



Βιοπληροφορική

Ενότητα 8:

Αναζήτηση Ομοιοτήτων σε Βάσεις Δεδομένων Ακολουθιών, 2 ΔΩ

Τμήμα: **Βιοτεχνολογίας**

Όνομα καθηγητή: **Τ. Θηραίου**



Ευρωπαϊκή Ένωση
Ευρωπαϊκό Κοινωνικό Ταμείο



ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ & ΘΡΗΣΚΕΥΜΑΤΩΝ, ΠΟΛΙΤΙΣΜΟΥ & ΑΘΛΗΤΙΣΜΟΥ
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης





Μαθησιακοί Στόχοι

- Κατανόηση της αναγκαιότητας των ευριστικών αλγορίθμων στην αναζήτηση ομοιοτήτων σε βάσεις δεδομένων ακολουθιών.
- Επεξήγηση των μεθόδων FASTA και BLAST.



Λέξεις Κλειδιά

- Λέξεις κλειδιά: Αναζήτηση ομοιοτήτων σε βάσεις δεδομένων ακολουθιών, Εφαρμογή φίλτρων.
- Key words: Sequence database similarity searching, FASTA, BLAST, Two-hit BLAST, Gapped BLAST, Filtering.



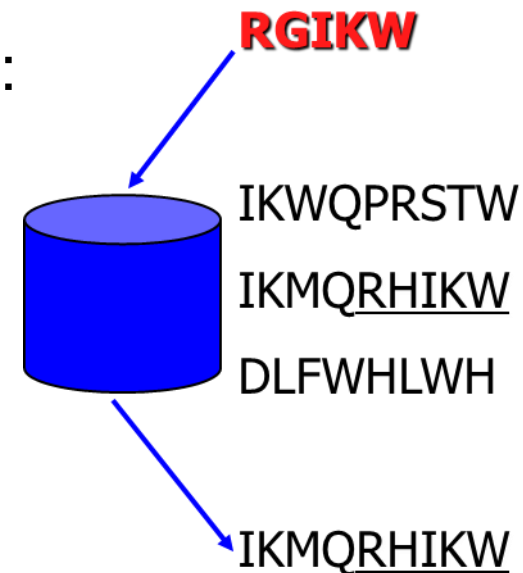
Αναζήτηση Ομοιοτήτων

● Δεδομένα.

- Ακολουθία επερώτησης (query sequence).
- Ακολουθίες στη Βάση Δεδομένων (subject sequences).

● Αναζήτηση.

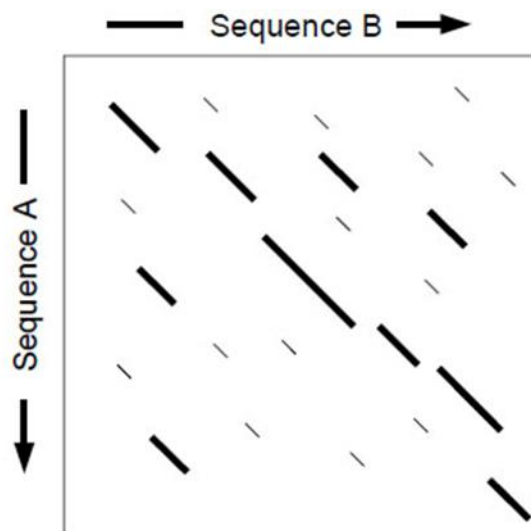
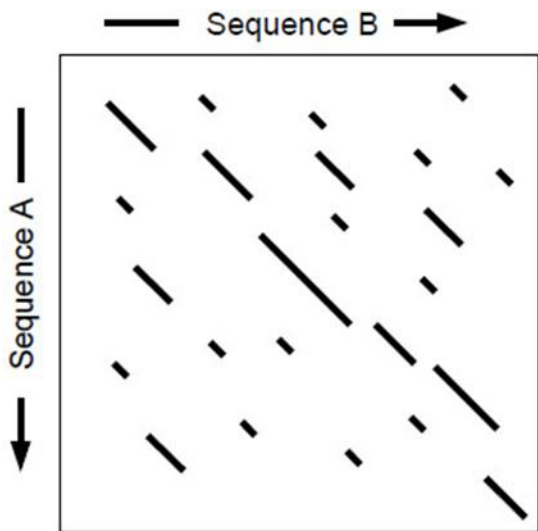
- Μέθοδοι δυναμικού προγραμματισμού:
 - Πρακτικοί μόνο για "μικρά" προβλήματα.
- Ευριστικοί Αλγόριθμοι:
 - Γρήγορη αναζήτηση.
 - Δεν εγγυώνται την βέλτιστη στοίχιση.
 - FASTA (Lipman and Pearson, 1985).
 - BLAST (Altschul et al, 1990).





