



ΠΑΝΕΠΙΣΤΗΜΙΟ
ΠΑΤΡΩΝ
UNIVERSITY OF PATRAS

ΑΝΟΙΚΤΑ ακαδημαϊκά
μαθήματα ΠΠ

Εισαγωγή στη Βιοπληροφορική

Ενότητα 3: Τεχνικές Ανάλυσης και Σύγκρισης
Ακολουθιών Βιολογικών Δεδομένων I

Μακρής Χρήστος, Τσακαλίδης Αθανάσιος,
Περδικούρη Αικατερίνη

Πολυτεχνική Σχολή

Τμήμα Μηχανικών Η/Υ και Πληροφορικής

Σκοποί ενότητας

- Σκοπός της ενότητας είναι η παρουσίαση του Δέντρο Επιθεμάτων (Suffix Tree) και οι εφαρμογές του



Περιεχόμενα ενότητας

- Δέντρο επιθεμάτων – Γενικευμένο Δέντρο Επιθεμάτων
- Εφαρμογές στην ανάλυση ακολουθιών βιολογικών δεδομένων



Βασική Βιβλιογραφική Πηγή στην οποία βασίζονται οι διαφάνειες

- Dan Gusfield , Algorithms on Strings, Trees and Sequences,, Cambridge University Press, 10th edition 2007

Τεχνικές Ανάλυσης και Σύγκρισης Ακολουθιών Βιολογικών Δεδομένων I

Τεχνικές Ανάλυσης και Σύγκρισης Ακολουθιών Βιολογικών Δεδομένων

- Δέντρο Επιθεμάτων- Suffix Tree
- Γενικευμένο Δέντρο Επιθεμάτων - Generalized Suffix Tree
- Εφαρμογές σε Προβλήματα Μοριακής Βιολογίας



Βασικοί Ορισμοί

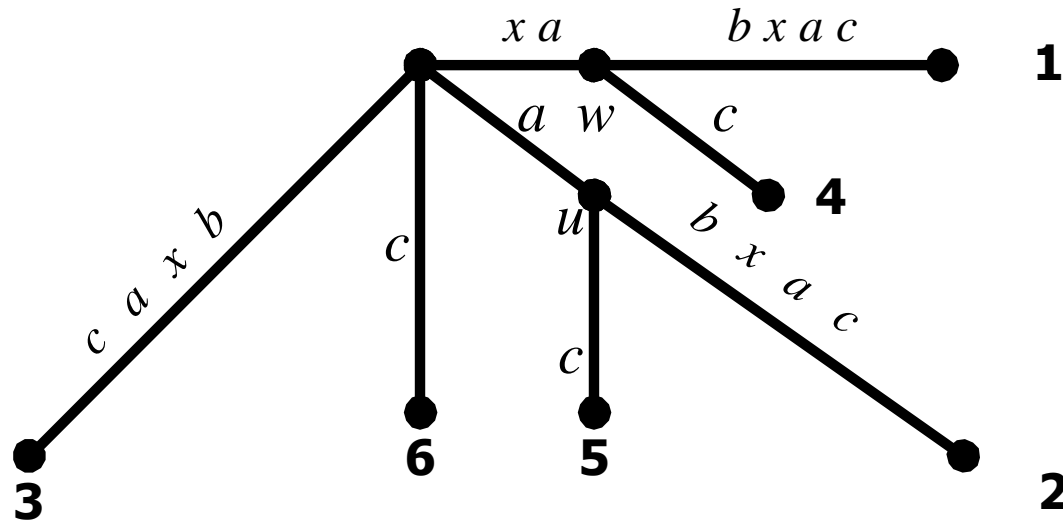
- Συμβολοσειρά-string: $x=x[1]x[2].....x[n]$, $x[i] \in \Sigma$ & $|x|=n$
 $x=acgttaaaca$, $|x|=10$ & $\Sigma=\{a,c,g,t\}$
- Κενή συμβολοσειρά: ϵ
- Υπο-συμβολοσειρά-substring w : $x=uwv$
- Πρόθεμα –Prefix w : $x=wu$
- Επίθεμα-Suffix w : $x=uw$
- Κάθε συμβολοσειρά S , μήκους $|S|=m$, έχει m δυνατά μη κενά επιθέματα που είναι τα ακόλουθα: $S[1...m]$, $S[2...m]$, $S[m-1...m]$ και $S[m]$.
- Παράδειγμα "sequence" : *sequence, equence, quence, uence, ence, nce, ce, e.*



Το Δέντρο Επιθεμάτων Suffix Tree

Ορισμός: «αποθηκεύει όλα τα δυνατά επιθέματα μιας συμβολοσειράς S ».

$x = \text{xabxac}$



Κατασκευή του Δέντρου Επιθεμάτων

- Μια απλοϊκή θεώρηση:
 - Ένθεση μιας πλευράς στο δέντρο για το επίθεμα $S[1\dots m]$,
 - Διαδοχική ένθεση των επιθεμάτων $S[i\dots m]$, για $i=2\rightarrow m$.

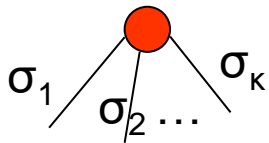


Trie

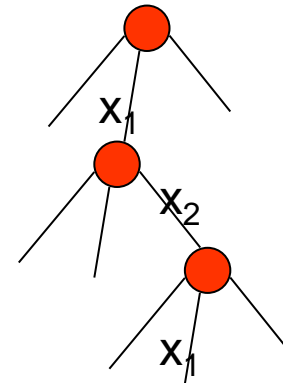
Ορισμός: Έστω σύμπαν $U = \Sigma^l$ για αλφάβητο Σ και $l > 0$.

$$(x \in U : x = d_1 d_2 \cdots d_l) \quad S \subseteq U :$$

Trie καλείται το κ -δικό δένδρο ($\kappa = |\Sigma|$) το οποίο περιέχει όλα τα προθέματα των στοιχείων του S . Κάθε επίπεδο του δένδρου αντιστοιχίζεται και σε ένα d_i ($d_1 \rightarrow$ ρίζα). Κάθε στοιχείο $x = x_1 x_2 \dots x_l$ τοποθετείται στο υποδένδρο $x_1 \rightarrow x_2 \rightarrow \dots$.

