



# Εφαρμογές Τεχνολογιών Γλωσσικής Επεξεργασίας στα Συστήματα Αναζήτησης των Ελληνικών Ακαδημαϊκών Βιβλιοθηκών

Άννα Μάστορα<sup>1</sup>, Μανόλης Πεπονάκης<sup>2</sup>, Σαράντος Καπιδάκης<sup>1</sup>

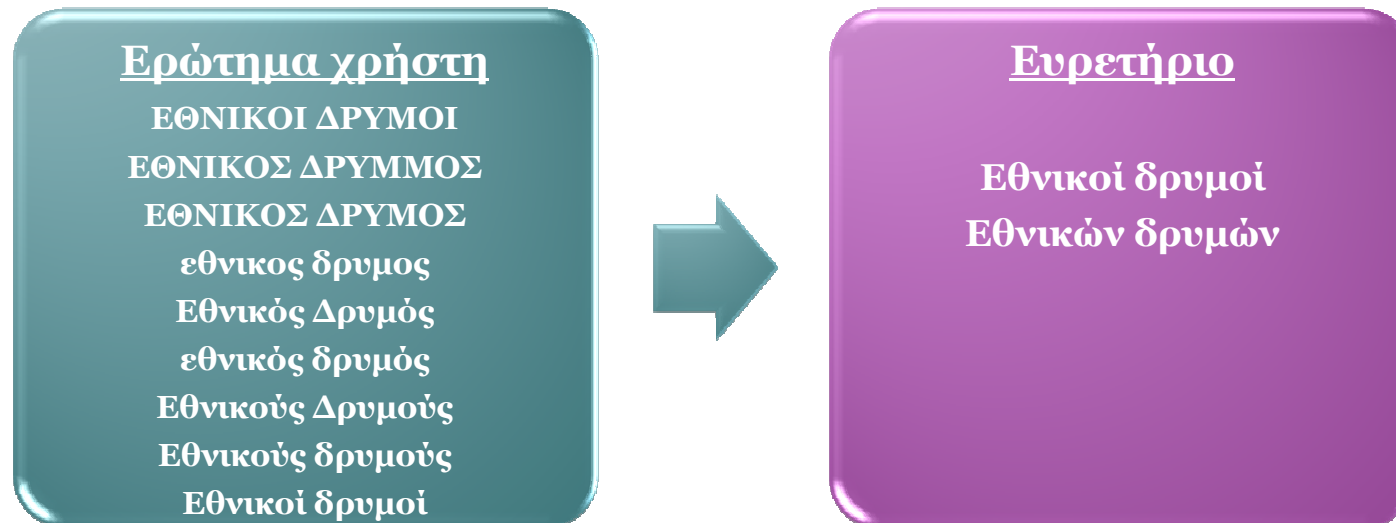
<sup>1</sup>Εργαστήριο Ψηφιακών Βιβλιοθηκών και Ηλεκτρονικής Δημοσίευσης,  
Τμήμα Αρχειονομίας - Βιβλιοθηκονομίας, Ιόνιο Πανεπιστήμιο  
{mastora, sarantos}@ionio.gr

<sup>2</sup>Εθνικό Κέντρο Τεκμηρίωσης, Εθνικό Ίδρυμα Ερευνών  
epepo@ekt.gr

21<sup>ο</sup> Πανελλήνιο Συνέδριο Ακαδημαϊκών Βιβλιοθηκών  
18-19 Οκτωβρίου 2012, Πειραιάς

# Κίνητρο της έρευνας

- Έντονη μορφολογία ελληνικής γλώσσας
- Αυξημένη χρήση της φυσικής γλώσσας στη διατύπωση ερωτημάτων
- Περιορισμένη χρήση ελεγχόμενου λεξιλογίου
- Ποικιλομορφία στη διατύπωση του ερωτήματος από τους χρήστες



# Γλωσσικές τεχνολογίες

- Στοχεύουν
  - Αυτόματη ανάλυση (και ίσως κατανόηση;) & παραγωγή γραπτών ή προφορικών εκφράσεων της φυσικής γλώσσας
- Πεδία εφαρμογής
  - Διόρθωση ορθογραφικών λαθών,
  - Εκμάθηση γλώσσας υποβοηθούμενη από Η/Υ,
  - Εξαγωγή πληροφορίας,
  - Αυτόματη περίληψη κειμένου,
  - **Ανάκτηση πληροφορίας**
  - κτλ.