FASTA 1/3

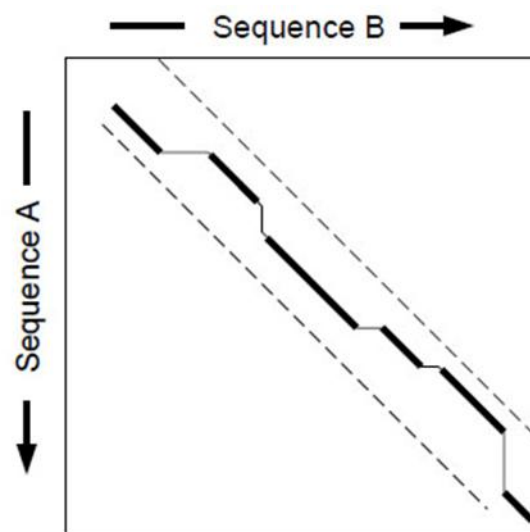
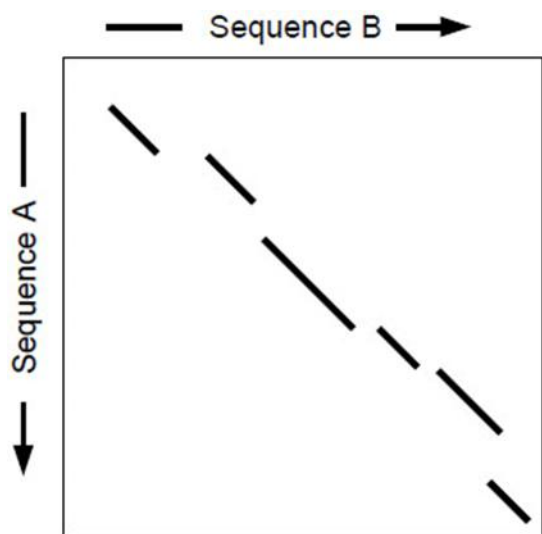
- **K-tuples (K-tup)**: λέξεις μεγέθους k .
- Για κάθε ακολουθία της ΒΔ:
 - αναγνώριση διαγωνίων ταύτισης.
 - βαθμολόγηση με τη χρήση πίνακα αντικατάστασης και επιλογή των διαγωνίων με το μεγαλύτερο score.





FASTA 2/3

- Για κάθε ακολουθία της ΒΔ:
 - Ένωση των διαγωνίων με score μεγαλύτερο ενός κατωφλίου.
 - Χρήση δυναμικού προγραμματισμού (banded Smith-Waterman) για τη βελτιστοποίηση της στοίχισης.





FASTA 3/3

- K-tuples (K-tup):
 - πρωτεΐνες: 2-tuples.
 - DNA: 6-tuples.
- η **τιμή του k** επηρεάζει την **ευαισθησία** και την **ταχύτητα** της αναζήτησης.
 - μεγάλη τιμή k:
 - μείωση λέξεων που ταυτίζονται αλλά δεν αντιστοιχούν σε πραγματικές στοιχίσεις (background word hits).
 - αύξηση ταχύτητας αλλά μείωση ευαισθησίας (sensitivity).



BLAST 1/12

Query: NKCKT**PQG**QRLVN, W=3, T=13

neighborhood words

word	score
PQG	18
PEG	15
PRG	14
PKG	14
PNG	13
PDG	13
PHG	13
PMG	13
PSG	13
PQA	12
PQN	12
...	...



BLAST 2/12

Expanded Word List

HEA
HDA
HQA
HKA
HES
HEC
HET
HEG
HEV
...

