



Dialectal Corpora Building (for oral and written sources)

by

Nikitas N. Karanikolas



An International Lecture, Based on the AMiGre project, Thalis framework.



European Union
European Social Fund



OPERATIONAL PROGRAMME
EDUCATION AND LIFELONG LEARNING
investing in knowledge society
MINISTRY OF EDUCATION & RELIGIOUS AFFAIRS, CULTURE & SPORTS
MANAGING AUTHORITY

Co-financed by Greece and the European Union



NSRF
2007-2013
programme for development
EUROPEAN SOCIAL FUND

This research has been co-financed by the European Union (European Social Fund - ESF) and Greek national funds through the Operational Program "Education and Lifelong Learning" of the National Strategic Reference Framework (NSRF) - Research Funding Program: Thalis. Investing in knowledge society through the European Social Fund.



AMiGre

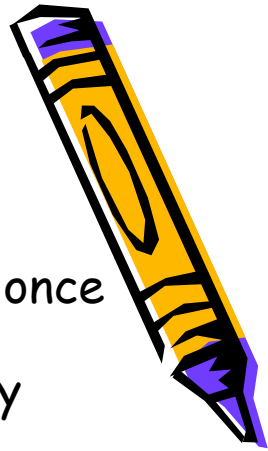
- Introduction
- Sources
- Applications Overview
- Design Overview



Introduction

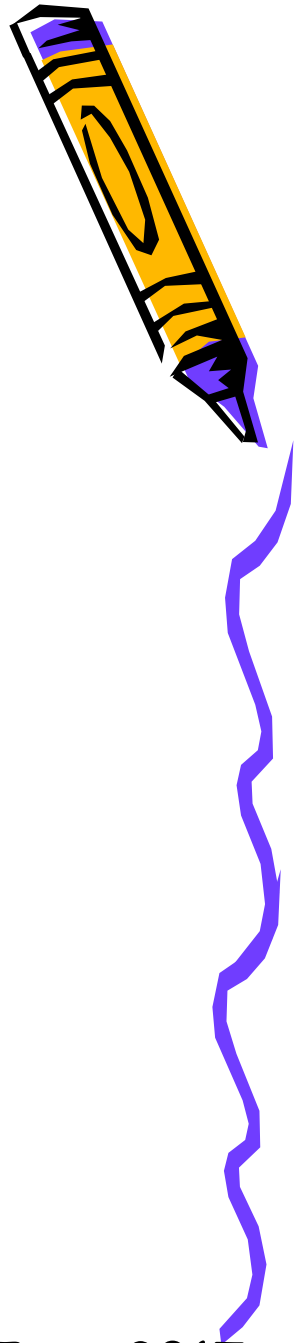
Pontus, Cappadocia, Aivali: In search of Asia Minor Greek

- Cappadocian, Pontic and Aivaliot are three varieties of Greek once spoken in Asia Minor.
- Following the population exchange between Greece and Turkey enforced by the Lausanne Treaty (1923), speakers of these varieties were relocated in various parts of Greece, leaving behind their lands and material possessions and carrying along with them only their history, mores and customs, traditions, and language.
- Due to their long-lasting contact with neighbouring Turkish and their isolation from the other Greek dialects, the Asia Minor Greek varieties (especially Pontic and Cappadocian) are regarded as ideal case-studies for shedding light on the linguistic evolution of Greek as well as on various language contact phenomena.
- Crucially, the three dialects channel a rich cultural and linguistic heritage but face a severe danger of extinction: the number of first generation refugees is shrinking rapidly and the next generations are being gradually absorbed by adstratal Modern Greek, both culturally and linguistically.
- Thus, the necessity of describing and preserving this precious piece of heritage is vital.



AMiGre

- Introduction
- **Sources**
- Applications Overview
- Design Overview



Sources

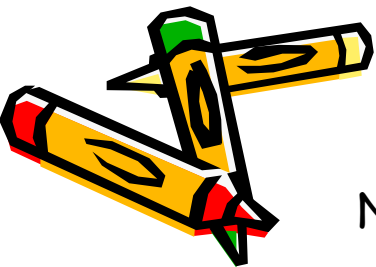
- Lexical data for 3 Asia Minor Greece dialects (see lecture *Dialectal lexicon building: requirements and technical specifications*)
- Digital Audio files (WAV) and Annotations (TextGrid - Praat - files)
- Written sources (digitized)
- their homogenized Transcriptions
- Morphological annotations
- Syntactic and Semantic annotations



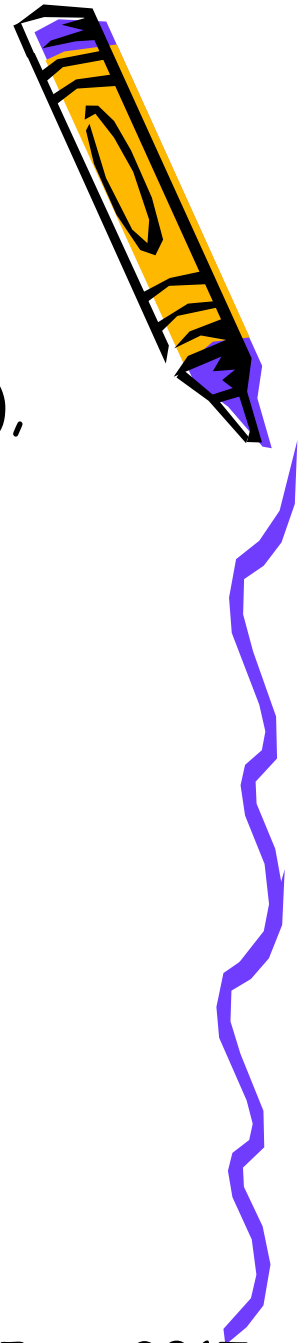
Digital Audio files (WAV) and Annotations (TextGrid, ELAN files)



- There are digitized audio files (WAV) of dialectal conversations
- Annotated with: sentences, phrases, syllables, segments, vowels, consonants, tones, phenomena, etc (TextGrid - Praat - files, ELAN files)

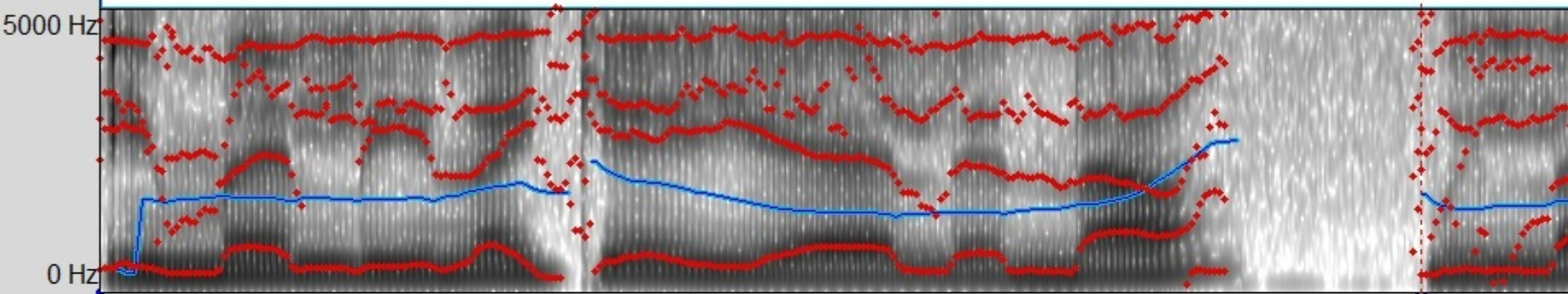


Items to annotate and some phenomena



- Tone / accent
(is it a question, an imperative statement, etc),
- Pragmatological phenomena,
- Intonations,
- Morphological words,
- Phonological words,
- Intonation phrases,
- Intonation sentences,
- Phonemes,
- Voices,
- Accent phenomena (stress, unstress).





1 (sentence) Η Μελίνα και η Έλενα μιλάμε με τον Μαν

2 Η Μελίνα και η Έλενα (intonation phrase)

3 Η Μελίνα και η Έλενα (intonation word)

4 i m e l i n a c e i e l e n a (phoneme)

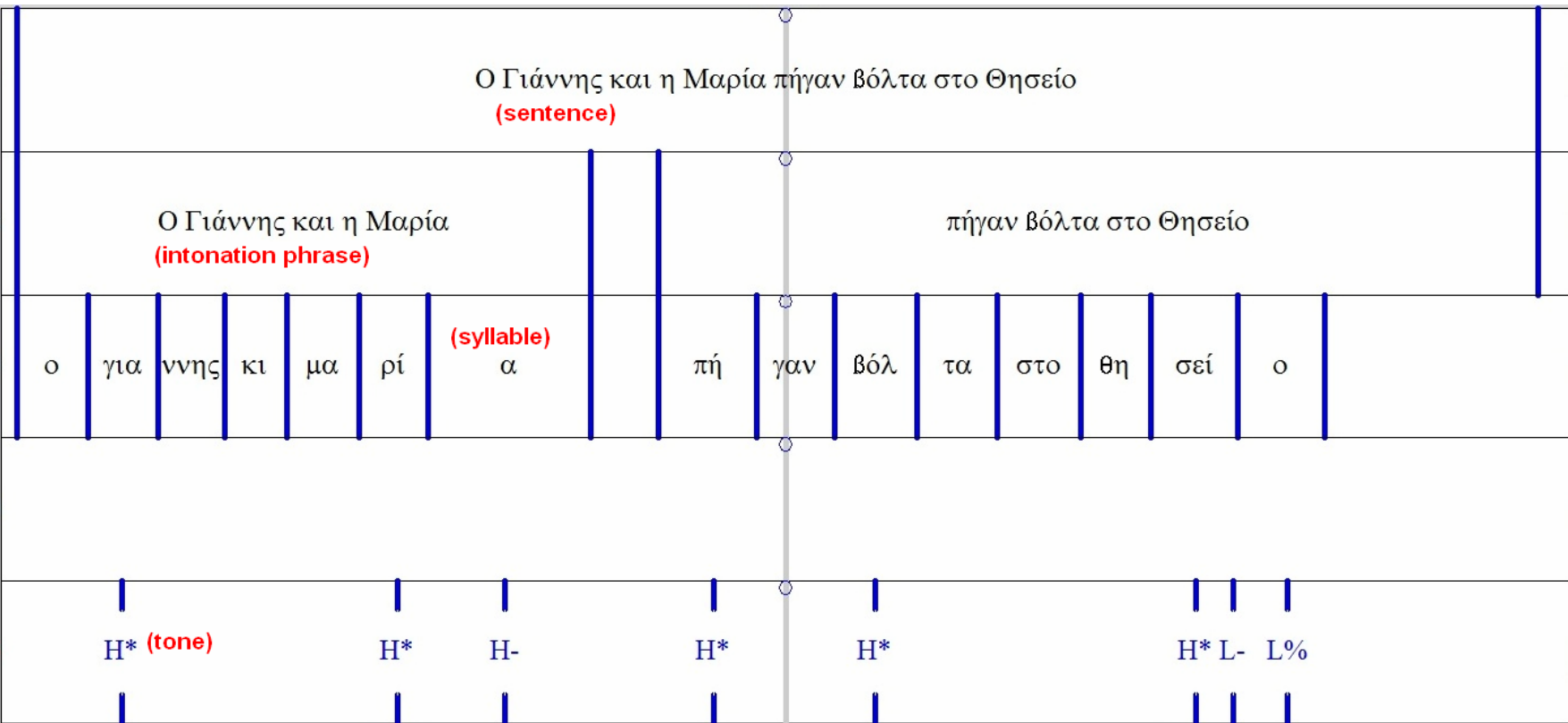
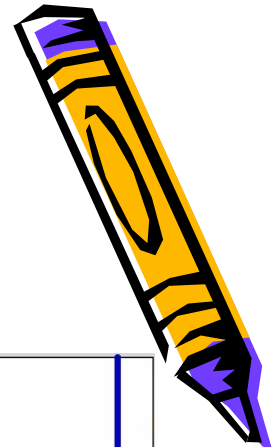
5 L*+H L* (tone) H-

1.717434

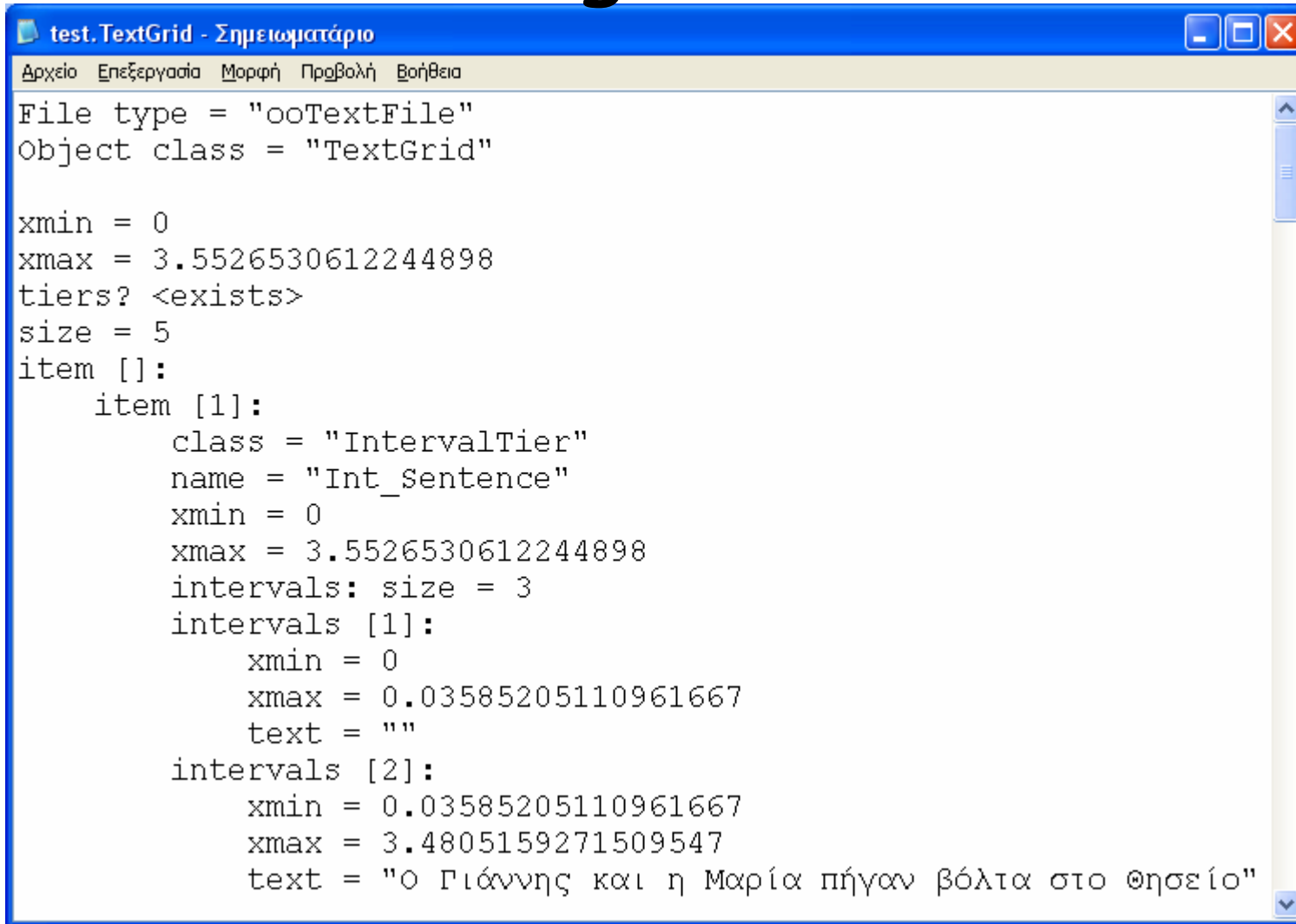
14521 0.414521 Visible part 2.958813 seconds