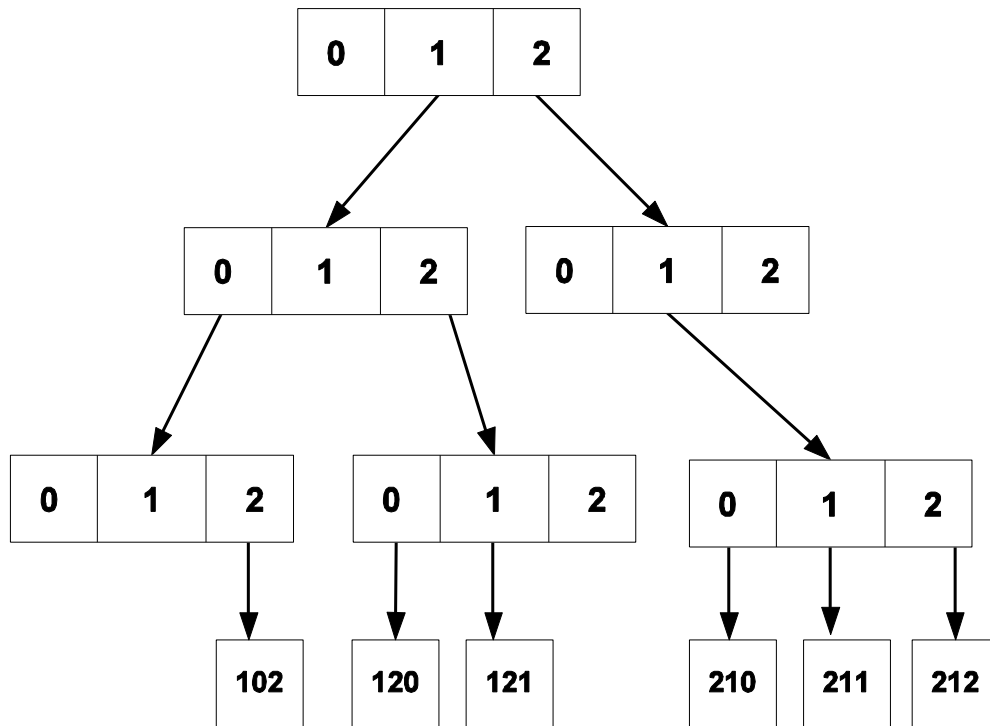


$$(\sigma_i \in \Sigma)$$



Trie - παράδειγμα

$S = \{102, 120, 121, 212, 211, 120\}$, $\Sigma = \{0, 1, 2\}$



digit 1

Χρόνος ins/del/search :
 $O(l)$

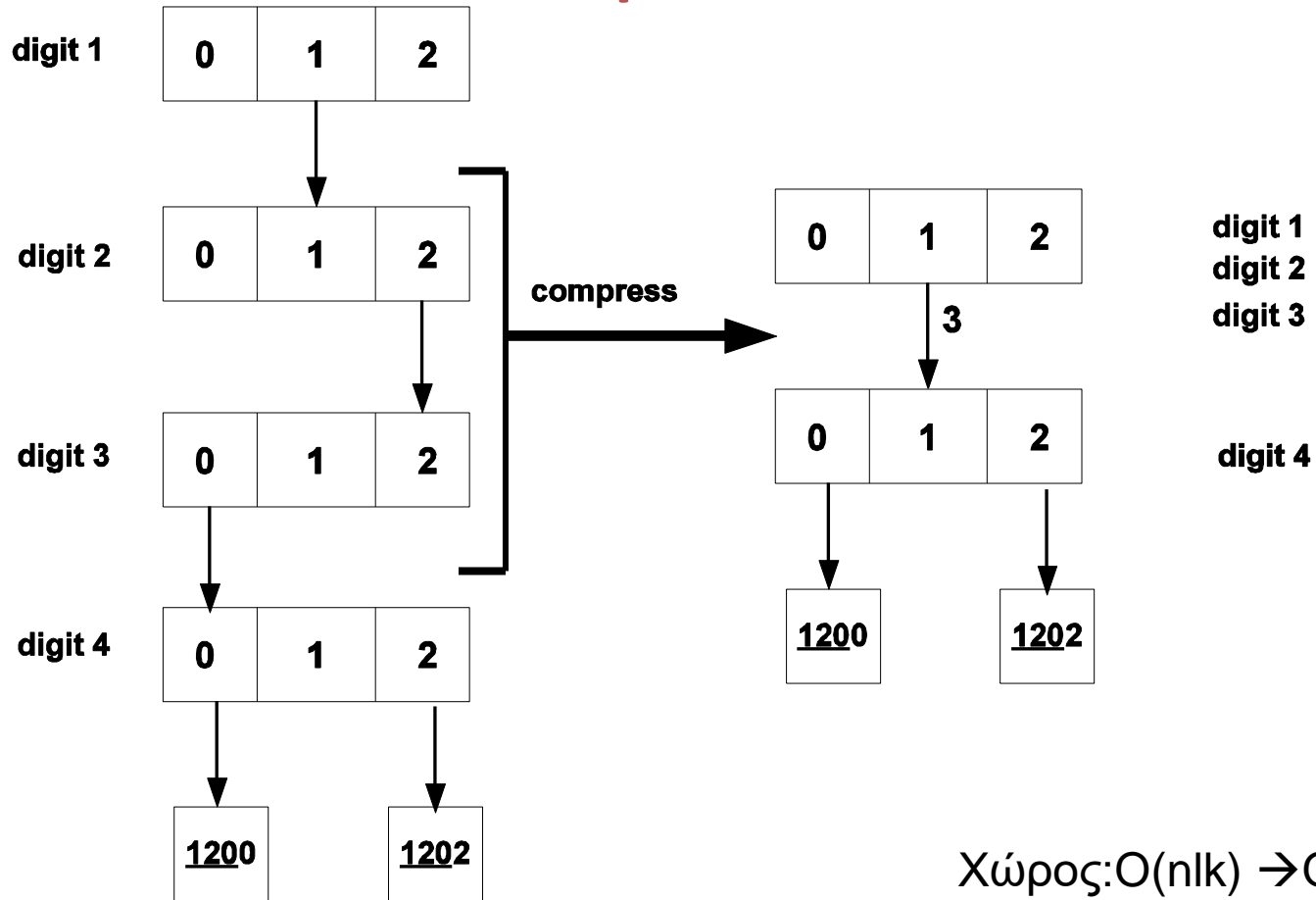
digit 2

Χώρος:
 $O(nlk)$

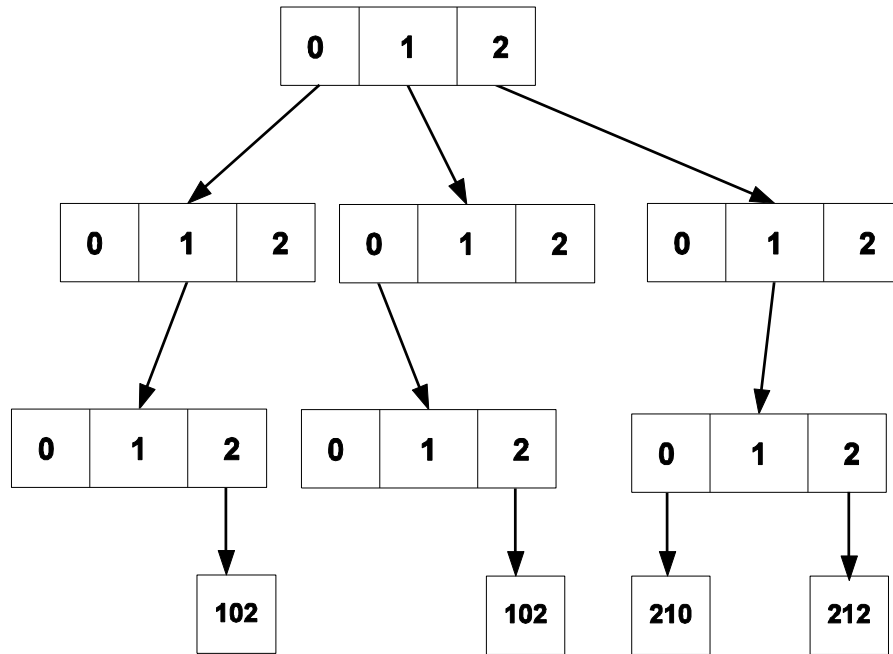
digit 3



Compressed Trie



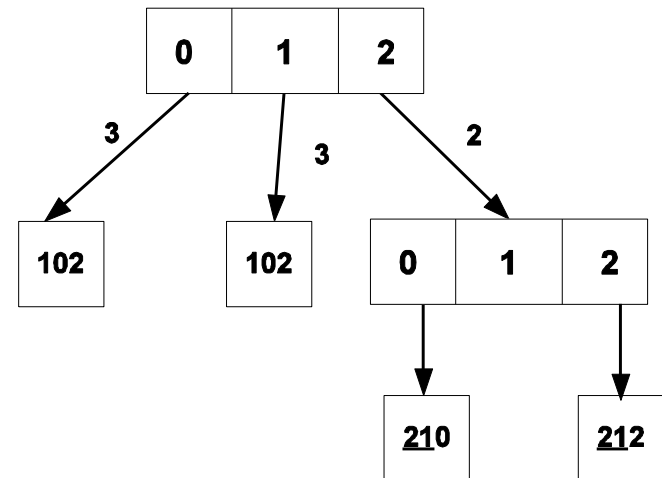
Compressed Trie - example



digit 1

digit 2

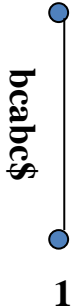
digit 3



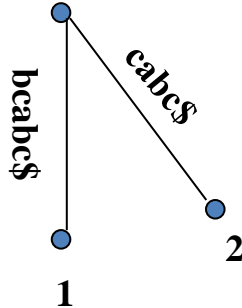
Naïve Κατασκευή

$S = bcabc\$$

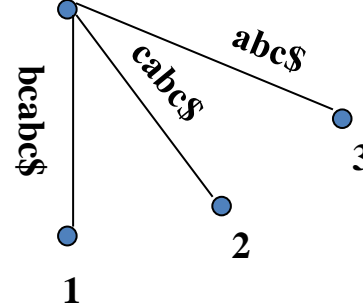
1: $bcabc\$$



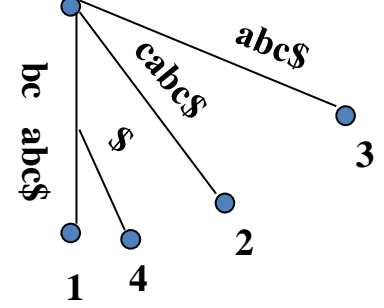
2: $cabc\$$



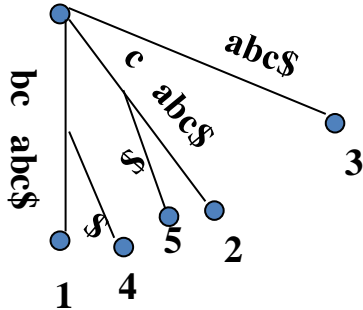
3: $abc\$$



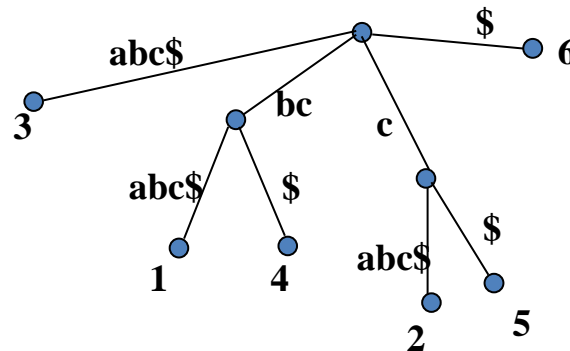
4: $bc\$$



5: $c\$$



6: $\$$

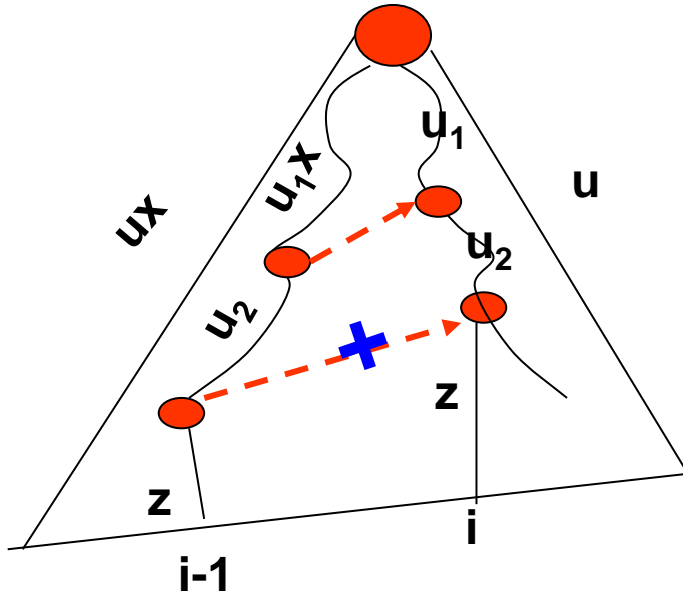


Χρόνος: $O(n^2)$