Sequence Database

```
>seq1
GTKCCEFQAKLCQFLAGKPEHPMTRETLNASHQAVRI ILEHLLEQVGFGIATAVASGAA
>seq2
FQAKLFGIATAVASEHLLGAGTKCCEAPMCQFLAGKPEQVASVRI ILEHGTRETLN
>seq3
QAKLFGFQAKLFGIAHECASEVAGKPVASRETLEAPMCQGKPVASVRI ILETRETLNEQI ILEHG
>seq4
VASVRGAGTKCCEAPMCQFLAGKPVASVRI ILETRETLNEQI ILEHGHL LFLAGKPVTL
>seq5
HGHL LFLAGKPVASRETLEAGAGTKCCEHGHL LFLAGKPVASVRI ILETAGKPVASVRI
>seq6
ILETRETLNGKPVASRETLEAPMCQEQI IHEVVAGKPVASRETLEAPMCQG
>seq7
LETHGHL LFLAGPMCQFLVRI IRI ILETAGKPVASVRIKPVASVHGHL LFLAGKPETLEARET
...
```



BLAST 3/12

```
>query
ASGDAAGVSEQTPKLAQYLADKPEHPLNRQRLDAKHEAFKIVVLQALTENCQDAPACNELRELAE
>seq1
GTKCCEFQAKLCQFLAGKPEHPMTRETLNASHQAVRIILEHLLEQVGFGIATAVASGAA
```



Align at seed word.

```
ASGDAAGVSEQTPKLAQYLADKPEHPLNRQRLDAKHEAFKIVVLQALTENCQDAPACNELRELAE
                                H+A
GTKCCEFQAKLCQFLAGKPEHPMTRETLNASHQAVRIILEHLLEQVGFGIATAVASGAA
```



Extend sequence to identify HSP.

```
ASGDAAGVSEQTPKLAQYLADKPEHPLNRQRLDAKHEAFKIVVLQALTENCQDAPACNELRELAE
      KL Q+LA KPEHP+ R+ L+A H+A +I++ L +
GTKCCEFQAKLCQFLAGKPEHPMTRETLNASHQAVRIILEHLLEQVGFGIATAVASGAA
```



Stop extension of HSP when quality of the alignment reaches a threshold value.
Calculate significance.

```
KLAQYLADKPEHPLNRQRLDAKHEAFKIVVLQALTE
KL Q+LA KPEHP+ R+ L+A H+A +I++ L +
KLCQFLAGKPEHPMTRETLNASHQAVRIILEHLLEQ
```

Score = 42.0 bits (97), Expect = 0.004
Identities = 17/36 (47%), Positives = 26/36 (72%)



BLAST 4/12

- **Words**: όλες οι δυνατές λέξεις μεγέθους W :
 - πρωτεΐνες: 3 κατάλοιπα.
 - DNA: 11 κατάλοιπα.
- Δημιουργία ευρετηρίου με όλες τις λέξεις στις ακολουθίες της ΒΔ.
- Για κάθε λέξη της ακολουθίας επερώτησης
 - εύρεση όλων των λέξεων (neighborhood words) που στοιχίζονται μαζί της με score μεγαλύτερο ενός κατωφλίου (**neighborhood score threshold, T**) χρησιμοποιώντας έναν πίνακα αντικατάστασης.



BLAST 5/12

- Αναζήτηση των neighborhood words έναντι του προκατασκευασμένου πίνακα με όλες τις λέξεις των ακολουθιών της ΒΔ.
- Επέκταση της στοίχισης μεταξύ της ακολουθίας επερώτησης και των ακολουθιών της ΒΔ για την αναγνώριση ενός **High-scoring Segment Pair (HSP)**.
- Τερματισμός της επέκτασης όταν το score γίνει μικρότερο ενός προκαθορισμένου κατωφλίου.
- **Μεγαλύτερο T:**
 - μείωση των προσπαθειών εύρεσης HSPs.
 - **αύξηση ταχύτητας αλλά μείωση ευαισθησίας.**



BLAST 6/12

● two-hit BLAST:

– Παρατηρήσεις:

- Η επέκταση γύρω από τα seed words αποτελεί το 90% του χρόνου εκτέλεσης του BLAST.
- Τα HSP έχουν μεγαλύτερο μήκος και περιέχουν πολλές λέξεις μήκους W.