## Γλωσσικές τεχνολογίες: Ανάκτηση πληροφορίας

- Αποκατάληξη (stemming)
- Λημματοποίηση (lemmatisation)
- Διαχωρισμός λέξεων (tokenisation)
- Διαχείριση σημείων στίξης & ανεπιθύμητων λέξεων
- Ορθογραφικός έλεγχος
- Εντοπισμός συνωνύμων
- Μορφολογική & συντακτική ανάλυση
- Αναγνώριση ονοματικών οντοτήτων (Named Entity Recognition)
- Διαχείριση χαρακτήρων: κεφαλαίων – πεζών, τονούμενων – άτονων

# Στόχος έρευνας

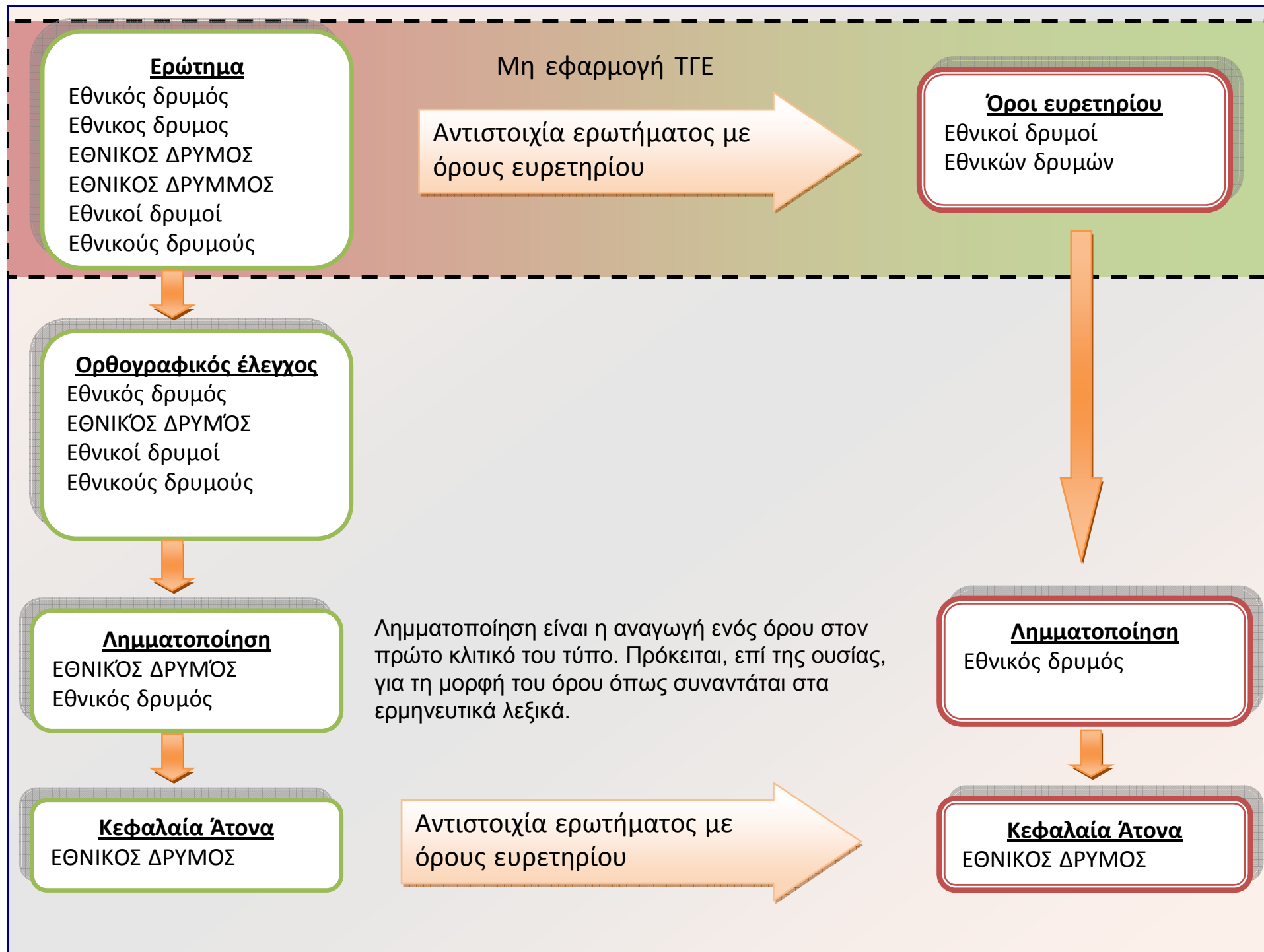
- Ανάδειξη πλεονεκτημάτων εφαρμογής γλωσσικών τεχνολογιών
  - Με έμφαση στην ελληνική γλώσσα
- Καταγραφή της υφιστάμενης κατάστασης αλλά και των προβλημάτων που προκύπτουν από την έλλειψη διαλειτουργικότητας στις εφαρμογές
  - Με έμφαση στα συστήματα αναζήτησης των Ελληνικών Ακαδημαϊκών Βιβλιοθηκών