$S = aaaaaa... \$$



Suffix Links – Speed Up



	i-1	i	
	x	u	z
	x	u₁	u₂
			z

$$\text{res}(i) = |u_2| + |z|$$

$$\text{res}(i+1) \leq \text{res}(i) - \text{int}_i$$

$$\text{Scan cost} = \text{length}(\text{head}_i) - \text{length}(\text{head}_{i-1})$$



Κατασκευή του Γενικευμένου Δέντρου Επιθεμάτων

- Μια απλοϊκή θεώρηση:
 - Διαδοχική ένθεση όλων των επιθεμάτων των συμβολοσειρών εισόδου.



Εφαρμογές στην ανάλυση ακολουθιών βιολογικών δεδομένων

- Ακριβής Εύρεση Προτύπου - Exact pattern matching problem
- Ακριβής Εύρεση Πολλαπλών Προτύπων- Σύγκριση με το Αυτόματο Aho- Corasick
- Μέγιστη κοινή υποσυμβολοσειρά 2 ακολουθιών- Longest Common Substring Problem
- DNA Contamination Problem
- Εύρεση Κοινών Μοτίβων σε 2 ή περισσότερες Βιολογικές Ακολουθίες
- Εύρεση Επαναλήψεων σε Βιολογικές Ακολουθίες



Ακριβής Εύρεση Προτύπου - Exact pattern matching problem

