



Εθνικόν και Καποδιστριακόν
Πανεπιστήμιον Αθηνών

Τμήμα Πληροφορικής και Τηλεπικοινωνιών

Φωνητικές Διεπαφές Χρήστη-Τεχνολογίες Φωνής

Ενότητα 5: Διαχείριση και Έλεγχος
Φωνητικού Διαλόγου
Γεώργιος Κουρουπέτρογλου

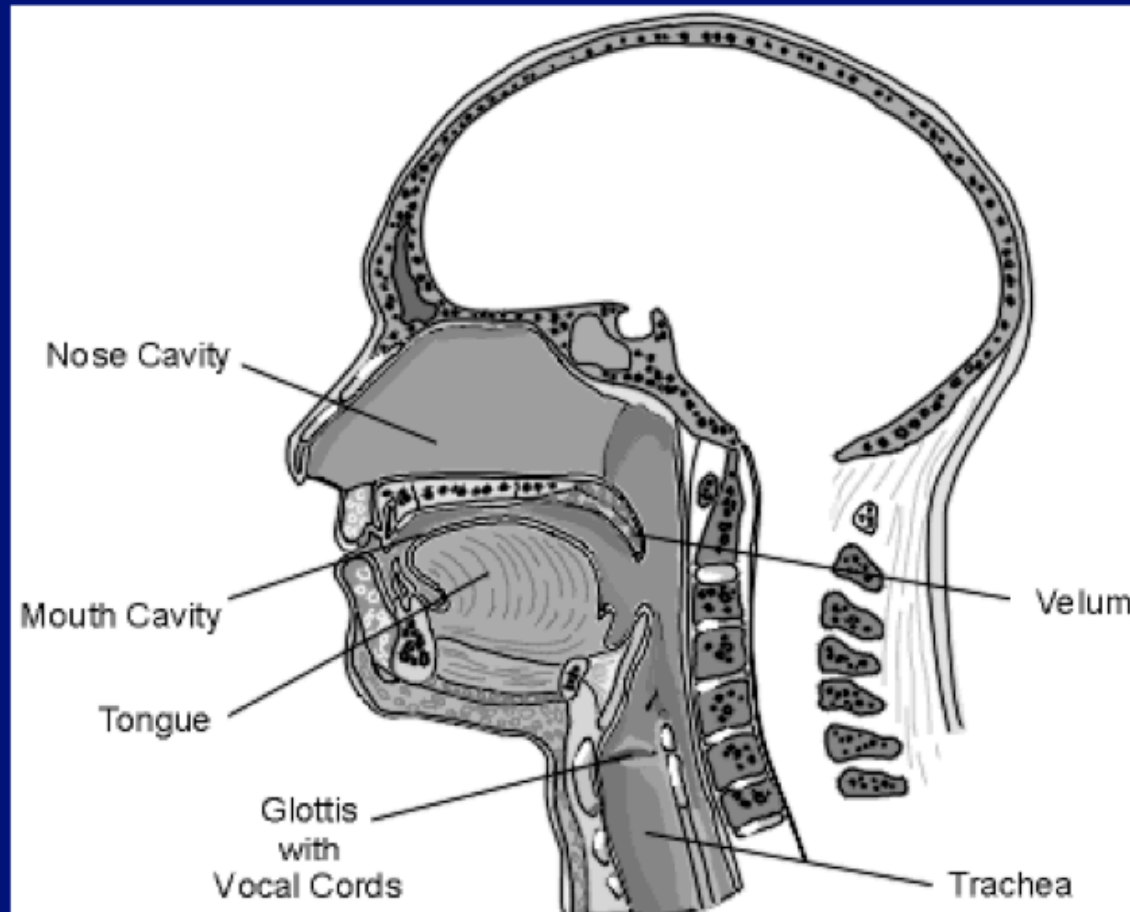
koupe@di.uoa.gr



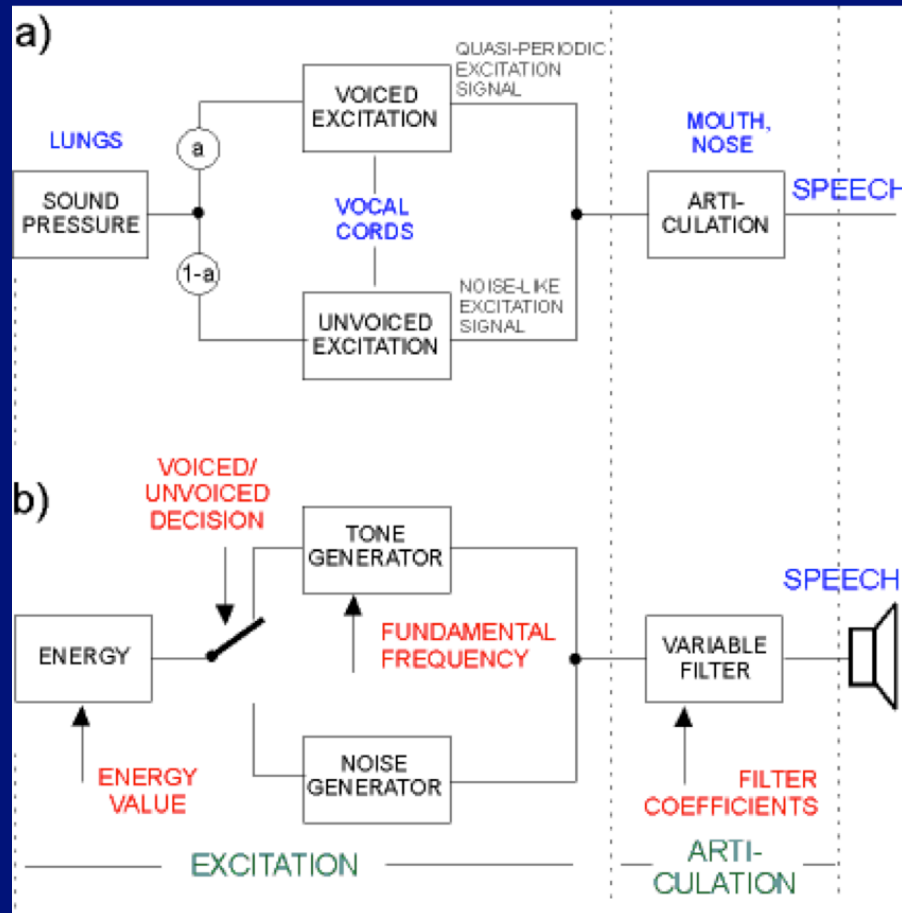
Περιεχόμενα ενότητας

Μοντέλα και εργαλεία διαχείρισης και ελέγχου
φωνητικού διαλόγου

Μοντέλο παραγωγής ομιλίας (1 από 2)



Μοντέλο παραγωγής ομιλίας (2 από 2)



Σύνθεση φωνοσυντονισμών

- Ελέγχεται από έναν πίνακα 40 παραμέτρων που ανανεώνουν την συμπεριφορά πηγών και φίλτρων κάθε 5 msec.
- «Είμαι ο πρώτος συνθέτης ομιλίας του Πανεπιστημίου Αθηνών» (1998)
 - Μέτρια καταληπτότητα
