**Slide 1:**

Spring 2024

# BIOS 477/877

*Bioinformatics and Molecular Evolution*

## Lecture 19

BIOS477/877 L19 - 1

**1**

**Slide 2:**

## TODAY'S TOPICS

➢ **Protein family/domain databases**
(InterPro, SMART, etc.)

➢ **Distance estimation**
• **Nucleotide substitutions**

BIOS477/877 L19 - 2

**2**

**Slide 3:**

InterPro — Classification of protein families

https://www.ebi.ac.uk/interpro/

**[Member databases]**

CATH-Gene3D — 7k entries
CDD — 19k entries
HAMAP — 2k entries
NCBIfam — 7k entries
PANTHER — 16k entries
Pfam — 21k entries
PIRSF — 3k entries
PRINTS — 2k entries
PROSITE profiles — 1k entries
PROSITE patterns — 1k entries
SFLD — 303 entries
SMART — 1k entries
SUPERFAMILY — 2k entries

**[Other information]**

Other sequence features: Phobius, SignalP, Coils, MobiDBLite, TMHMM

Other category: SignalP_EUK, SignalP_GRAM_POSITIVE, SignalP_GRAM_NEGATIVE, AntiFam, FunFam, PIRSR

Domains will be searched against all these databases

Transmembrane, signal peptide, coils, disordered regions will be predicted

BIOS477/877 L19 - 3

**3**

**Slide 4:**

InterPro — Classification of protein families

https://www.ebi.ac.uk/interpro/

To browse each member database (Pfam, etc.)

Home | Search | Browse | Results | Release notes | Download | Help | About

By InterPro / By Member DB / By Protein / By Structure / By Taxonomy / By Proteome / By Set

InterPro is the new home of Pfam. The Pfam website was shut down on October 5th, but InterPro offers the same functionality and data. A legacy version is available at https://pfam-legacy.xfam.org/ but will not receive any updates and will be decommissioned in Spring 2023.

Search by sequence | Search by Domain Architecture

Sequence, in FASTA format

Enter your sequence

Choose file | Example protein sequence

▶▶ Advanced options

Submit | Clear

Powered by InterProScan

BIOS477/877 L19 - 4

**4**

**Slide 5:**

InterPro — Classification of protein families

**[Pfam database entries]**

https://www.ebi.ac.uk/interpro/

Home | Search | Browse | Results | Release notes | Download | Help | About | Contact us

Browse / By Entry / Pfam

Select your database:

AntiFam 263
CATH-Gene3D 7k
CDD 19k
HAMAP 2k
NCBIfam 7k
PANTHER 16k
Pfam 21k
PIRSF 3k
PRINTS 2k
PROSITE profiles 1k
PROSITE patterns 1k
SFLD 303
SMART 1k
SUPERFAMILY 2k

1 - 20 of 21k entries in Pfam

| ACCESSION | NAME | PFAM TYPE | DB | INTEGRATED INTO |
|---|---|---|---|---|
| PF00001 | 7 transmembrane receptor (rhodopsin family) | family | | IPR000276 |
| PF00002 | 7 transmembrane receptor (Secretin family) | family | | IPR000832 |
| PF00003 | 7 transmembrane sweet-taste receptor of 3 GCPR | domain | | IPR017978 |
| PF00004 | ATPase family associated with various cellular activities (AAA) | domain | | IPR003959 |
| PF00005 | ABC transporter | domain | | IPR003439 |

BIOS477/877 L19 - 5

**5**

**Slide 6:**

InterPro — Classification of protein families

https://www.ebi.ac.uk/interpro/

Home | Search | Browse | Results | Release notes | Download | Help | About

InterPro is the new home of Pfam. The Pfam website (pfam.xfam.org) was shut down on October 5th, but InterPro offers the same functionality and data. A legacy version of Pfam is available at https://pfam-legacy.xfam.org/ but will not receive any updates and will be decommissioned in Spring 2023.