# Αντικειμενικοί σκοποί

- Ανάδειξη αποτελεσμάτων από πειραματικά δεδομένα με εφαρμογή γλωσσικών τεχνολογιών
- Έρευνα σε Ελληνικές Ακαδημαϊκές Βιβλιοθήκες και Ιδρυματικά Αποθετήρια
  - Συγκριτική προσέγγιση ώστε να διαπιστωθεί η διαλειτουργικότητα

# Δεδομένα & εργαλεία πειράματος

- Διαπιστώθηκε η χρήση διαφορετικών μορφολογικών τύπων της ελληνικής γλώσσας στην υποβολή των ερωτημάτων
  - Διαπιστώθηκε αναντιστοιχία μεταξύ όρων ερωτήματος και όρων ευρετηρίου
  - Αναπτύχθηκε σουίτα εργαλείων
    - Εφαρμογή ad-hoc τεχνολογιών γλωσσικής επεξεργασίας
      - Διαχείριση τονούμενων-άτονων, διαχείριση σημείων στίξης, μετατροπή χαρακτήρων σε πεζούς ή κεφαλαίους, κτλ
    - Δυνατότητα χρήσης εξωτερικών εργαλείων
      - Ορθογραφικός έλεγχος (ASpell) -> Βελτίωση\* 9,75%
      - Λημματοποίηση (ilsp\_nlp, από το ΙΕΛ) -> Βελτίωση\* 16,7%
- \* Βελτίωση θεωρήθηκε η δυνατότητα ανάκτησης αποτελεσμάτων σε ερωτήματα που αρχικά είχαν μηδενικά αποτελέσματα





