

# 10: Biological Applications for HMMs

## Machine Learning and Real-world Data (MLRD)

Ann Copestake  
(based on slides created by Simone Teufel)

Lent 2018

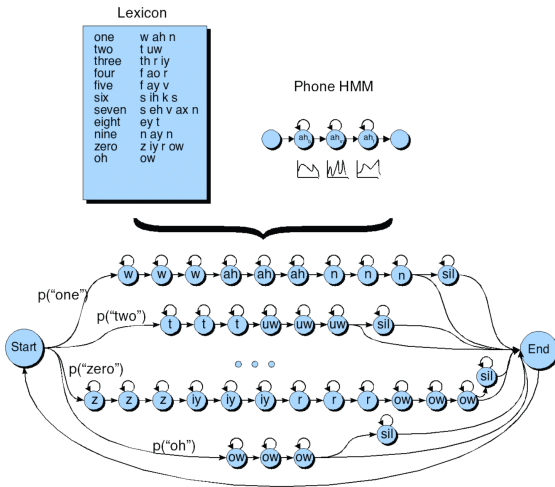
# Last session: dice world and HMM decoding

- You may by now have written a decoder, i.e., an algorithm that can determine the most likely state sequence of an HMM.
- From the task before that, you also have code that can estimate the parameters from a labelled HMM sequence.
- But the dice world is very simple/artificial.

# Sequence Learning in the real world

- HMMs for speech recognition
  - Goal: determine from signal which words were said
  - States: words
  - Observations: acoustic inputs from signal
- HMMs for parts of speech tagging
  - Goal: determine the parts of speech for text
  - States: parts of speech
  - Observations: words
- HMM for protein analysis
  - Goal: Find which sections of proteins are in cell membranes
  - States: zones relating to cells
  - Observations: amino acids

# HMMs in Automatic Speech Recognition (ASR)



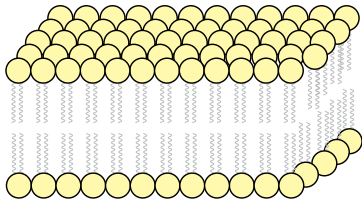
## A biological application

```
#MNQGKIWTVVNPAIGIPALLGSVTVIAILVHLAILSHTTWFPAYWQGGVKKAA  
iiiiiiiiiMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMM
```

- top line records the amino acid sequence (one character per amino acid)
- bottom line shows the states:
  - i: inside the cell
  - M: within the cell membrane
  - o: outside the cell
- Ignoring the start and end sequence states/labels for simplicity.

# Eight minutes about biology of cells

- living organisms are made up of cells
- multicellular organisms have lots of cells
- cells are surrounded by a cell membrane
- cell membranes are lipid bilayers: inside the membrane is hydrophobic (water-hating), the two sides are hydrophilic (water-loving)

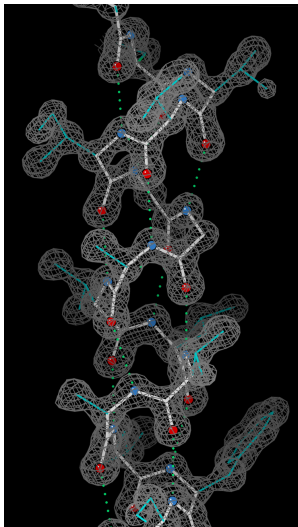


Jerome Walker - Own work, CC BY 2.5, <https://commons.wikimedia.org/w/index.php?curid=915557>

# Proteins

- in cell metabolism: proteins make sure the right thing happens in the right place at the right time
- proteins are made up of amino acid sequences
- all amino acids have amine and carboxyl groups, but they have very different **side chains**
- 20 amino acids are coded for directly by DNA
- amino acid sequences fold into very complex 3-D protein structure

