Marius Audenis, Gabriel Chiche, Leana Ortolani
Supervisor: Carlos Pulido Quetglas

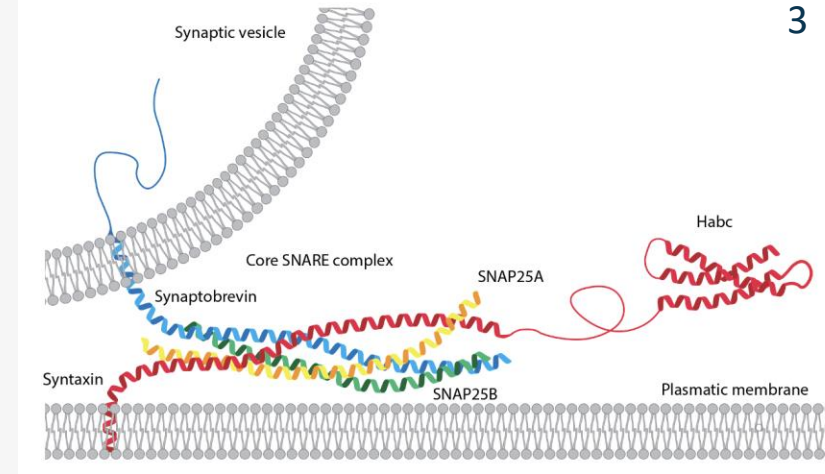# HMM based classification for conserved protein domains (SNARE)
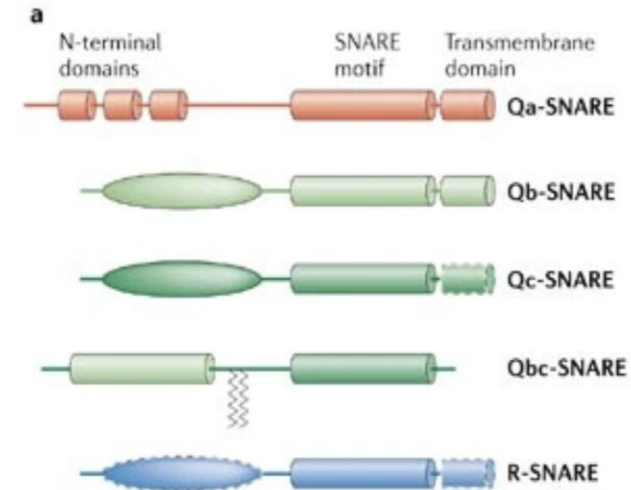
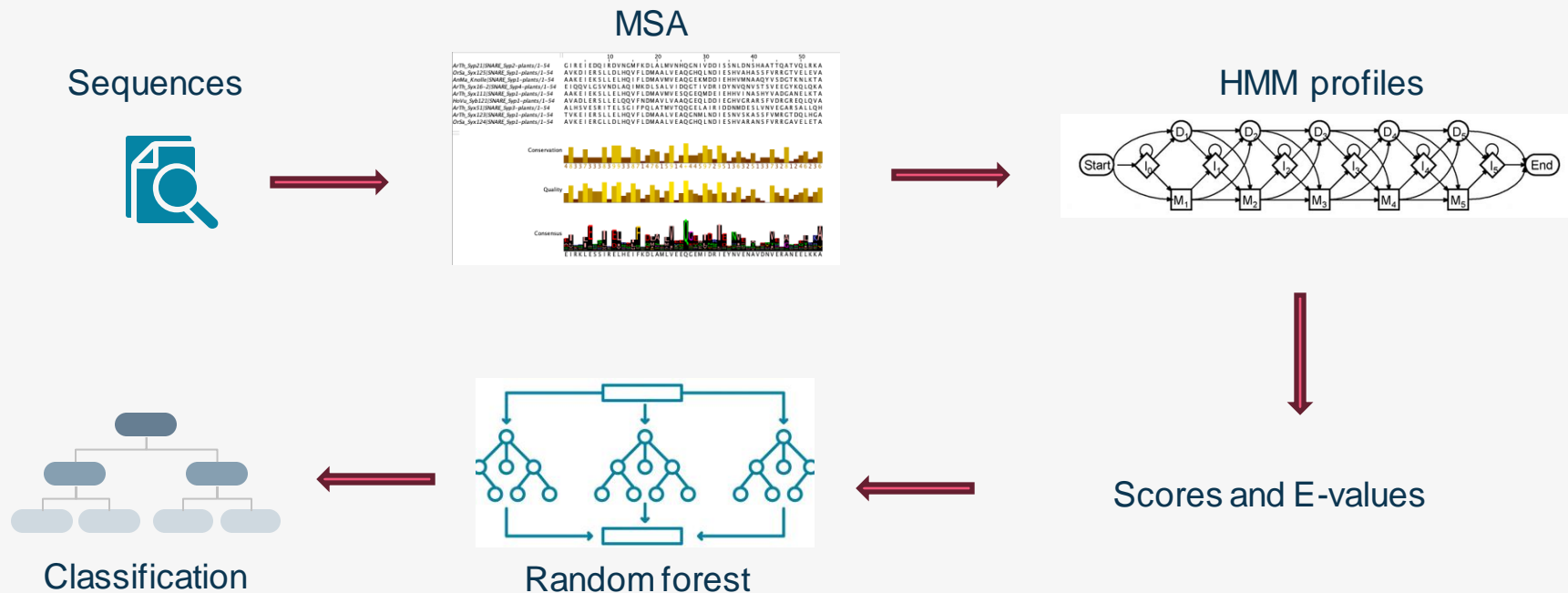# 01