Total duration 24.879728 seconds

Annotations in various levels with PRAAT



Textgrid file

A screenshot of a Praat TextGrid file window. The window title is "test.TextGrid - Σημειωματάριο". The menu bar includes "Αρχείο", "Επεξεργασία", "Μορφή", "Προβολή", and "Βοήθεια". The main text area contains the following content:

```
File type = "ooTextFile"
Object class = "TextGrid"

xmin = 0
xmax = 3.5526530612244898
tiers? <exists>
size = 5
item []:
  item [1]:
    class = "IntervalTier"
    name = "Int_Sentence"
    xmin = 0
    xmax = 3.5526530612244898
    intervals: size = 3
    intervals [1]:
      xmin = 0
      xmax = 0.03585205110961667
      text = ""
    intervals [2]:
      xmin = 0.03585205110961667
      xmax = 3.4805159271509547
      text = "Ο Γιάννης και η Μαρία πήγαν βόλτα στο Θησείο"
```

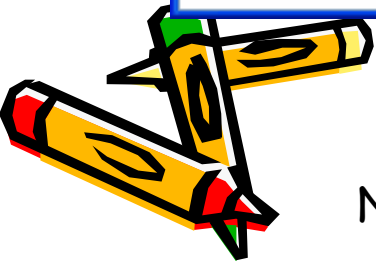
The Sentence tier of Praat (textgrid) file

Textgrid file



```
test.TextGrid - Σημειωματάριο
Αρχείο Επεξεργασία Μορφή Προβολή Βοήθεια

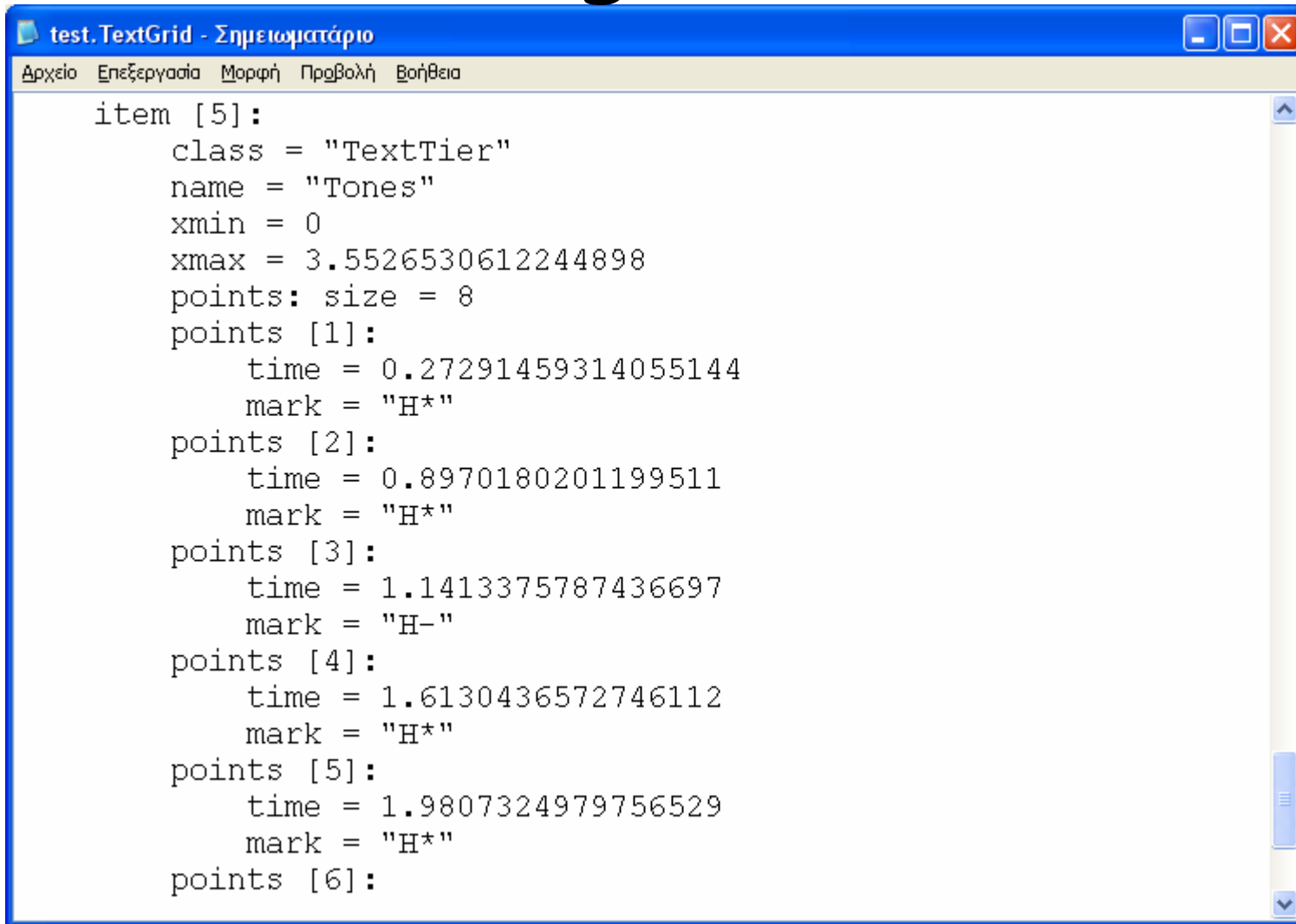
item [2]:
  class = "IntervalTier"
  name = "Int_phrase"
  xmin = 0
  xmax = 3.5526530612244898
  intervals: size = 5
  intervals [1]:
    xmin = 0
    xmax = 0.03585205110961667
    text = ""
  intervals [2]:
    xmin = 0.03585205110961667
    xmax = 1.3348580212179022
    text = "Ο Γιάννης και η Μαρία"
  intervals [3]:
    xmin = 1.3348580212179022
    xmax = 1.489674375197288
    text = ""
  intervals [4]:
    xmin = 1.489674375197288
    xmax = 3.4805159271509547
    text = "πήγον βόλτα στο Θησείο"
```



The Intonation Phrase tier of Praat (textgrid) file

Nikitas N. Karanikolas - Dialectal Corpora Building - June 2017

Textgrid file



```
test.TextGrid - Σημειωματάριο
Αρχείο Επεξεργασία Μορφή Προβολή Βοήθεια

item [5]:
  class = "TextTier"
  name = "Tones"
  xmin = 0
  xmax = 3.5526530612244898
  points: size = 8
  points [1]:
    time = 0.27291459314055144
    mark = "H*"
  points [2]:
    time = 0.8970180201199511
    mark = "H*"
  points [3]:
    time = 1.1413375787436697
    mark = "H-"
  points [4]:
    time = 1.6130436572746112
    mark = "H*"
  points [5]:
    time = 1.9807324979756529
    mark = "H*"
  points [6]:
```

The Tone point tier of Praat (textgrid) file

Written sources (digitized) Pontic (Ποντιακά)



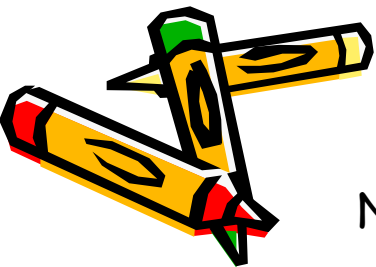
Ἔτον ἕνας πολλά πλούσιος καὶ εἶχεν ἕναν παιδὶν καὶ τὸ παιδὶν ἀτ' ἐπέγ'νεν ἔς σὸ σχολεῖον. Τ' ἄλλ' τὰ παιδία, π' ἐκράτ'ναν μει' ἐκεῖνον, εἶχαν βαγγέλον κ' ἐκεῖνος ἔκ' εἶχεν. Εἶπεν ἕναν ἡμέραν τῇ μάννῃν ἀχτε «μάννα, τὰ παιδία, ὅλα ποὺ κρατοῦν μετ' ἐμέν, ἔχουν βαγγέλα κ' ἐγὼ ἔκ' ἔχω, καὶ ἔκ' ἰλές τὸν κύρη μ' καὶ παίρ' κ' ἐμέν ἕναν βαγγέλον;». Ἡ μάννα ἔτ' πα εἶπεν ἀτο τὸν κύρ'ν ἀτ' κ' ἐπῆρεν κ' ἐδέκεν ἄ κ' ἐδέστεν. Ἄς σὸ ἐδέστεν κ' ὕστερον, ἔγκεν ἄ ἔς ἕναν κουῖμτῶην κ' ἐντῶκεν ἀπάν' ἐκὰν πεντακόσα φηριλία κ' ἐδέκεν ἀτο τὸ γιόν ἀτ'. Κι ἀτὸς πάει κ' ἔρται ἔς σὸ σχολεῖον. Ἐρθεν ἕναν ἡμέραν ἕνας καλόγερος ἔς σοῦ πλούσιονος καὶ ἐρώτεσαν ἀτον «ἀπόθεν ἔρχεσαι καὶ ποῦ πᾶς;» Εἶπεν ἀτ'ς κ' ἐκεῖνος «ἄς σ' Ἄγιον Ὅρος ἔρχουμαι καὶ ἔς σὸν Ἄιν Τάφον πάγω». Εἶπεν ἀτον τὸν πουρνόν ὁ Γιαννίτσης τοῦ



Their homogenized Transcriptions



έτον ένας πολλά πλούσιος και είΣεν έναν παιδίν και το παιδίν ατ' επέγνεν σο
σχολείον. τ' άλλ' τα παιδιά, π' εκράτναν μετ' εκείνον, είχαν βαγγέλΟν κι εκείνος
'κ' είΣεν. είπεν έναν ημέραν την μάναν αχτε «μάνα, τα παιδιά, όλα που κρατούν
μετ' εμέν, έχουν βαγγέλα κι εγώ 'κ' έχω, κά 'κι λες τον κύρη μ' και παίρ' κι εμέν
έναν βαγγέλΟν;». Η μάνα 'τ' πα είπεν ατο τον κύρ'ν ατ κι επήρεν κι εδέκεν α κι
εδέστεν. ασο εδέστεν κι ύστερον έγκεν ας έναν κουιμτΣήν κι εντώκεν απάν' εκάν
πεντακόΣα φηριλία κι εδέκεν ατο τον γιόν ατ. κι ατος πάει κι έρται σο σχολείον.
έρθεν έναν ημέραν ένας καλόγερος σου πλούσιονος και ερώτεσαν ατον «απόθεν
έρΣεσαι και πού πας;». είπεν ατ'ς κι εκείνος «ας άγιον όρος έρχουμαι και σον
άιν τάφον πάγω». είπεν ατον τον πουρνόν ο γιαννίτσης του

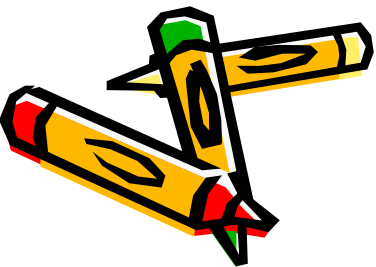


2. Ἀτματσᾶς καὶ ἡ ποθίκα¹.

(Ἔπος)

Ἀτματσᾶς² καὶ ἡ ποθίκα³ ἐποίκανε
καρτασλόυκ⁴. Ἀτματσᾶς ἐποίκε
πουλία ἐπὶ σὴν πέτρα⁵ καὶ ἡ πο-
θίκα ἐποίκε ἐπουκά σὴν πέτρα. Ὑ-
στὲρ ἀτματσᾶς ἐκατέβε καὶ ἔφραε τσῆ

ποθίκα τὰ πουλία. Ἡ ποθίκα οὐκ ἐ-
πόρευε ν' ἀνιβαίν ἐπὶ σὴν πέτρα
νὰ τρώει τ' ἀτματσᾶ τὰ πουλία καὶ ἐ-
πορπάτενε ἐπουκαῖ⁶ καὶ ἐκλαιε. Ἐβηθε
βακοίτ⁷ καὶ ἀτματσᾶς ἐποίησε νὰ βρῖσκ
φαῖ γιὰ τὰ πουλία τ' ἐδέαβε σ' ἕνα
κάμπο μερέα. Ἐκεῖ σὸν κάμπο ἐκά-
θουσανε ἀργάτ καὶ ἔψηνανε κρέας ἀ-
πανκῆ⁸ σ' ἄψιμο. Ἀτματσᾶς παλ
ἐκατέβε ἐπῆρε ἕναν παρτσᾶ⁹ κρέας ἀσ-
σ' ἄψιμο ἀπὸν καὶ ἔφερεν ἀ⁴ στή φω-
λέα τ. Μικερ⁵ ἐκρατενε ἐκεῖ ἄψιμο.
Ἐκολλίε⁶ ἡ φωλέαν ἀτ, ἐρώξανε⁷ ἐ-
πουκά τὰ πουλία τ' καὶ ἔφραεν ἀτα ἡ
ποθίκα.





ατμασας τσε η ποθίκα εποίκανε καρτασλουκ. ατμασας εποίτσε πουλία
επάν σην πέτρα τσε η ποθίκα εποίτσε επουκά σην πέτρα. υστέρ ατμασάς
εκατέβε τσε έφασε τση ποθίκας τα πουλία. η ποθίκα ουτσ επόρενε ν' ανιβαίν
επάν σην πέτρα να τρώει τ' ατμασά τα πουλία τσε επορπάτενε επουκατσέσ
τσε έκλαιε. έρθε βακΙτ τσε ατμασάς εποίε να βρίσκ φαΐ για τα πουλία 'τ
τσ' εδΑβε σ' ένα κάμπο μερέα. ετσεί σον κάμπο εκάθουσανε αργάτ και
έψηνανε κρέας απαντσέσ σ' άψιμο. ατμασάς παλ εκατέβε επήρε έναν
παρτσά κρέας ας άψιμο. εκολλίε η φωλέαν ατ', ερώξανε επουκά τα πουλία 'τ
τσε έφασεν ατα η ποθίκα.



