

# n-gram analysis of biological sequences in R

---

Michał Burdukiewicz<sup>1</sup>, Piotr Sobczyk<sup>2</sup>, Małgorzata Kotulska<sup>3</sup>, Paweł Mackiewicz<sup>1</sup>

<sup>1</sup>University of Wrocław, Department of Genomics,

<sup>2</sup>Wrocław University of Science and Technology, Faculty of Pure and Applied Mathematics,

<sup>3</sup>Wrocław University of Science and Technology, Department of Biomedical Engineering

# Introduction

---

# Biological sequences

Long chains of amino acids (proteins) or nucleotides (RNA or DNA).

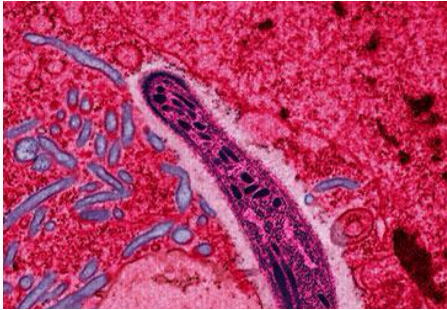
Sample protein sequence:

MRKLYCVLLLSAFEFTYMINFGRGQNYWEHPYQKSDVYHP

INEHREHPKEYQYPLHQEHTYQQEDSGEDENTLQHAYPID

HEGAEPAPQEQNLFSSIEIV...

# Biological sequences



*Plasmodium falciparum*. Source: <http://www.protists.ensembl.org>

Protein of *Plasmodium falciparum*:

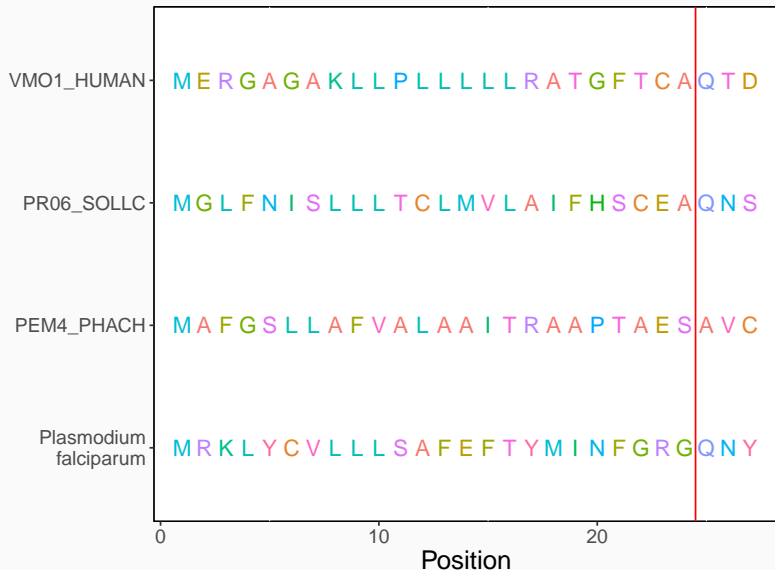
```
MRKLYCVLLLSAFEFTYMINFGRGQNYWEHPYQKSDVYHP  
INEHREHPKEYQYPLHQEHTYQQEDSGEDENTLQHAYPID  
HEGAEPAPQEQLFSSIEIV...
```

# Biological sequences

Signal peptide (red): n-terminal amino acid sequence directing proteins to the endomembrane system and next to extracellular localizations.

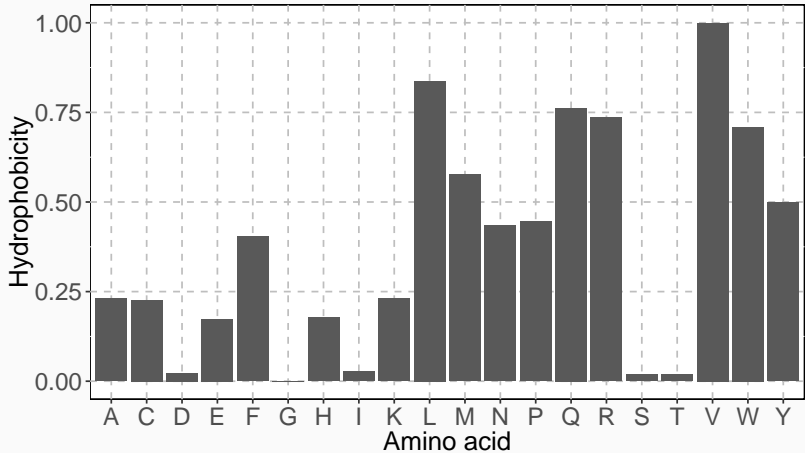
MRKLYCVLLLSAFEFTYMINFGRGQNYWEHPYQKSDVYHP  
INEHREHPKEYQYPLHQEHTYQQEDSGEDENTLQHAYPID  
HEGAEPAPQEQNLFSSIEIV...

# Biological sequences



n-grams (k-tuples) are vectors of  $n$  characters derived from input sequence(s).

# Properties of amino acids



Amino acids may be described using their physicochemical properties.



# Biological sequences