 - Χαμηλή φυσικότητα

Σύνθεση συρραφής διφώνων

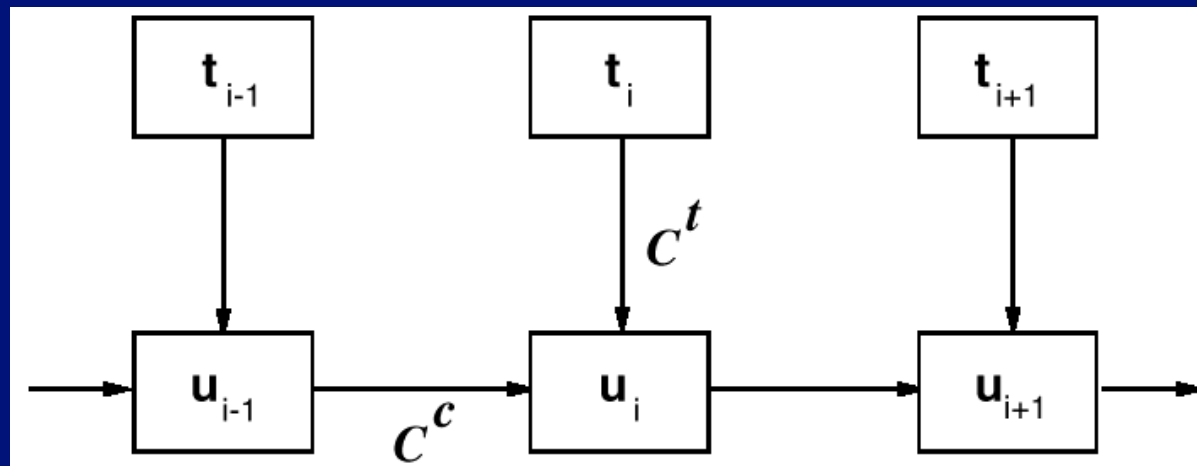
- Οι παράμετροι της ομιλίας εμπεριέχονται σε μικρά ηχογραφημένα τμήματα (φωνήματα, δίφωνα, λέξεις, φράσεις).
- «Είμαι η πρώτη φωνητική βάση διφώνων του Πανεπιστημίου Αθηνών» (2001)
 - Υψηλή καταληπτότητα
 - Μέτρια φυσικότητα

Προσωδία

- «Η θέλησή του Δημοσθένη ήταν τόσο μεγάλη, ώστε, όπως μας αναφέρει ο Πλούταρχος, έβαζε στο στόμα του μικρά χαλίκια την ώρα που απήγγειλε λόγους, προκειμένου να βελτιώσει την άρθρωσή του. Οι προσπάθειες του αυτές, απέδωσαν καρπούς και εξελίχθηκε σε σπουδαίο ρήτορα και πολιτικό.»
- Από εμπειρικά μοντέλα... σε εκπαιδευόμενα...

Σύνθεση από σώμα ηχογραφήσεων...