Cappadocian (Καππαδοκικά)



· Τότε πιάσαν βουδαχχήρε να κόψουν το φαβάχ. Κόφτουν το φαβάχ. Δέν πλερούται· πλεμνίσκει λιγόδικο. Το παλτά σακουται. Τότε πιάνουν ένα jadé φαρά· έδωκάν δο ένα πολά σταφίρες νά τα πλύν. Τα καλά επέτανέν da, και τα κōτία βαήνεν da. Το κορίσ λέχ το, “Όί ζάεις; τα καλά πετάνεις τα, και τα κōτία στέγγουν.” “Όί να ποίκω; Δέ χιωρῶ.” Σόνγρα πιάνουν ένα βασκά jadé φαρά, και δίνουν δο, να ζυμῶς ζυμάρ. Ζύμωνέν δο μέ το πράϊ τ. Όί ζάεις;” λέχ το κορίσ. “Μέ το πράχ ζυμοῦται ζυμάρ μί;” λέχ. Τότε το κορίσ κατέβη και ζύμωσέν δο. Σόνγρα νανέβη. Δέν δο βάκε· πιάσεν da ἄς τα μαλιά τ. Τότε ήρτε πατισαχιού το παιρί· πήρεν δο. Και σεράνδα μέρες έπκαν γάμος.



Cappadocian transcription



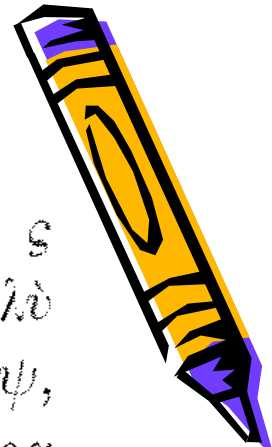
τότε πιάσαν μπουδαχτσήρε να κόψουν το γαβάχ. Κόφτουν το γαβάχ. δεν πλερούται. πλεμνίσκει λιγότησικο. το παλτά σακούται. τότε πιάνουν ένα dZadé γαρά. έδωκάν do ένα πολά σταφίρες να τα πλύν. τα καλά επέτανέν da, και τα κΟτία βαήνεν da. το κορίΣ λέχ το, «τΣί ζάεις; τα καλά πετάνεις τα, και τα κΟτία στέγνουν». «τΣί να ποίκω; δε χιωρώ». σόνγρα πιάνουν ένα βασκά dZadé γαρά, και δίνουν do, να ΖυμώΣ Ζυμάρ. Ζύμωνέν do με το πράι τ'. «τΣί ζάεις;» λέχ το κορίΣ. «με το πράχ Ζυμούται Ζυμάρ μι;» λέχ. τότε το κορίΣ κατέβη και Ζύμωσέν do. σόνγρα νανέβη. δεν το βάκε. πιάσεν da ας τα μαλλιά τ. τότε ήρτε πατιΣαχιού το παιρί. πήρεν do. και σεράνδα μέρες έπκαν γάμος.



Αϊναλιότ (Αϊβαλιώτικα)

Μν'ὰ φουρά η̄δαν ένας βασίλης τσ' εἶχι τς
τοῦ τσιφάλ' ένα τσιρατέλ' τσι τοῦ εἶχι πουλὸ
ἀκουφά. "Οποιοὺν βιρβέρο ἐπιονι νὰ τοῦ γουρέψ,
τοὺν ἔκανι τιβίχ¹⁾ νὰ μὴ τοῦ λέγ ὄξου. Τώρα
οὐλ' οἱ βιρβέροδισ δὲν ἰβουροῦσαν νὰ τοῦ βαστά-
ξιν ἀκουφά· ἵ' ἀφτὸ τς ἔσφαξι.

Πίσου πίσου πῆρι ένα βιρβέρο, τσι σὰ δοῦ
ἀπουκούριψι, τ εἶπι, νὰ μὴ τοῦ πῆ σὶ κανέναν,
ποῦς ἔχ τσέρατου, ἵατὶ θὰ πάρ τοῦ τσιφάλ' τ.
'Ἡ βιρβέρος δὲν ἰβόροσι νὰ βαστάξ, πῆρι, ἔστουψι
μὲς ένα πγάδ τσι φώναξι μ' οὐλ' τ γαρδιά τ:
„Ἡ βασίλης ἔχ τσιρατέλ'." Τώρα τοῦ πγάδ ξι-
ράθτσι, φύτροουσι μέσα μν'ὰ καλαμν'ά. Μιγάλ'νι
ἢ καλαμν'ά. Πέρα μν'ὰ μέρα ένας δζουβάν'ς,
ἔκουψι ἄ γαλαμν'ά τσ' ἔκανι μν'ὰ τσαβούνα τσι
την ἔπιξι. 'Ἡ τσαβούνα ἤλιγι: „Βί! ἰ βασίλης
ἔχ τσιρατέλ'." Τοῦ ἤκσαν, τοῦ εἶπαν τ βασίλέ.



Aivaliot transcription



μια φουρά ήδαν ένας βασιλέσ τσ' είχι στου τσιφάΛ ένα τσιρατέΛ τσι του είχι πουλύ ακρυφά. όποιοιun бирβέρ έπιρνι να του γουρέψ, τουν έκανι τιβίχ να μη του λέΓ' όξου. τώρα ούΛ οι бирβέρδισ δεν ιμπουρούσαν να του βαστάξιν ακρυφά. γί αυτό τς έσφαξι. πίσου πίσου πήρι ένα бирβέρ, τσι σα δου απουκούριψι, τ' είπι να μη του πη σι κανέναν, πους έχ' τσέρατου, γιατί θα πάρ' του τσιφάΛ τ. ι бирβέρς δεν ιμπόρσι να βαστάξ, πήγι, έστουψι μες ένα πγάδ τσι φώναξι μ'ούΛ τ' γαρδιά τ': «ι βασιλέσ έχ' τσιρατέΛ». τώρα του πγάδ ξιράθτσι, φύτρουσι μέσα μια καλαμιά. μιγάΛνι η καλαμνιά. πέρνα μια μέρα ένας τΖουβάΝς, έκουψι d' γαλαμιά τσ' ' έκανι μια τσαβούνα τσι τ'ν έπιζι. Η τσαβούνα ήλιγι: «Βί! ι βασιλέσ έχ' τσιρατέΛ». Του ήκσαν, του είπαν τ' βασιλέ.



Morphological annotations

Morphological categories

Word

Grammatical category

(noun, verb, gerund, particle,
adjective, pronominal, adverb, ...)

Special characteristics

Loan word

Origin

Archaism

Gender alteration

Other

Simple

Structured

Declinable

Production

Composition

Merging

Other



More Morphol. annotations

Morphological process

Declension

Noun

Number (αριθμός)

Gender (γένος)

Case (πτώση)

...

Verb

Person (πρόσωπο)

Number (αριθμός)

Tense (χρόνος)

Mood (έγκλιση)

Voice (φωνή)

...

Production

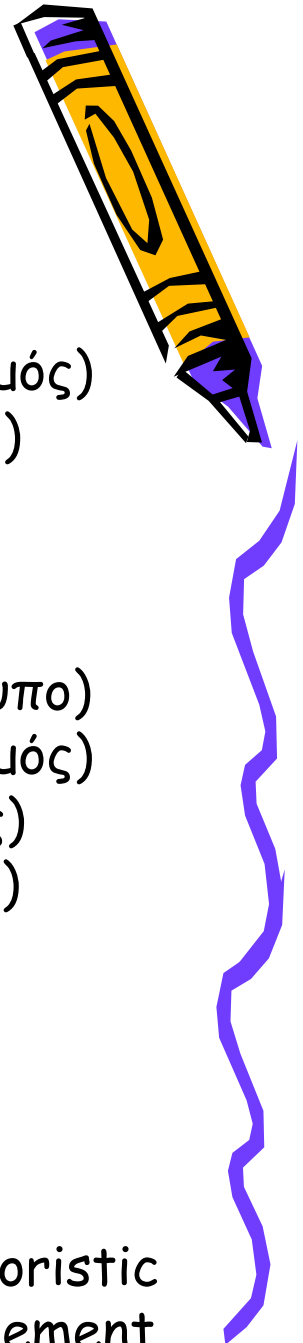
With Postfix

Noun

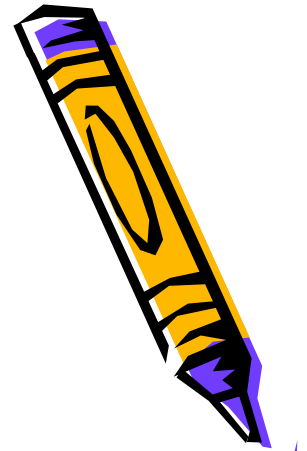
Hypocoristic

Enlargement

Verb



Syntactic and Semantic annotations

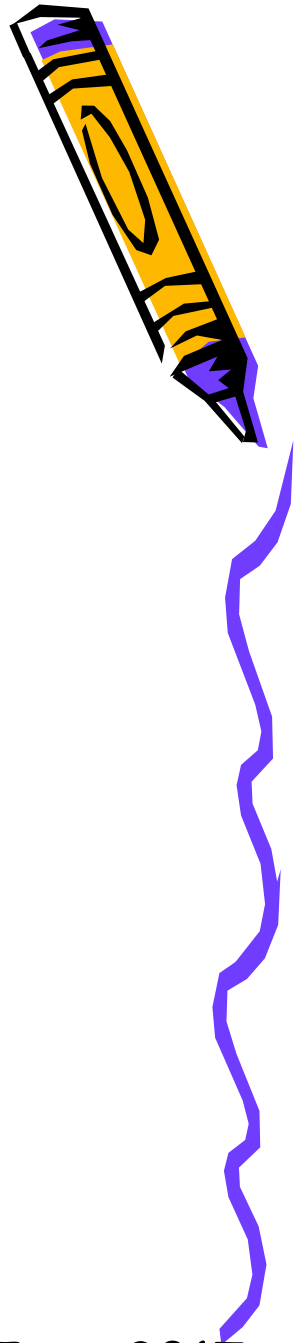


- Intentionally left blank
- It is of less interest in the context of AMiGre
- But, it is implemented



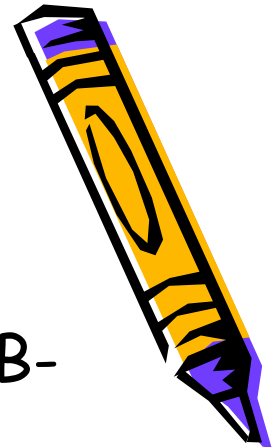
AMiGre

- Introduction
- Sources
- **Applications Overview**
- Design Overview



Why our own system

- Initially, advanced software tools such as LaBB-CAT seemed as the perfect solution to the problem of storing and processing of oral and written dialectal resources.
- But, they proved to be unable to meet all the requirements of the AMiGre project:
 - Annotations at many different linguistic levels
 - Combined search at different levels of representation (phonology, morphology, metadata and, potentially, syntax and semantics)
 - Combined search on both oral and written resources
- So we designed and implemented our own system in order to meet all requirements



Applications Overview

- Oral + Written modules
- Written Sources GUI
- Oral Sources GUI
- Oral & Written Sources Retrieval GUI
- Oral & Written Browsing Web Interface
- Dissemination of effort & Results Web site



Oral + Written - modules



Name	Interpetation	Oral	Written
<u>B. Oral</u>	<u>Browse Oral</u>	X	
<u>B. Written</u>	<u>Browse Written</u>		X
<u>I. Oral</u>	<u>Import Oral</u>	X	
<u>Meta Oral</u>	<u>Metadata for Oral</u>	X	
<u>G. Oral</u>	<u>Oral triptych</u>	X	
<u>G. Written</u>	<u>Written triptych</u>		X
<u>Meta Written</u>	<u>Metadata for Written</u>		X
<u>P.I. Written</u>	<u>Page (or Part) Import Written</u>		X
<u>Morph. Tag</u>	<u>Morphological Tagging</u>	X	X
<u>Syn. Tag</u>	<u>Syntactic Tagging</u>	X	X
<u>Sem. Tag</u>	<u>Semantic Tagging</u>	X	X
<u>Ph. Tagger</u>	<u>Phonological Tagging</u>	X	X
<u>T. Imaging</u>	<u>Text Imaging</u>		X
<u>T. Transcription</u>	<u>Text Transcription</u>		X
<u>M.I. Oral</u>	<u>Massive Import Oral</u>	X	
<u>M.I. Written</u>	<u>Massive Import Written</u>		X



Written Sources GUI - 3fold

display the transcription of document together with attributes of the selected word



Μορφολογική Ανάλυση

Εφαρμογή Εργαλεία

Προβολή σελίδας 1 από 3

Προβολή Εικόνας Προβολή ορίων λέξεων Προσθήκη νέας σελίδας Στην αρχή του κειμένου Προσθήκη Επεξεργασία Διαγραφή Επιστροφή στη λίστα

έσανε

δύς Σεράντ, είχανε απ' ένα παιδί. έστειλαν α στην κΣεντεία. ο ένας είπε το παιδίν ατς: «άδεια μη κάσαι, έναν παρά πάλ αν ευρίσκεις, δουλέψο». η άλλη είπε το παιδίν ατς: «ασά είκοσ παράδΑΣ εξ ούκ μη δουλεύεις». επήγανε εδούλεψανε. ε κείνος οπ εδούλευε σ' έναν παρά αργατικό ουκ εχάσε. ο-γι-άλλο ουκ εύρε δουλεία να δουλεύ σα είκοσ παράδΑΣ. απάν σο χρόνο εκλώστανε να πάνε σ' οσπίτ. σο δρόμο είπαν «ας μετρούμε τα παράδΑΣ μουνα». εμέ τρεσανε. ένας είσε ελίγα, ο-γι-άλλο είσε πολλά. εκείνος οπ είσε ελίγα είπε τον άλλονα «αδά σον κόσμο ποίο κυριεύ, η ψευτία γιόξα η αληθεία;» ο-γι-άλλο είπε «η αληθεία». εκείνος είπε «η ψευτία». εποίκα νε κάβλ. ο είς είπεν «αν κυριεύ η ψευτία, εγώ να δίγω σε τα παράδΑΣ». ο-γι-άλλο πάλ είπε «αν κυριεύ η αληθεία, εγώ πάλ να να δίγω σε τα παράδΑΣ». είπανε «ατάρα σάτι πάμε, ό,τινα τσατεύομε, ερωτούμε, να τρερούμε ποίο κυριεύ».

