

BIOS 477/877 Bioinformatics and Molecular Evolution

Instructor: Etsuko Moriyama
(School of Biological Sciences)

Spring 2026 Lecture 18

BIOS477/877 L18 - 1

1

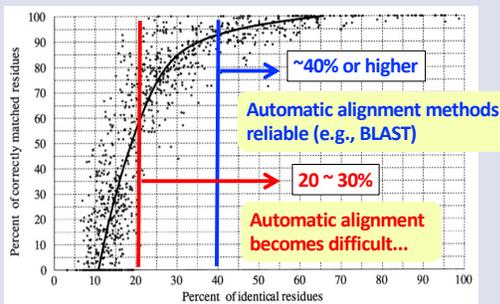
Today's topics

- Pattern and profile (PSSM)
- Hidden Markov model and profile HMM
- PSI-BLAST

BIOS477/877 L18 - 2

2

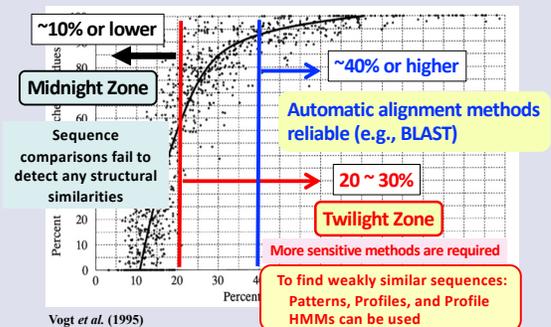
Sequence similarity and search sensitivity



BIOS477/877 L18 - 3

3

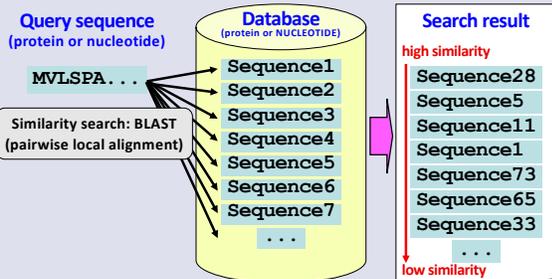
Sequence similarity and search sensitivity