- Κατασκευή του Δ.Ε. για την ακολουθία εισόδου T σε $O(|T|)$ χρόνο,
- Ξεκινώντας από τη ρίζα, σύγκρινε έναν προς έναν τους χαρακτήρες του P , ακολουθώντας το κατάλληλο μονοπάτι.
 - Εάν εμφανιστεί κάποιο μη-ταίριασμα, τότε το πρότυπο δεν εμφανίζεται στην ακολουθία,
 - διαφορετικά ανέφερε ως απάντηση όλα τα φύλλα που βρίσκονται κάτω από τον κόμβο του τελευταίου χαρακτήρα του P .
- Η αναζήτηση στοιχίζει $O(n+k)$ χρόνο, $|P|=n$ & k : το πλήθος εμφανίσεων του P στο T .



Ακριβής Εύρεση Πολλαπλών Προτύπων

- Κατασκευή του Δ.Ε. για την ακολουθία εισόδου T σε $O(|T|)$ χρόνο,
- Ξεκινώντας από τη ρίζα, αναζητούμε διαδοχικά όπως και στην προηγούμενη εφαρμογή το σύνολο των προτύπων $P=\{P_1, P_2, \dots\}$
- Η αναζήτηση στοιχίζει $O(n+|P|+k_p)$ χρόνο, $|P|=|P_1|+|P_2|+\dots$ & k_p : το πλήθος εμφανίσεων των προτύπων στο T .



