



CEITEC

Central European Institute of Technology
BRNO | CZECH REPUBLIC

**Introduction to Bioinformatics
(LF:DSIB01)**

Week 3 : Pattern Recognition



Sequence Patterns

We will learn:

1. How to define a pattern
2. How to identify the presence of a pattern in a sequence
3. How to count pattern occurrences
4. Calculation of overrepresentation
5. Creation of elaborate queries

Defining pattern

There are several ways to define a pattern

- Deterministic Patterns
- Patterns with mismatches
- Position Weight Matrices
- Stochastic Models

Deterministic Patterns

- Defined Alphabet $\Sigma=[A,T,G,C]$
- Simple Sequence: e.g. TATAAAA
- Ambiguous character: e.g. TAT[AT]AAA : [AT] = either A or T
- Wildcard: TAT . AAA : . = any character
- Flexible gap: TAT. **{1,3}**AAA : **{1,3}** = one to three times any character

Patterns with mismatches

- Allow exact matching of Deterministic Pattern + a certain number of mistakes
- This category will be covered in depth over the next 2 lectures (Week 4,5)

Position Probability Matrices (PPM)

- Ambiguous symbol [AT] gives the same % on both symbols
- PPM is a table Position x Alphabet containing probability (%) scores
- Using this PPM one can score the probability that this pattern produces this sequence.

$$M = \begin{matrix} A \\ C \\ G \\ T \end{matrix} \begin{bmatrix} 0.3 & 0.6 & 0.1 & 0.0 & 0.0 & 0.6 & 0.7 & 0.2 & 0.1 \\ 0.2 & 0.2 & 0.1 & 0.0 & 0.0 & 0.2 & 0.1 & 0.1 & 0.2 \\ 0.1 & 0.1 & 0.7 & 1.0 & 0.0 & 0.1 & 0.1 & 0.5 & 0.1 \\ 0.4 & 0.1 & 0.1 & 0.0 & 1.0 & 0.1 & 0.1 & 0.2 & 0.6 \end{bmatrix} \cdot \prod_{i=1}^k \frac{M[x_i, i]}{f(x_i)}$$

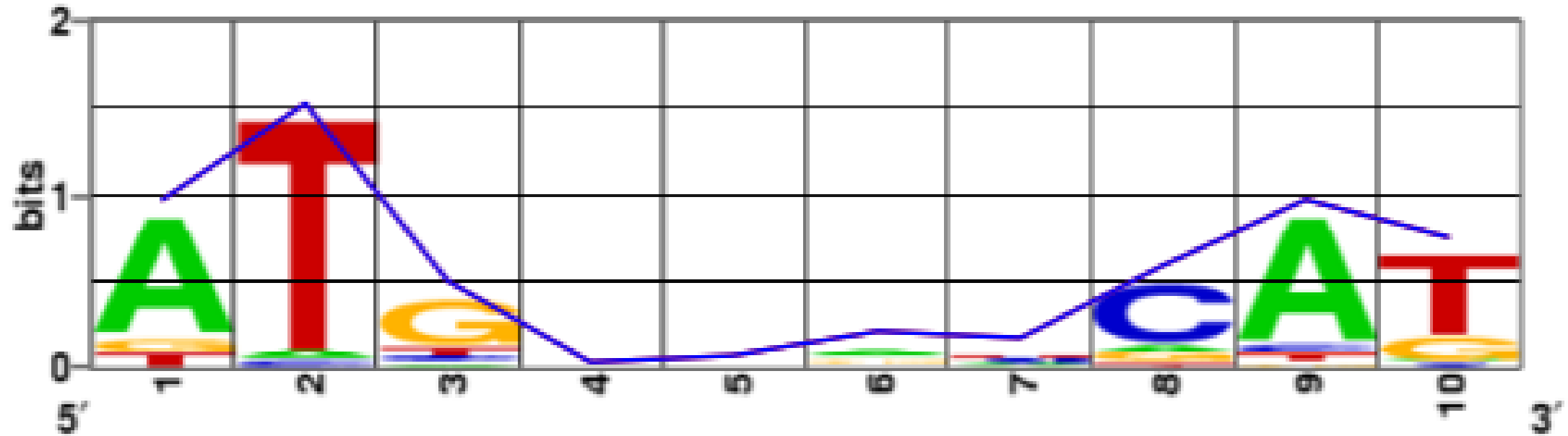
Position Weight Matrices (PWM)