Search by sequence | Search by text | Search by Domain Architecture

Sequence, in FASTA format

Enter your sequence

**InterProScan**

Choose file | Example protein sequence

▶▶ Advanced options

Submit | Clear

Powered by InterProScan

BIOS477/877 L19 - 6

**6**

**7**



**8**



**9**



**10**



**11**



**12**

**13**



**14**



**15**



**16**



**17**



**18**

**19**



**20**



**21**



**22**



**23**



**24**

**25**



**26**



**27**



**28**



**29**



**30**

**31**



**HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment**

Michael Remmert, Andreas Biegert, Andreas Hauser & Johannes Söding ✉

*Nature Methods* **9**, 173–175 (2012) | Cite this article

https://toolkit.tuebingen.mpg.de/hhblits

Query sequence → ① Query HMM → ② Accurate HMM-HMM search → Prefilter

Add sequences to query HMM

Database HMMs

③ Accepted HMMs

Rejected HMMs

① **Convert the query sequence to a profile HMM by adding context-specific pseudocounts**

② **HMM-HMM search against an HMM database (e.g., generated from clustered UniProt at 20% similarity; UniProt20)**

③ **Sequences from the accepted HMMs are iteratively added to the query sequences**

BIOS477/877 L19 - 31

---

**32**



**HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment**

Michael Remmert, Andreas Biegert, Andreas Hauser & Johannes Söding ✉

*Nature Methods* **9**, 173–175 (2012) | Cite this article

https://toolkit.tuebingen.mpg.de/hhblits

Query sequence → Query HMM → 219-letter profile → Fast prefilter

Fast search comparing sequences coded into 219 letters

Add sequences to query HMM

Accurate HMM-HMM search

Database HMMs → Select

Accepted HMMs

Rejected HMMs

Profile database as 219-letter sequences

**219 profile codes**

Each position of MSA is coded into single letter

MSA + pseudocounts

Query/DB MSA

BIOS477/877 L19 - 32

---

**33**



Steinegger *et al. BMC Bioinformatics* (2019) 20:473
https://doi.org/10.1186/s12859-019-3019-7

**BMC Bioinformatics**

SOFTWARE — Open Access

## HH-suite3 for fast remote homology detection and deep protein annotation

Martin Steinegger[1,2], Markus Meier[1], Milot Mirdita[1], Harald Vöhringer[1,3], Stephan J. Haunsberger[4] and Johannes Söding[1*]

HHblits3 (AVX2), HHblits3 (SSE2), HHblits2.0.16, HMMer, PSI–BLAST

**Faster** — Time (sec)

Number of iterations

HHblits 3, HHblits 2.0.16, HMMer, PSI–BLAST
1 iteration, 2 iterations, 3 iterations

Fraction of queries

**More accurate**

AUC up to the first 1st false positve

https://toolkit.tuebingen.mpg.de/hhblits

BIOS477/877 L19 - 33

---

**34**

### Distance estimation

**Ancestral sequence ?**

ACTGTAGGAATCGC ⟷ AATGAAAGAATCGC

S1 ←— **Amount of evolution** —→ S2

**Distance**

S1 A**C**TG**T**A**G**GAATCGC
  :X::X:X::::::::
S2 A**A**TG**A**A**A**GAATCGC

(No. of differences = 3)

BIOS477/877 L19 - 34

---

**35**

### Distance estimation

➢ **The simplest method ($p$-distance)**

➜ **Number of substitutions per site ($p$) or degree of divergence**

$$p = \frac{n_d}{L}$$

$$V(p) = \frac{p(1-p)}{L} \text{ or } \sigma(p) = \sqrt{\frac{p(1-p)}{L}}$$

$n_d$: Number of differences between the two sequences
$L$: Number of nucleotides (or amino acids) compared
$V(p)$, $\sigma(p)$: Variance or standard error of the mean ($p$) for binomial distribution

➜ **Can be used for both nucleotide and amino acid substitutions**

A**C**TG**T**A**G**GAATCGC
:X::X:X::::::::
A**A**TG**A**A**A**GAATCGC

$n_d = 3, L = 14$
$p = 3/14 = 0.214$
$\sigma_p = \sqrt{\{0.214 \times (1-0.214)/14\}} = 0.110$
$p = 0.214 \pm 0.110$

BIOS477/877 L19 - 35

---

**36**

### Distance estimation

**Ancestral sequence ?**

ACTGTAGGAATCGC ⟷ AATGAAAGAATCGC

S1 A**C**TG**T**A**G**GAATCGC
  :X::X:X::::::::
S2 A**A**TG**A**A**A**GAATCGC

(No. of differences = 3)

**Were there only three changes during the evolution?**

BIOS477/877 L19 - 36

## Slide 37

**Distance estimation**

AATGTAGGAATCGC  [Ancestral]