έσανε
δύς
Σεράντ
είχανε
απ'
ένα
παιδί
έστειλαν
α
σην
κΣεντεία
ο
ένας
είπε
το
παιδίν
ατς
άδεια
μη
κάσαι
έναν
παρά
πάλ

Αποθήκευση Αναίρεση αλλαγών

Βασικές Πληροφορίες

Λήμμα είμαι

Σημασία

Μορφολογική Διαδικασία Κλίση-Ακλισία

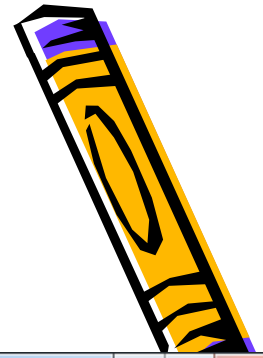
Γραμματική Κατηγορία Ρήμα

Γένος -

Κλιτική Τάξη -

Πρωτότυπη Δάνεια Λέξη

Display transcription with word borders



Μορφολογική Ανάλυση

Εφαρμογή Εργασία

Προβολή σελίδας 1 από 3

Προβολή Εικόνας Προβολή ορίων λέξεων Προσθήκη νέας σελίδας Στην αρχή του κειμένου Προσθήκη Επεξεργασία Διαγραφή Επιστροφή στη λίστα

έσανε δύς Σεράντ, είχανε απ' ένα παιδί. έστειλαν α σην κΣεντεία. ο ένας είπε το παιδίν ατς: «άδεια μη κάσαι, έναν παρά πάλ αν ευρίσκεις, δουλέψο». η άλλη είπε το παιδίν ατς: «ασά είκοσ παράδας εξ ούκ μη δουλεύεις». επήγανε εδούλεψανε. εκείνος οπ εδούλεψε σ' έναν παρά αργατικό ουκ εχάσε. ο-γι-άλλο ουκ εύρε δουλεία να δουλεύ σα είκοσ παράδας. απάν σο χρόνο εκλώσανε να πάνε σ' οσπίτ. σο δρόμο είπανε «ας μετρούμε τα παράδας μouna». εμέ τρεσανε. ένας είσε ελίγα, ο-γι-άλλο είσε πολλά. εκείνος οπ είσε ελίγα είπε τον άλλονα «αδά σον κόσμο ποίο κυριεύ, η ψευτία γιόξα η αληθεία;» ο-γι-άλλο είπε «η αληθεία». εκείνος είπε «η ψευτία». εποίκα νε κάβλ. ο είς είπεν «αν κυριεύ η ψευτία, εγώ να δίγω σε τα παράδας». ο-γι-άλλο πάλ είπε «αν κυριεύ η αληθεία, εγώ πάλ να να δίγω σε τα παράδας». είπανε «ατώρα σάτι πάμε, ό,τινα τσατεύομε, ερωτούμε, ν

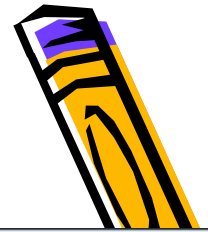
έσανε
δύς
Σεράντ
είχανε
απ'
ένα
παιδί
έστειλαν
α
σην
κΣεντεία
ο
ένας
είπε
το
παιδίν
ατς
άδεια
μη
κάσαι
έναν
παρά
πάλ

Αποθήκευση Αναίρεση αλλαγών

Βασικές Πληροφορίες

Λήμμα	είμαι
Σημασία	
Μορφολογική Διαδικασία	Κλίση-Ακλισία
Γραμματική Κατηγορία	Ρήμα
Γένος	-
Κλιτική Τάξη	-
Πρωτότυπη Δάνεια Λέξη	

Transcription together with the original image scan



Μορφολογική Ανάλυση

Εφαρμογή Εργαλεία

Προβολή λέξεων Προβολή ορίων λέξεων Προσθήκη νέας σελίδας Στην αρχή του κειμένου

Προβολή σελίδας 1 από 3

ΠΑΡΑΜΥΘΙΑ ΟΦΕΩΣ

1

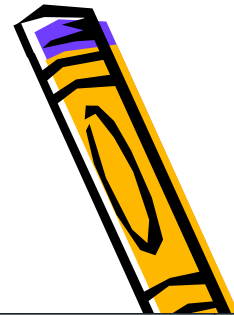
Ἔσαν δύο ἄνθρωποι, εἶχανε ἀπὸ ἑνὸς παιδίου. Ἐστειλαν ἕνα σὺν κτενεῖα. Ὁ ἕνας εἶπε τὸ παιδίν αὐτῷ εὐχόμενος, ἕνα παρὰ πάλιν εὐρίσκεις, δούλεψο». Ἡ ἄλλη εἶπε τὸ παιδίν αὐτῷ: «ἀσὰ εἴκοσ παράδας ἐξούκ μη δουλεύεις». ἐπήγαγε ἐδούλεψαν. ἐκεῖνος ὅπου ἐδούλευε σὲ ἕναν παρὰ ἀργατικὸν οὐκ εἶχε. ο-γι-ἄλλο οὐκ εἶχε δουλείαν νὰ δουλεύῃ σὰ εἴκοσ παράδας. ἀπὸν σὸ χρόνον ἐκλώσταν νὰ πάνε σὲ ὄσπιτον. σὸ δρόμον εἶπαν «ἀς μετροῦμε τὰ παράδας μὴ οὐκ». ἐμέτρησαν. ἕνας εἶπε εὐχόμενος, ο-γι-ἄλλο εἶπε πολλὰ. ἐκεῖνος ὅπου εἶπε εὐχόμενος εἶπε τὸν ἄλλον «ἀδὰ σὸν κόσμον ποῖον κυριεύῃ, ἡ ψευτιά γιόξα ἡ ἀληθεία;» ο-γι-ἄλλο εἶπε «ἡ ἀληθεία». ἐκεῖνος εἶπε «ἡ ψευτιά». ἐποίησαν κάβλον. ὁ εἷς εἶπεν «ἀν κυριεύῃ ἡ ψευτιά, ἐγὼ νὰ δίδω σοὶ τὰ παράδας». ο-γι-ἄλλο πάλιν εἶπε «ἀν κυριεύῃ ἡ ἀληθεία, ἐγὼ πάλιν νὰ δίδω σοὶ τὰ παράδας». εἶπεν «ἀτώρα σὰτι πάμε, ὅτινα τσατεύομε, ἐρωτοῦμε, νὰ τερῶμε ποῖον κυριεύῃ».

ἐπήγαγε ἐποίησαν ἕναν ποπά νέον. ἐκεῖνος εἶπε «ἡ ψευτιά κυριεύῃ». εἶπεν ἐκεῖν «νὰ ρωτοῦμε δύο νοματοῦς καὶ ἄλλο». ἐρώτησαν ἕνα μεσοκαιρίτη ποπά. ἐκεῖνος πάλιν εἶπεν «ἡ ψευτιά». τὸ ὑστερὸν ἐρώτησαν ἕνα γέρον ποπά. ἐκεῖνος πάλιν εἶπε «ἡ ψευτιά κυριεύῃ». ἔτι τὸ παιδί ἐδῶκε τὰ παράδας αὐτῷ τὸν ἄλλον τὸν τε(μ)πέλ καὶ αὐτὸς ἐκλώσταν ὀπίσ, τσοῦγκ οὐκ εἶπε παράδας ν' ἐπέγινε σὲ ὄσπιτον. ο-γι-ἄλλο ἐπῆγε τὰ παράδας αὐτῷ καὶ ἐπῆγε αὐτῷ τὴν μάννα ἡ».

Ἔσαν δύο ἄνθρωποι, εἶχανε ἀπὸ ἑνὸς παιδίου. Ἐστειλαν ἕνα σὺν κτενεῖα. Ὁ ἕνας εἶπε τὸ παιδίν αὐτῷ εὐχόμενος, ἕνα παρὰ πάλιν εὐρίσκεις, δούλεψο». Ἡ ἄλλη εἶπε τὸ παιδίν αὐτῷ: «ἀσὰ εἴκοσ παράδας ἐξούκ μη δουλεύεις». ἐπήγαγε ἐδούλεψαν. ἐκεῖνος ὅπου ἐδούλευε σὲ ἕναν παρὰ ἀργατικὸν οὐκ εἶχε. ο-γι-ἄλλο οὐκ εἶχε δουλείαν νὰ δουλεύῃ σὰ εἴκοσ παράδας. ἀπὸν σὸ χρόνον ἐκλώσταν νὰ πάνε σὲ ὄσπιτον. σὸ δρόμον εἶπαν «ἀς μετροῦμε τὰ παράδας μὴ οὐκ». ἐμέτρησαν. ἕνας εἶπε εὐχόμενος, ο-γι-ἄλλο εἶπε πολλὰ. ἐκεῖνος ὅπου εἶπε εὐχόμενος εἶπε τὸν ἄλλον «ἀδὰ σὸν κόσμον ποῖον κυριεύῃ, ἡ ψευτιά γιόξα ἡ ἀληθεία;» ο-γι-ἄλλο εἶπε «ἡ ἀληθεία». ἐκεῖνος εἶπε «ἡ ψευτιά». ἐποίησαν κάβλον. ὁ εἷς εἶπεν «ἀν κυριεύῃ ἡ ψευτιά, ἐγὼ νὰ δίδω σοὶ τὰ παράδας». ο-γι-ἄλλο πάλιν εἶπε «ἀν κυριεύῃ ἡ ἀληθεία, ἐγὼ πάλιν νὰ δίδω σοὶ τὰ παράδας». εἶπεν «ἀτώρα σὰτι πάμε, ὅτινα τσατεύομε, ἐρωτοῦμε, νὰ τερῶμε ποῖον κυριεύῃ».

ἐπήγαγε ἐποίησαν ἕναν ποπά νέον. ἐκεῖνος εἶπε «ἡ ψευτιά κυριεύῃ». εἶπεν ἐκεῖν «νὰ ρωτοῦμε δύο νοματοῦς καὶ ἄλλο». ἐρώτησαν ἕνα μεσοκαιρίτη ποπά. ἐκεῖνος πάλιν εἶπεν «ἡ ψευτιά». τὸ ὑστερὸν ἐρώτησαν ἕνα γέρον ποπά. ἐκεῖνος πάλιν εἶπε «ἡ ψευτιά κυριεύῃ». ἔτι τὸ παιδί ἐδῶκε τὰ παράδας αὐτῷ τὸν ἄλλον τὸν τε(μ)πέλ καὶ αὐτὸς ἐκλώσταν ὀπίσ, τσοῦγκ οὐκ εἶπε παράδας ν' ἐπέγινε σὲ ὄσπιτον. ο-γι-ἄλλο ἐπῆγε τὰ παράδας αὐτῷ καὶ ἐπῆγε αὐτῷ τὴν μάννα ἡ».

Insertion (selection) of an image for a written source



Προσθήκη Σελίδας

Βήμα 1: Επιλογή Εικόνας & Εισαγωγή Κειμένου

Εικόνα Σελίδας:

* Κείμενο Σελίδας:

τ' όρωμαν για την πίταν.

τρει νομάτ συντρόφ εβραδασταν σ' έναν μέρος. έψεσαν έναν πίταν, έφααν τ' ημψόν και για τ' άλλο πα είπαν «θα τρώει ατο ήντισα ν ελέπ πολλά αdζαϊπκον όρωμαν. το πουρνόν εκάτσαν και λέγνε τ' ορώματα τουν.

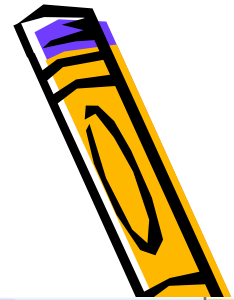
-εγώ, είπεν ο είνας, είδα εξέβα ψηλά ψηλά απάν σ' ουρανού την τΑπΑν.

-ατό τιδέν 'κι έν, είπεν άλλος, εγώ είδα εκατήβα αφκά σον κόλον τη ής.

-εγώ, είπεν ο τρίτον, όνταν είδα σας να πάτε ατόσον μακρά, ελογάριασα 'κι θα έρχουζνε και έφαα την πίταν.



Symbols' selection in order to separate words of a transcription



Προσθήκη Σελίδας

Βήμα 2: Εξαγωγή Λέξεων

* Κείμενο Σελίδας:

τ' όρωμαν για την πίταν.

τρεί νομάτ συντρόφ εβραδᾶσαν σ' έναν μέρος. έψεσαν έναν πίταν, έφααν τ' ημψόν και για τ' άλλο πα είπαν «θα τρώει ατο ήντισαν ελέπ πολλά αδΖαίπκον όρωμαν. το πουρνόν εκάτισαν και λέγνε τ' ορώματα τουν.

-εγώ, είπεν ο είνας, είδα εξέβα ψηλά ψηλά απάν σ' ουρανού την τᾶπᾶν.

-ατό τιδέν 'κι έν, είπεν άλλος, εγώ είδα εκατήβα αφκά σον κόλον τη ής.

-εγώ, είπεν ο τρίτον, όνταν είδα σας να πάτε ατόσον μακρά, ελογάριασα 'κι θα έρχουζνε και έφαα την πίταν.