- Κατασκευή του Αυτομάτου σε χρόνο $O(|P|)$
- Αναζήτηση: $O(m+k)$



Μέγιστη κοινή υποσυμβολοσειρά 2 ακολουθιών- Longest Common Substring Problem

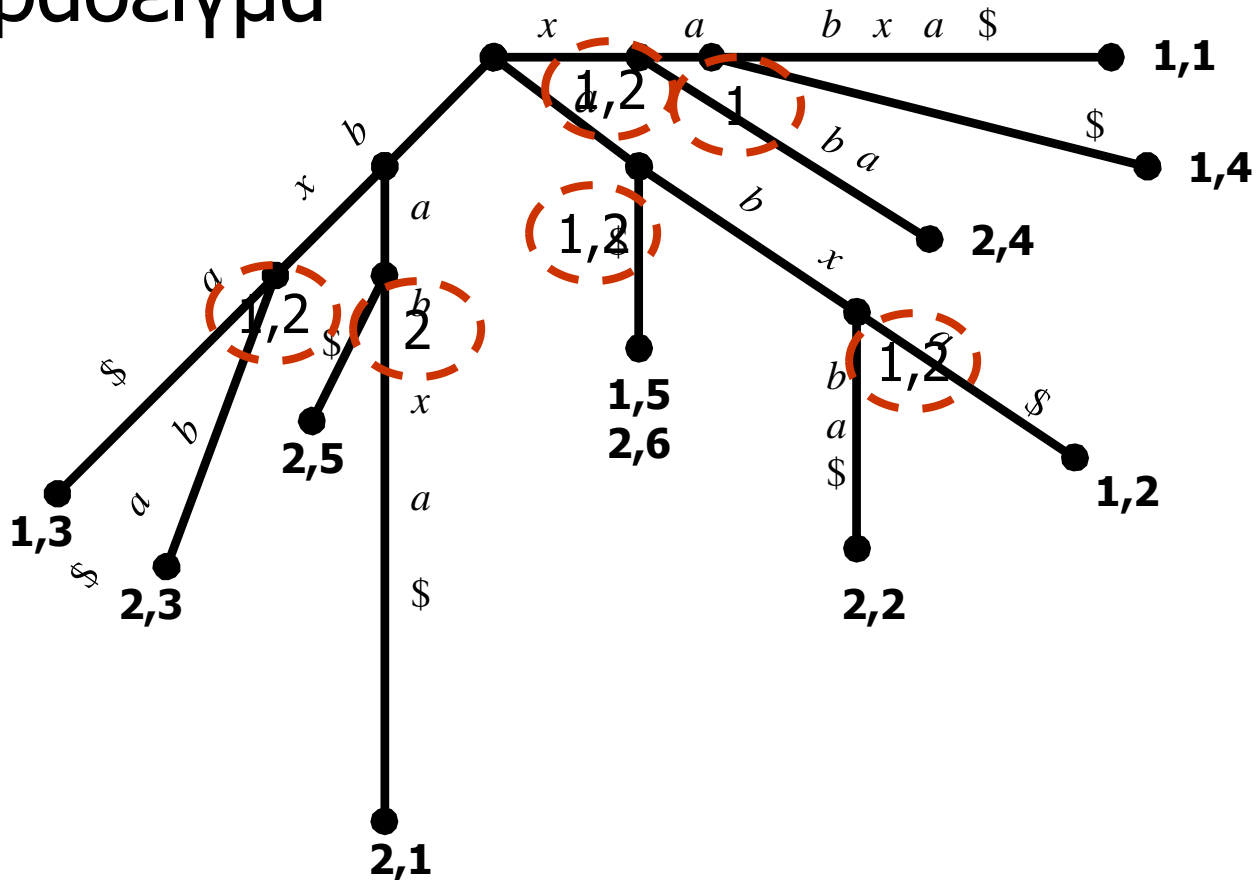
```
stringmatchersaniachers  
teacherA
```

- Κατασκευή του Γ.Δ.Ε. για τις ακολουθίες εισόδου S_1, S_2, \dots
- Σημειώνουμε κάθε εσωτερικό κόμβο του δέντρου u , με "1" ή "2", αν εμπεριέχει στο υπόδεντρο του u , κάποιο φύλλο που αναπαριστά κάποιο επίθεμα της ακολουθίας S_1 ή S_2 .
- Η ετικέτα μονοπατιού - path label, κάθε εσωτερικού κόμβου που σημειώνεται ταυτόχρονα με "1" και "2", αποτελεί μια κοινή υποσυμβολοσειρά των δυο ακολουθιών S_1 και S_2 .



Μέγιστη κοινή υποσυμβολοσειρά 2 ακολουθιών- Longest Common Substring Problem

- Παράδειγμα



DNA Contamination Problem

- **DNA Contamination Problem:** Για μια δοσμένη ακολουθία DNA S_1 , που έχει πρόσφατα απομονωθεί και ταυτοποιηθεί και μια ήδη γνωστή ακολουθία S_2 , (επιμέρους τμήματα που πιθανά έχουν μολυνθεί), αναζητούμε όλες τις υπο-συμβολοσειρές της S_2 που εμφανίζονται στην S_1 , με μήκος μεγαλύτερο από λ .
- Κατασκευάζουμε το **Γενικευμένο Δέντρο Επιθεμάτων** για τις ακολουθίες S_1 και S_2 .
- Ακολουθούμε τη μεθοδολογία και αναφέρουμε όλους τους κόμβους με βάθος $string-depth(u) \geq \lambda$.



Εύρεση Κοινών Μοτίβων σε 2 ή περισσότερες Βιολογικές Ακολουθίες

- **Το Πρόβλημα της Εύρεσης κοινών μοτίβων:** Για ένα σύνολο K ακολουθιών με συνολικό μήκος $\Sigma(|K|) = n$, και έναν ακέραιο k , ($2 < k < K$), ορίζουμε ως $\lambda(k)$, το μήκος του μέγιστου μοτίβου που εμφανίζεται σε τουλάχιστον k υποσυμβολοσειρές. Το πρόβλημα ανάγεται στον υπολογισμό όλων των δυνατών τιμών του $\lambda(k)$ και λύνεται σε γραμμικό χρόνο $O(n)$, ως προς το μήκος των ακολουθιών εισόδου.



Εύρεση Κοινών Μοτίβων σε 2 ή περισσότερες Βιολογικές Ακολουθίες

- Παράδειγμα
- Έστω $K = \{sandollar, sandlot, handler, grand, pantry\}$.

k	$l(k)$	μοτίβο
2	4	<i>sand</i>
3	3	<i>and</i>
4	3	<i>and</i>
5	2	<i>an</i>



Εύρεση Επαναλήψεων σε Βιολογικές Ακολουθίες

- Οι επαναλήψεις σε βιολογικές ακολουθίες κατηγοριοποιούνται στις εξής 3 βασικές κατηγορίες:
 - επαναλήψεις περιορισμένου μήκους που εμφανίζονται σε τοπικό επίπεδο, και των οποίων η λειτουργία είναι γνωστή,
 - επαναλήψεις περιορισμένου μήκους που εμφανίζονται σε όλο το μήκος της ακολουθίας, και των οποίων η λειτουργία δεν είναι απόλυτα γνωστή,
 - δομημένες επαναλήψεις μεγάλου μήκους των οποίων η λειτουργία δεν έχει προσδιοριστεί.