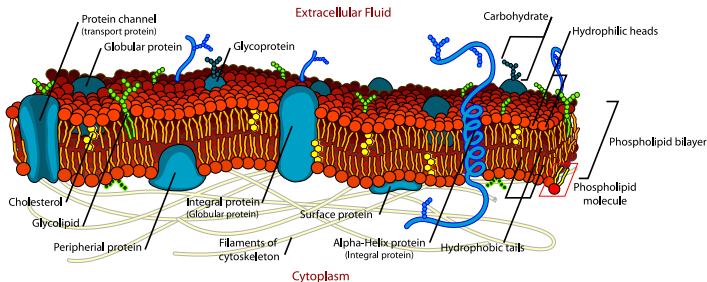
# Alpha helix





# Cell membranes and proteins

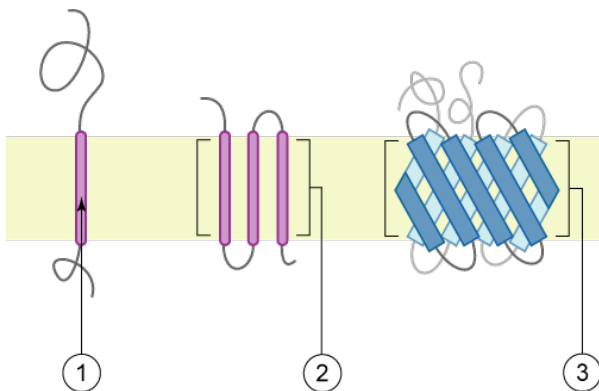
- cell membranes have to let things in and out of the cell (e.g., water, glucose, sodium ions, calcium ions)
- proteins which are part of the cell membrane allow this (membrane proteins do other things too)



# Transmembrane proteins

- transmembrane proteins go through the membrane one or more times
- the regions of the protein which lie inside and outside the cell tend to have more hydrophilic amino acids
- the regions inside the membranes tend to have more hydrophobic amino acids
- many transmembrane proteins involve one or more  $\alpha$ -helices in the membrane
- the channels formed by the protein allow ions and molecules through, in a controlled way

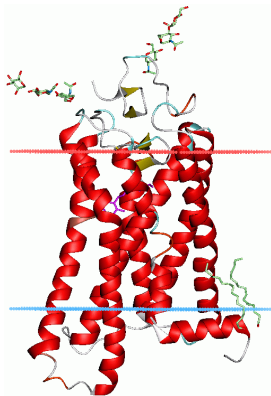
# Transmembrane protein: schematic diagram



1. a single transmembrane  $\alpha$ -helix (bitopic membrane protein)
2. a polytopic transmembrane  $\alpha$ -helical protein
3. a polytopic transmembrane  $\beta$ -sheet protein

By Foobar - self-made by Foobar, CC BY 2.5, <https://commons.wikimedia.org/w/index.php?curid=802476>

# Transmembrane protein: Bovine rhodopsin



- one of the visual pigments
- accurate structure via x-ray crystallography: difficult and time-consuming, membrane location undetermined

## A biological application

```
#MNQGKIWTVVNPAIGIPALLGSVTVIAILVHLAILSHTTWFPAYWQGGVKKAA  
iiiiiiiiiMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMM
```

- HMM-based modelling: much, much easier and quicker than x-ray crystallography
- distinguish interior of membrane from inside/outside of cell
- simple HMM in practical, but could be improved: more discussion in practical notes

# Your Task

## Task 9:

- Download the biological dataset and familiarise yourself with it.
- Modify your code so that your HMM parameter estimation from Task 7 and decoder from Task 8 works with this data format.
- Use 10-fold cross validation.
- Evaluate.

For Task 10 (next week), you will need to download gephi (graph visualization): please do this in advance of the scheduled session if possible.

# Strike action

- As things stand, the mini-lectures on Feb 23, Feb 26, March 5 and March 12 are cancelled due to strike action.
- The digital timetable will be updated appropriately: check this! There is a possibility the strike action will be called off.
- At least some demonstrators will probably continue to work, so practical sessions and ticking are expected to continue.