- Οι παράμετροι της ομιλίας εμπεριέχονται σε μικρά ηχογραφημένα τμήματα (φωνήματα, δίφωνα, λέξεις, φράσεις).
 - Υψηλή καταληπτότητα
 - Υψηλή φυσικότητα



...γενικών και περιορισμένων θεματικών πεδίων

- «Αυτό το έκθεμα είναι ένας αμφορέας και σήμερα βρίσκεται στο μουσείο της Αθήνας»
 - Rhetorical – rVoice 1 (γενικό)
 - Rhetorical – rVoice 2 (γενικό)
 - Loquendo – Actor 1 (γενικό)
 - ΔΗΜΟΣΘΕΝΗΣ – (περιορισμένο)

Τι είναι «κείμενο»;

- «Ο Νίκος παίζει πιάνο.»
- «Στις 21/12 ο Νίκος θα παίξει 3 από τις 4 μπαλάντες του στο ΚΨΜ.»
- «Ο Νίκος (ο αδελφός του κ. Πέτρου) δουλεύει στο ΚΨΜ.»
- «Ο Νίκος είναι:
 - Ψηλός
 - Έξυπνος
 - Αφηρημένος»
- Ο **Νίκος** έρχεται με τα πόδια, ενώ η Μαρία με **αυτοκίνητο**.
- | | |
|---------|--------|
| Μοντέλο | Κυβικά |
| Ibiza | 1600 |
| Amazon | 2000 |

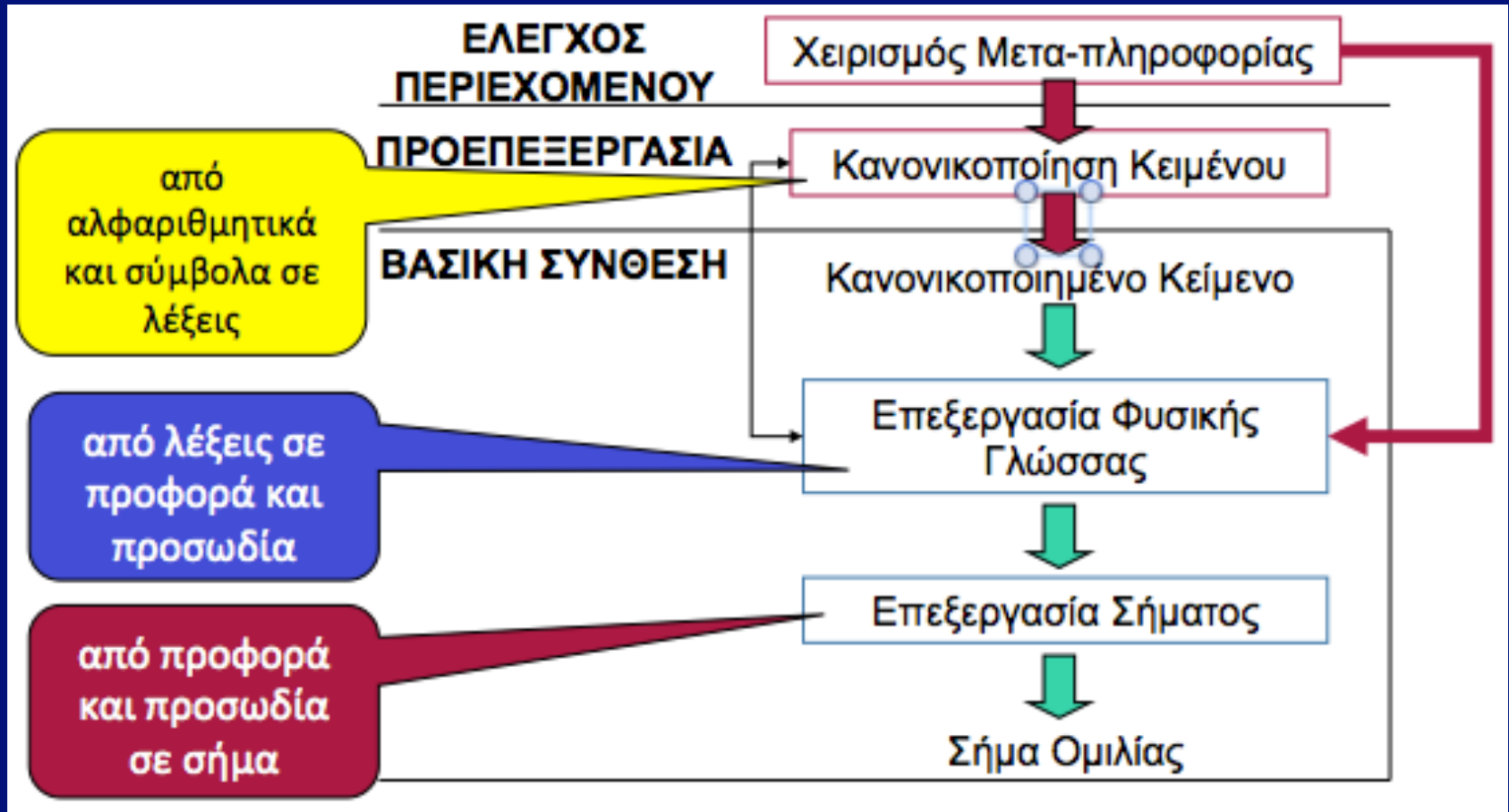
Τι είναι «συνθετική ομιλία»;

- Ομιλία που παράγεται με αυτόματες διαδικασίες από έναν ηλεκτρονικό υπολογιστή και προσομοιάζει την συμπεριφορά της ανθρώπινης ομιλίας.
- Σαν τεχνολογία υφίσταται πειραματικά από τα τέλη '50 και εμπορικά από το '70.
- Διεπαφή με τον χρήστη τελευταίας γενιάς.
- Από hardware (κουτάκια, χαμηλή ρομποτική ποιότητα) σε software (ευέλικτα, υψηλή ποιότητα)