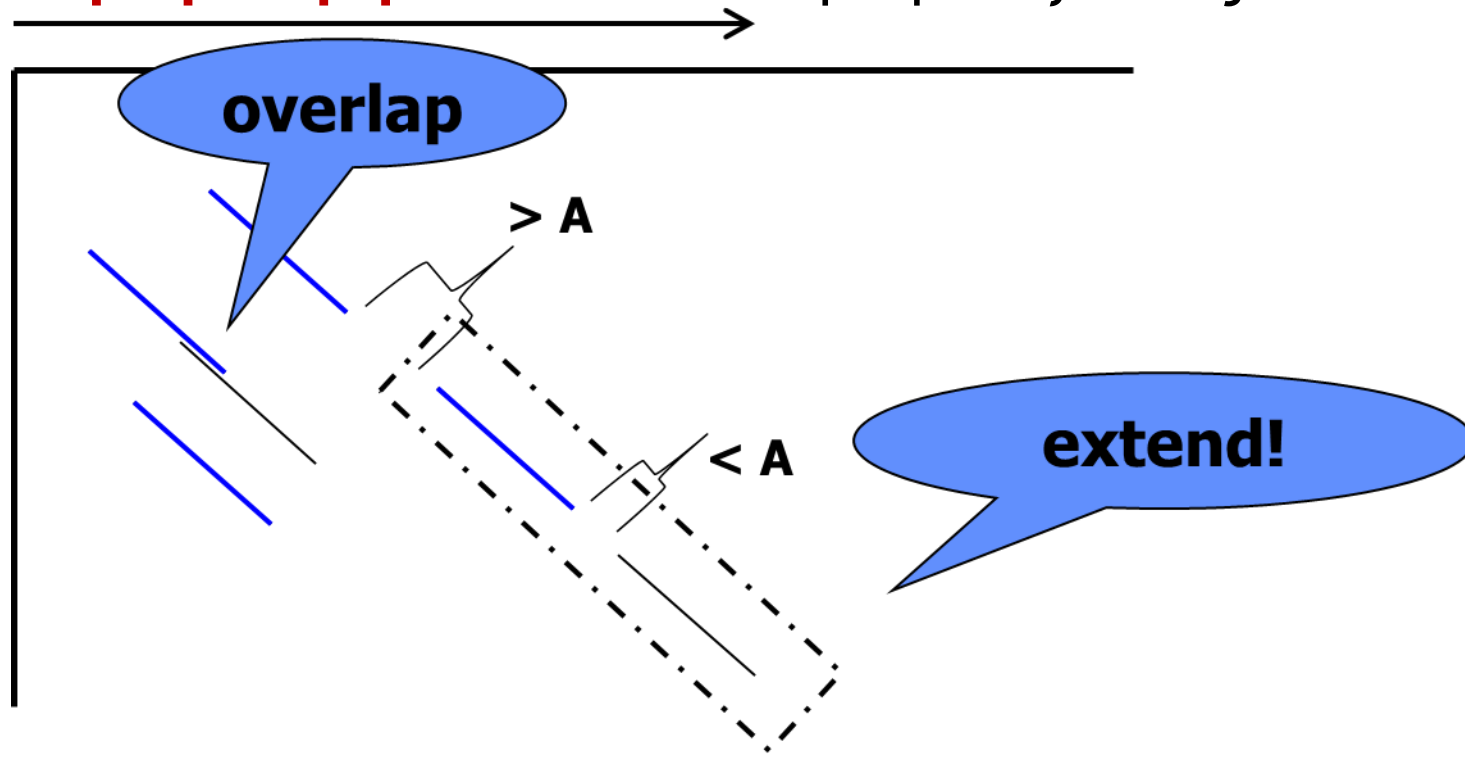
```
FLKDSAGVVDSPKLGAHAEKVFGMVRDSAVQLRATGEVV--LDGK
| : | : || : || : || | : || | ||
DLSNPGAVMGNPKVKAHGKKV-----LHSEGEGVKHLNLDNL
```



BLAST 7/12

● two-hit BLAST:

- Επέκταση μόνο όταν υπάρχουν **δύο μη επικαλυπτόμενα** ζεύγη λέξεων σε **απόσταση μικρότερη** από ένα κατώφλι μεταξύ τους.