- More useful is the log-odds score (weight)
- Odds Score = \log_2 (frequency / background frequency)
- Sum: **How different** is the scored seq from a background random seq
- 0 => equal prob of being M or Background, + => M more prob, - => M less prob

$$M = \begin{matrix} A \\ C \\ G \\ T \end{matrix} \begin{bmatrix} 0.3 & 0.6 & 0.1 & 0.0 & 0.0 & 0.6 & 0.7 & 0.2 & 0.1 \\ 0.2 & 0.2 & 0.1 & 0.0 & 0.0 & 0.2 & 0.1 & 0.1 & 0.2 \\ 0.1 & 0.1 & 0.7 & 1.0 & 0.0 & 0.1 & 0.1 & 0.5 & 0.1 \\ 0.4 & 0.1 & 0.1 & 0.0 & 1.0 & 0.1 & 0.1 & 0.2 & 0.6 \end{bmatrix} \cdot \prod_{i=1}^k \frac{M[x_i, i]}{f(x_i)}$$

$$M = \begin{matrix} A \\ C \\ G \\ T \end{matrix} \begin{bmatrix} 0.26 & 1.26 & -1.32 & -\infty & -\infty & 1.26 & 1.49 & -0.32 & -1.32 \\ -0.32 & -0.32 & -1.32 & -\infty & -\infty & -0.32 & -1.32 & -1.32 & -0.32 \\ -1.32 & -1.32 & 1.49 & 2.0 & -\infty & -1.32 & -1.32 & 1.0 & -1.32 \\ 0.68 & -1.32 & -1.32 & -\infty & 2.0 & -1.32 & -1.32 & -0.32 & 1.26 \end{bmatrix} \cdot \sum_{i=1}^k M'[x_i, i].$$

Information Content



Information Content

Height = Information content (bits)

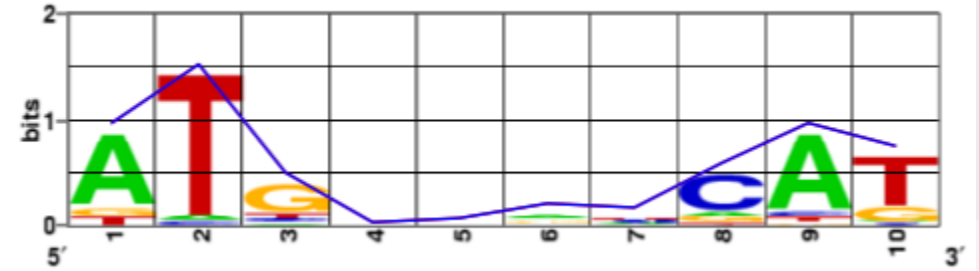
Information = degree of decrease in uncertainty

Hartley 1928: $I(N) = \log(N)$

Shannon 1948: $I(a_i) = \log\left(\frac{1}{P(a_i)}\right) = \log(1) - \log(P(a_i)) = -\log(P(a_i)), \quad P(a_i) \in [0,1]$

$$I(a_i) = -\log\left(\frac{1}{n}\right) = -[\log(1) - \log(n)] = \log(n), \quad P(a_i) = \frac{1}{n}$$

$$I(a_1) = I(a_2) = -\log\left(\frac{1}{2}\right) = \log(2) \quad < \text{Definition of a bit}$$