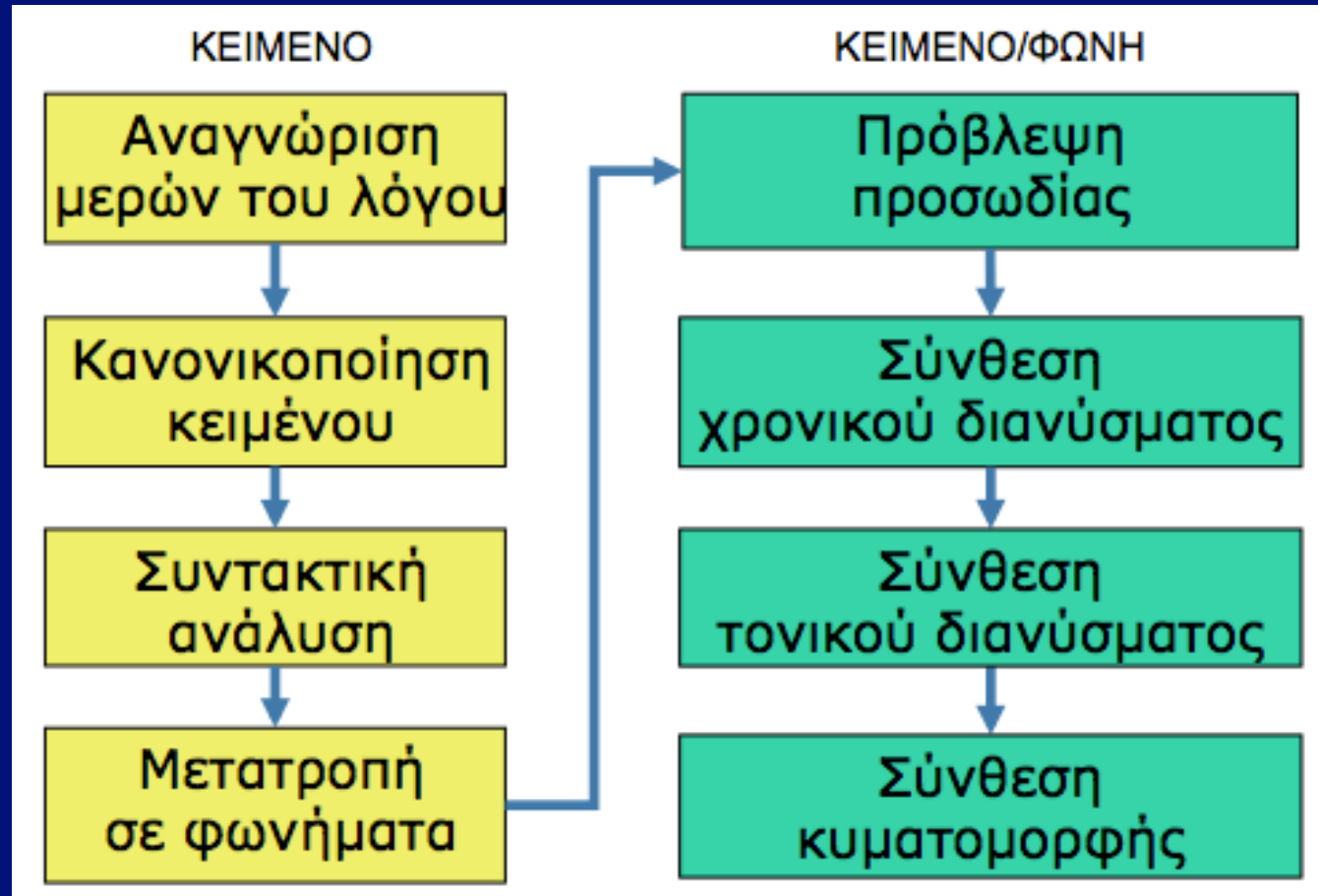
Τι είναι ένα «Σύστημα Μετατροπής Κειμένου σε Ομιλία»;

- Ένα σύστημα το οποίο δέχεται μία ακολουθία **συμβόλων** και με τεχνικές νοημοσύνης συνθέτει μία ή περισσότερες **εκδοχές** από αντίστοιχα **ακουστικά σήματα ομιλίας**.
- Ποιό είναι το πεδίο ορισμού των συμβόλων;
- Πόσο εξελιγμένες είναι οι τεχνικές νοημοσύνης;
- Πως επιλέγεται η κατάλληλη εκδοχή;
- Πόσο φυσικά είναι τα ακουστικά σήματα που παράγονται;

Διαδικασία σύνθεσης



Κλασική Μετατροπή ΚΣΟ



Κανονικοποίηση κειμένου (1 από 3)

- «Μένω Πατησίων και 3ης Σεπτεμβρίου»
- «Πάρε με στο 210 7275320»
- «Ήταν 13 Ιανουαρίου»
- «Ήταν 13 του μηνός»
- «Ήταν 1 γυναίκα, 1 άντρας και 1 παιδί»
- «Εξετάστηκαν οι αιτήσεις 1654 φοιτητών»
- «Ελάτε μεταξύ 3 και 4 η ώρα το απόγευμα»
- «1500 λουλούδια ήταν οι υποψήφιοι»
- «1500 ήταν οι υποψήφιοι»