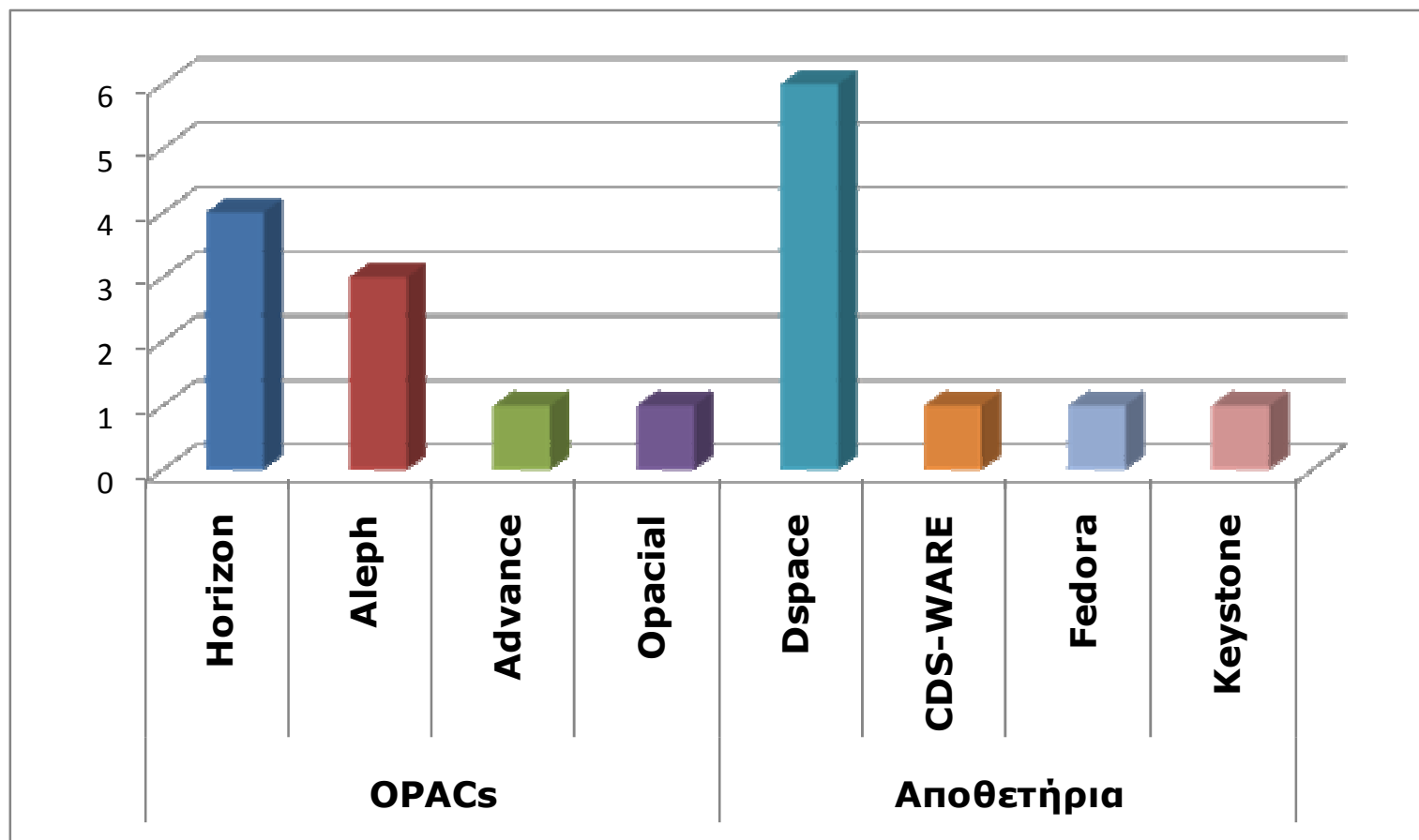
# Συμπεράσματα εφαρμογής ΤΓΕ

- Οι διαφορετικοί κλιτικοί τύποι ενός όρου μπορεί να επηρεάσουν σημαντικά την απόδοση της ανάκτησης πληροφορίας
- Εκ πρώτης όψης οι ΕΤΓΕ φαίνεται να είναι κάτι που μπορεί να βοηθήσει πολύ, ωστόσο, δεν υπάρχει ουσιαστική χρήση τους από τα συστήματα αναζήτησης ελληνικού περιεχομένου
- Δεδομένης της ύπαρξης εργαλείων και τεχνικών που θα μπορούσαν να συμβάλουν αποφασιστικά στην ανάπτυξη των υπηρεσιών αναζήτησης είναι σημαντικό να διερευνηθούν οι δυνατότητες συνέργειας στο πλαίσιο διεπιστημονικής προσέγγισης του θέματος

## Οι γλωσσικές τεχνολογίες στα συστήματα αναζήτησης των Ακαδημαϊκών Βιβλιοθηκών

- Επιλέξαμε: εκείνα τα ιδρύματα που είχαν δημόσιο κατάλογο (OPAC) και διέθεταν επίσης Ιδρυματικό Αποθετήριο
- Εξετάσαμε: τα συστήματα αναζήτησης των παραπάνω ιδρυμάτων ώστε να καταγραφεί η κατάσταση σχετικά με την εφαρμογή Τεχνολογιών Γλωσσικής Επεξεργασίας

# Κατανομή συστημάτων που εξετάστηκαν



## ΟΡΑCs και Εφαρμογές Τεχνολογιών Γλωσσικής Επεξεργασίας

ΟΡΑCs				
Βιβ/θήκη	Σύστημα	Διαφοροποιεί τονούμενα-άτονα	Ορθογραφικός έλεγχος ή προτεινόμενοι όροι	Αποκατάληξη ή λημματοποίηση
A	Horizon	Όχι	Όχι	Όχι
B	Aleph	Όχι	Ναι	Όχι
Γ	Aleph	Όχι	Ναι	Όχι
Δ	Horizon	Όχι	Όχι	Όχι
E	Advance	Όχι	Όχι	Όχι
ΣΤ	Horizon	Όχι	Όχι	Όχι
Z	Opacial	Όχι	Όχι	Όχι
H	Horizon	Όχι	Όχι	Όχι
Θ	Aleph	Όχι	Ναι	Όχι

## Ιδρυματικά Αποθετήρια και Εφαρμογές Τεχνολογιών Γλωσσικής Επεξεργασίας

Συστήματα Αποθετηρίων				
Βιβ/θηκη	Σύστημα	Διαφοροποιεί τονούμενα-άτονα	Ορθογραφικός έλεγχος ή προτεινόμενοι όροι	Αποκατάληξη ή λημματοποίηση
A	CDS-WARE	Όχι	<u>Ναι</u>	Όχι
B	DSpace	Όχι	Όχι	Όχι
Γ	Keystone	Όχι	Όχι	Όχι
Δ	DSpace	<u>Ναι</u>	Όχι	Όχι
E	DSpace	<u>Ναι</u>	Όχι	Όχι
ΣΤ	DSpace	<u>Ναι</u>	Όχι	Όχι
Z	Fedora	Όχι	Όχι	Όχι
H	DSpace	Όχι	Όχι	Όχι
Θ	DSpace	Όχι	Όχι	Όχι