A→C
Single substitution

ACTGTAGGAATCGC    AATGAAAGAATCGC

ACTGTAGGAATCGC
:X::X:X:::::::
AATGAAAGAATCGC

BIOS477/877 L19 - 37

**37**

## Slide 38

**Distance estimation**

AATGTAGGAATCGC  [Ancestral]

A→C
Single substitution

T→A
Single substitution

G→A
Single substitution

ACTGTAGGAATCGC    AATGAAAGAATCGC

ACTGTAGGAATCGC
:X::X:X:::::::
AATGAAAGAATCGC

**3 substitutions: All substitutions found**

BIOS477/877 L19 - 38

**38**

## Slide 39

**Distance estimation**

AATGTAGGAATCGC  [Ancestral]

Reality?

ACTGTAGGAATCGC    AATGAAAGAATCGC

ACTGTAGGAATCGC
:X::X:X:::::::
AATGAAAGAATCGC

Observed number of differences = 3

BIOS477/877 L19 - 39

**39**

## Slide 40

**Distance estimation**

AATGTAGGAATCGC  [Ancestral]

A→G + G→C
2 substitutions
?

G

Reality?

C

T→C + C→T
2 substitutions
?

ACTGTAGGAATCGC    AATGAAAGAATCGC

ACTGTAGGAATCGC
:X::X:X:::::::
AATGAAAGAATCGC

Observed number of differences = 3

Actual number of substitutions = 6
[3 hidden substitutions]

Each site has only a limited number of possibilities (4 nucleotides or 20 amino acids)
More hidden substitutions can happen with longer evolution

BIOS477/877 L19 - 40

**40**

## Slide 41

**Effect of multiple substitutions**

Saturation plot

Divergence

Actual divergence
(Actual number of substitutions)

Actual >> Observed

Observed divergence
(Some substitutions are hidden, cannot be observed)

Observed number of substitutions reaches a plateau after a certain time
→ Saturation of substitutions

Actual = Observed

Time

BIOS477/877 L19 - 41

**41**

## Slide 42

**Effect of multiple substitutions**

➢ When the degree of divergence between two sequences is small,
  ➜ the chance of having more than one substitution at any site is negligible
  ➜ Observed divergence ≈ actual divergence
➢ When the degree of divergence becomes larger,
  ➜ more than one substitution could happen at any site [multiple substitutions or multiple hits]
  ➜ Observed divergence << actual divergence [Saturation effect]
➢ Effect of multiple hits is larger for nucleotide substitutions
➢ Methods to uncover the number of hidden substitutions need to be used [Multiple hit correction]
  ➜ Actual divergence level is estimated based on the observed degree of divergence

BIOS477/877 L19 - 42

**42**

7

## Distance estimation for nucleotide substitutions

➢ **Jukes-Cantor (one-parameter) method**

**Jukes and Cantor (1969)**

|   | A | C | G | T |
|---|---|---|---|---|
| **A** | – | α | α | α |
| **C** | α | – | α | α |
| **G** | α | α | – | α |
| **T** | α | α | α | – |

**All substitutions occur with equal probability**
[**Jukes-Cantor model of nucleotide substitutions**]

(Derivation of the JC equation: a note on Canvas)

$$k = -\frac{3}{4}\ln\left(1 - \frac{4}{3}p\right)$$

$k$: Expected number of nucleotide substitutions per site or **Distance**
$p$: Proportion of nucleotide differences (observed)

$$V(k) = \frac{9p(1-p)}{(3-4p)^2 L}$$

$$\sigma(k) = \frac{3}{(3-4p)}\sqrt{\frac{p(1-p)}{L}}$$

$L$: number of nucleotide positions compared

BIOS477/877 L19 - 43

**43**

---

## Distance estimation for nucleotide substitutions

```
ACTGTAGGAATCGC
:X::X:X:::::::
AATGCAAGAATCGC
```

Number of differences = 3
Sequence length = 14

- **Without multiple-hit correction (*p*-distance):**

$$p = \frac{n_d}{L}$$

$n_d$: number of differences, $L$: number of nucleotides compared
$p = 3/14 = 0.214$

- **With multiple-hit correction by Jukes-Cantor method:**

$$k = -\frac{3}{4}\ln\left(1 - \frac{4}{3}P\right)$$

$p$: (Observed) proportion of nucleotide differences = 0.214
$k = -3/4\ln(1 - 4 \times 0.214/3) = 0.252$