Παραδείγματα Επαναλήψεων

- 1η κατηγορία:
 - τα **συμπληρωματικά παλίνδρομα** σε ακολουθίες DNA & RNA, που ρυθμίζουν τη μετεγγραφή του DNA,
 - τα **εμφωλευμένα συμπληρωματικά παλίνδρομα** σε ακολουθίες tRNA
- 2η κατηγορία:
 - **συνεχόμενες επαναλήψεις- tandem repeats,**
 - **δορυφορικά τμήματα DNA- satellite DNA, (micro & mini satellite DNA)**
- 3η κατηγορία:
 - **SINE-Short Interspersed Nuclear Sequences (π.χ.: *Alu family*)**
 - **LINE-Long Interspersed Nuclear Sequences.**

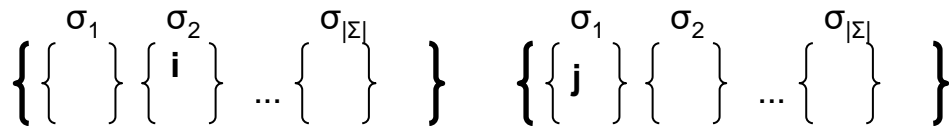
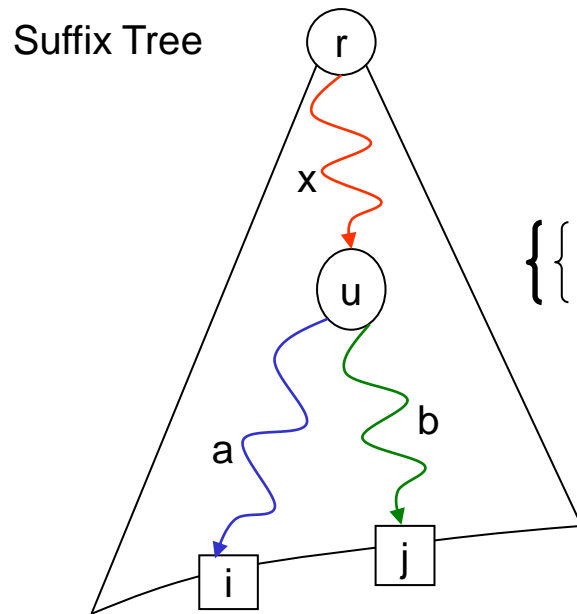
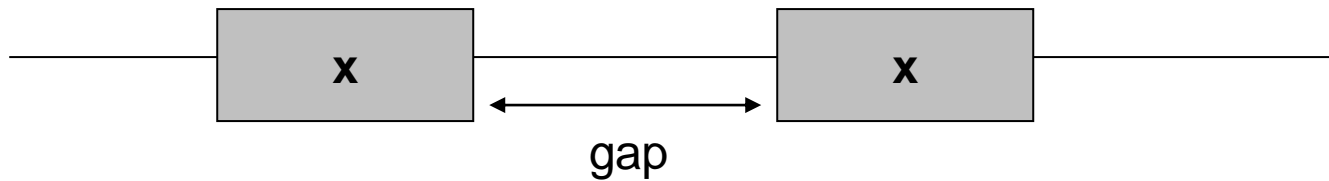


Παλίνδρομα

- **Ορισμός:** Ένα **παλίνδρομο- *palindrome*** αποτελεί την επαναλαμβανόμενη εμφάνιση της υπο-συμβολοσειράς που διαβάζεται ως ίδια και προς τις 2 κατευθύνσεις (από αριστερά προς τα δεξιά και από δεξιά προς τα αριστερά): *xyaayx*
- **Ορισμός:** Ένα παλίνδρομο σε μια ακολουθία DNA ή RNA, ονομάζεται **συμπληρωματικό παλίνδρομο- *complemented palindrome***, αν προκύπτει από την αντικατάσταση όλων των χαρακτήρων από την αρχή έως τη μέση με τις αντίστοιχες συμπληρωματικές βάσεις: *agctcgcgagct*



Maximal Pairs

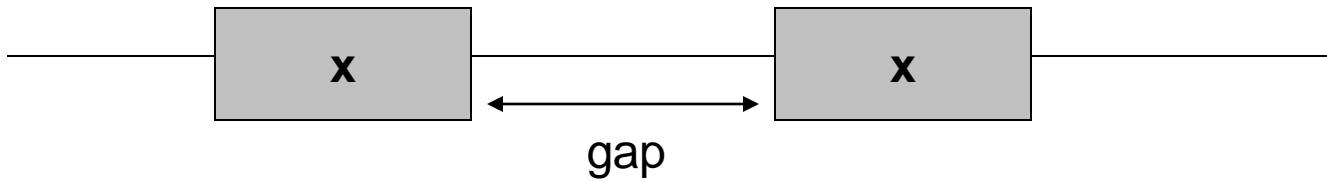
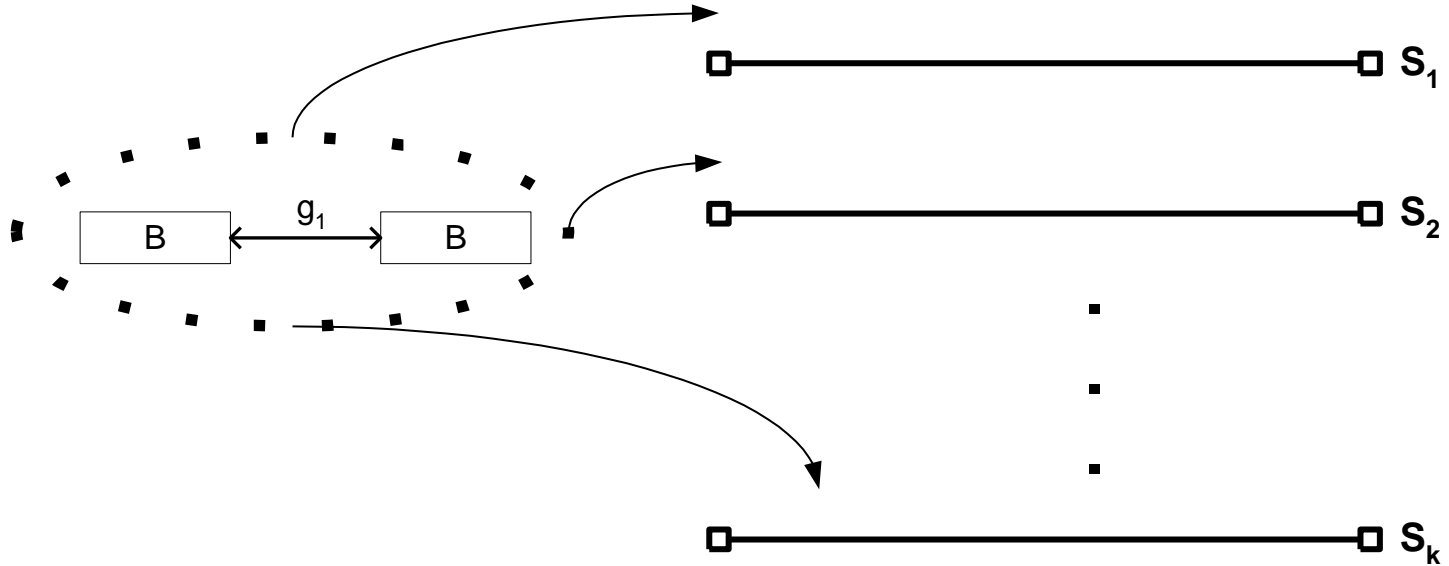