Κανονικοποίηση κειμένου (2 από 3)

alpha	EXPN	abbreviation	<i>adv, N.Y, mph, gov't</i>
	LSEQ	letter sequence	<i>CIA, D.C, CDs</i>
	ASWD	read as word	<i>CAT, proper names</i>
	MSPL	misspelling	<i>geogaphy</i>
N U M B E R S	NUM	number (cardinal)	<i>12, 45, 1/2, 0.6</i>
	NORD	number (ordinal)	<i>May 7, 3rd, Bill Gates III</i>
	NTEL	telephone (or part of)	<i>212 555-4523</i>
	NDIG	number as digits	<i>Room 101</i>
	NIDE	identifier	<i>747, 386, I5, pc110, 3A</i>
	NADDR	number as street address	<i>5000 Pennsylvania, 4523 Forbes</i>
	NZIP	zip code or PO Box	<i>91020</i>
	NTIME	a (compound) time	<i>3.20, 11:45</i>
	NDATE	a (compound) date	<i>2/2/99, 14/03/87 (or US) 03/14/87</i>
	NYER	year(s)	<i>1998, 80s, 1900s, 2003</i>
S	MONEY	money (US or other)	<i>\$3.45, HK\$300, Y20,000, \$200K</i>
	BMONEY	money tr/m/billions	<i>\$3.45 billion</i>
	PRCT	percentage	<i>75%, 3.4%</i>
M I S C	SPLT	mixed or "split"	<i>WS99, x220, 2-car</i> (see also SLNT and PUNC examples)
	SLNT	not spoken, word boundary	<i>M.bath, KENT*RLTY, _really_</i>
	PUNC	not spoken, phrase boundary	non-standard punctuation: "****" in <i>\$99,9K***Whites, "..."</i> in <i>DECIDE...Year</i>
	FNSP	funny spelling	<i>sllooooooww, sh*t</i>
	URL	url, pathname or email	<i>http://apj.co.uk, /usr/local, phj@tpt.com</i>
	NONE	should be ignored	ascii art, formatting junk

Κανονικοποίηση κειμένου (3 από 3)

- «Μένω Πατησίων και 3ης Σεπτεμβρίου»
- «Ήταν 13 Ιανουαρίου»
- «Ήταν 13 του μηνός»
- «Πάρε με στο 210 7275320»
- «Ήταν 1 γυναίκα, 1 άντρας και 1 παιδί»
- «Εξετάστηκαν οι αιτήσεις 1654 φοιτητών»
- «Ελάτε μεταξύ 3 και 4 η ώρα το απόγευμα»
- «1500 λουλούδια ήταν οι υποψήφιοι»
- «1500 ήταν οι υποψήφιοι»

ΣΥΝΤΑΚΤΙΚΕΣ ΣΥΜΦΩΝΙΕΣ

- Οι κλιτές (inflected) γλώσσες έχουν επιπλέον το πρόβλημα της κλίσης των Non Standard Words (NSW).
- Για την αντιμετώπιση του προβλήματος χρησιμοποιούμε (ΔΗΜΟΣΘΕΝΗΣ) γλωσσολογική επεξεργασία που στηρίζεται στη Γραμματική της Νέας Ελληνικής (Μπαμπινιώτη – Χρήστου).
- Φροντίζουν για τη γραμματική συνέπεια ανάμεσα σε αριθμητικά και ουσιαστικά.
- 36,34% των περιπτώσεων βασίζονται σε αυτές τις συμφωνίες για την ορθή κανονικοποίηση τους.
- Συμφωνία ονοματικών συνόλων: **1636 φοιτητών**
- Αντικείμενο στην αιτιατική: **Το μουσείο δέχεται καθημερινά 1500 επισκέπτες**
- Κατηγορούμενο στην ονομαστική: **Οι επιτυχόντες είναι 1501**
- ...

Μετατροπή σε φωνήματα

- Ταύτιση της προφορικής αναπαράστασης με τη γραπτή.
- Εξαρτάται από ιδιώματα κάθε περιοχής.

«Τα παιδιά είδαν για μιά στιγμή τον Κώστα και την Χαρά σε μία ταβέρνα.»

«ta pEDÉÜ βDan Ña MÉa stiãmβ ton gýsta KE tíí xarÜ
sE mβa tanÝrna »

- «χαρά» «xarÜ»
- «χέρι» «÷Ýri»

Προσωδία (1 από 4)

- Τονικές (κινέζικα, αφρικάνικα) & επιτονικές γλώσσες (ευρωπαϊκές)
- Τονικές: διαφορετικά επίπεδα τόνου αλλάζουν το νόημα μίας λέξης.
- Επιτονικές γλώσσες: free-stress (π.χ. Αγγλικά) και fixed-stress (π.χ. Γαλλικά).
- Ο μελωδικός τονισμός και η καμπύλη επιτονισμού εξαρτώνται από την οργάνωση των λέξεων σε μονάδες υψηλότερου επιπέδου (προσωδιακές λέξεις, ονοματικά σύνολα, ενδιάμεσες φράσεις κλπ)

Προσωδία (2 από 4)

- Συντακτική ανάλυση σχετικά εύκολη σε σχέση με σημασιολογική και πραγματολογική
- Είναι ένα σημαντικό θέμα, σε ένα πρώτο επίπεδο, η απόδοση της προσωδίας σε συνάρτηση των συντακτικών σχέσεων ανάμεσα στις λέξεις των φράσεων.
 - Μελωδικές φράσεις συνήθως συναντώνται μέσα σε μορφο-συντακτικές ομάδες
 - Είναι σπάνια η εμφάνιση παύσεων ανάμεσα σε στενά συνδεδεμένες (συντακτικά) λέξεις.
- Είναι όμως δυνατόν να διαβαστεί ένα κείμενο χωρίς καταρχήν να γίνεται κατανοητό;

Προσωδία (3 από 4)

- Εστίαση (focus): «Ο Νίκος ήρθε στη Νάξο με πλοίο»
 - Σωστό μήκος συλλαβών και επιτονισμός
 - Λάθος εστίαση οδηγεί σε παρεξήγηση
 - Το TtS πρέπει να καταλαβαίνει όχι απομονωμένες προτάσεις αλλά όλο το κείμενο...
- Κοινός τόπος εκφωνητή (E) και ακροατή (A):
 - Ένα προσωδιακό πρότυπο από τον E -> A εξαρτάται από το τι γνωρίζει ο A για την κατάσταση που του περιγράφεται
 - Σταμάτα να ρωτάς αν ο «Νίκος ήρθε στη Νάξο με πλοίο».
Είναι η τρίτη φορά που σου λέω ότι ΔΕΝ ήρθε.
 - Η ιδέα ότι ο Νίκος θα ερχότανε στη Νάξο με πλοίο είναι αστεία: σιχαίνεται τα πλοία.

