GD αβ()γδ`;

or:

- a line-initial instance, after a line-breaking hyphen, of [] (with optional digits qualifying the brackets: [1] 1 (), [2] 2 ⟨ ⟩ etc.):
 - $$\begin{array}{l} \text{AB- } \alpha\beta- \\ [] \quad [] \end{array}$$
- or] preceded by:
 - nothing:
 - $$\begin{array}{l} \text{AB- } \alpha\beta- \\]^{4G} \quad] \gamma \end{array}$$
 - non-spacing beta escapes (\$&<>{}): font shifts, quasi-brackets (typographical formatting), non-text):
 - $$\begin{array}{l} \text{AB- } \alpha\beta- \\ \$10]^G \gamma \end{array}$$
 - or [followed by at least one space (but not an equivalent beta-escape, since these indicate *indentation* rather than lacunae):
 - $$\begin{array}{l} \text{AB- } \alpha\beta- \\ [\text{ AB}]^G [\alpha\beta] \gamma \end{array}$$
 - but not

$$\begin{array}{l} @ @ *) / \text{AB- } \text{''}A- \\ [@ @]^{\text{DHRA}} [] \beta \delta \eta \rho \alpha \end{array}$$
- a beta-escape which can stand in for omitted text:
 - the asterisk %2 (but only when followed immediately by another asterisk in the text);
 - the hyphen %19 (but only at the beginning or end of a word or word continuation, not word-internally);

- the diacritics without letters %30-%39, %132-%134;
- the metrical signs %40-%49, %140, %141, %144, %145 ;
- the partial letter sign #7;
- the unintelligible letter sign #99;
- the blot/dot sign #459.

Thus,

- Any one of the following has an initial partial word boundary:

!A . α
 @]A] α
]A] α
]\$10A] α
 #7A \forall α

- Any one of the following has a final partial word boundary:

A! α .
 A[@ α [
 A[α [
 A%2%2 α^{**}
 (but not A%2 α^*)
 A#459 α ■

- Any one of the following contains a word boundary partial on both sides:

A!B $\alpha.\beta$,
 A[]B $\alpha[]\beta$,
 A%40%40%41B $\alpha^{\sim}\sim\beta$.

- And in the following hyphenated words, the hyphen becomes a partial word boundary:

A- []B	α - [] β
A- [&`12 spaces\$]B	α - [12 spaces] β
A-]B	α -] β
A- %19B	α - - β
A- [%40]B-]	α - [~] β -]

However,

A- α -
 @@B] β]

does not have the hyphen as a partial word boundary: the @ act merely to indent, and not to indicate a lacuna. Likewise, the bracket pair in AB[1%1]1 $\alpha\beta(?)$ does not indicate a lacuna: the presence of any code between brackets is taken to indicate the bracket pair is no lacuna --- unless the code itself denotes a lacuna (e.g. AB[1%40]1G $\alpha\beta(\sim)\gamma$).

Furthermore,

A- α-
\$13]B β

contains a partial word boundary; the \$13 is a font shift which can be ignored as a non-spacing Beta escape, not representing any actual text intervening between the hyphen and the right bracket. However,

A- α-
#13]B ※]β

does not contain a partial word boundary, since the *asteriskos* is a spacing Beta escape, and is deemed to start the current text line.

By this reckoning, the following are incomplete words:

- ο(MO! in 0059.005:

A)NTI\ TOU= "3O(MO%19"3, "3A)%19"3
ἀντὶ τοῦ ‘όμο-’, ‘ά-’

- !NEMON in 0006.029:

[*U*Y.] %19NEMON A)/GAGE/ POTE/
[ΥΨ.] -νεμον ἄγαγε ποτέ

- SAP[F! in 0009.001:

*SAP- Σαπ-
[F%19]GXANON φ-]γχανον

The following are complete words, by contrast:

- XILIONTA%19ETHRI/DAN in 2866.001 (χιλιοντα-ετηρίδαν),
- YALLO/%19[4MENON]4 in 2703.001 (ψαλλό- [[μενον]]).

And the alternative readings A)MAMA/CUD[1%19OS, %19ES]1 (0009.001: ἀμαμάξυδ(-ος, -εξ) = A)MAMA/CUDOS *vel* A)MAMA/CUDES) are indexed as A)MAMA/CUDOS, !ES.

If a beta escape denoting a lacuna occurs at the beginning of a word, it is deemed part of the word unless a blank, dash, or new line follows it: the words %40A ~α and [%40]A [~]α contain the breve, and thus are words with incomplete beginnings; but the word [%40] A [~] α is not incomplete, since the breve is considered a separate word.

Exceptionally, as of November 2002, if a lacuna is followed immediately by a vowel with a breathing mark, the word is deemed complete. This is because the breathing mark is deemed to mark the beginning of a word rather than a coronis, and the source texts are rarely consistent in delimiting words from lacunae. For instance, !!!A)NH\R ...ἀνήρ is treated as the complete word A)NH/R ἀνήρ .

5. Special hyphen rules

If a beta escape denoting a lacuna occurs on a new line after a hyphen, it is deemed part of the hyphenated word automatically; if a blank, dash or new line ensues, the word is then terminated. Thus,

A- α-
%40B ~β

is indexed as A!B = A!, !B, but

A- α-
%40 B ~β

is indexed as A!, B.

A hyphenated word is also terminated immediately if a space follows a left bracket in a hyphenated word continuation; what follows that space is not deemed part of the word. For instance,

I (STORI /- ιστορί-
[A)NAGKAI /AN] [άναγκαίαν]

is indexed as I (STORI /!, A)NAGKAI /AN—not I (STORI /A)NAGKAI /AN *ιστορίάναγκαίαν*.

If a hyphen follows a lower case letter within a portion of a word containing an unclosed left bracket, and the continuation of the word in the next line consists of a capital letter not followed by at least one more capital letter, then it is deemed impossible for the two fragments to be part of the same word, and the hyphen is considered a final partial word boundary. Thus, in an instance like

F[ILE- φ[ιλε-
)AS]KLHP IOU= 'Ασ]κληπιου̃

(0020.004), the program refuses to join the two fragments, and indexes them as FILE!, *)ASKLHP IOU=. This results because the bracket in F[ILE- explicitly flags what follows in the word as fragmentary. The casing of *)ASKLHP IOU=, on the other hand, with its initial capital, indicates that it must constitute a new word. Without that casing, the algorithm would still join the two word fragments:

F[ILE- φ[ιλε-
S]KLHP IOU= σ]κληπιου̃

is interpreted as FILESKLHP IOU= *φιλεσκληπιου̃*, and

F[ILE- φ[ιλε-
*S]*K*L*H*P*I*O*U Σ]ΚΛΗΠΙΟΥ

as FILE*S*K*L*H*P*I*O*U *φιλεΣΚΛΗΠΙΟΥ*.

If a hyphen preceded by an unclosed left bracket is followed by a continuation word with a *breathing mark*, as distinct from a coronis (for which the heuristic is that the vowel with the breathing mark occur before any consonants in the line), it is likewise deemed impossible for the two words to be joined. Thus,

*KLEANAK[TID- Κλεανακ[τιδ-
H(ἥ

(0009.003) is indexed as KLEANAKTID!, H(: by having a breathing mark, H(cannot be anything but the beginning of a new word.

The necessity of the bracket in such checks is shown by the instance of

in the non-fragmentary text 2734.013.

If a hyphen is followed in the same line by a letter or punctuation, the word is deemed to be incomplete, and is not joined with the next line. Thus, if a line ends in FILE-., that text is not joined with the next line, but is indexed as FILE! .

6. Ellipses & Abbreviations

An instance of more than one contiguous dot, either within or on the boundary of a word, is treated as a lacuna: thus, ..A, A.., and A..B are analyzed as !A, A!, and A!B = A!, !B. If a space intervenes between the multiple dots and the word, the dots are not treated as part of the word, which thus remains complete; this is how ellipses as punctuation are distinguished from lacunae. Thus, A)NH/R ... ἀνὴρ ... is analyzed as A)NH/R, not A)NH/R!.

Multiple dots in the continuation of a hyphenated word are treated like the other lacunae: A- B... is indexed as AB!, and A- ...B as A!B. A single dot, on the other hand, terminates a word: A.B is indexed as A followed by B, and not as AB or A!B = A!, !B. If the dot was to represent a missing letter, it would have been encoded as ! in the Beta code instead.

Abbreviations are thus treated as separate words rather than a single word: e.g. κ.ο.κ. *κ.ο.κ.* (= *καὶ οὕτω καθεξῆς* 'and so forth') is indexed as κ, ο, κ.

The one exception to this (new to CD ROM #E) is κ.τ.λ. *κ.τ.λ.* ('etc.') which is indexed as the single word κτλ *κτλ* (with no periods.)

The foregoing is an algorithmic approach to determining what constitutes a complete or partial word. It is by no means infallible, and several words will inevitably be indexed wrongly. Substantial further improvement of the process, however, can only be done manually.