Gusfield : $O(n+a)$

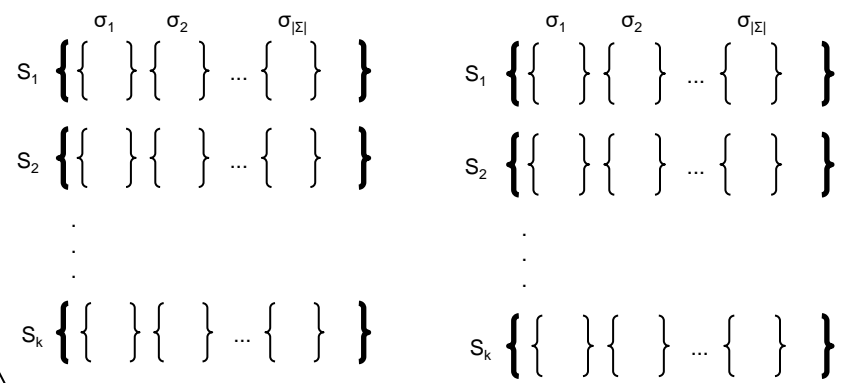
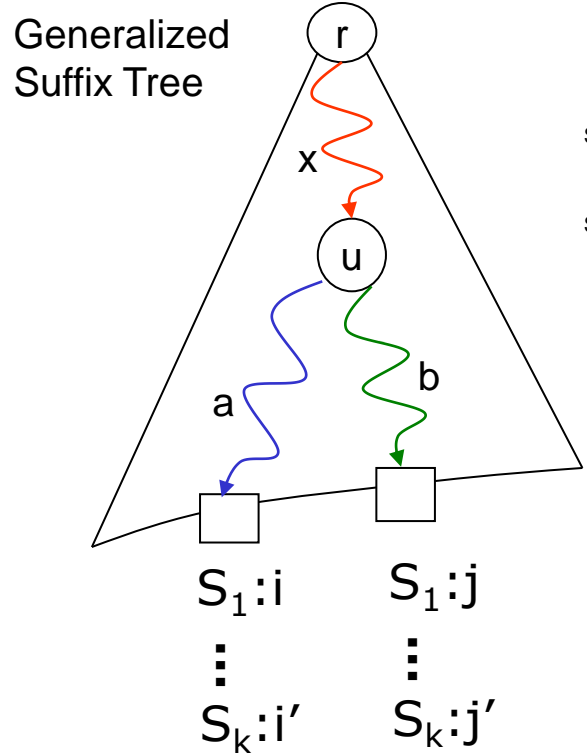
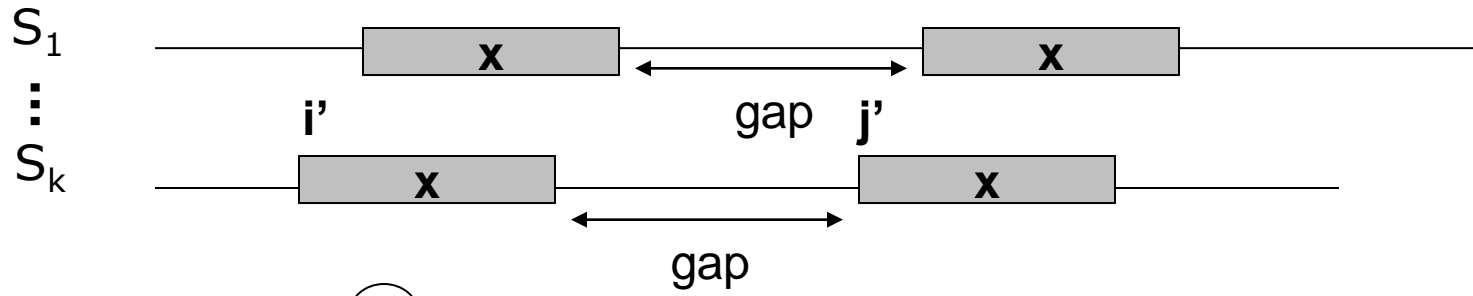
Brodal : $O(n \log n + a)$, $t_1 \leq \text{gap} \leq t_2$
 $O(n+a)$, $t_1 \leq \text{gap}$