Entropy

Entropy is a measure of the *unpredictability* of a state (*average information content*)

$$H(X) = E[I(X)] = E[-\log(P(X))].$$

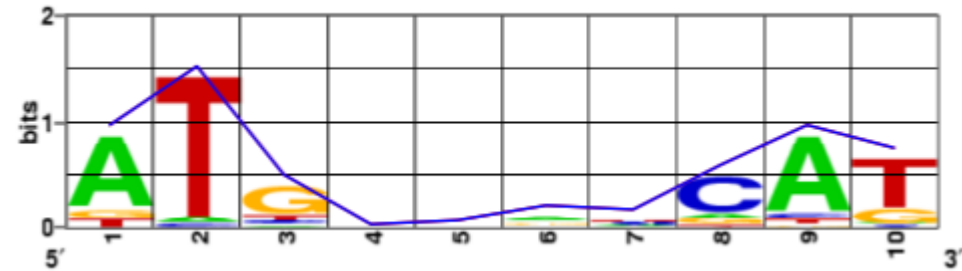
	1	2	3	4	5	6	7	8	9	10
A	0.76	0.04	0.08	0.28	0.12	0.44	0.24	0.12	0.80	0.04
C	0.00	0.04	0.12	0.32	0.28	0.12	0.28	0.68	0.08	0.04
T	0.12	0.92	0.16	0.16	0.28	0.12	0.40	0.08	0.08	0.68
G	0.12	0.00	0.64	0.24	0.32	0.32	0.08	0.12	0.04	0.24

$$H(X) = H_{\text{before}}(l) = -\sum_{S \in \Omega} [f(S) \cdot (\log_2(f(S)))]$$

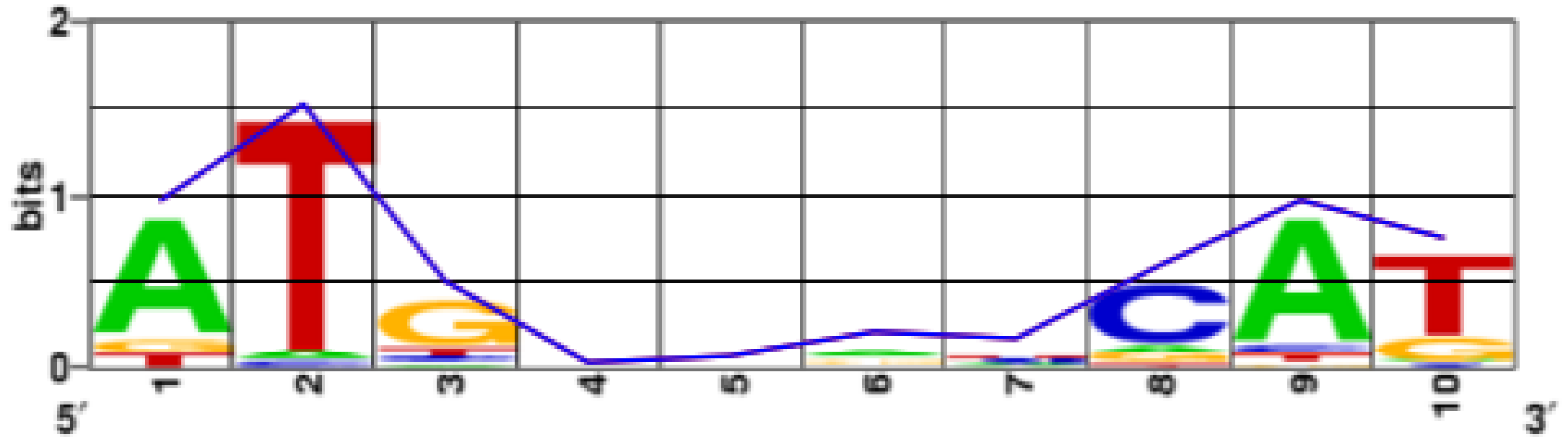
$$H(X | Y) = H_{\text{after}}(l) = -\sum_{S_l \in \Omega} (p(S_l) \cdot \log_2(p(S_l)))$$

$$I(l) = H_{\text{before}}(l) - H_{\text{after}}(l) = \left[-\sum_{S \in \Omega} (f(S) \cdot (\log_2(f(S)))) \right] - \left[-\sum_{S_l \in \Omega} (p(S_l) \cdot \log_2(p(S_l))) \right]$$