Σύμβολα Διαχωρισμού:

- (Regular Expression: '\s')
- , (Regular Expression: ',')
- . (Regular Expression: '\.')
- (Regular Expression: '\-')
- ? (Regular Expression: '\?')
- ; (Regular Expression: ';')
- ' (Regular Expression: '\')
- " (Regular Expression: '\"')
- / (Regular Expression: '/')

Regular Expression:

(?!\\V[\\s,\\-;!«»:]

✘ Ακυρο

⏪ Προηγούμενο Βήμα: Εισαγωγή Καμένου & Εικόνας

✔ Επόμενο: Προεπισκόπηση

Improved 3fold presentation



Μορφολογική Ανάλυση

Εφαρμογή Εργαλεία

Προβολή Εικόνας Προβολή ορίων λέξεων Προσθήκη νέας σελίδας Στην αρχή του κειμένου Προσθήκη Επεξεργασία Διαγραφή Επιστροφή στη λίστα

Προβολή σελίδας 1 από 2

είς εφτωχός σίτε **έρτον** ασήν χαμελέτεν με το παιδί ν ατ, ενεγκάστεν και εκάτσεν κά ν' αναπάγετον κι ενεστέναξεν και είπεν «ωφ!» ευτύς εξέβεν ας έναν σπέλον κέσ εις όφισ και είπεν «γιατί κουίεις με;» κι ο φτωχόν είπεν «εγώ εσέν 'κι κουίζω». κι επεκεί είπεν ο όφισ «γιατί 'κι δίς μ' από το παιδί σ', ας μαθίζ στο τέχνης;» και με το λόγον εκείνον εδέκεν α κεί και ο όφισ επήρεν ατον κι επήγεν αφκά σην νηγήν απέσ σ' έναν σπέλον. εκΑπέσ ο όφισ είξεν έναν κ ουτσήν κι ατέ εγάπεσεν τον παιδάν. ύστερ από κα μπόσον καιρόν η κουτσή είπεν στον παιδάν «ο κύρ η μ' αν λέει σε, έμαθες τέχνην; εσύ πέ, 'κι έμαθα, ήνταν λέει σε, έμαθες; εσύ, 'κι έμαθα, πέ». ύστερ ον ερώτεσεν ο όφισ «έμαθες ακομάν τέχνης;» ο παιδίας είπεν «'κι έμαθα». ατότεσ ο όφισ εχτύπεσεν στον έναν σιλέν κι εχάταμεν ατον. άμαν ο παιδίας έμαθεν έτον τέχνης. έρθεν στον κύρν ατ και είπεν « τ'ΑΤΑ, εγώ ας ίνομαι μουλάρ κι εσύ πούλτσον με, άμαν τηνάν πουλείς με το δουκάλι μ' μη δίς στον ». έντον μουλάρ και έρθεν εκείνος ο όφισ ν' αγορ άζ Ατον. ο κύρτσ τη παιδά το μουλάρ εδέκεν και το

είς
εφτωχός
σίτε
έρτον
ασήν
χαμελέτεν
με
το
παιδί
ν
ατ
ενεγκάστεν
και
εκάτσεν
κά
ν'
αναπάγετον
κι
ενεστέναξεν
και
είπεν
ωφ
ευτύς
εξέβεν
ας
έναν
σπέλον
κέσ
είς
όφισ
και
είπεν
γιατί

Μορφολογική Επισημείωση

ΒΑΣΙΚΕΣ ΠΛΗΡΟΦΟΡΙΕΣ

Λήμμα: **έρχομαι**
Μορφολογική Διαδικασία: **Κλίση-Ακλισία**
Γραμματική Κατηγορία: **Ρήμα**
Καταγωγή Λήμματος: **Ελληνική**

ΚΛΙΣΗ

Χρόνος: **Ενεστώτας**
Αριθμός: **Πληθυντικός**

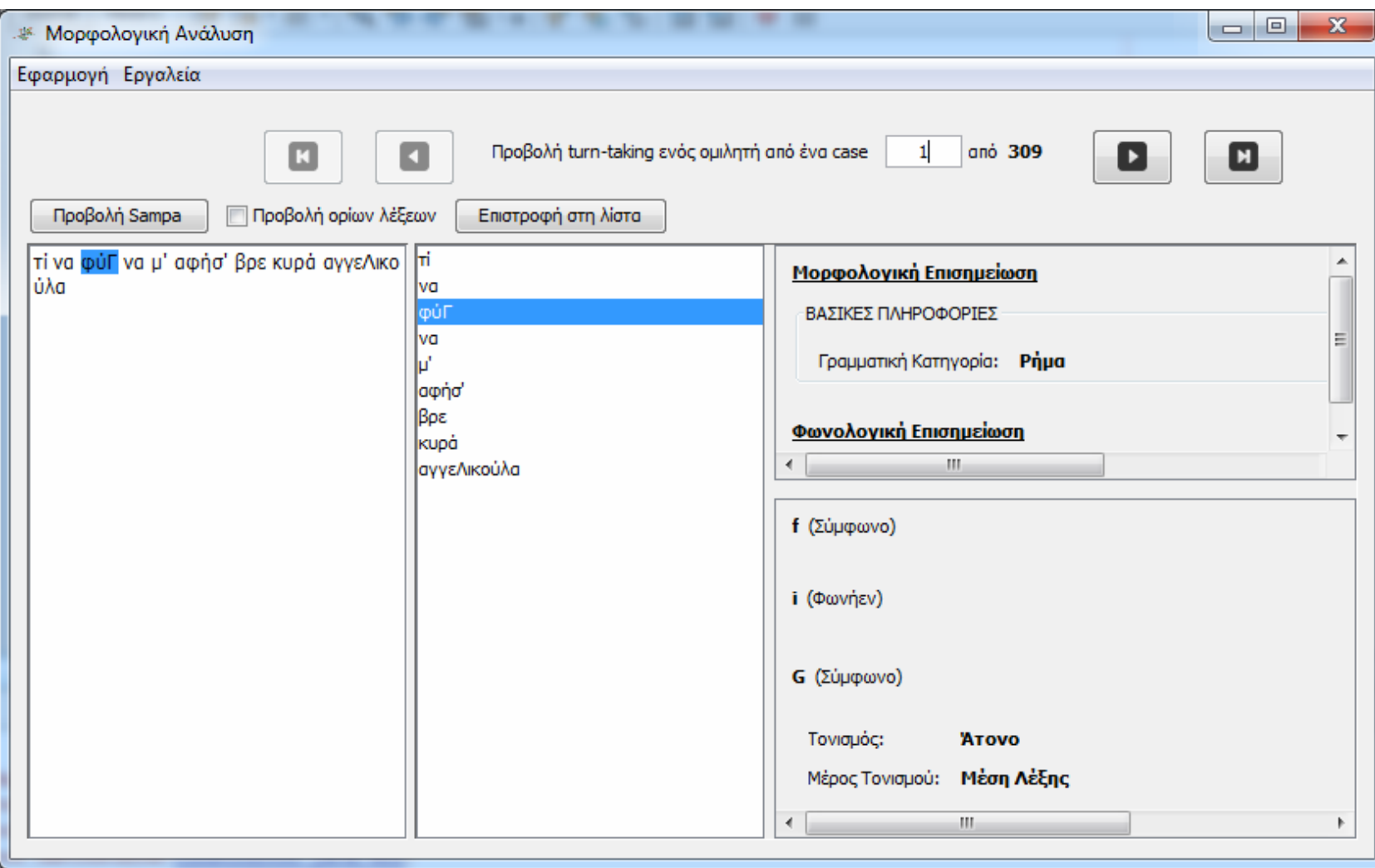
ΜΕΤΑ-ΠΛΗΡΟΦΟΡΙΕΣ

Περιοχή: **ΤΡΑΠΕΖΟΥΝΤΙΑΚΑ (Τραπεζούντα, Κερασούντα, Ριζούντα, Σούρμενα, Όφισ, Λιβερά, Τρίπολις, Ματσούκα)**

Φωνολογική Επισημείωση

Φαινόμενα: **Ανομοίωση Συμφώνου**

Oral Sources GUI - 3fold



The screenshot shows a software window titled "Μορφολογική Ανάλυση" (Morphological Analysis). The interface includes a menu bar with "Εφαρμογή" and "Εργαλεία", a toolbar with navigation buttons, and a main display area. The main display area is divided into three sections: a text input field on the left, a word list in the middle, and a detailed morphological analysis on the right. The word "φύΓ" is selected in the list, and its analysis is shown in the right panel.

Μορφολογική Ανάλυση

Εφαρμογή Εργαλεία

Προβολή turn-taking ενός ομιλητή από ένα case 1| από 309

Προβολή Sampra Προβολή ορίων λέξεων Επιστροφή στη λίστα

τί να φύΓ να μ' αφήσ' βρε κυρά αγγελικο ύλα

τί
να
φύΓ
να
μ'
αφήσ'
βρε
κυρά
αγγελικούλα

Μορφολογική Επισημείωση

ΒΑΣΙΚΕΣ ΠΛΗΡΟΦΟΡΙΕΣ

Γραμματική Κατηγορία: **Ρήμα**

Φωνολογική Επισημείωση

f (Σύμφωνο)

i (Φωνήεν)

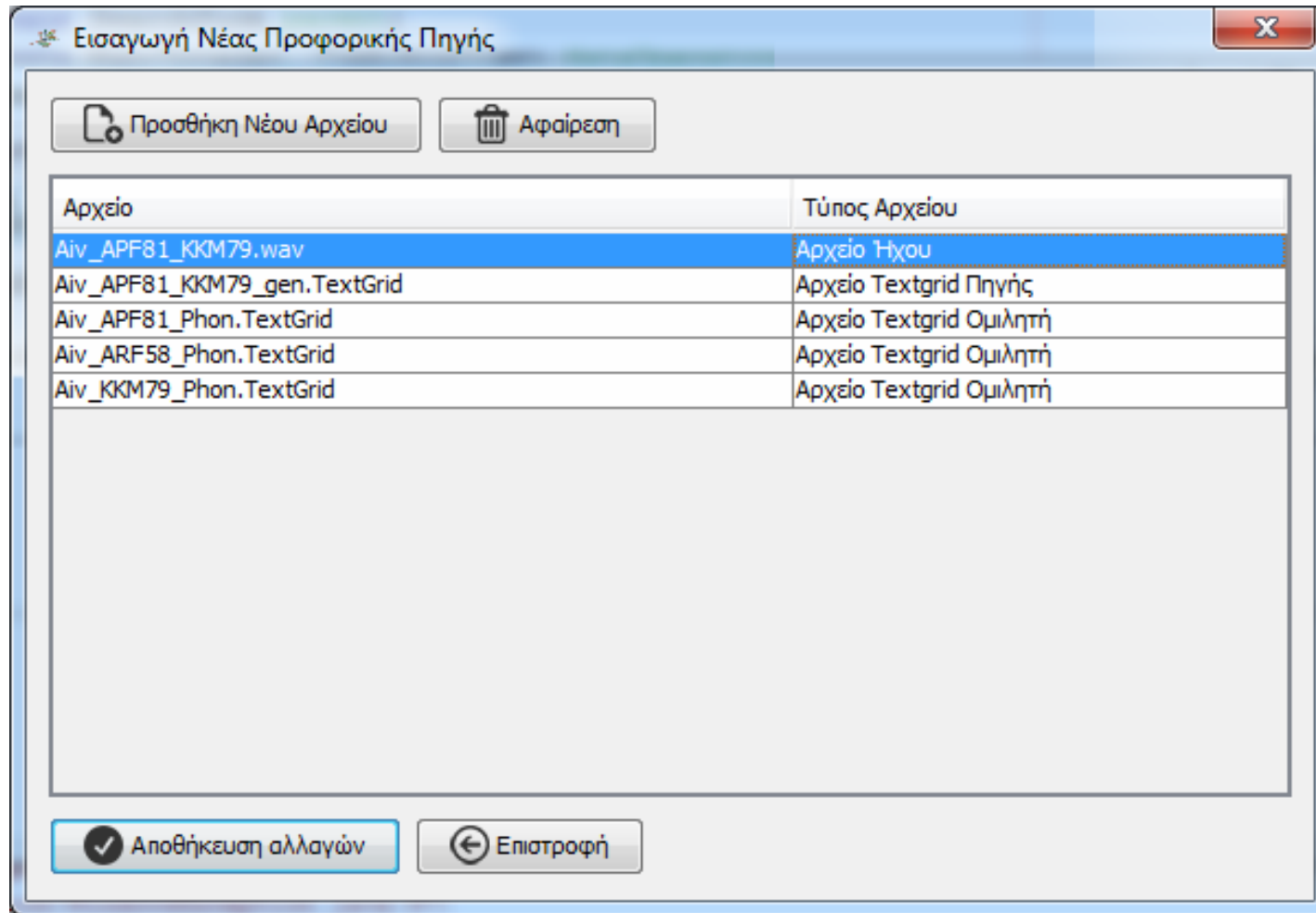
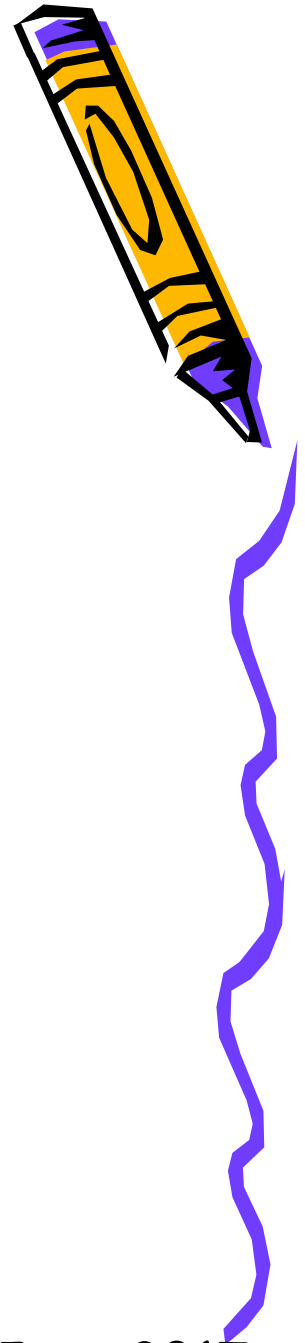
G (Σύμφωνο)

Τονισμός: **Άτονο**

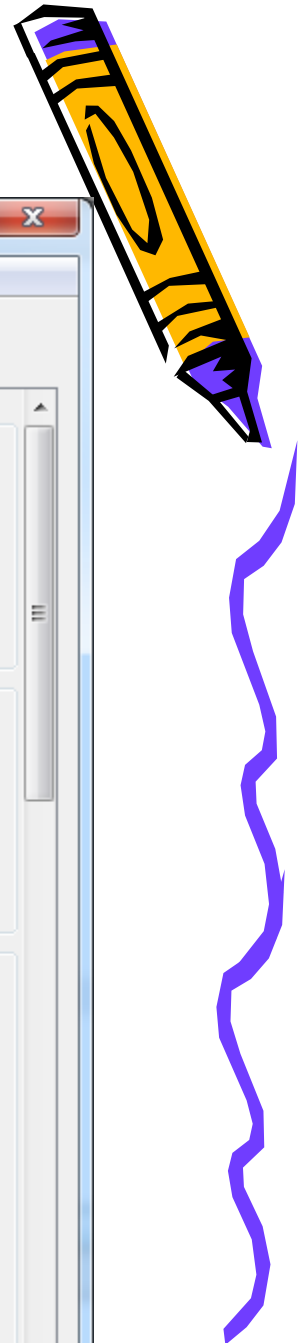
Μέρος Τονισμού: **Μέση Λέξης**



Oral Sources - Import



Oral Sources - Metadata



Μεταδεδομένα: Επεξεργασία Ιδιοτήτων

Εφαρμογή Εργαλεία

Αποθήκευση αλλαγών Επιστροφή

ΣΤΟΙΧΕΙΑ ΑΡΧΕΙΟΥ

Αύξων Αριθμός Αρχείου: 1

Όνομα Αρχείου: Αίν_ΚΚΜ80

Θέση του αρχείου:

Ελεύθερο/Κλαδωμένο: Ελεύθερο

ΔΙΑΛΕΚΤΟΣ

Όνομα Διαλέκτου: Αιβαλιώτικα

Γεωγραφικός Προσδιορισμός Διαλέκτου: Αιβαλί Μικράς Ασίας

Τόπος Γέννησης:

Τόπος ανατροφής (πού μεγάλωσαν):

ΕΡΕΥΝΗΤΙΚΟ ΠΡΟΓΡΑΜΜΑ

Όνομα: AMIGRe

Πηγή Χρηματοδότησης: -

Διάρκεια ερευνητικού προγράμματος:

Επιστημονικός Υπεύθυνος: Αγγελική Ράλλη

Υπεύθυνος Έρευνας Πεδίου: Αγγελική Ράλλη

Ερευνητής Πεδίου:

- Αγγελική Ράλλη
- Δημήτρης Παπαζαχαρίου
- Δήμητρα Μελισσαροπούλου



AMiGre

- Introduction
- Sources
- Applications Overview
- **Design Overview**

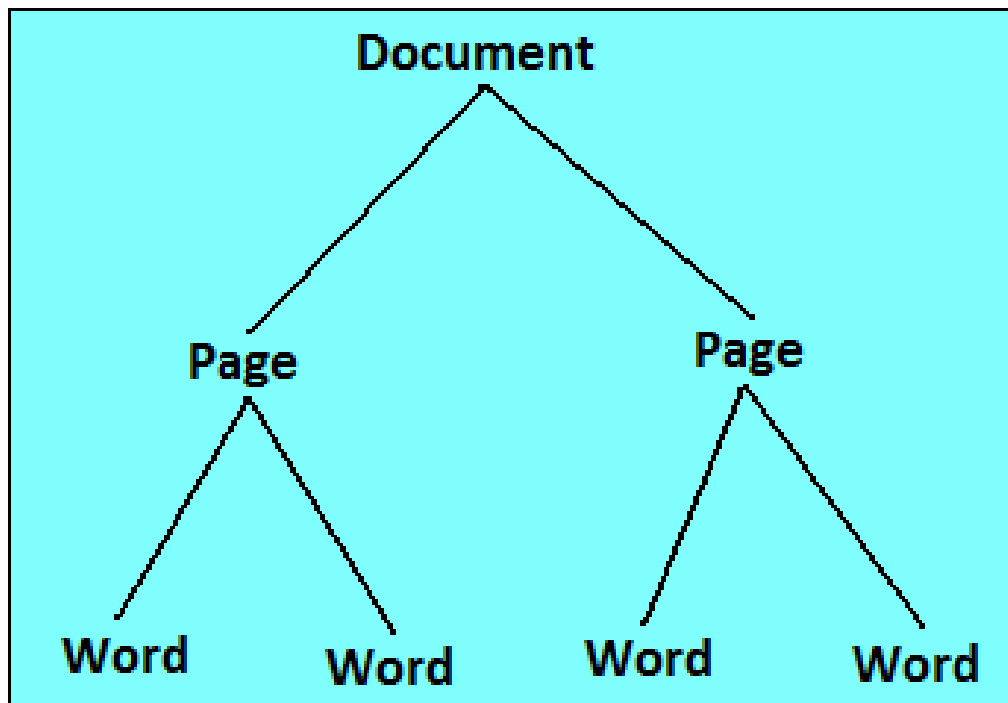


Design Overview

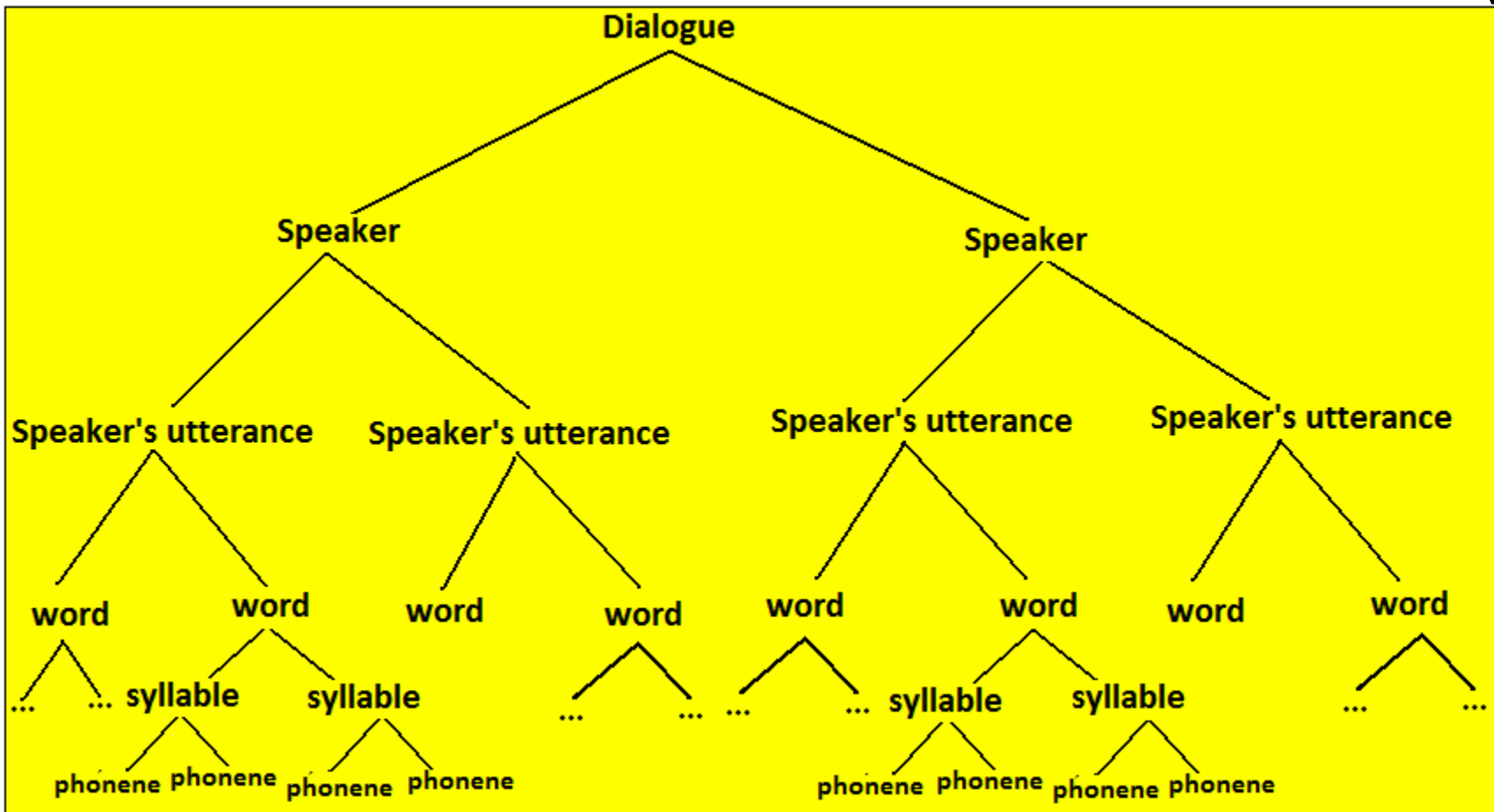
- Alignment of Oral and Written data
- Oral & Written - System overview
- Struct (relational) databases
- EAV data structures



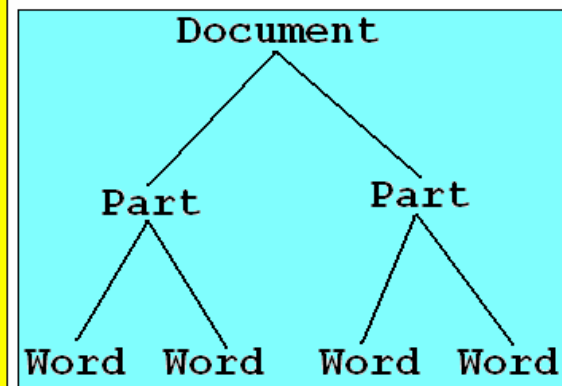
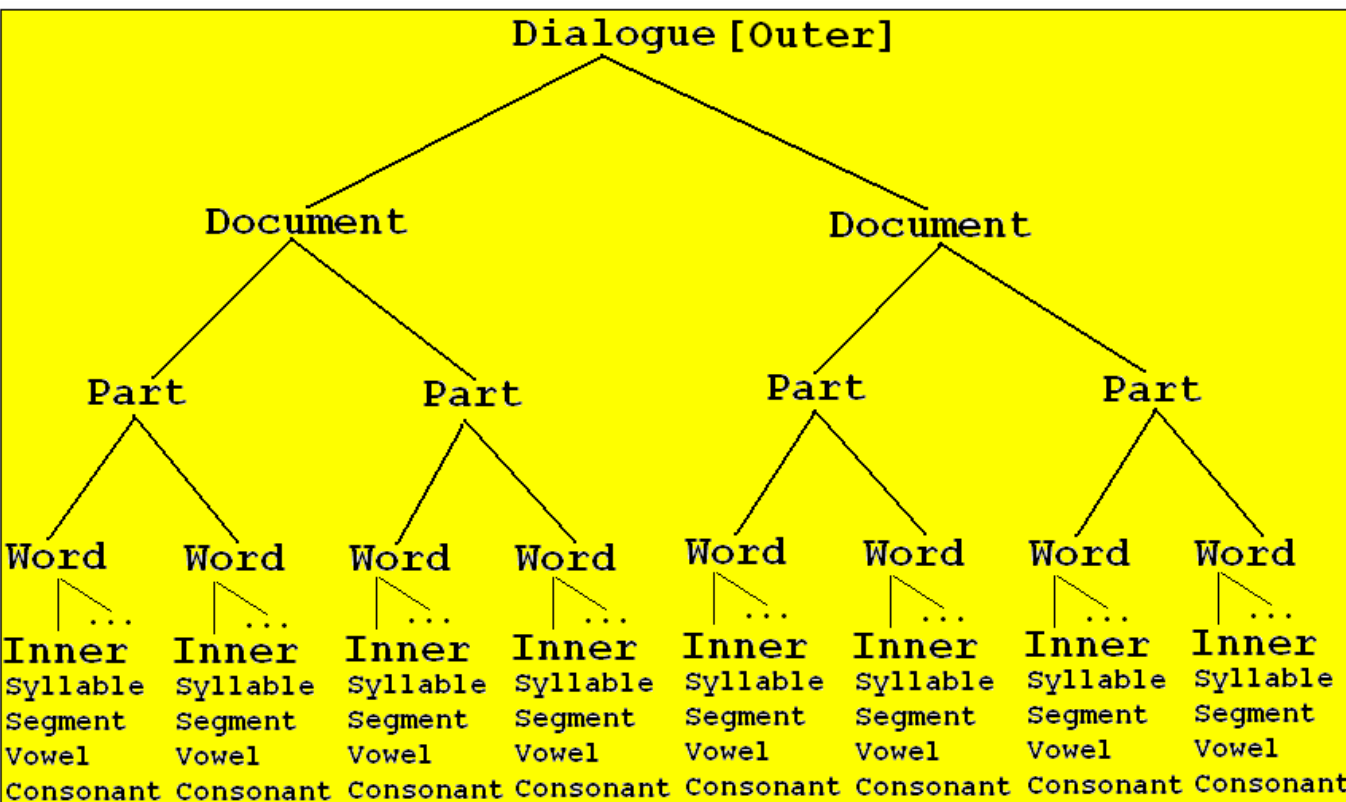
Alignment of Oral and Written - structure of Written data