BIOS477/877 L18 - 4

4

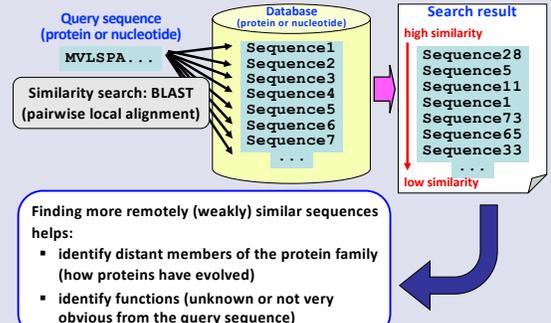
Similarity search



BIOS477/877 L18 - 5

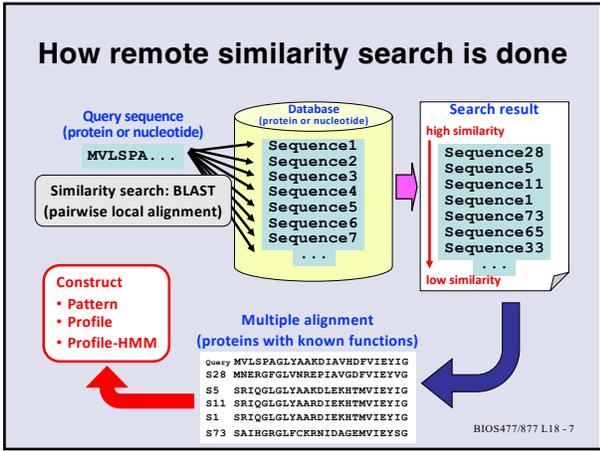
5

Finding remote similarity

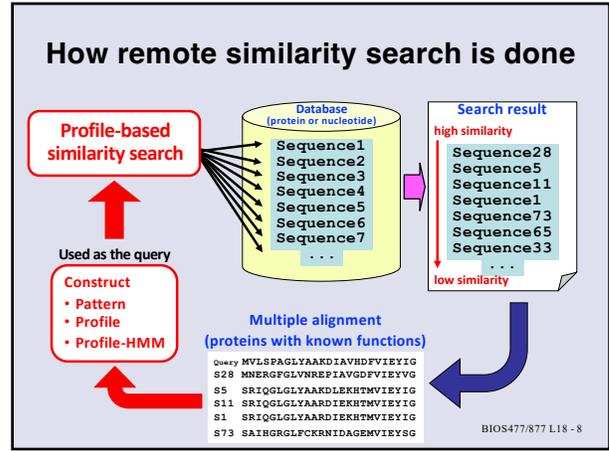


BIOS477/877 L18 - 6

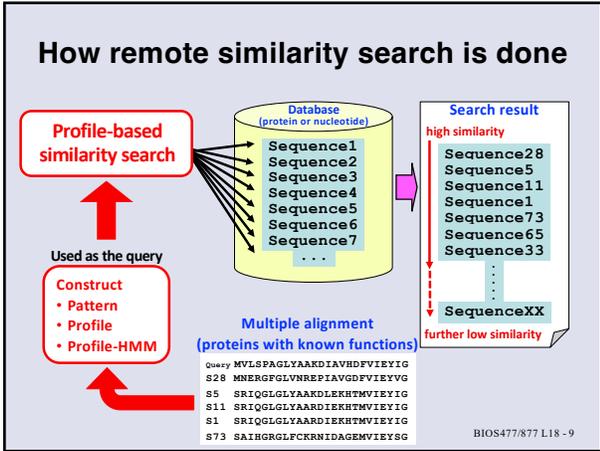
6



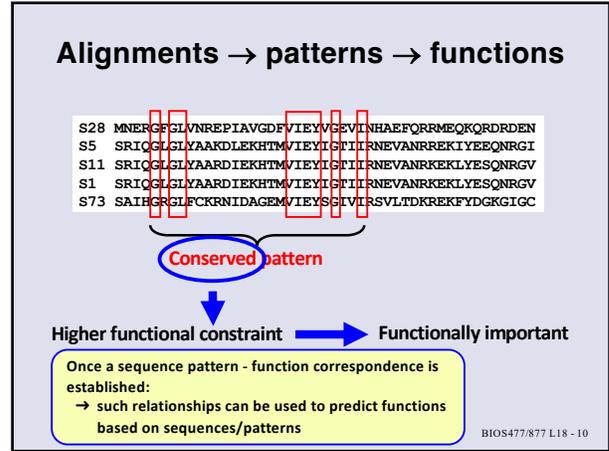
7



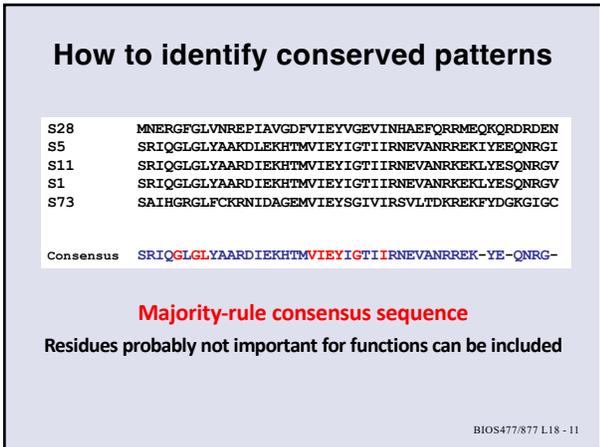
8



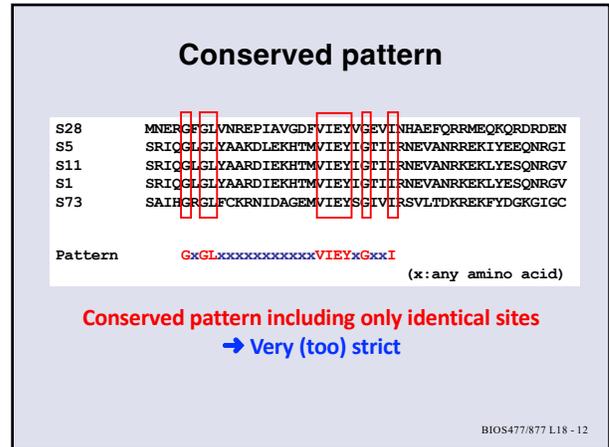
9



10



11



12

Regular expression pattern

```

S28 MNERGEGSLVNREPIAVGDFVIEYVGEVINHAEFQRMEQKQRDRDEN
S5  SRICQLGLIYAAKDLEKHTMVEYIGTIIIRNEVANRREKLYEEQNRGI
S11 SRICQLGLIYAAKDIEKHTMVEYIGTIIIRNEVANRREKLYEYQNRGV
S1  SRICQLGLIYAAKDIEKHTMVEYIGTIIIRNEVANRREKLYEYQNRGV
S73 SAIHGRGLFCKRNIDAGEMVIEYSGIVLRISVLTDKREKFDYDGGKIGC
  
```

G-[FLR]-G-L-X10-[FM]-V-I-E-Y-[VIS]-G-[ETI]-[VI]-I
 (10 any amino acids)

Regular expression
 → More flexible than strict conserved pattern

BIOS477/877 L18 - 13

13

PROSITE:

Database of protein domains, families and functional sites



<https://prosite.expasy.org/>

- PROSITE Documentation (entry)
 - Biological information about protein families, domains, etc.
 - PROSITE documentation PDOC00210
 - G-protein coupled receptors family 1 signature and profile
- Two types of motif descriptors
 - **Patterns (regular expressions)**
 G_PROTEIN_RECEP_F1_1, PS00237; G-protein coupled receptors family 1 signature (PATTERN)
 - **Profiles (position-specific scoring matrix)**
 G_PROTEIN_RECEP_F1_2, PS50262; G-protein coupled receptors family 1 profile (MATRIX)

BIOS477/877 L18 - 14

14

PROSITE pattern

PROSITE: PS00237 (G-protein coupled receptors family 1 signature)

[GSTALIVMFYWC]-[GSTANCPDE]-[EDPKRH]-x(2)-[LIVMNOGA]-x(2)-[LIVMFT]-[GSTANC]-[LIVMFYWSTAC]-[DENH]-R-[FYWCSH]-x(2)-[LIVM]

SH1A_FUGRU/131-147	SSIIh	caIALR	rwaI
SH1A_HUMAN/122-138	SSIIh	caIALR	rwaI
SH1A_MOUSE/122-138	SSIIh	caIALR	rwaI
SH1A_PANTR/122-138	SSIIh	caIALR	rwaI
SH1B_CAVPO/134-150	ASImh	cvIALR	rwaI
SH1B_CRIGR/131-147	ASImh	cvIALR	rwaI
SH1B_DIDMA/134-150	ASImh	cvIALR	rwaI
SH1B_FUGRU/119-135	SSIIh	cvIALR	rwaI
SH1B_HUMAN/135-151	ASImh	cvIALR	rwaI
SH1B_MOUSE/131-147	ASImh	cvIALR	rwaI
SH1B_RABIT/135-151	ASImh	cvIALR	rwaI
SH1B_RAT/131-147	ASImh	cvIALR	rwaI
SH1B_SPAEH/131-147	ASImh	cvIALR	rwaI
SH1D_CANFA/124-140	ASIIh	cvIALR	rwaI
SH1D_CAVPO/124-140	ASIIh	cvIALR	rwaI
SH1D_FUGRU/122-138	ASIIh	cvIALR	rwaI
SH1D_HUMAN/124-140	ASIIh	cvIALR	rwaI
SH1D_MOUSE/121-137	ASIIh	cvIALR	rwaI
SH1D_PIG/44-60	ASIIh	cvIALR	rwaI
SH1D_RABIT/124-140	ASIIh	cvIALR	rwaI
SH1D_RAT/121-137	ASIIh	cvIALR	rwaI
SH1E_HUMAN/108-124	CSIIh	cvIALR	rwaI
SH1E_PANTR/108-124	CSIIh	cvIALR	rwaI
SH1E_PTC/55-71	CSIIh	cvIALR	rwaI

PROSITE pattern syntax: <https://prosite.expasy.org/prosuser.html#meth1>

BIOS477/877 L18 - 15

15

PROSITE pattern

PROSITE: PS00237 (G-protein coupled receptors family 1 signature)

[GSTALIVMFYWC]-[GSTANCPDE]-[EDPKRH]-x(2)-[LIVMNOGA]-x(2)-[LIVMFT]-[GSTANC]-[LIVMFYWSTAC]-[DENH]-R-[FYWCSH]-x(2)-[LIVM]

SH1A_FUGRU/131-147	SSIIh	caIALR	rwaI
SH1A_HUMAN/122-138	SSIIh	caIALR	rwaI
SH1A_MOUSE/122-138	SSIIh	caIALR	rwaI
SH1A_PANTR/122-138	SSIIh	caIALR	rwaI
SH1B_CAVPO/134-150	ASImh	cvIALR	rwaI
SH1B_CRIGR/131-147	ASImh	cvIALR	rwaI
SH1B_DIDMA/134-150	ASImh	cvIALR	rwaI
SH1B_FUGRU/119-135	SSIIh	cvIALR	rwaI
SH1B_HUMAN/135-151	ASImh	cvIALR	rwaI
SH1B_MOUSE/131-147	ASImh	cvIALR	rwaI
SH1B_RABIT/135-151	ASImh	cvIALR	rwaI
SH1B_RAT/131-147	ASImh	cvIALR	rwaI
SH1B_SPAEH/131-147	ASImh	cvIALR	rwaI
SH1D_CANFA/124-140	ASIIh	cvIALR	rwaI
SH1D_CAVPO/124-140	ASIIh	cvIALR	rwaI
SH1D_FUGRU/122-138	ASIIh	cvIALR	rwaI
SH1D_HUMAN/124-140	ASIIh	cvIALR	rwaI
SH1D_MOUSE/121-137	ASIIh	cvIALR	rwaI
SH1D_PIG/44-60	ASIIh	cvIALR	rwaI
SH1D_RABIT/124-140	ASIIh	cvIALR	rwaI
SH1D_RAT/121-137	ASIIh	cvIALR	rwaI
SH1E_HUMAN/108-124	CSIIh	cvIALR	rwaI
SH1E_PANTR/108-124	CSIIh	cvIALR	rwaI
SH1E_PTC/55-71	CSIIh	cvIALR	rwaI