$R_{\text{sequence}}(l)$



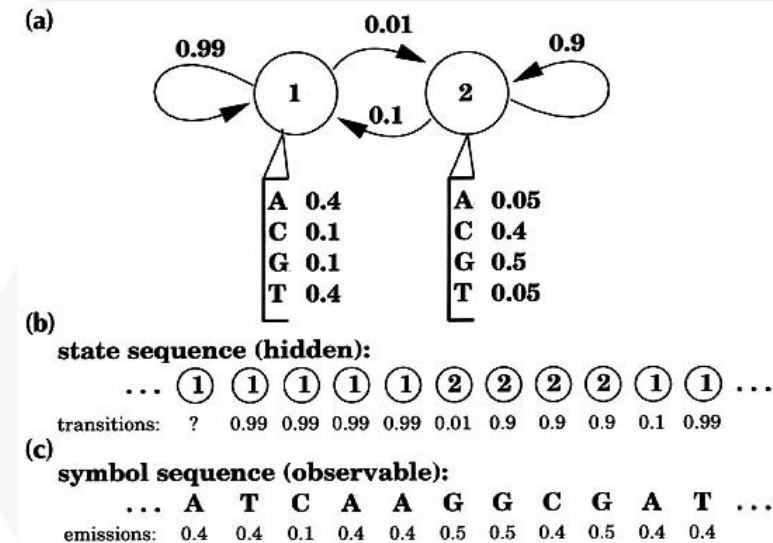
How to read a Seq Logo



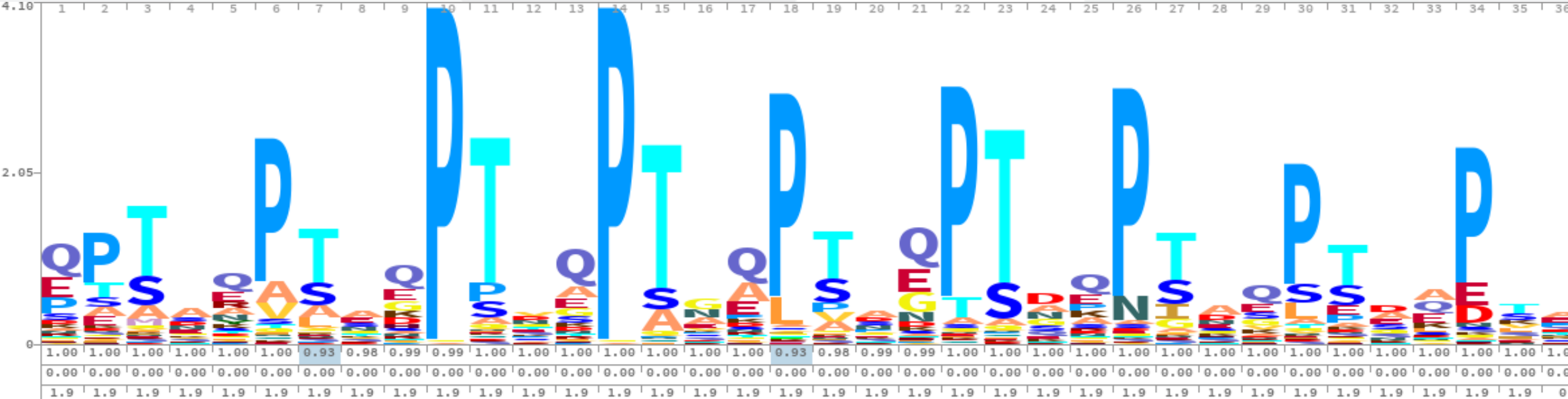
Stack height indicates the Information Content per position (Rseq(I))
Letter height indicates the Base Frequency per position

Stochastic models

- A set of rules or a machine learning method
- Must be able to discriminate / classify / score sequence
- Commonly used: Hidden Markov Models
- We will not talk about stochastic models in this course



HMM Logo (PFAM)