Προσωδία (4 από 4)

Τα βήματα που ακολουθεί ένα TtS προκειμένου να γεννήσει τα προσωδιακά διανύσματα μίας φράσης:

- Αναγνώριση και κατηγοριοποίηση φράσεων.
- Πρόβλεψη προσωδιακής δομής φράσεων
 - Θέση, διάρκεια και τύπος παύσεων (break indices)
 - Θέση και τύπος μελωδικών τόνων (pitch accents)
 - Θέση και τύπος τόνων στα όρια (endtones – prosodic phrase tones + boundary tones)
- Υπολογισμός διανύσματος διάρκειας
- Απόδοση καμπύλης F0.
- Υπολογισμός διανύσματος έντασης

Αναγνώριση μερών του λόγου

- Απαραίτητη διαδικασία για:
 - Κανονικοποίηση κειμένου σε κλιτές γλώσσες (π.χ. Ελληνικά)
 - Μετατροπή σε φωνήματα (π.χ. «χρόνια», «record»)
 - Πρόβλεψη προσωδιακής δομής μίας πρότασης:
 - Προσδιορισμός προσωδιακών φράσεων
 - Προσδιορισμός παύσεων στις προσωδιακές φράσεις
 - Προσδιορισμός μελωδικού τονισμού (pitch accent) σε φράσεις
 - Προσδιορισμός τόνων ορίων (boundary tones) σε φράσεις
- Ποσοστό επιτυχίας (ελληνικά): ~95%

Θεματικά πεδία (1 από 2)

- Πεδία εφαρμογής TtS στα οποία η θεματολογία είναι συγκεκριμένη.
 - Περιορισμένο λεξιλόγιο.
 - Περιορισμένα γλωσσολογικά φαινόμενα.
 - Περιορισμένα προσωδιακά φαινόμενα (από επαγγελματίες εκφωνητές του πεδίου).
 - Άρα, μικρότερο πεδίο προσομοίωσης και λιγότερα λάθη κατά την επεξεργασία
- Ευκολύνεται η δημιουργία στατιστικών μοντέλων από πραγματικές μετρήσεις (π.χ. Μορφοσυντακτική ανάλυση, προσωδία -> ~95%)
- Ακουστικό σήμα: γίνεται προσιτή η χρήση μεγαλύτερων ακουστικών μονάδων συρραφής -> μεγαλύτερη φυσικότητα

Θεματικά πεδία (2 από 2)

- Παραδείγματα:
 - Υπηρεσίες καταλόγου (131, σινεμά, βενζινάδικα κλπ)
 - Ανακοινώσεις (δελτίο καιρού, οικονομικό δελτίο κλπ)
 - Περιήγηση σε μουσεία
 - Τεχνικά ή νομικά κείμενα
 - Ειδήσεις
- Με βάση τους περιορισμούς που πρέπει να θέτει ένα θεματικό πεδίο, δεν θεωρούνται θεματικά πεδία:
 - Λογοτεχνικά κείμενα
 - Ποίηση

«Φωνή»

- Ο καθένας θέλει μία φωνή με προσωποποιημένα χαρακτηριστικά.
- Τα ομιλούντα προϊόντα (θα) θέλουν να έχουν την δική τους ξεχωριστή φωνή.
- Τι είναι φωνή;
 - Καταρχήν, ένα συγκεκριμένο ηχόχρωμα
 - Μία ακολουθία από διαδικασίες επεξεργασίας φυσικής γλώσσας: διαφορετικές τοπολαλίες, ιδιώματα, τρόποι προφοράς συνεπτηγμένων μορφών, προσωδιακή συμπεριφορά.