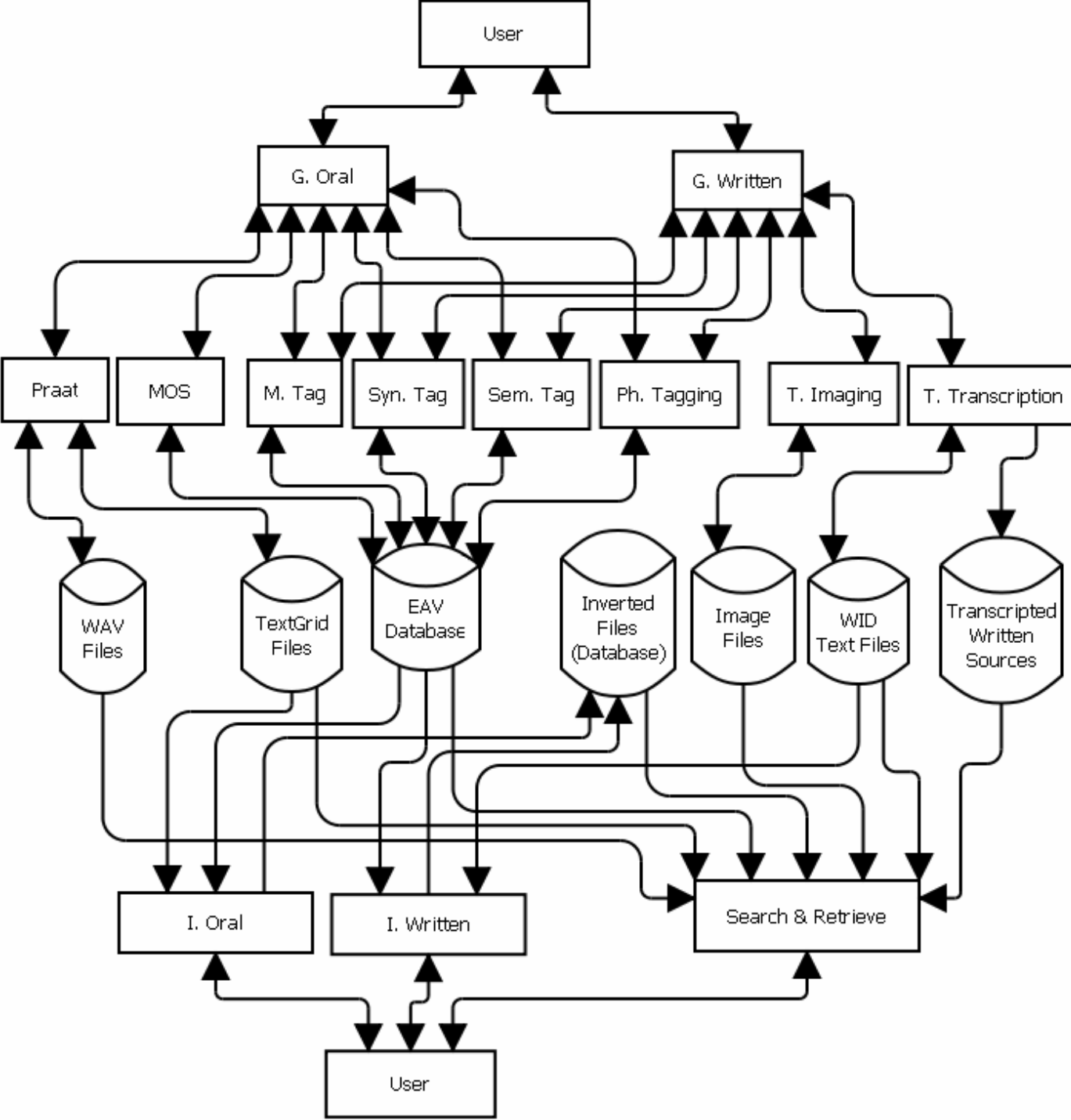


Alignment of Oral and Written - structure of Oral data

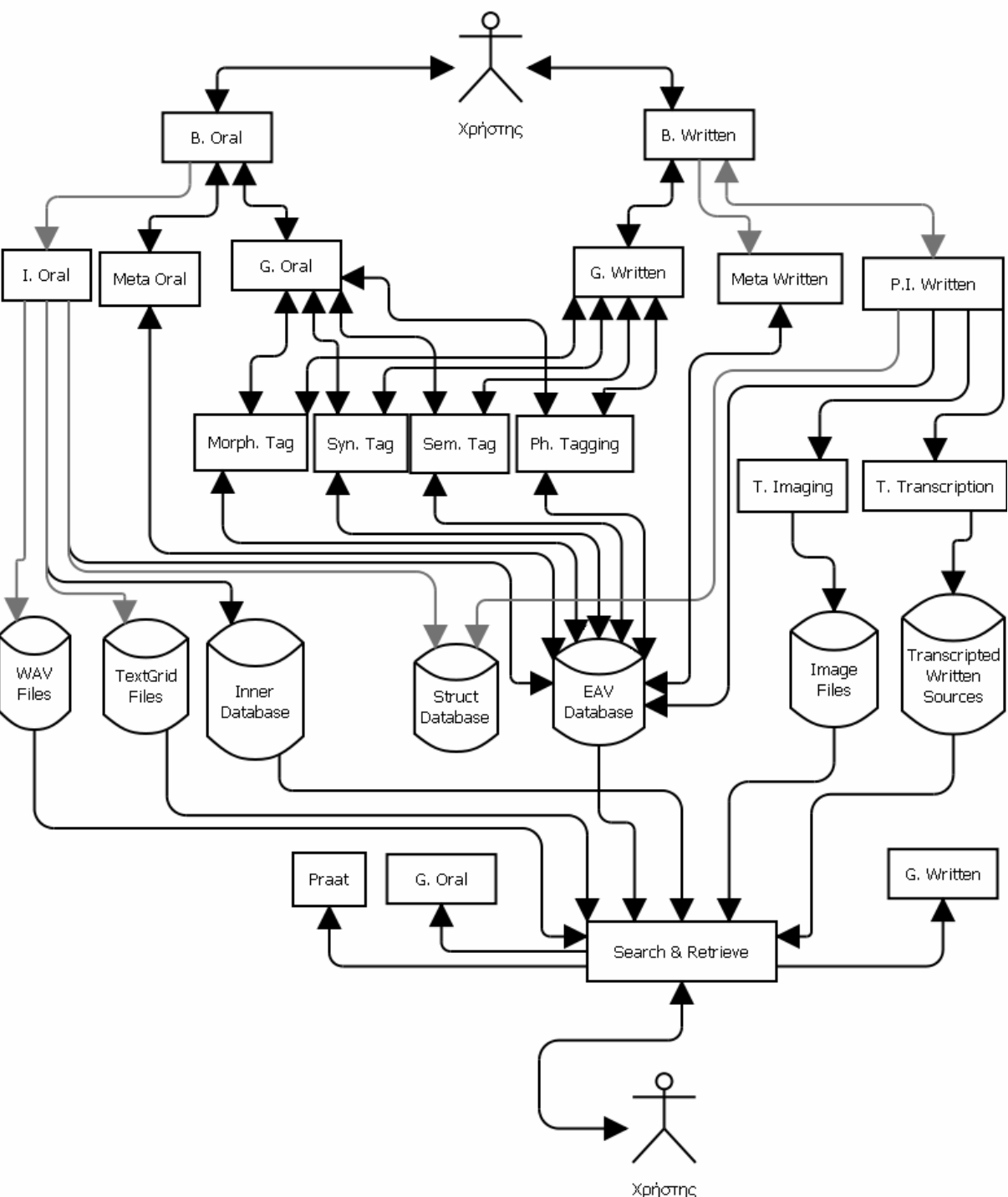


Alignment of Oral and Written - common structure





Oral &
Written
Sources -
System
overview
- OLD



Oral &
Written
Sources -
System
overview -
updated

Struct database



- Within the *Struct* database, the components of the documents are organized on consecutive levels of refinement which will be annotated with the help of the *EAV* database.
- The implementation of the abstract structure of the *Struct* database uses two quasi similar relational schemas.
- The only difference is that the implementation for the oral documents *Struct oral* is composed of all 5 levels, while the implementation for the written documents *Struct written* contains only the 3 intermediate levels.

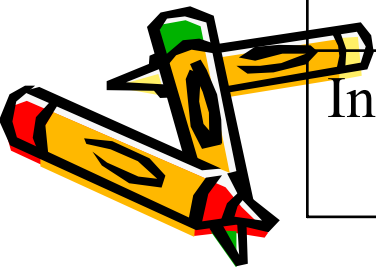


Alignments of Struct (oral and written) databases



- Next table define the data alignments between the two Struct databases (*Struct oral* and *Struct written*):

Abstract name	Oral	Written
Dialog	Overall oral document	--
Document	Speaker / interlocutor	Overall written document
Part	Speaker's utterance	Page of written document
Word	Morphological word	Morphological word
Inner	Syllables, vowels, consonants etc,	--



Struct db for oral documents



oral_sources [Dialogue]	
*OralSourceId	INT(11)
°Title	VARCHAR(500)
°TextGridFilePath	VARCHAR(200)
°TextGridUploadedOn	DATETIME
°MetadataFilePath	VARCHAR(200)
°MetadataUploadedOn	DATETIME
°Notes	TEXT
°IsDeleted	BIT(1)
°CreatedOn	DATETIME

oral_speakers [Document]	
*SpeakerId	INT(11)
*Code	VARCHAR(100)
*Name	VARCHAR(200)
°TextGridFilePath	VARCHAR(200)
°TextGridUploadedOn	DATETIME
°OrderIndex	INT(11)
°IsDeleted	BIT(1)
*OralSourceId	INT(11)

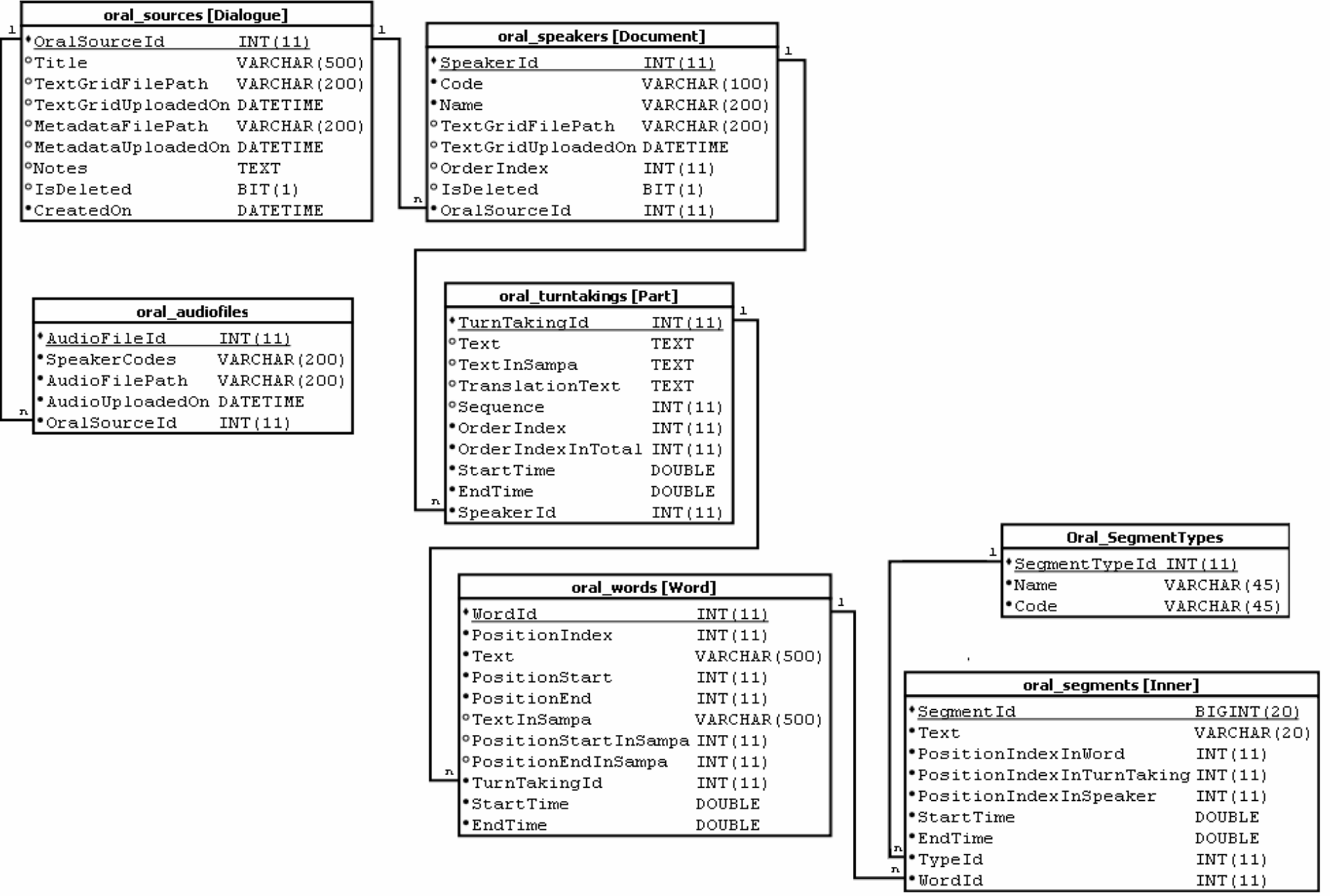
oral_audiofiles	
*AudioFileId	INT(11)
*SpeakerCodes	VARCHAR(200)
*AudioFilePath	VARCHAR(200)
*AudioUploadedOn	DATETIME
*OralSourceId	INT(11)

oral_turntakings [Part]	
*TurnTakingId	INT(11)
°Text	TEXT
°TextInSampa	TEXT
°TranslationText	TEXT
°Sequence	INT(11)
*OrderIndex	INT(11)
*OrderIndexInTotal	INT(11)
*StartTime	DOUBLE
*EndTime	DOUBLE
*SpeakerId	INT(11)

oral_words [Word]	
*WordId	INT(11)
*PositionIndex	INT(11)
*Text	VARCHAR(500)
*PositionStart	INT(11)
*PositionEnd	INT(11)
°TextInSampa	VARCHAR(500)
°PositionStartInSampa	INT(11)
°PositionEndInSampa	INT(11)
*TurnTakingId	INT(11)
*StartTime	DOUBLE
*EndTime	DOUBLE

Oral_SegmentTypes	
*SegmentTypeId	INT(11)
*Name	VARCHAR(45)
*Code	VARCHAR(45)

oral_segments [Inner]	
*SegmentId	BIGINT(20)
*Text	VARCHAR(20)
*PositionIndexInWord	INT(11)
*PositionIndexInTurnTaking	INT(11)
*PositionIndexInSpeaker	INT(11)
*StartTime	DOUBLE
*EndTime	DOUBLE
*TypeId	INT(11)
*WordId	INT(11)



Struct db for written documents



wrt_documents [Document]	
*DocumentId	INT(11)
◦CoverImageFilePath	VARCHAR(200)
◦CoverImageUploadedOn	DATETIME
◦Notes	TEXT
◦IsDeleted	BIT(1)
◦CreatedOn	DATETIME

wrt_wordextractconfigurations	
*ConfigurationId	INT(11)
◦Configuration	VARCHAR(2000)
◦PageId	INT(11)

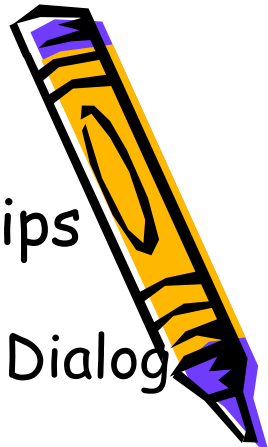
wrt_pages [Part]	
*PageId	INT(11)
•Text	TEXT
◦Notes	TEXT
◦MainImageFilePath	VARCHAR(200)
◦MainImageUploadedOn	DATETIME
◦TranslationImageFilePath	VARCHAR(200)
◦TranslationImageUploadedOn	DATETIME
◦IsDeleted	BIT(1)
•OrderIndex	INT(11)
•DocumentId	INT(11)

wrt_words [Word]	
*WordId	INT(11)
•Text	VARCHAR(200)
•PositionIndex	INT(11)
•PositionStart	INT(11)
•PositionEnd	INT(11)
•PageId	INT(11)