BIOS477/877 L19 - 44

**44**

---

## Distance estimation for nucleotide substitutions

$$k = -\frac{3}{4}\ln\left(1 - \frac{4}{3}P\right)$$

**If $p \geq 0.75$,
JC distance cannot be estimated due to arithmetic violation**

Expected
$k$ (Jukes-Cantor distance)

$(k = p)$

$p = 0.75$

$p$ (uncorrected nucleotide difference)
**Observed**

BIOS477/877 L19 - 45

**45**

---

## Distance estimation for nucleotide substitutions

➢ **Kimura two-parameter method** Kimura (1980)

|   | A | C | G | T |
|---|---|---|---|---|
| **A** | – | β | α | β |
| **C** | β | – | β | α |
| **G** | α | β | – | β |
| **T** | β | α | β | – |

**Difference in Ts and Tv substitutions (usually Ts > Tv) can be considered**
[**Kimura 2-parameter model of nucleotide substitutions**]

$$k = \frac{1}{2}\ln\left[\frac{1}{(1-2P-Q)}\right] + \frac{1}{4}\ln\left[\frac{1}{(1-2Q)}\right]$$

$P$: Proportion of transitional (Ts) differences
$Q$: Proportion of transversional (Tv) differences

$$V(k) = \frac{1}{L}\left[P\left\{\frac{1}{(1-2P-Q)}\right\}^2 + Q\left\{\frac{1}{(2-4P-2Q)} + \frac{1}{(2-4Q)}\right\}^2 \right.$$
$$\left. - \left\{\frac{P}{(1-2P-Q)} + \frac{Q}{(2-4P-2Q)} + \frac{Q}{(2-4Q)}\right\}^2\right]$$

$L$: number of nucleotide positions compared

BIOS477/877 L19 - 46

**46**

---

## Distance estimation for nucleotide substitutions

```
ACTGTAGGAATCGC
:X::X:X:::::::
AATGCAAGAATCGC
```

Number of differences = 3
Ts = 2, Tv = 1
Alignment length = 14

- **Without multiple-hit correction (*p*-distance):**

$$p = \frac{n_d}{L} \qquad V(p) = \frac{p(1-p)}{L}$$

$p = 0.214 \pm 0.110$

- **Jukes-Cantor distance:**

$$k = -\frac{3}{4}\ln\left(1 - \frac{4}{3}p\right) \qquad V(k) = \frac{9p(1-p)}{(3-4p)^2 L}$$

$k = -3/4\ln(1 - 4 \times 0.214/3)$
$\quad = 0.252 \pm 0.154$

- **Kimura 2-parameter distance:**

$$k = \frac{1}{2}\ln\left[\frac{1}{(1-2P-Q)}\right] + \frac{1}{4}\ln\left[\frac{1}{(1-2Q)}\right]$$

$P = 2/14, Q = 1/14$
$k = 1/2\ln[1/(1-4/14-1/14)] + 1/4\ln[1/(1-2/14)]$
$\quad = 0.259 \pm 0.143$

$$V(k) = \frac{1}{L}\left[P\left\{\frac{1}{(1-2P-Q)}\right\}^2 + Q\left\{\frac{1}{(2-4P-2Q)} + \frac{1}{(2-4Q)}\right\}^2 - \left\{\frac{P}{(1-2P-Q)} + \frac{Q}{(2-4P-2Q)} + \frac{Q}{(2-4Q)}\right\}^2\right]$$

BIOS477/877 L19 - 47

**47**

---

## Sequence evolution as Markov process

**Transition probability**

| [Time] | T1 | T2 | T3 | T4 |
|--------|----|----|----|----|

$A \xrightarrow{P_1} G \xrightarrow{P_2} G \xrightarrow{P_3} A$

**Markov Chain: a discrete-time stochastic process**

**In more general continuous-time scale,**
➜ **Markov Process**

BIOS477/877 L19 - 48

**48**

8

**Slide 49: Sequence evolution as Markov process**

Transition probability

$$A \xrightarrow{P_1} G \xrightarrow{P_2} G \xrightarrow{P_3} A$$

[Time]   T1    T2    T3    T4

|   | A | C | G | T |
|---|---|---|---|---|
| A | - | α | α | α |
| C | α | - | α | α |
| G | α | α | - | α |
| T | α | α | α | - |