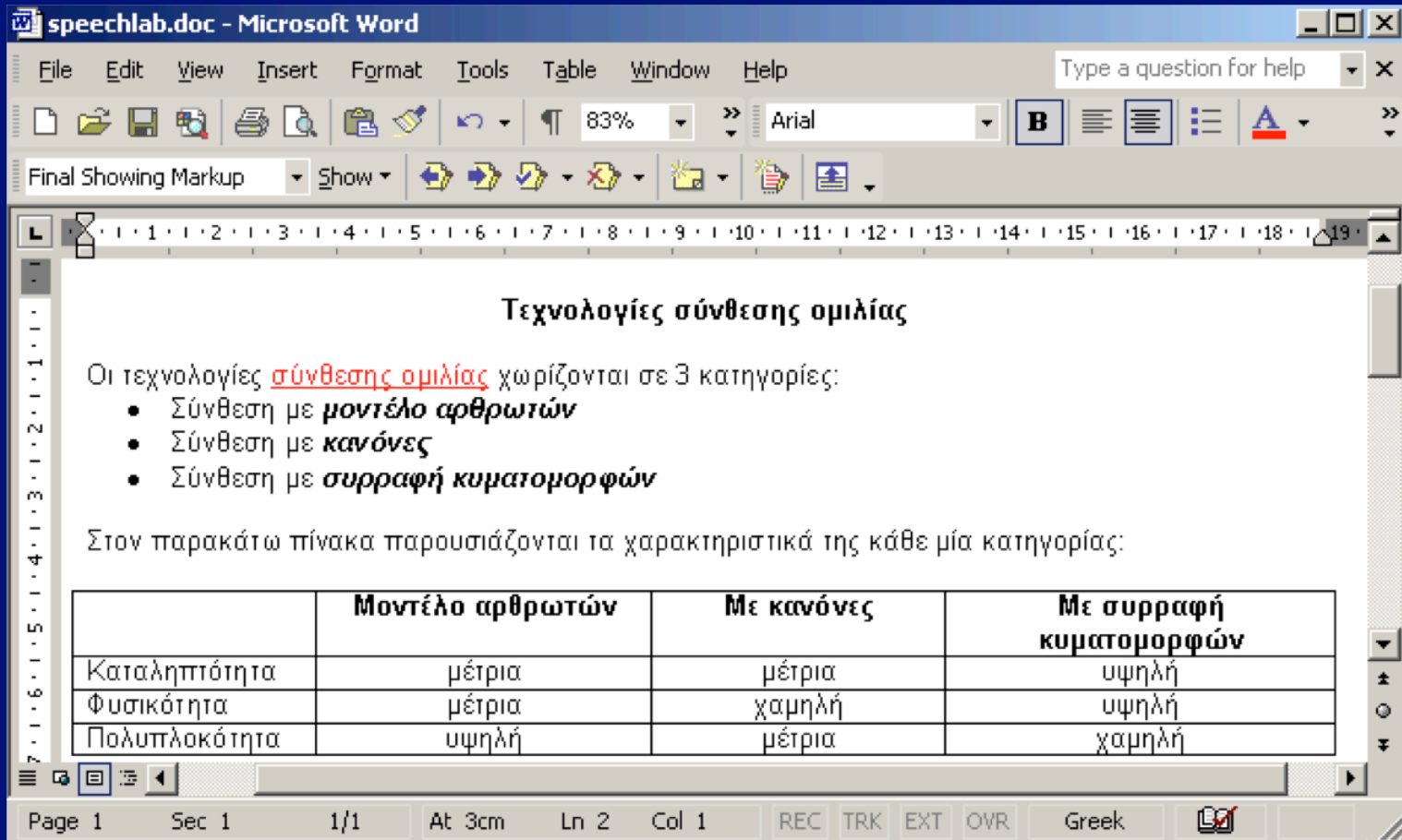
Ποιότητα συνθετικής ομιλίας

- **Καταληπτότητα:** ήταν κατανοητό το περιεχόμενο της ομιλίας;
- **Φυσικότητα:** πόσο κοντά στην φυσική χροιά ήταν η ομιλία;
- Επιπλέον, μεταδόθηκε σωστά η **πληροφορία**;

Μετα-πληροφορία (1 από 2)

- **Οδηγίες οπτικοποίησης:** bold, italics, tables, bullets, big-small letters κλπ (π.χ. HTML, MS-Word)
- **Οδηγίες δομής:** header, title, section, record content κλπ (π.χ. HTML, XML, SQL)
- **Δομή κειμένου:** παρενθέσεις, σημάδια, κλπ
- **Γλωσσολογική πληροφορία:** ρητορικές σχέσεις, σύνταξη, γραμματική, μορφολογία (π.χ. SOLE, plain κείμενο)
- **Οδηγίες ομιλίας:** prosody, emp, rate, pitch κλπ (π.χ. SABLE, VoiceXML, SSML, ACSS)

Μετα-πληροφορία (2 από 2)



The screenshot shows a Microsoft Word window titled "speechlab.doc". The document content is as follows:

Τεχνολογίες σύνθεσης ομιλίας

Οι τεχνολογίες **σύνθεσης ομιλίας** χωρίζονται σε 3 κατηγορίες:

- Σύνθεση με **μοντέλο αρθρωτών**
- Σύνθεση με **κανόνες**
- Σύνθεση με **συρραφή κυματομορφών**

Στον παρακάτω πίνακα παρουσιάζονται τα χαρακτηριστικά της κάθε μία κατηγορίας:

	Μοντέλο αρθρωτών	Με κανόνες	Με συρραφή κυματομορφών
Καταληπτικότητα	μέτρια	μέτρια	υψηλή
Φυσικότητα	μέτρια	χαμηλή	υψηλή
Πολυπλοκότητα	υψηλή	μέτρια	χαμηλή

Page 1 Sec 1 1/1 At 3cm Ln 2 Col 1 REC TRK EXT OVR Greek



Άσκηση 5.1

Α) Εγκαταστήστε στον υπολογιστή σας το σύστημα Μετατροπής Κειμένου σε Ομιλία ΔΗΜΟΣΘΕΝΗΣ <http://demosthenes.di.uoa.gr/>

Β) Γράψτε 10 ελληνικές προτάσεις που να περιέχει κάθε μία τουλάχιστον δύο αριθμητικά από τις κατηγορίες: τηλέφωνο, ημερομηνία, ποσό (να παραδοθούν και οι προτάσεις). Δώστε τις προτάσεις αυτές να τις εκφωνήσει ο ΔΗΜΟΣΘΕΝΗΣ και σημειώστε (και αναφέρετε) τα τυχόν λάθη στο τμήμα Κανονικοποίησής του.