Struct dbs some explanations

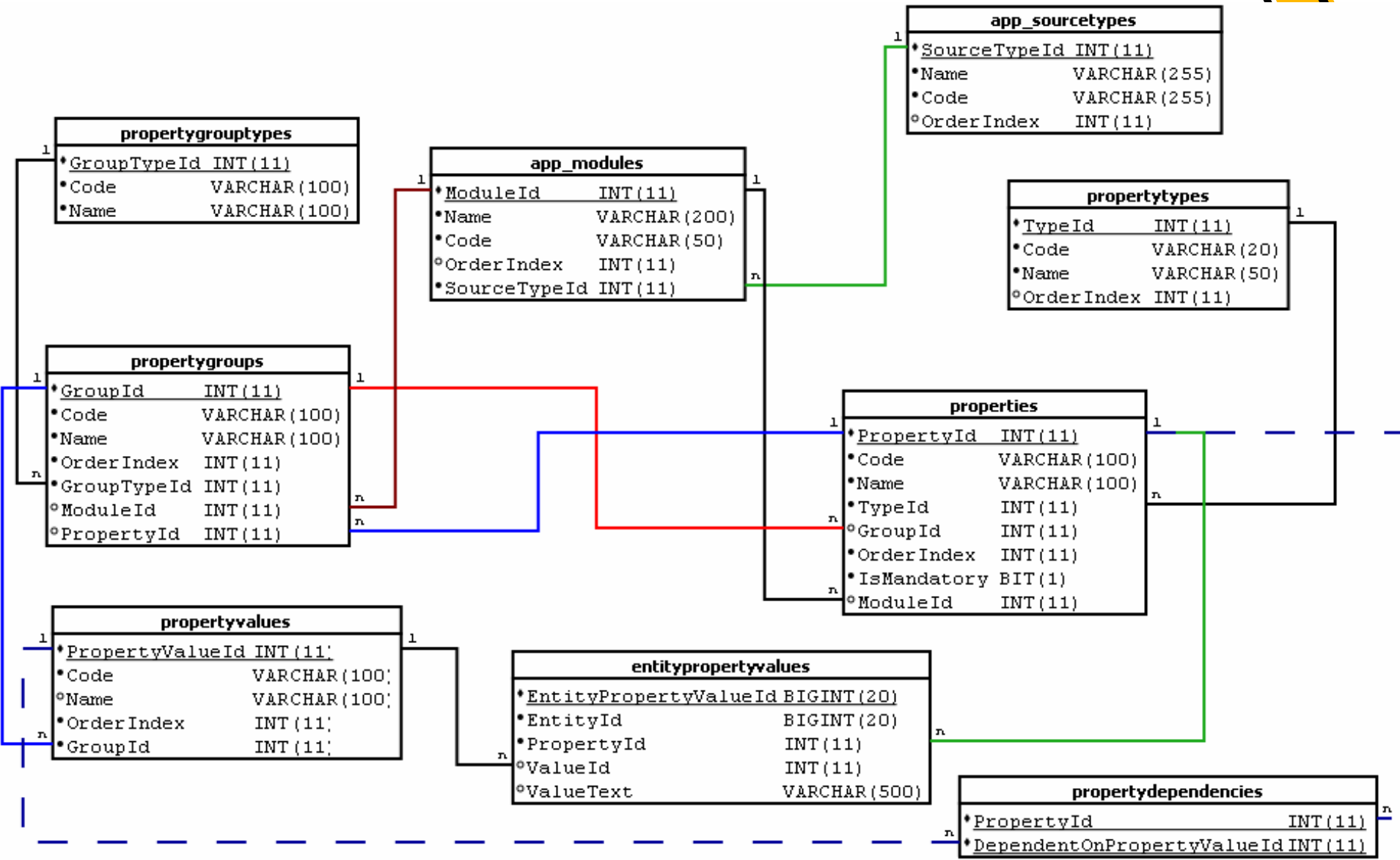
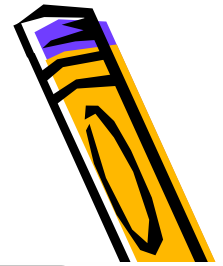
- In both db schemas, one-to-many relationships exist for each pairs of levels. Examples are:
 - many speakers/interlocutors participate in a Dialog (the association between *oral_sources* and *oral_speakers* is 1: n)
 - a written document is composed of many pages (the association between *wrt_documents* and *wrt_pages* is 1:n).
- A number of auxiliary tables are also included in the diagrams:
 - Table *oral_audiofiles* is used for storing one or more digital audio files associated to an oral document.
 - Table *oral_SegmentTypes* contains segments of a morphological word (syllables, phonemes etc).
 - Table *wrt_wordextractconfigurations* is used for storing the tokenization configuration of a written document page (separators, regular expressions etc)



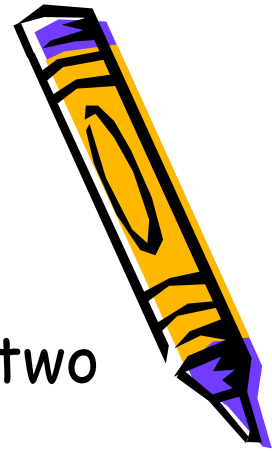
EAV database

- The EAV database keeps record of all kinds of annotation.
- It is based on the Entity-Attribute-Value representation aiming at the exemption from the continual Schema evolution problem and the extended use of null values.
- Annotations concern entities (tuples) of the database tables *oral_speakers*, *oral_turntakings*, *oral_words*, *oral_segments*, *wrt_documents*, *wrt_pages* and *wrt_words*.
- In other words, the entity (E of EAV) takes its values from the primary keys of the 7 abovementioned tables.
- As all annotations have the same functionality all annotation modules could be merged into one which would update a different set of attributes.
- In addition, meta-information could be managed by the same module since they share the same functionality but they update different sets of attributes.
- So, all 11 modules, i.e. 9 annotation modules (word part annotation for oral documents, and phonological, morphological, syntactic and semantic for both types of documents) and 2 meta-information modules could be merged into one which is applied to the set of attributes it applies to.

EAV schema which supports all the requirements



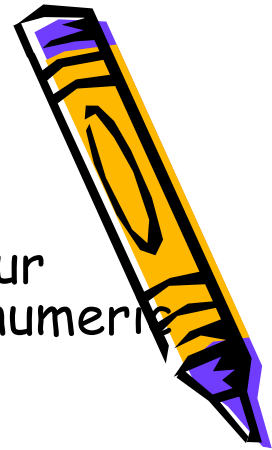
Brief description of the 9 tables consisting the EAV



1. The *app_sourcetypes* table contains the two document types (oral / written).
2. The *app_modules* table defines the 11 modules needed for the processing.
3. The *propertygroups* table is used for two reasons: (a) the definition of thematic subsets of attributes and (b) the definition of predefined values of attributes (lookups).
4. The auxiliary *propertygrouptypes* table contains the two reasons for which the *propertygroups* table is used.
5. The *properties* table contains all properties used by the 11 modules.



Brief description of the 9 tables consisting the EAV



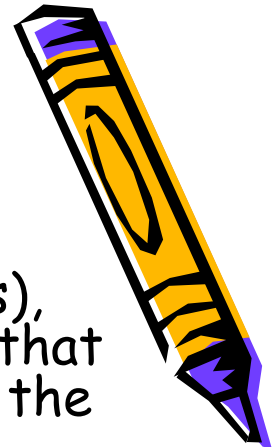
6. The auxiliary *propertytypes* table contains the four possible types of a property (alphanumeric, alphanumeric with multiple values, predefined value, multiple predefined value).
7. The *propertyvalues* table contains all acceptable values of lookups.
8. The *entitypropertyvalues* table is the main table of the EAV database. The *EntityId* field contains the Entities and takes its values among the primary keys of the 7 primary tables of the Struct database, the *PropertyId* contains the Attributes and takes its values from the primary keys of the properties table. The *ValueId* or the *ValueText* field contains the Values. The value domain of the *ValueId* is the primary key of the *propertyvalues* table. In the case where Attribute defined by the *PropertyId* field is an alphanumeric the *ValueText* is filled instead of the *ValueId*.



9. The *propertydependencies* table contains the properties which appear under the constraint that another property has a particular value. If an property does not appear in the *propertydependencies* table, then it is constraint free and it always appear in the module it is assigned to.

Need for a relation for the Inner level

- Annotations defining Syllables, Phonemes (Segments), Vowels and Consonants are the results of a process that imports TextGrid (Praat output) files. The way that the imported data are encapsulated should aim at:
 - a) Defining criteria for retrieving items at the three main levels (document, part, word) and the inner level (syllables, phonemes, etc). For example, we would like to be able to formulate a criterion such as seeking words ending with a stressed [u]. Obviously, this criterion should combine with other criteria (for example, metadata-based criteria such as that the speaker should be at least 75 years old and originates from Trabzon (Greek "Τραπεζούντα", [trape`zunta], Turkish "Trabzon" [ˈtrabzon])).
 - b) We should be able to create (on the fly) an artifact TextGrid (praat-like output) file with all the relevant annotations, from the information extracted from inner database. In the previous example (seeking words that end with a stressed [u]), our system should be able to create a Textgrid (praat-like output) file representing the word and all of its annotations, i.e. word, syllables, segments, vowels, consonants.



Tentative Structure of the Inner level relation



Phenomenon	WID	start	stop	Level	interval_no
u_s_e	30852	4.1234	4.2345	Vowel (11)	22
t	30852	3.9876	4.1233	Consonant (12)	27
u	30852	4.1234	4.2345	Segment (10)	28
tu	30852	3.9876	4.2345	Syllable (6)	12



This is our first approach. We have switched to EAV. However we keep it because it is easy for explanation.

inner level relation explanations



- All retrieved intervals of the same tier can be ordered based on the property `interval_no` which emanates from the original TextGrid file.
- In this way, we can formulate queries concerning the distance between segments at a certain level.
- The values of the attributes `phenomenon`, `start` and `stop` also emanate from the original TextGrid file.
- This structure enables the hierarchical reproduction of the data from the lower levels of fig. 1 and, at the same time, the serial access to the elements of the lower level.



Interface requirements

- Intuitive usage,
- Support Multi valued fields. As a consequence, the "And" operator is introduced for the values of a single criterion. That means that a demand for two or more values in a single record (item of a level of the data hierarchy) should be met, in addition to classical data demands (Exact, Range, Disjunction),
- 2 kinds of criteria (*main criteria* and *distance criteria*),
- Conjunction between main criteria (implicit use of *And* between rows of conditions),
- Expression of Retrieval requirements for: actual data, data aggregations, artifacts (on the fly created data),
- Expression of distance conditions (distance criteria) between items which are compatible with the main criteria. Therefore, the interface should support three different distance conditions (*Part, Word, Inner*).



Interface template



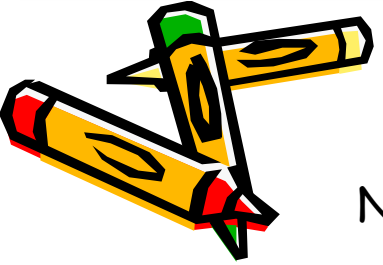
Word/ token /phenomenon		Location					
<Value>	{ Between, And, Or, -- }	<Value>	<DB>	<At- tribute >	<Part distanc es>	<Word distances>	<Interval_no distances >



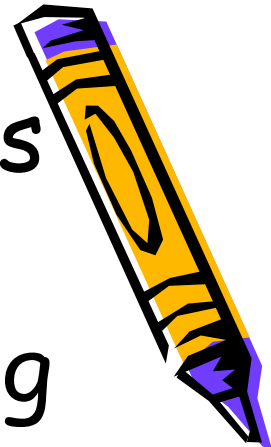


search of parts (pages for written resources) that contain the phenomenon of *Vowel Archaism*, followed by an adjective which is a loan word with a *Noun* Part of Speech and a *Masculin* Gender

Word/ token /phenomenon		Location				
vowel archaism	--	EAV Phon	↓	-	X	-
Adjective	--	EAV Morpho	PART OF SPEECH	-	Y in (X+1, X+10)	-
Noun	--	EAV Morpho	PART OF SPEECH OF LOAN WORD	-	Y	-
Masculin	--	EAV Morpho	GENDER OF LOAN WORD	-	Y	-
Output		Part	-			



search of parts (intonation phrases
in case of oral resources) ending
with an unstressed vowel, appearing
in the (oral) collection



Word/ token /phenomenon			Location				
?_u_f	--		detailed Phon	Vowel	-	-	-
Output			Document		count_part		



search based on metadata of participants



Word/ token /phenomenon			Location					
<u>Ifigenia Zisi</u>	Or	Mary <u>Karra</u>	EAV (O)	Meta	Annotator	-	-	-
Male	--		EAV (O)	Meta	Inf. Sex	-	-	-
75	Between	100	EAV (O)	Meta	Inf. Age	-	-	-
<u>cappadocians</u>	--		EAV (O)	Meta	Inf. Origin	-	-	-
Output			Document			-		



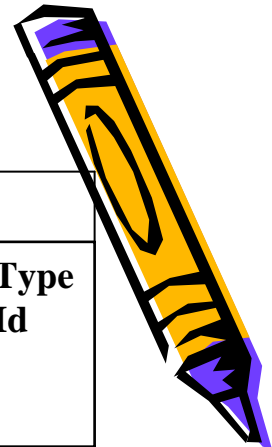
EAV implementation of Inner level annotations



- The tentative relational implementation for the inner level of annotations has constraints as for the expansion of annotations caused by the schema evolution problem.
- In order to solve this problem, we have decided to store information on the stress of the vowels in the EAV database.
- Instead of using the previous table (the Inner level Relation), we have placed the *oral_segments* table in the *Struct* database and also created the *oral_SegmentTypes* auxiliary table.
- Actually, the *oral_segments* table contains only the transcriptions in the lower levels (syllable, phoneme, etc) while, the eventual lower level element attributes are stored in the EAV database where the *Properties* and *PropertyValues* tables are updated with information on stress and stress position.



Inner Annotations in EAV



Oral_Segments								
Segment Id	Text	Word Id	Position Index In Word	Position Index In Turn Taking	Position Index In Speaker	Start Time	End Time	Type Id
8501	u	30852	3	13	22	4.1234	4.2345	1
8502	t	30852	4	14	27	3.9876	4.1233	2
8503	tu	30852	1	5	12	3.9876	4.2345	3

oral_SegmentTypes		
SegmentTypeId	Name	Code
1	Φωνήεν	VOWEL
2	Σύμφωνο	CONSONANT
3	Συλλαβή	SYLLABLE

EntityPropertyValues				
EAV_Id	EntityId	Property Id	Value Id	ValueText
11240	8501	148	1446	NULL
11241	8501	149	1450	NULL



Definitions of attributes and values in EAV



Properties				
Property Id	Name	TypeId	Is Mandatory	ModuleId
148	accent	3	0	11
149	Accent location	3	0	11

PropertyValues			
PropertyValueId	Code	Name	Property Id
1445	Unstressed	Άτονο	148
1446	Stressed	Τονισμένο	148
1447	Accented	Εστιασμένο	148
1448	Beginning of word	Αρχή Λέξης	149
1449	Middle of word	Μέση Λέξης	149
1450	End of word	Τέλος Λέξης	149
1451	End of phrase	Τέλος Φράσης	149



See also

- (in Greek) Νικήτας Ν. Καρανικόλας, Ελένη Γαλιώτου, Κωνσταντίνος Αθανασάκος & Γεώργιος Κορωνάκης. Ένα πολυτροπικό σύστημα αρχειοθέτησης και διαχείρισης γραπτών και προφορικών πηγών μελέτης της γλώσσας και των γλωσσικών ιδιωμάτων. Στο Αγγελική Ράλλη, Πρόγραμμα Θαλής: Πόντος, Κατπαδοκία, Αϊβαλί: Στα Χνάρια της Μικρασιατικής Ελληνικής, ISBN 978-960-99426-2-1.

http://users.teiath.gr/nnk/papers/C03_CR.pdf

http://users.teiath.gr/nnk/papers/C03_extended.pdf



Closing



- Thank you for your attention!
- Questions can be asked.
- nnk@teiath.gr
- <http://users.teiath.gr/nnk/>