Jukes-Cantor model
(α: substitution rate)

Transition probability matrix

$$\mathbf{P}(t) = \begin{bmatrix} r_t & s_t & s_t & s_t \\ s_t & r_t & s_t & s_t \\ s_t & s_t & r_t & s_t \\ s_t & s_t & s_t & r_t \end{bmatrix}$$

$r_t$: Prob. of no change
$s_t$: Prob. of changes

where $r_t + 3s_t = 1$ (row sum)
thus $r_t = 1 - 3s_t$

BIOS477/877 L19 - 49

**Slide 50: Jukes-Cantor model of sequence evolution**

|   | A | C | G | T |
|---|---|---|---|---|
| A | - | α | α | α |
| C | α | - | α | α |
| G | α | α | - | α |
| T | α | α | α | - |

(e.g., α = 5x10⁻⁹ substitutions/site/year)

$r_t$: Prob. of no change
0.25
$s_t$: Prob. of changes

Time (million years)

Transition probability matrix

$$\mathbf{P}(t) = \begin{bmatrix} r_t & s_t & s_t & s_t \\ s_t & r_t & s_t & s_t \\ s_t & s_t & r_t & s_t \\ s_t & s_t & s_t & r_t \end{bmatrix}, \text{ where } \begin{cases} r_t = \dfrac{1}{4} + \dfrac{3}{4}e^{-4\alpha t} \\ s_t = \dfrac{1}{4} - \dfrac{1}{4}e^{-4\alpha t} \end{cases}$$

$t=0$: $r_0=1$ & $s_0=0$
$t \rightarrow \infty$:
$r_\infty=0.25$ & $s_\infty=0.25$

(Derivation of $r_t$, $s_t$, and J-C distance equations, read "Derivation of the JC equation" on Canvas.)

$$k = -\frac{3}{4}\ln(1 - \frac{4}{3}p)$$

BIOS477/877 L19 - 50

**Slide 51: Nucleotide substitution models**

Jukes-Cantor (JC)
Equal base frequency
($f_A=f_T=f_G=f_C=0.25$)

|   | A | C | G | T |
|---|---|---|---|---|
| A | – | α | α | α |
| C | α | – | α | α |
| G | α | α | – | α |
| T | α | α | α | – |

Felsenstein (F81)
Unequal base frequency

|   | A | C | G | T |
|---|---|---|---|---|
| A | – | $\pi_C\alpha$ | $\pi_G\alpha$ | $\pi_T\alpha$ |
| C | $\pi_A\alpha$ | – | $\pi_G\alpha$ | $\pi_T\alpha$ |
| G | $\pi_A\alpha$ | $\pi_C\alpha$ | – | $\pi_T\alpha$ |
| T | $\pi_A\alpha$ | $\pi_C\alpha$ | $\pi_G\alpha$ | – |

|   | A | C | G | T |
|---|---|---|---|---|
| A | – | β | α | β |
| C | β | – | β | α |
| G | α | β | – | β |
| T | β | α | β | – |

Kimura 2-parameter (K2P)
Equal base frequency
($f_A=f_T=f_G=f_C=0.25$)

|   | A | C | G | T |
|---|---|---|---|---|
| A | – | $\pi_C\beta$ | $\pi_G\alpha$ | $\pi_T\beta$ |
| C | $\pi_A\beta$ | – | $\pi_G\beta$ | $\pi_T\alpha$ |
| G | $\pi_A\alpha$ | $\pi_C\beta$ | – | $\pi_T\beta$ |
| T | $\pi_A\beta$ | $\pi_C\alpha$ | $\pi_G\beta$ | – |

Hasegawa et al. (HKY85)
Unequal base frequency

3-parameter model:
Ts(AG), Ts(TC), and Tv

|   | A | C | G | T |
|---|---|---|---|---|
| A | – | $\pi_C a$ | $\pi_G b$ | $\pi_T c$ |
| C | $\pi_A a$ | – | $\pi_G d$ | $\pi_T e$ |
| G | $\pi_A b$ | $\pi_C d$ | – | $\pi_T f$ |
| T | $\pi_A c$ | $\pi_C e$ | $\pi_G f$ | – |

There are many more!

General reversible (REV)
Unequal base frequency

BIOS477/877 L19 - 51

9