Maximal Pairs in Multiple Strings



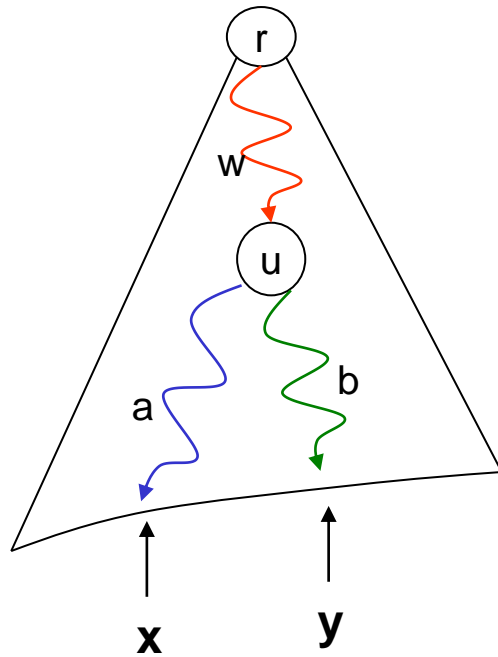
Maximal Pairs in Multiple Strings



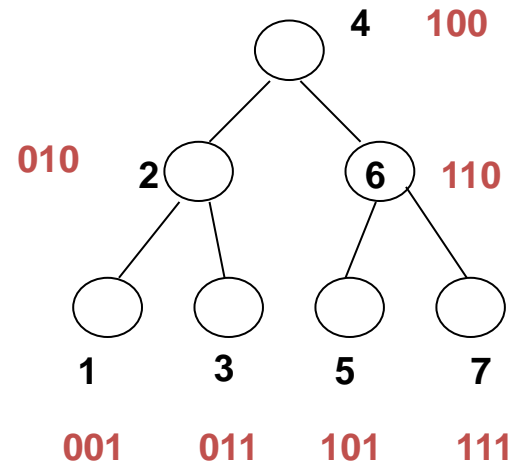
$O(n \log n + ak)$, $\text{gap} \leq t_1$
 $O(n+a)$, gap unbounded



Nearest Common Ancestor & Suffix Tree



$nca(x,y)=u$ σε χρόνο $O(1)$



$nca(\mathbf{001}, \mathbf{101}) = \text{leftmost}_1(\text{XOR}(\mathbf{001}, \mathbf{101})) = \mathbf{100} = 100$
 $nca(\mathbf{001}, \mathbf{111}) = \text{leftmost}_1(\text{XOR}(\mathbf{001}, \mathbf{111})) = \mathbf{110} = 100$
 $nca(\mathbf{011}, \mathbf{010}) = \text{leftmost}_1(\text{XOR}(\mathbf{011}, \mathbf{010})) = \mathbf{010}$



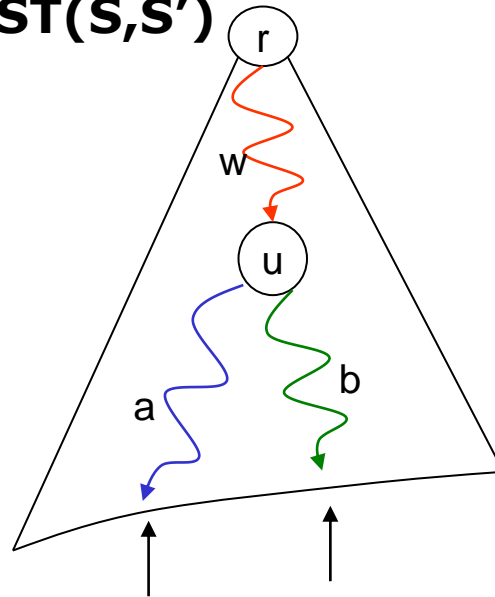
Maximal Palindromes

ávva

abactgaaccaat

taaccaagtcabaa

GST(S,S')



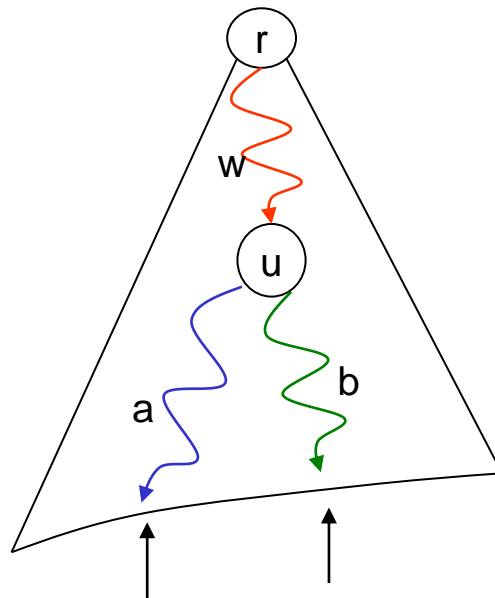
Exact Matching with wild cards

text _____

acgtttaacctttgagttgggcv

pattern * * * * * * *

a**t



Τέλος Ενότητας

Χρηματοδότηση

- Το παρόν εκπαιδευτικό υλικό έχει αναπτυχθεί στο πλαίσιο του εκπαιδευτικού έργου του διδάσκοντα.
- Το έργο «**Ανοικτά Ακαδημαϊκά Μαθήματα στο Πανεπιστήμιο Πατρών**» 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 Πανεπιστήμιο Πατρών, Μακρής Χρήστος, Περδικούρη Αικατερίνη.
«Εισαγωγή στη Βιοπληροφορική. Τεχνικές Ανάλυσης και Σύγκρισης
Ακολουθιών Βιολογικών Δεδομένων Ι». Έκδοση: 1.0. Πάτρα 2015. Όλες οι
εικόνες έχουν δημιουργηθεί από την κυρία Περδικούρη Αικατερίνη, εκτός αν
αναφέρεται διαφορετικά. Διαθέσιμο από τη δικτυακή διεύθυνση:

<https://eclass.upatras.gr/courses/CEID1047/>



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.