<https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-15-7>



Hunting for Models

Motif Discovery with MEME

▼ Motif Discovery

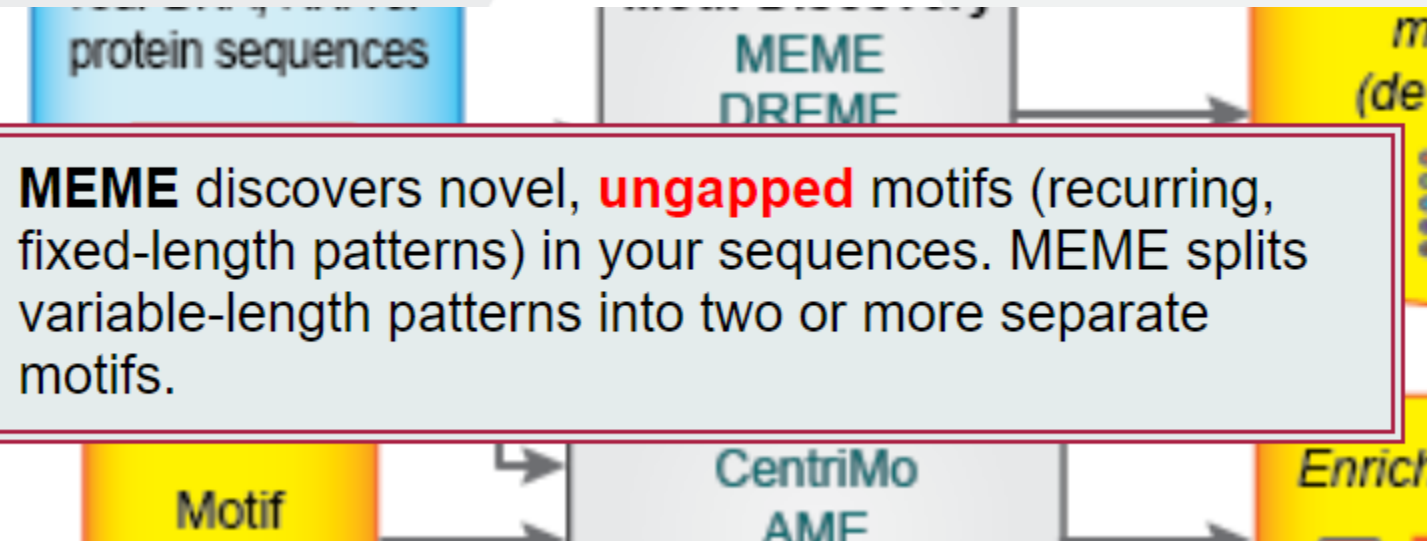
MEME

DREME

MEME-ChIP

GLAM2

MoMo



MEME discovers novel, **ungapped** motifs (recurring, fixed-length patterns) in your sequences. MEME splits variable-length patterns into two or more separate motifs.

Motif Discovery with MEME

```
>SEQ1
TCAGTGCAGTCATGCACATGCATGCATGCATGCATGCATGCATGCATG
>SEQ2
TGTGCTGACTGCATGACTCTATCTGCATGACTGTTTCTGCGCGGC
>SEQ3
TGCATGCATGCACTGAAAAAAAAAATGCATGCATGCACTGACATGCTGACTGA
```




MEME discovers novel, (recurring, fixed-length) sequences (sample output MEME splits variable-length or more separate motifs. See more information.


Data Submission Form

Perform motif discovery on DNA, RNA, protein or custom alphabet datasets.

Select the motif discovery mode


Classic mode Discriminative mode Differential Enrichment mode 


Select the sequence alphabet

Use sequences with a standard alphabet or specify a custom alphabet. 

DNA, RNA or Protein Custom


Input the primary sequences

Enter sequences in which you want to find motifs. 


Type in sequences 

```
>SEQ1
TCAGTGCAGTCATGCACATGCATGCATGCATGCATGCATGCATGCATG
>SEQ2
```

Select the site distribution

How do you expect motif sites to be distributed in sequences? 