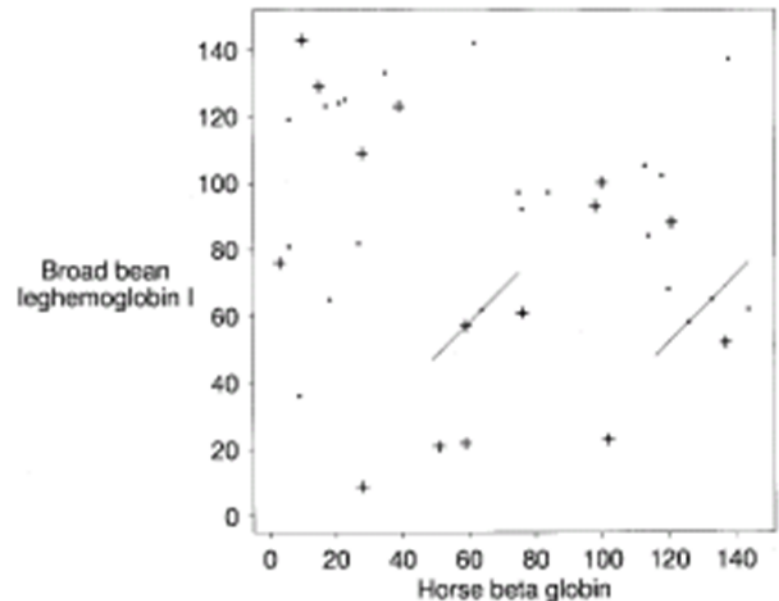




BLAST 8/12

● two-hit BLAST:

- για να διατηρηθεί η ευαισθησία της αναζήτησης, απαιτείται **μικρότερο T** (neighborhood score threshold).
 - περισσότερα hits.
 - αλλά μόνο ένα μικρό ποσοστό από αυτά σχετίζονται με δεύτερο hit.
- αύξηση της ταχύτητας.





BLAST 9/12

● Gapped BLAST:

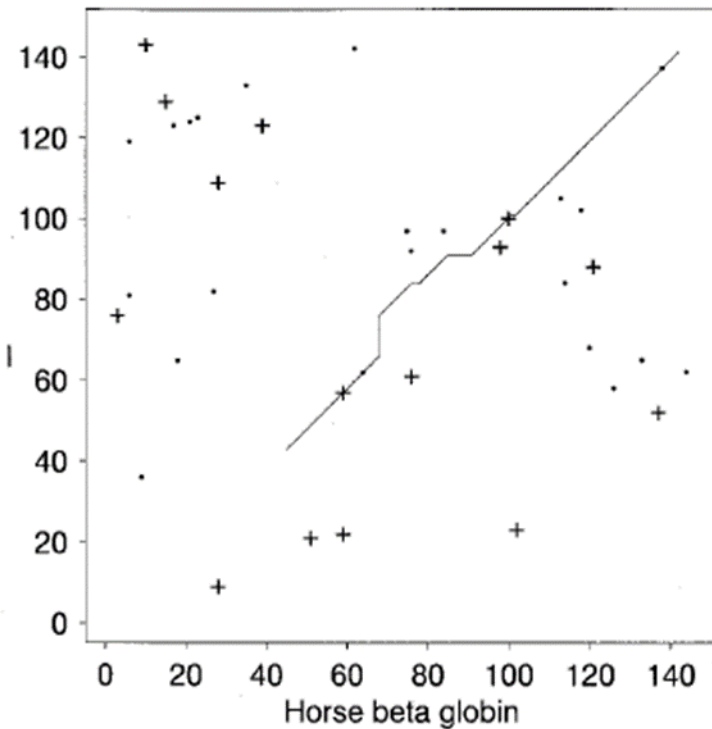
- two-hit BLAST για τη δημιουργία ενός HSP.
- επέκταση του HSP με κενά, αν έχει score μεγαλύτερο ενός κατωφλίου S_g .



BLAST 10/12

● Gapped BLAST.

Broad bean
leghemoglobin I



```

Leghemoglobin 43 FSFLKDSAGVVDSPKLGAHAEKVFGMVRDSAVQLRATGEVV--LDGKDGS----- 90
                  F L +   V+ +PK+ AH +KV                L + GE V LD  G+
Beta globin    45 FGDLSNPGAVMGNPKVKAHGKKV-----LHSFGEGVHHLNCLKGTFAALSE 90

Leghemoglobin 91 IHIQKGVLDP-HFVVVKEALKTIKEASGDKWSEELSAAWAVAYDGLATAI 140
                  +H K +DP +F ++  L+ +   G ++ EL A+++  G+A A+
Beta globin    91 LHCDKLHVDPENFRLGLGNVLVVVLARHFGKDFTPPELQASYQKVVAGVANAL 141

```




BLAST 12/12

- Εφαρμογή φίλτρων:
 - επαναλήψεις.
 - περιοχές χαμηλής πολυπλοκότητας.
 - δημιουργία στατιστικώς σημαντικών στοιχίσεων, αλλά χωρίς βιολογικό νόημα.
 - **Soft filtering:**
 - φιλτράρισμα μόνο στη φάση αναζήτησης.
 - **Hard filtering:**
 - φιλτράρισμα στη φάση αναζήτησης και στη φάση τελικής στοίχισης.



Servers

● FASTA:

- http://fasta.bioch.virginia.edu/fasta_www2/fasta_list2.shtml
- <http://www.ebi.ac.uk/Tools/sss/fasta/>

● BLAST:

- <http://blast.ncbi.nlm.nih.gov/>
- <http://www.ebi.ac.uk/Tools/sss/ncbiblast/>

● Sequence Similarity Searching:

- <http://www.ebi.ac.uk/Tools/sss/>



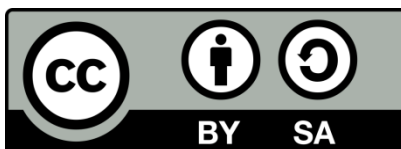
Βιβλιογραφία

- David Mount, "Bioinformatics: Sequence and Genome Analysis", Cold Spring Harbor Laboratory Press; 2nd edition (March 12, 2013).
- Jonathan Pevsner, "Bioinformatics and Functional Genomics", Wiley-Blackwell; 2nd edition (May 4, 2009).
- Andreas D. Baxevanis, B. F. Francis Ouellette, "Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins", Wiley-Interscience; 3rd edition (October 29, 2004).



Άδειες Χρήσης

- Το παρόν εκπαιδευτικό υλικό υπόκειται σε άδειες χρήσης Creative Commons.
- Για εκπαιδευτικό υλικό, όπως εικόνες, που υπόκειται σε άλλου τύπου άδεια χρήσης, η άδεια χρήσης αναφέρεται ρητώς.





Χρηματοδότηση

- Το παρόν εκπαιδευτικό υλικό έχει αναπτυχθεί στο πλαίσιο του εκπαιδευτικού έργου του διδάσκοντα.
- Το έργο «**Ανοικτά Ακαδημαϊκά Μαθήματα Γεωπονικού Πανεπιστημίου 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.



Ευρωπαϊκή Ένωση
Ευρωπαϊκό Κοινωνικό Ταμείο



ΕΠΙΧΕΙΡΗΣΙΑΚΟ ΠΡΟΓΡΑΜΜΑ
ΕΚΠΑΙΔΕΥΣΗ ΚΑΙ ΔΙΑ ΒΙΟΥ ΜΑΘΗΣΗ
επένδυση στην κοινωνία της γνώσης
ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ ΚΑΙ ΘΡΗΣΚΕΥΜΑΤΩΝ
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



ΕΣΠΑ
2007-2013
πρόγραμμα για την ανάπτυξη
ΕΥΡΩΠΑΪΚΟ ΚΟΙΝΩΝΙΚΟ ΤΑΜΕΙΟ



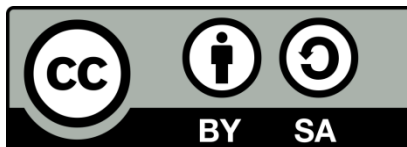
Σημείωμα Αναφοράς

Copyright Γεωπονικό Πανεπιστήμιο Αθηνών 2015. Τμήμα Βιοτεχνολογίας, Θηραίου Τριάς. «Βιοπληροφορική». Έκδοση: 1.0. Αθήνα 2015. Διαθέσιμο από τη δικτυακή διεύθυνση:
<https://mediasrv.aua.gr/eclass/courses/OCDB100/>



Σημείωμα Αδειοδότησης

Το παρόν υλικό διατίθεται με τους όρους της άδειας χρήσης Creative Commons Αναφορά, 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».



Η άδεια αυτή ανήκει στις άδειες που ακολουθούν τις προδιαγραφές του Ορισμού Ανοικτής Γνώσης [2], είναι ανοικτό πολιτιστικό έργο [3] και για το λόγο αυτό αποτελεί ανοικτό περιεχόμενο [4].

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

[2] <http://opendefinition.org/okd/ellinika/>

[3] <http://freedomdefined.org/Definition/EI>

[4] <http://opendefinition.org/buttons/>



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.