## Συμπεράσματα ΕΤΓΕ σε συστήματα Ελληνικών Ακαδημαϊκών Βιβλιοθηκών

- Τα συστήματα αναζήτησης των Ελληνικών Ακαδημαϊκών Βιβλιοθηκών εφαρμόζουν σε περιορισμένη έκταση γλωσσικές τεχνολογίες
  - **τονούμενα-άτονα**
    - διαφορετικές πολιτικές τόσο μεταξύ των συστημάτων κάθε ιδρύματος όσο και των ιδρυμάτων μεταξύ τους
  - **κεφαλαία-πεζά**
    - όλα τα συστήματα αναζήτησης (OPACs & Αποθετήρια) εξισώνουν κεφαλαία-πεζά
      - Εξαίρεση: το τελικό σίγμα «ς»: μη εξίσωση του «ς» με το «Σ» ή το «σ» για τους OPACs των βιβλιοθηκών Α, Δ, ΣΤ, Η και στο Αποθετήριο Ε

# Επισήμανση!

- Η εξίσωση κεφαλαίων-πεζών ή τονούμενων-άτονων χαρακτήρων δεν αποτελεί (πάντοτε) ζητούμενο (τονικά παρώνυμα κτλ.)
  - Αθ<sup>ή</sup>να – Αθην<sup>ά</sup>
  - τσι<sup>ί</sup>πουρα – τσιπου<sup>ύ</sup>ρα
  - γέ<sup>ε</sup>ρος – γερό<sup>ο</sup>s
  - νό<sup>ο</sup>μος - νομό<sup>ο</sup>s
  - Μα<sup>ρ</sup>γαρίτα – μα<sup>ρ</sup>γαρίτα

## Συμπεράσματα: σύνοψη

- Οι γλωσσικές τεχνολογίες μπορούν να έχουν θετικά αποτελέσματα για την ανάκτηση πληροφοριών από τα συστήματα αναζήτησης
  - Όστε να διευκολύνουν το χρήστη στην αναζήτηση
    - Όπως έδειξαν τα αποτελέσματα της έρευνάς μας, η ληματοποίηση μπορεί να παίξει καθοριστικό ρόλο στην ανάκτηση ελληνικού περιεχομένου
- Είναι πολύ σημαντική η χρήση κοινών πολιτικών ώστε να επιτυγχάνεται διαλειτουργικότητα
  - Διαφορετικά, ο χρήστης βρίσκεται στην ίδια (δύσκολη) θέση όπως και χωρίς την εφαρμογή ΤΓΕ



# Μελλοντική έρευνα

- Είναι αναγκαία η χρήση (δοκιμή) των εργαλείων και τεχνικών ΕΤΓΕ σε μεγάλο όγκο δεδομένων ελληνικού περιεχομένου
  - ❖ Ανάγκη απόκτησης δεδομένων και από ελληνικούς φορείς για επιβεβαίωση των πειραματικών δεδομένων
- Δοκιμή περισσότερων εργαλείων και πιο εκλεπτυσμένων τεχνικών για πιο στοχευμένη ερμηνεία της πρόθεσης του χρήστη (query intent)
  - Μορφοσυντακτική ανάλυση
  - Αναγνώριση ονοματικών οντοτήτων
  - Εντοπισμός συνωνύμων, κτλ.
- Μελέτες συμπεριφοράς χρηστών για εξαγωγή-επιβεβαίωση μοντέλων αναζήτησης πληροφορίας

# ΕΥΧΑΡΙΣΤΩ ΓΙΑ ΤΗΝ ΠΡΟΣΟΧΗ ΣΑΣ!

*Η παρούσα έρευνα έχει συγχρηματοδοτηθεί από την Ευρωπαϊκή Ένωση (Ευρωπαϊκό Κοινωνικό Ταμείο - ΕΚΤ) και από εθνικούς πόρους μέσω του Επιχειρησιακού Προγράμματος Εκπαίδευση και Δια Βίου Μάθηση» του Εθνικού Στρατηγικού Πλαισίου Αναφοράς (ΕΣΠΑ) – Ερευνητικό Χρηματοδοτούμενο Έργο: Ηράκλειτος II. Επένδυση στην κοινωνία της γνώσης μέσω του Ευρωπαϊκού Κοινωνικού Ταμείου.*



Ευρωπαϊκή Ένωση  
Ευρωπαϊκό Κοινωνικό Ταμείο



ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ & ΘΡΗΣΚΕΥΜΑΤΩΝ, ΠΟΛΙΤΙΣΜΟΥ & ΑΘΛΗΤΙΣΜΟΥ  
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης