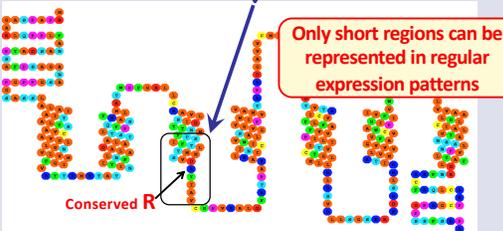
BIOS477/877 L18 - 16

16

PROSITE pattern

PROSITE: PS00237 (G-protein coupled receptors family 1 signature)

[GSTALIVMFYWC]-[GSTANCPDE]-[EDPKRH]-x(2)-[LIVMNOGA]-x(2)-[LIVMFT]-[GSTANC]-[LIVMFYWSTAC]-[DENH]-R-[FYWCSH]-x(2)-[LIVM]



Only short regions can be represented in regular expression patterns

Conserved R

[OPRD_HUMAN]

BIOS477/877 L18 - 17

17

Profile

- **Position Specific Scoring Matrix (PSSM)**
 - Constructed from multiple alignments
 - Short conserved domains (BLOCKS, PRINTS)
 - Protein families (PROSITE)
 - Results of similarity search (PSI-BLAST)
- More flexible than simple patterns
 - PSSM lists amino acid frequencies for each alignment position
- Profiles (PSSMs) can be used for database search to identify remote similarities

BIOS477/877 L18 - 18

18

How to build a profile: (a simple method)

➤ **EMBOSS "Protein Profile" tools** EMBOSS_website

- PROPHECY: creates profiles from multiple alignment [option]
- Simple amino acid "frequency"

```
# Columns are amino acid counts A->Z
# Rows are alignment positions 1->n
Simple
Name          mymatrix
Length        5
Maximum score  11
Thresh        75
Consensus     RCEGH
1 0 0 0 0 0 0 0 0 1 0 ... 0 ... 0 2 ...
2 0 0 3 0 0 0 0 0 0 0 ... 0 ... 0 0 ...
3 0 0 0 0 2 0 0 0 0 0 ... 0 ... 0 1 0 ...
4 1 0 0 0 0 0 0 2 0 0 ... 0 ... 0 0 ...
5 0 0 0 0 0 0 0 0 2 0 ... 1 ... 0 0 ...
A B C D E F G H I ... N ... Q R ...
```

Seq1 RCQAH
Seq2 HCEGH
Seq3 RCEGN

BIOS477/877 L18 - 19

19

How to build a profile: (a simple method)

➤ **EMBOSS "Protein Profile" tools** EMBOSS_website

- PROPHECY: creates profiles from multiple alignment [option]
- Simple amino acid "frequency"

```
# Columns are amino acid counts A->Z
# Rows are alignment positions 1->n
Simple
Name          mymatrix
Length        5
Maximum score  11
Thresh        75
Consensus     RCEGH
1 0 0 0 0 0 0 0 0 1 0 ... 0 ... 0 2 ...
2 0 0 3 0 0 0 0 0 0 0 ... 0 ... 0 0 ...
3 0 0 0 0 2 0 0 0 0 0 ... 0 ... 0 1 0 ...
4 1 0 0 0 0 0 0 2 0 0 ... 0 ... 0 0 ...
5 0 0 0 0 0 0 0 0 2 0 ... 1 ... 0 0 ...
A B C D E F G H I ... N ... Q R ...
```

Seq1 RCQAH
Seq2 HCEGH
Seq3 RCEGN

Lots of 0s!
Do it cause any problem?

BIOS477/877 L18 - 20

20

Over-fitting problem in profile

Seq1 RDA
Seq2 REA
Seq3 REG

↓

A simple profile
(frequency)

Pos	A	R	N	D	C	Q	E	G	H	...	K	T	...
1	0	1.0	0	0	0	0	0	0	0	...	0	0	...
2	0	0	0	0.3	0	0	0.6	0	0	...	0	0	...
3	0.6	0	0	0	0	0	0	0.3	0	...	0	0	...

Lots of 0s!

Position 1 has only Arg (R)
Position 2 has only Asp (D) & Glu (E)
Position 3 has only Ala (A) & Gly (G)

Can we find **KEA**?

BIOS477/877 L18 - 21

21

Over-fitting problem in profile

Seq1 RDA
Seq2 REA
Seq3 REG

↓

A simple profile
(frequency)

Pos	A	R	N	D	C	Q	E	G	H	...	K	T	...
1	0	1.0	0	0	0	0	0	0	0	...	0	0	...
2	0	0	0	0.3	0	0	0.6	0	0	...	0	0	...
3	0.6	0	0	0	0	0	0	0.3	0	...	0	0	...

Position 1 has only Arg (R)
Position 2 has only Asp (D) & Glu (E)
Position 3 has only Ala (A) & Gly (G)

Can we find **KEA**? No flexibility!
How about **REA**?

0 should be avoided → Use something else instead

BIOS477/877 L18 - 22

22

Pseudocount methods

➤ **Simplest method**

- Add a small constant (C) to all the counts:
 $W(b) = \{n(b) + C\} / (N + C)$

where $W(b)$: Frequency of amino acid b
 $n(b)$: Number of amino acid b
 N : Number of sequences

- Pseudocount (C): something small but not 0 (e.g., $C = 0.1$)

Without pseudocount:
 $W(b) = n(b) / N$

➤ **Substitution matrix dependent**

- Proportional to scores, S_{ij}

➤ **Dirichlet mixtures** (Sjölander *et al.* 1996)

- Mixture of different types of pseudocounts
- Representing various context of protein sequences (e.g., loop region, hydrophobic region)
- Used in **profile HMM**

BIOS477/877 L18 - 23

23

Gribskov profile (PSSM)

Gribskov *et al.* (1987)

- Used in **PROSITE profile**
- Weighted scoring matrix for each position

→ Generated from:

- Multiple alignment
- Scoring matrix (e.g., PAM250)
- Scoring matrix is weighted with amino acid frequencies at each position

$$M(p,a) = \sum_{b=1}^{20} W(p,b) * Y(a,b)$$

$W(p,b)$: Frequency of amino acid b in position p
 $Y(a,b)$: Value in the scoring matrix for AA pair a and b

BIOS477/877 L18 - 24

24

Markov models: A simple example (diagram)

▪ Coin tossing (fair coin)

At time t+1

	H	T
H	0.5	0.5
T	0.5	0.5

At time t

• States: **Head or Tail**
• Transition probability

BIOS477/877 L18 - 31

31

Markov models: A simple example (diagram)

▪ Coin tossing (loaded coin: more Tail than Head)

At time t+1

	H	T
H	0.4	0.6
T	0.4	0.6

At time t

• States: **Head or Tail**
• Transition probability

BIOS477/877 L18 - 32

32

Markov models: A simple example (diagram)

▪ Coin tossing (general model)

At time t+1

	H	T
H	1-a	a
T	b	1-b

At time t

• States: **Head or Tail**
• Transition probability

BIOS477/877 L18 - 33

33

Markov models: State paths

▪ Loaded coin tossing [states: H and T]

Transition probability matrix

$$P_L = \begin{matrix} & \begin{matrix} H & T \end{matrix} \\ \begin{matrix} H \\ T \end{matrix} & \begin{bmatrix} 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix} \end{matrix}$$

Prob(H @ 2nd Toss) = 0.16 + 0.24 = 0.4

BIOS477/877 L18 - 34

34

Markov models: State paths

▪ Loaded coin tossing [states: H and T]

Transition probability matrix

$$P_L = \begin{matrix} & \begin{matrix} H & T \end{matrix} \\ \begin{matrix} H \\ T \end{matrix} & \begin{bmatrix} 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix} \end{matrix}$$

Prob(H @ 2nd Toss) = 0.16 + 0.24 = 0.4

BIOS477/877 L18 - 35

35

Markov models: State paths

▪ Loaded coin tossing [states: H and T]

Transition probability matrix

$$P_L = \begin{matrix} & \begin{matrix} H & T \end{matrix} \\ \begin{matrix} H \\ T \end{matrix} & \begin{bmatrix} 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix} \end{matrix}$$

Prob(T @ 2nd Toss) = 0.24 + 0.36 = 0.6

BIOS477/877 L18 - 36

36

Markov models: State paths

▪ Loaded coin tossing [states: H and T]

Transition probability matrix

$$P_L = \begin{matrix} & \begin{matrix} H & T \end{matrix} \\ \begin{matrix} H \\ T \end{matrix} & \begin{bmatrix} 0.4 & 0.6 \\ 0.4 & 0.6 \end{bmatrix} \end{matrix}$$

Initial probability

$$P_0 = \begin{matrix} H & T \\ 0.4 & 0.6 \end{matrix}$$

Probabilities for 1st Toss:

- Prob(HH) = $0.4 \times 0.4 = 0.16$
- Prob(HT) = $0.4 \times 0.6 = 0.24$
- Prob(TH) = $0.6 \times 0.4 = 0.24$
- Prob(TT) = $0.6 \times 0.6 = 0.36$

Equation: $P_1 = P_0 \times P_L$

BIOS477/877 L18 - 37

37

Markov models: State paths

▪ Loaded coin tossing [states: H and T]

Initial probability:

$$P_0 = \begin{matrix} H & T \\ 0.4 & 0.6 \end{matrix}$$

Transition probability:

$$P_L = \begin{matrix} H & T \\ 0.4 & 0.6 \end{matrix}$$

Equation: $P_n = P_0 \times (P_L)^n$

P_n : Probability of the state H or T after n times of tossing

BIOS477/877 L18 - 38

38

What is hidden Markov model?

▪ Fair vs. loaded coin [states: H and T]

Transition matrices:

$$P_{F0} = \begin{matrix} H & T \\ 0.5 & 0.5 \end{matrix}$$

$$P_F = \begin{matrix} H & T \\ 0.5 & 0.5 \end{matrix}$$

$$P_{L0} = \begin{matrix} H & T \\ 0.4 & 0.6 \end{matrix}$$

$$P_L = \begin{matrix} H & T \\ 0.4 & 0.6 \end{matrix}$$

BIOS477/877 L18 - 39

39

What is hidden Markov model?

▪ Fair vs. loaded coin [states: H and T]

Transition matrices:

$$P_{F0} = \begin{matrix} H & T \\ 0.5 & 0.5 \end{matrix}$$

$$P_F = \begin{matrix} H & T \\ 0.5 & 0.5 \end{matrix}$$

$$P_{L0} = \begin{matrix} H & T \\ 0.5 & 0.5 \end{matrix}$$

$$P_L = \begin{matrix} H & T \\ 0.4 & 0.6 \end{matrix}$$

Which coin do we have?

Try tossing the coin several times. Observe what happens!

BIOS477/877 L18 - 40

40

Hidden Markov model (observed states)

▪ Fair vs. loaded coin [states: H and T]

Transition matrices:

$$P_{F0} = \begin{matrix} H & T \\ 0.5 & 0.5 \end{matrix}$$

$$P_F = \begin{matrix} H & T \\ 0.5 & 0.5 \end{matrix}$$

$$P_{L0} = \begin{matrix} H & T \\ 0.5 & 0.5 \end{matrix}$$

$$P_L = \begin{matrix} H & T \\ 0.4 & 0.6 \end{matrix}$$

Which coin do we have?

Observation: H, H, H

$P(\text{HHH} | \text{Fair coin}) = ?$ $P(\text{HHH} | \text{Loaded coin}) = ?$

($P_{F0}=P_{L0}$, so we can ignore the initial probability)

BIOS477/877 L18 - 41

41

Hidden Markov model (observed states)

▪ Fair vs. loaded coin [states: H and T]

Transition matrices:

$$P_{F0} = \begin{matrix} H & T \\ 0.5 & 0.5 \end{matrix}$$

$$P_F = \begin{matrix} H & T \\ 0.5 & 0.5 \end{matrix}$$

$$P_{L0} = \begin{matrix} H & T \\ 0.5 & 0.5 \end{matrix}$$

$$P_L = \begin{matrix} H & T \\ 0.4 & 0.6 \end{matrix}$$

Which coin do we have?

Observation: H, H, H

$P(\text{HHH} | \text{Fair coin}) > P(\text{HHH} | \text{Loaded coin})$

Based on the observed sequence: HHH → Fair coin is more likely

BIOS477/877 L18 - 44

44

Hidden Markov model (observed vs. hidden states)

▪ If fair and loaded coins are mixed!

$P_{F0} = \begin{bmatrix} H & T \\ 0.5 & 0.5 \end{bmatrix}$
 $P_F = \begin{bmatrix} H & T \\ 0.5 & 0.5 \end{bmatrix}$
 $P_{L0} = \begin{bmatrix} H & T \\ 0.5 & 0.5 \end{bmatrix}$
 $P_L = \begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$

We have to consider these probabilities, too!

If the coin can be changed?

For observed sequence: H, H, H;
 $P(H_1H_2H_3)$ or $P(H_1H_2H_3)$ or $P(H_1H_2H_3)$ or $P(H_1H_2H_3)$ or ... ?

BIOS477/877 L18 - 45

45

Hidden Markov model (observed vs. hidden states)

▪ Mixed coins: fair and loaded

States: Fair and Loaded

Transition probability matrix $\rightarrow \begin{bmatrix} F & L \\ F & 0.9 & 0.1 \\ L & 0.1 & 0.9 \end{bmatrix}$

Each state emits symbols H and T with certain emission probabilities

BIOS477/877 L18 - 46

46

Hidden Markov model (hidden state sequence)

➤ State sequence is unknown (hidden): Fair or Loaded

Observation: H T H T T T

Hidden state: ForL ForL ForL ForL ForL ForL

Can we guess the hidden state sequence?

What we know

- Transition probabilities: $\begin{bmatrix} F & L \\ F & 0.9 & 0.1 \\ L & 0.1 & 0.9 \end{bmatrix}$
- Initial probabilities: (0.5, 0.5)
- Emission probabilities: Fair: $\begin{bmatrix} H & T \\ 0.5 & 0.5 \end{bmatrix}$; Loaded: $\begin{bmatrix} H & T \\ 0.4 & 0.6 \end{bmatrix}$

BIOS477/877 L18 - 47

47

Hidden Markov model (hidden state sequence)

➤ State sequence is unknown (hidden): Fair or Loaded

Observation: H T H T T T

Fair: $\begin{matrix} (emission) & 0.5 & 0.5 & 0.5 & 0.5 & 0.5 & 0.5 \\ F_H & F_T & F_H & F_T & F_T & F_T & F_T \end{matrix}$

Loaded: $\begin{matrix} L_H & L_T & L_H & L_T & L_T & L_T & L_T \\ (emission) & 0.4 & 0.6 & 0.4 & 0.6 & 0.6 & 0.6 \end{matrix}$

Probability of this path? $Prob=0.0046$

[Prob(path)] = [initial prob. X emission prob. X transition prob.]

Prob = $0.5 \times 0.5 \times 0.9 \times 0.5 \times 0.9 \times 0.5 \times 0.9 \times 0.5 \times 0.9 \times 0.5$
 $= (0.5)^7 \times (0.9)^5 = 0.0046$

BIOS477/877 L18 - 48

48

Hidden Markov model (hidden state sequence)

➤ State sequence is unknown (hidden): Fair or Loaded

Observation: H T H T T T

Fair: $\begin{matrix} (emission) & 0.5 & 0.5 & 0.5 & 0.5 & 0.5 & 0.5 \\ F_H & F_T & F_H & F_T & F_T & F_T & F_T \end{matrix}$

Loaded: $\begin{matrix} L_H & L_T & L_H & L_T & L_T & L_T & L_T \\ (emission) & 0.4 & 0.6 & 0.4 & 0.6 & 0.6 & 0.6 \end{matrix}$

Probability of this path? $Prob=0.0046$

[Prob(path)] = [initial prob. X emission prob. X transition prob.]

Prob = $0.5 \times 0.4 \times 0.9 \times 0.6 \times 0.9 \times 0.4 \times 0.9 \times 0.6 \times 0.9 \times 0.6 \times 0.9 \times 0.6$
 $= 0.5 \times (0.4)^2 \times (0.6)^3 \times (0.9)^5 = 0.0061$

BIOS477/877 L18 - 49

49

Hidden Markov model (hidden state sequence)

➤ State sequence is unknown (hidden): Fair or Loaded

Observation: H T H T T T

Fair: $\begin{matrix} (emission) & 0.5 & 0.5 & 0.5 & 0.5 & 0.5 & 0.5 \\ F_H & F_T & F_H & F_T & F_T & F_T & F_T \end{matrix}$

Loaded: $\begin{matrix} L_H & L_T & L_H & L_T & L_T & L_T & L_T \\ (emission) & 0.4 & 0.6 & 0.4 & 0.6 & 0.6 & 0.6 \end{matrix}$

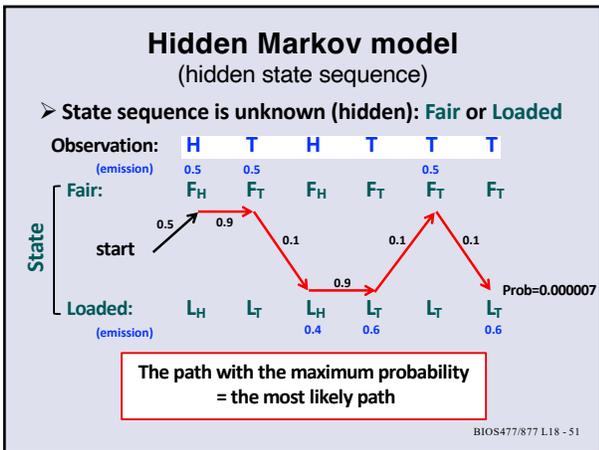
Probability of this path? $Prob=0.000007$

[Prob(path)] = [initial prob. X emission prob. X transition prob.]

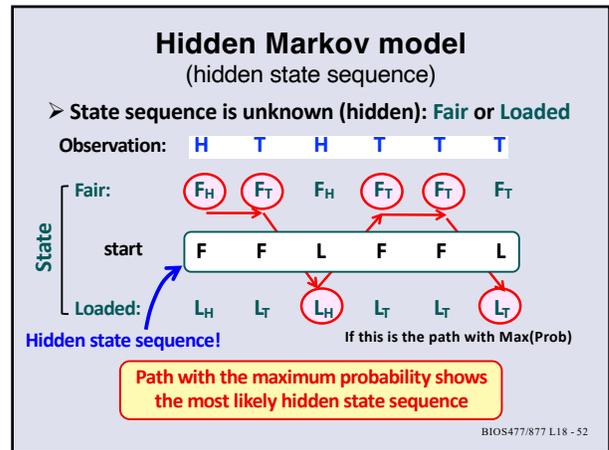
Prob = $0.5 \times 0.5 \times 0.9 \times 0.5 \times 0.1 \times 0.4 \times 0.9 \times 0.6 \times 0.1 \times 0.5 \times 0.1 \times 0.6$
 $= (0.5)^7 \times (0.9)^2 \times (0.1)^3 \times 0.4 \times (0.6)^2 = 0.000007$

BIOS477/877 L18 - 50

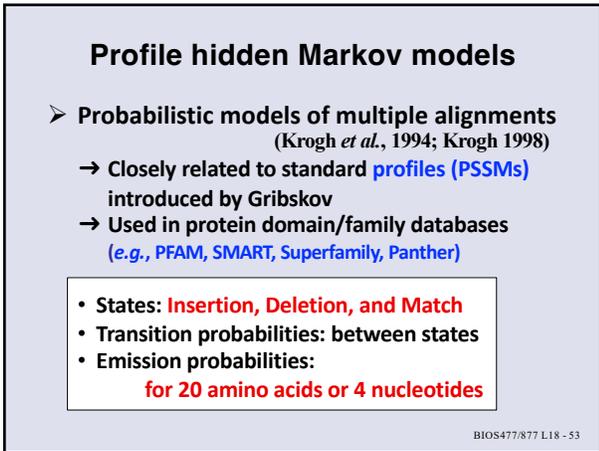
50



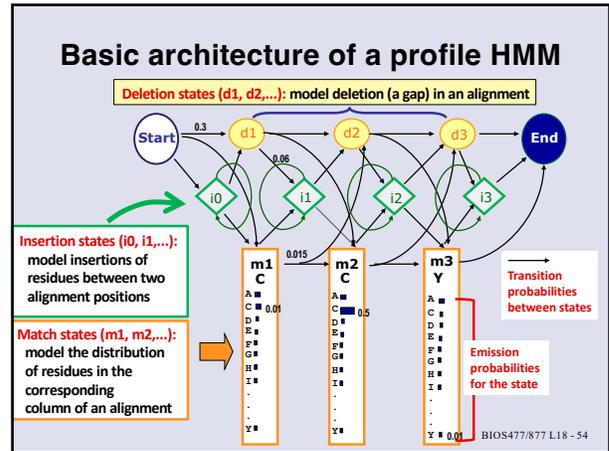
51



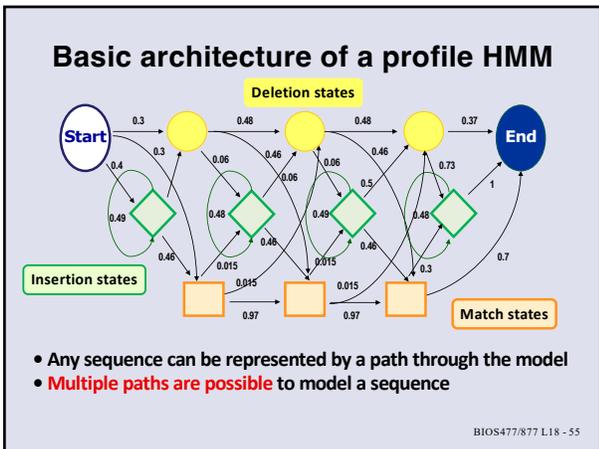
52



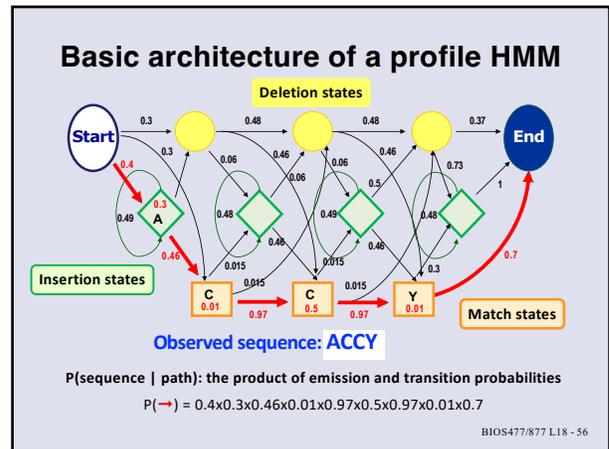
53



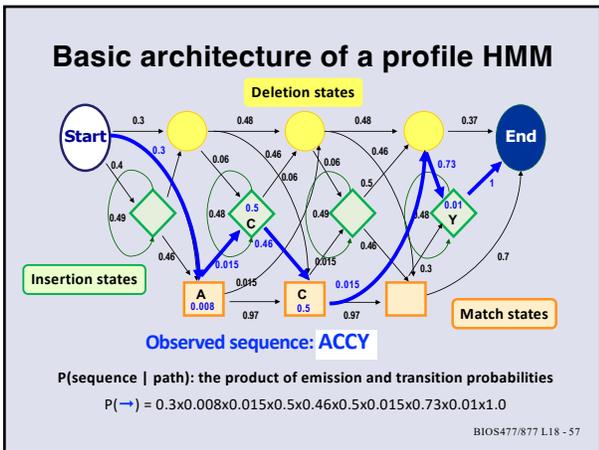
54



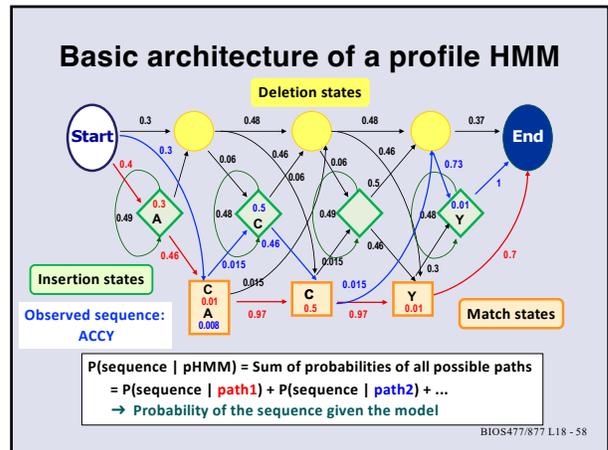
55



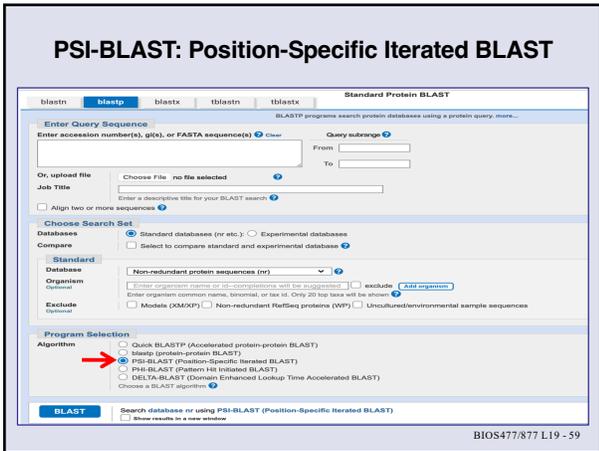
56



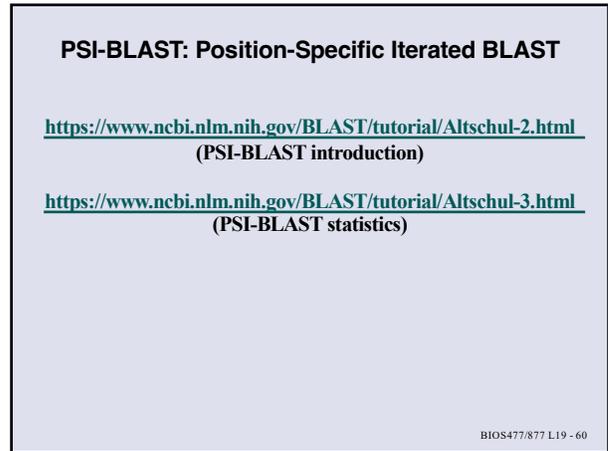
57



58



59



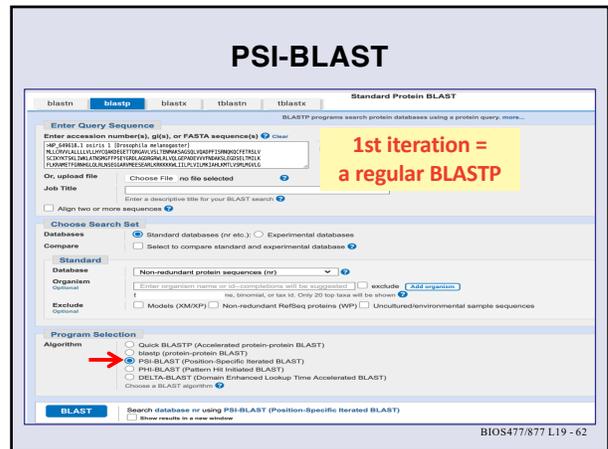
60

PSI-BLAST: Position-Specific Iterated BLAST

- **PSI-BLAST** <https://blast.ncbi.nlm.nih.gov/Blast.cgi>
- ➔ **1st iteration: a regular BLASTP search**
 - Uses a scoring matrix (e.g., BLOSUM62)
- ➔ **2nd iteration**
 - Multiple alignment is constructed from the highly similar hits
 - **Positive-specific scoring matrix (PSSM)** is constructed
 - Similarity search using the **PSSM** instead of the single query sequence
- ➔ **3rd, ..., iterations**
 - Stop when no more new hit (or anytime)

BIOS477/877 L19 - 61

61



62

PSI-BLAST

Two different E-value thresholds

Algorithm parameters

General Parameters

Max target sequences: 500

Short queries: Automatically reduce

Expect threshold: 0.05 **Threshold E-value for reporting matches**

Word size: 3

Max matches in a query range: 0

Scoring Parameters

Matrix: BLOSUM62

Gap Costs: Existence: 11 Extension: 1

Compositional adjustments: Conditional compositional score matrix adjustment

Filters and Masking

Filter: Low complexity regions

Mask: Mask for lookup table only Mask lower case letters

PSI/PHDELTA BLAST

Upload PSSM: Choose File (no file selected)

PSI-BLAST

Threshold: 0.005 **PSI-BLAST threshold: E-value threshold for sequences to be included to create the PSSM for the next iteration**

Postcount: 0

BIOS477/877 L19 - 63

63

PSI-BLAST

Sequences producing significant alignments

Sequences with E-value BETTER than threshold

select all 132 sequences selected

PSI-BLAST threshold (0.005)

Description	Scientific Name	Max Score	Total Score	Query Cover	E-value	Per. Ident.	Acc. Len.	Accession
seqs_1_Drosophila melanogaster	Drosophila melanogaster	638	638	100%	0.0	100.0%	308	NP_548818.1
uncharacterized protein LOC3726827_Drosophila simulans	Drosophila simulans	624	624	100%	0.0	97.4%	308	XP_020102229.1
uncharacterized protein LOC2282475_Drosophila obscura	Drosophila obscura	525	525	100%	0.0	95.2%	308	XP_043854955.1
uncharacterized protein LOC552548_Drosophila obscura	Drosophila obscura	523	523	100%	0.0	93.1%	308	XP_001879262.1
uncharacterized protein LOC0615469_Drosophila achaeae	Drosophila achaeae	520	520	100%	0.0	95.2%	308	XP_020382821.1
uncharacterized protein LOC3246180_Drosophila santomea	Drosophila santomea	513	513	100%	0.0	95.6%	308	XP_038495195.1
uncharacterized protein LOC5817079_Drosophila yakuba	Drosophila yakuba	508	508	100%	0.0	94.3%	308	XP_020388121.1
uncharacterized protein FB98A_001384_Drosophila obscura	Drosophila obscura	491	491	100%	0.0	94.7%	308	XP_043854955.1
uncharacterized protein LOC10954823_Drosophila yakuba	Drosophila yakuba	488	488	100%	0.0	94.7%	311	XP_018939523.1
hypothetical protein RN58_055423_Drosophila burnsi	Drosophila burnsi	479	479	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC119681019_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_037207270.1
hypothetical protein RN209_000924_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC112170668_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC10815288_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
hypothetical protein RN554_007721_Drosophila yakuba	Drosophila yakuba	474	474	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC108154648_Drosophila obscura	Drosophila obscura	474	474	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC108186281_Drosophila obscura	Drosophila obscura	463	463	100%	0.0	94.1%	311	XP_017133955.1

BIOS477/877 L19 - 64

64

PSI-BLAST

Sequences producing significant alignments

Sequences with E-value BETTER than threshold

select all 132 sequences selected

PSI-BLAST threshold (0.005)

Description	Scientific Name	Max Score	Total Score	Query Cover	E-value	Per. Ident.	Acc. Len.	Accession
seqs_1_Drosophila melanogaster	Drosophila melanogaster	638	638	100%	0.0	100.0%	308	NP_548818.1
uncharacterized protein LOC0615469_Drosophila achaeae	Drosophila achaeae	520	520	100%	0.0	95.2%	308	XP_020382821.1
uncharacterized protein LOC3246180_Drosophila santomea	Drosophila santomea	513	513	100%	0.0	95.6%	308	XP_038495195.1
uncharacterized protein LOC5817079_Drosophila yakuba	Drosophila yakuba	508	508	100%	0.0	94.3%	308	XP_020388121.1
uncharacterized protein FB98A_001384_Drosophila obscura	Drosophila obscura	491	491	100%	0.0	94.7%	308	XP_043854955.1
uncharacterized protein LOC10954823_Drosophila yakuba	Drosophila yakuba	488	488	100%	0.0	94.7%	311	XP_018939523.1
hypothetical protein RN58_055423_Drosophila burnsi	Drosophila burnsi	479	479	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC119681019_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_037207270.1
hypothetical protein RN209_000924_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC112170668_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC10815288_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
hypothetical protein RN554_007721_Drosophila yakuba	Drosophila yakuba	474	474	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC108154648_Drosophila obscura	Drosophila obscura	474	474	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC108186281_Drosophila obscura	Drosophila obscura	463	463	100%	0.0	94.1%	311	XP_017133955.1

Run PSI-BLAST Iteration 2 with max number of sequences: 500

Sequences with E-value WORSE than threshold

0.05 > E-value > 0.005

E-value < 0.005 (PSI-BLAST threshold)

E-value < 0.05 (reporting threshold)

BIOS477/877 L19 - 65

65

PSI-BLAST

Sequences producing significant alignments

Sequences with E-value BETTER than threshold

select all 132 sequences selected

PSI-BLAST threshold (0.005)

Description	Scientific Name	Max Score	Total Score	Query Cover	E-value	Per. Ident.	Acc. Len.	Accession
seqs_1_Drosophila melanogaster	Drosophila melanogaster	638	638	100%	0.0	100.0%	308	NP_548818.1
uncharacterized protein LOC0615469_Drosophila achaeae	Drosophila achaeae	520	520	100%	0.0	95.2%	308	XP_020382821.1
uncharacterized protein LOC3246180_Drosophila santomea	Drosophila santomea	513	513	100%	0.0	95.6%	308	XP_038495195.1
uncharacterized protein LOC5817079_Drosophila yakuba	Drosophila yakuba	508	508	100%	0.0	94.3%	308	XP_020388121.1
uncharacterized protein FB98A_001384_Drosophila obscura	Drosophila obscura	491	491	100%	0.0	94.7%	308	XP_043854955.1
uncharacterized protein LOC10954823_Drosophila yakuba	Drosophila yakuba	488	488	100%	0.0	94.7%	311	XP_018939523.1
hypothetical protein RN58_055423_Drosophila burnsi	Drosophila burnsi	479	479	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC119681019_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_037207270.1
hypothetical protein RN209_000924_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC112170668_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC10815288_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
hypothetical protein RN554_007721_Drosophila yakuba	Drosophila yakuba	474	474	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC108154648_Drosophila obscura	Drosophila obscura	474	474	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC108186281_Drosophila obscura	Drosophila obscura	463	463	100%	0.0	94.1%	311	XP_017133955.1

Run PSI-BLAST Iteration 2 with max number of sequences: 500

Sequences with E-value WORSE than threshold

0.05 > E-value > 0.005

Checked sequences are used to construct PSSM and used for the next search

Start the 2nd iteration of search

BIOS477/877 L19 - 66

66

PSI-BLAST

Sequences producing significant alignments

Sequences with E-value BETTER than threshold

select all 132 sequences selected

PSI-BLAST threshold (0.005)

Description	Scientific Name	Max Score	Total Score	Query Cover	E-value	Per. Ident.	Acc. Len.	Accession
seqs_1_Drosophila melanogaster	Drosophila melanogaster	638	638	100%	0.0	100.0%	308	NP_548818.1
uncharacterized protein LOC0615469_Drosophila achaeae	Drosophila achaeae	520	520	100%	0.0	95.2%	308	XP_020382821.1
uncharacterized protein LOC3246180_Drosophila santomea	Drosophila santomea	513	513	100%	0.0	95.6%	308	XP_038495195.1
uncharacterized protein LOC5817079_Drosophila yakuba	Drosophila yakuba	508	508	100%	0.0	94.3%	308	XP_020388121.1
uncharacterized protein FB98A_001384_Drosophila obscura	Drosophila obscura	491	491	100%	0.0	94.7%	308	XP_043854955.1
uncharacterized protein LOC10954823_Drosophila yakuba	Drosophila yakuba	488	488	100%	0.0	94.7%	311	XP_018939523.1
hypothetical protein RN58_055423_Drosophila burnsi	Drosophila burnsi	479	479	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC119681019_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_037207270.1
hypothetical protein RN209_000924_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC112170668_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC10815288_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
hypothetical protein RN554_007721_Drosophila yakuba	Drosophila yakuba	474	474	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC108154648_Drosophila obscura	Drosophila obscura	474	474	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC108186281_Drosophila obscura	Drosophila obscura	463	463	100%	0.0	94.1%	311	XP_017133955.1

Run PSI-BLAST Iteration 2 with max number of sequences: 500

Sequences with E-value WORSE than threshold

0.05 > E-value > 0.005

Checked sequences are used to construct PSSM and used for the next search

More sequences can be included. But very careful!

- If unrelated sequences are included, generated PSSM loses the specificity
- Errors can be amplified quickly with iterations
- Profile corruption problem

BIOS477/877 L19 - 67

67

PSI-BLAST

Sequences producing significant alignments

Sequences with E-value BETTER than threshold

select all 500 sequences selected

PSI-BLAST threshold (0.005)

Description	Scientific Name	Max Score	Total Score	Query Cover	E-value	Per. Ident.	Acc. Len.	Accession
seqs_1_Drosophila melanogaster	Drosophila melanogaster	638	638	100%	0.0	100.0%	308	NP_548818.1
uncharacterized protein LOC0615469_Drosophila achaeae	Drosophila achaeae	520	520	100%	0.0	95.2%	308	XP_020382821.1
uncharacterized protein LOC3246180_Drosophila santomea	Drosophila santomea	513	513	100%	0.0	95.6%	308	XP_038495195.1
uncharacterized protein LOC5817079_Drosophila yakuba	Drosophila yakuba	508	508	100%	0.0	94.3%	308	XP_020388121.1
uncharacterized protein FB98A_001384_Drosophila obscura	Drosophila obscura	491	491	100%	0.0	94.7%	308	XP_043854955.1
uncharacterized protein LOC10954823_Drosophila yakuba	Drosophila yakuba	488	488	100%	0.0	94.7%	311	XP_018939523.1
hypothetical protein RN58_055423_Drosophila burnsi	Drosophila burnsi	479	479	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC119681019_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_037207270.1
hypothetical protein RN209_000924_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC112170668_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC10815288_Drosophila obscura	Drosophila obscura	478	478	100%	0.0	94.8%	308	XP_020382821.1
hypothetical protein RN554_007721_Drosophila yakuba	Drosophila yakuba	474	474	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC108154648_Drosophila obscura	Drosophila obscura	474	474	100%	0.0	94.8%	308	XP_020382821.1
uncharacterized protein LOC108186281_Drosophila obscura	Drosophila obscura	463	463	100%	0.0	94.1%	311	XP_017133955.1

Run PSI-BLAST Iteration 2 with max number of sequences: 500

Sequences with E-value WORSE than threshold

0.05 > E-value > 0.005

Sequences newly added this iteration

Profile corruption problem

BIOS477/877 L19 - 68

68

PSI-BLAST

First iteration (BLASTP search)

XP_045473886.1: 0.001

Sequence	Score	E-value	Query	Accession
uncharacterized protein LOC128801113 (Hemorrhoea axyroidis)	62.4	52.4	49%	0.001
uncharacterized protein LOC128801113 (Hemorrhoea axyroidis)	62.4	52.4	49%	0.001
uncharacterized protein LOC113498740 (Isotria medeolae ssp.)	25.60%	268	XP_037967634.1	0.001
uncharacterized protein LOC12887560 (Panicum latifolium)	25.60%	268	XP_045780333.1	0.004
hypothetical protein L128_L12800546 (L. acidithiobacillus)	51.6	51.6	79%	0.002
uncharacterized protein LOC113498740 (Isotria medeolae ssp.)	51.6	51.6	81%	0.003
uncharacterized protein LOC105385659 (Isotria medeolae ssp.)	50.4	50.4	77%	0.004

Second iteration

XP_037967634.1: 6e-40

XP_037967632.1: 4e-39

XP_045473886.1: 9e-32

BIOS477/877 L19 - 69

- 1st iteration has the real E-values for the query
- After the 2nd iteration, E-values are for the PSSM

69

PSI-BLAST

First iteration (BLASTP search)

Description	Scientific Name	Max Score	Total Score	Query Cover	E	Per. Ident	Acc. Len	Accession	Select for PSI-blast	Used to build PSSM	Newly added
genus 1 (Drosophila melanogaster)	Drosophila melanogaster	638	638	100%	0.0	100.00%	308	XP_045473886.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
uncharacterized protein LOC128801113 (Drosophila simulans)	Drosophila simulans	624	624	100%	0.0	97.40%	308	XP_002102233.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
uncharacterized protein LOC128801113 (Drosophila teissleri)	Drosophila teissleri	530	530	94%	0.0	97.40%	308	XP_045658119.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
uncharacterized protein LOC850248 (Drosophila erecta)	Drosophila erecta	523	523	94%	0.0	95.21%	308	XP_001929306.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
uncharacterized protein LOC9814089 (Drosophila sechellii)	Drosophila sechellii	520	520	94%	0.0	95.21%	308	XP_002038523.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
uncharacterized protein LOC10454180 (Drosophila santomea)	Drosophila santomea	513	513	94%	0.0	95.95%	308	XP_039485195.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
uncharacterized protein LOC4801909 (Drosophila pseudoobscura)	Drosophila sp.	421	421	100%	9e-144	97.73%	430	XP_001398203.6	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
uncharacterized protein LOC128801113 (Drosophila simulans)	Drosophila sp.	406	406	94%	4e-140	97.40%	308	XP_002102233.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
uncharacterized protein LOC128801113 (Drosophila teissleri)	Drosophila sp.	407	407	94%	9e-140	97.40%	308	XP_045658119.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
genus 1 (Drosophila melanogaster)	Drosophila sp.	406	406	94%	3e-139	100.00%	308	XP_045473886.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
uncharacterized protein LOC850248 (Drosophila erecta)	Drosophila sp.	404	404	94%	1e-139	97.40%	308	XP_001929306.1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	

E-values can become lower or higher with more iterations depending on the sequences included to build PSSM!!

BIOS477/877 L19 - 70

70