Select the number of motifs

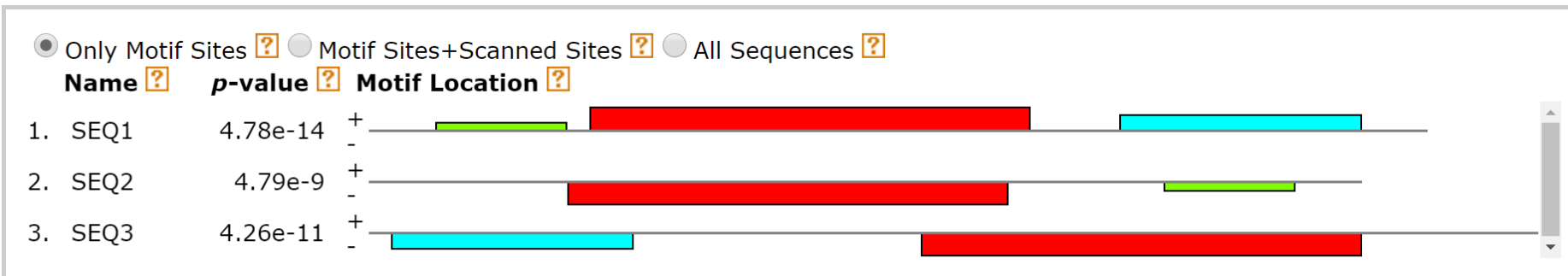
How many motifs should MEME find? 

Motif Discovery with MEME

DISCOVERED MOTIFS



MOTIF LOCATIONS

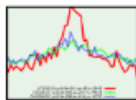


Motif Discovery with MEME



MEME

Multiple Em for Motif Elicitation



CentriMo

Local Motif Enrichment Analysis



FIMO

Find Individual Motif Occurrences



DREME

Discriminative Regular Expression Motif Elicitation



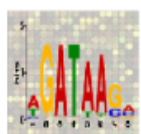
AME

Analysis of Motif Enrichment



MAST

Motif Alignment & Search Tool



MEME-ChIP

Motif Analysis of Large Nucleotide Datasets



SpaMo

Spaced Motif Analysis Tool



MCAST

Motif Cluster Alignment and Search Tool



GLAM2

Gapped Local Alignment of Motifs



GOMo

Gene Ontology for Motifs



GLAM2Scan

Scanning with Gapped Motifs



MoMo

Modification Motifs



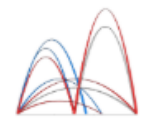
Tomtom

Motif Comparison Tool



GT-Scan

Identifying Unique Genomic Targets



CisMapper


Predicting Regulatory Links



Regular Expressions

- For Deterministic Patterns
 - Regex is a commonly used language for definition of Deterministic Patterns
 - Extremely powerful syntax – but ATTENTION for unintended results
1. Meta-characters
 2. Special Characters
 3. Sets

Meta-characters

Character	Description	Example
[]	A set of characters	"[a-m]"
\	Signals a special sequence (can also be used to escape special characters)	"\d"
.	Any character (except newline character)	"he..o"
^	Starts with	"^hello" 
\$	Ends with	"world\$"
*	Zero or more occurrences	"aix*"
+	One or more occurrences	"aix+"
{}	Exactly the specified number of occurrences	"al{2}"
	Either or	"falls stays"
()	Capture and group	

```
import re
str = "hello world"
#Check if the string starts with 'hello':
x = re.findall("^hello", str)
if (x):
    print("Yes, the string starts with 'hello'")
else:
    print("No match")
```

Special Characters

\ Signals a special sequence (can also be used to escape special characters)

```
import re

str = "The rain in Spain"

#Check if the string starts with "The":

x = re.findall("\AThe", str)

print(x)

if (x):
    print("Yes, there is a match!")
else:
    print("No match")
```