Γ) Πάρτε ένα τυχαίο ελληνικό κείμενο από το διαδίκτυο τουλάχιστον 100 λέξεων (να παραδοθεί το κείμενο αυτό) και δώστε το να το εκφωνήσει ο ΔΗΜΟΣΘΕΝΗΣ. Ελέγξατε αν το τμήμα Μετατροπής σε Φωνήματα του ΔΗΜΟΣΘΕΝΗΣ λειτουργεί σωστά σε όλες τις περιπτώσεις φωνημάτων και αλλοφώνων. Σημειώστε (και αναφέρετε) τα προβλήματα που βρήκατε.

Τέλος Ενότητας

Διαχείριση και Έλεγχος Φωνητικού Διαλόγου

Χρηματοδότηση

- Το παρόν εκπαιδευτικό υλικό έχει αναπτυχθεί στο πλαίσιο του εκπαιδευτικού έργου του διδάσκοντα.
- Το έργο «**Ανοικτά Ακαδημαϊκά Μαθήματα στο Πανεπιστήμιο Αθηνών**» έχει χρηματοδοτήσει μόνο την αναδιαμόρφωση του εκπαιδευτικού υλικού.
- Το έργο υλοποιείται στο πλαίσιο του Επιχειρησιακού Προγράμματος «**Εκπαίδευση και Δια Βίου Μάθηση**» και συγχρηματοδοτείται από την Ευρωπαϊκή Ένωση (Ευρωπαϊκό Κοινωνικό Ταμείο) και από εθνικούς πόρους.



Σημειώματα

Σημείωμα Ιστορικού Εκδόσεων Έργου

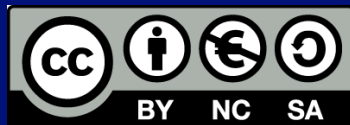
Το παρόν έργο αποτελεί την έκδοση 1.0.

Σημείωμα Αναφοράς

Copyright Εθνικών και Καποδιστριακών Πανεπιστημίων Αθηνών 2015, Γεώργιος Κουρουπέτρογλου 2015. Γεώργιος Κουρουπέτρογλου.
«Φωνητικές Διεπαφές Χρήστη-Τεχνολογίες Φωνής. Διαχείριση και Έλεγχος Φωνητικού Διαλόγου». Έκδοση: 1.0. Αθήνα 2015. Διαθέσιμο από τη δικτυακή διεύθυνση: <http://opencourses.uoa.gr/courses/DI37/>.

Σημείωμα Αδειοδότησης

Το παρόν υλικό διατίθεται με τους όρους της άδειας χρήσης Creative Commons Αναφορά, Μη Εμπορική Χρήση Παρόμοια Διανομή 4.0 [1] ή μεταγενέστερη, Διεθνής Έκδοση. Εξαιρούνται τα αυτοτελή έργα τρίτων π.χ. φωτογραφίες, διαγράμματα κ.λ.π., τα οποία εμπεριέχονται σε αυτό και τα οποία αναφέρονται μαζί με τους όρους χρήσης τους στο «Σημείωμα Χρήσης Έργων Τρίτων».



[1] <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Ως **Μη Εμπορική** ορίζεται η χρήση:

- που δεν περιλαμβάνει άμεσο ή έμμεσο οικονομικό όφελος από την χρήση του έργου, για το διανομέα του έργου και αδειοδόχο
- που δεν περιλαμβάνει οικονομική συναλλαγή ως προϋπόθεση για τη χρήση ή πρόσβαση στο έργο
- που δεν προσπορίζει στο διανομέα του έργου και αδειοδόχο έμμεσο οικονομικό όφελος (π.χ. διαφημίσεις) από την προβολή του έργου σε διαδικτυακό τόπο

Ο δικαιούχος μπορεί να παρέχει στον αδειοδόχο ξεχωριστή άδεια να χρησιμοποιεί το έργο για εμπορική χρήση, εφόσον αυτό του ζητηθεί.

Διατήρηση Σημειωμάτων

Οποιαδήποτε αναπαραγωγή ή διασκευή του υλικού θα πρέπει να συμπεριλαμβάνει:

- το Σημείωμα Αναφοράς
- το Σημείωμα Αδειοδότησης
- τη δήλωση Διατήρησης Σημειωμάτων
- το Σημείωμα Χρήσης Έργων Τρίτων (εφόσον υπάρχει)

μαζί με τους συνοδευόμενους υπερσυνδέσμους.

Σημείωμα Χρήσης Έργων Τρίτων (1 από 2)

Οι φωτογραφίες που περιέχονται στην παρουσίαση αποτελούν πνευματική ιδιοκτησία τρίτων. Απαγορεύεται η αναπαραγωγή, αναδημοσίευση και διάθεσή τους στο κοινό με οποιονδήποτε τρόπο χωρίς τη λήψη άδειας από τους δικαιούχους.

Σημείωμα Χρήσης Έργων Τρίτων (2 από 2)

Η δομή και οργάνωση της παρουσίασης, καθώς και το υπόλοιπο περιεχόμενο, αποτελούν πνευματική ιδιοκτησία του συγγραφέα και του Πανεπιστημίου Αθηνών και διατίθενται με άδεια Creative Commons Αναφορά Μη Εμπορική Χρήση Παρόμοια Διανομή Έκδοση 4.0 ή μεταγενέστερη.