Character	Description	Example
\A	Returns a match if the specified characters are at the beginning of the string	"\AThe"
\b	Returns a match where the specified characters are at the beginning or at the end of a word	r"\bain" r"ain\b"
\B	Returns a match where the specified characters are present, but NOT at the beginning (or at the end) of a word	r"\Bain" r"ain\B"
\d	Returns a match where the string contains digits (numbers from 0-9)	"\d"
\D	Returns a match where the string DOES NOT contain digits	"\D"
\s	Returns a match where the string contains a white space character	"\s"
\S	Returns a match where the string DOES NOT contain a white space character	"\S"
\w	Returns a match where the string contains any word characters (characters from a to Z, digits from 0-9, and the underscore _ character)	"\w"
\W	Returns a match where the string DOES NOT contain any word characters	"\W"
\Z	Returns a match if the specified characters are at the end of the string	"Spain\Z"

\A vs ^ : ^ matches start of LINE while \A start of string

Sets

Set	Description
[arn]	Returns a match where one of the specified characters (<code>a</code> , <code>r</code> , or <code>n</code>) are present
[a-n]	Returns a match for any lower case character, alphabetically between <code>a</code> and <code>n</code>
[^arn]	Returns a match for any character EXCEPT <code>a</code> , <code>r</code> , and <code>n</code>
[0123]	Returns a match where any of the specified digits (<code>0</code> , <code>1</code> , <code>2</code> , or <code>3</code>) are present
[0-9]	Returns a match for any digit between <code>0</code> and <code>9</code>
[0-5][0-9]	Returns a match for any two-digit numbers from <code>00</code> and <code>59</code>
[a-zA-Z]	Returns a match for any character alphabetically between <code>a</code> and <code>z</code> , lower case OR upper case
[+]	In sets, <code>+</code> , <code>*</code> , <code>.</code> , <code> </code> , <code>()</code> , <code>\$</code> , <code>{}</code> has no special meaning, so <code>[+]</code> means: return a match for any <code>+</code> character in the string

Python regex

Function Description

<code>findall</code>	Returns a list containing all matches
<code>search</code>	Returns a <u>Match object</u> if there is a match anywhere in the string
<code>split</code>	Returns a list where the string has been split at each match
<code>sub</code>	Replaces one or many matches with a string

Python Regex Cheatsheet

Regular Expression Basics		Regular Expression Character Classes		Regular Expression Flags	
.	Any character except newline	[ab-d]	One character of: a, b, c, d	i	Ignore case
a	The character a	[^ab-d]	One character except: a, b, c, d	m	^ and \$ match start and end of line
ab	The string ab	[b]	Backspace character	s	. matches newline as well
a b	a or b	\d	One digit	x	Allow spaces and comments
a*	0 or more a's	\D	One non-digit	L	Locale character classes
\	Escapes a special character	\s	One whitespace	u	Unicode character classes
		\S	One non-whitespace	(?iLmsux)	Set flags within regex
		\w	One word character		
		\W	One non-word character		
Regular Expression Quantifiers		Regular Expression Assertions		Regular Expression Special Characters	
*	0 or more	^	Start of string	\n	Newline
+	1 or more	\A	Start of string, ignores m flag	\r	Carriage return
?	0 or 1	\$	End of string	\t	Tab
{2}	Exactly 2	\Z	End of string, ignores m flag	\YYY	Octal character YYY
{2, 5}	Between 2 and 5	\b	Word boundary	\xYY	Hexadecimal character YY
{2,}	2 or more	\B	Non-word boundary		
{,5}	Up to 5	(?=...)	Positive lookahead	Regular Expression Replacement	
Default is greedy. Append ? for reluctant.		(?!...)	Negative lookahead	\g<0>	Insert entire match
		(?<=...)	Positive lookbehind	\g<Y>	Insert match Y (name or number)
		(?!<...)	Negative lookbehind	\Y	Insert group numbered Y
		(?())	Conditional		
Regular Expression Groups					
(...)	Capturing group				
(?P<Y>...)	Capturing group named Y				
(?...)	Non-capturing group				
\Y	Match the Y'th captured group				
(?P=Y)	Match the named group Y				
(?#...)	Comment				

New to Debuggex? Check out the regex tester!



CEITEC



@CEITEC_Brno

Thank you for your attention!
60 minutes lunch break.



www.ceitec.eu