# Background

# SNAREs

- Role : membrane fusion

- Conserved domain : SNARE motif

- Classification : Qa, Qb, Qc, R, SNAP

- **Goal** : automatise the classification with a model

- Classification can be used to hypothesize important informations about a protein (function, location, etc...)

# Classification – broad picture



Sequences

MSA

HMM profiles

Scores and E-values

Random forest

Classification

Goal : input = sequence | output = classification of the sequence (group & subgroup)
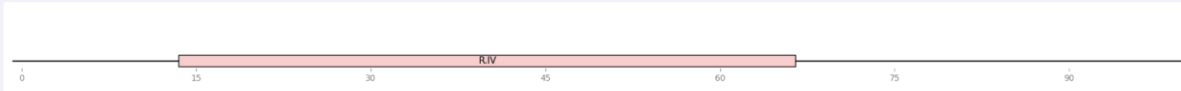
# 02

# MSA

Multiple Sequence Alignment

# Data Collection

- Tracey database : SNAREs protein collection

- Motif sequences and full sequences

- Sequences from specific taxa :
  - Archaeplastida --> Viridiplantae
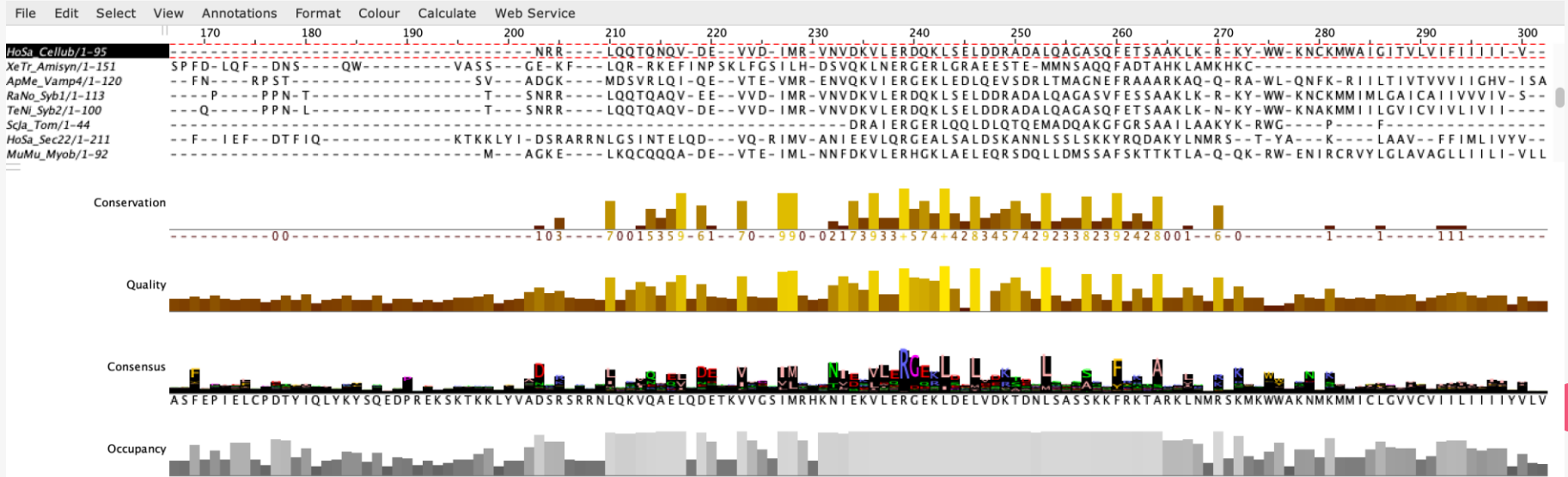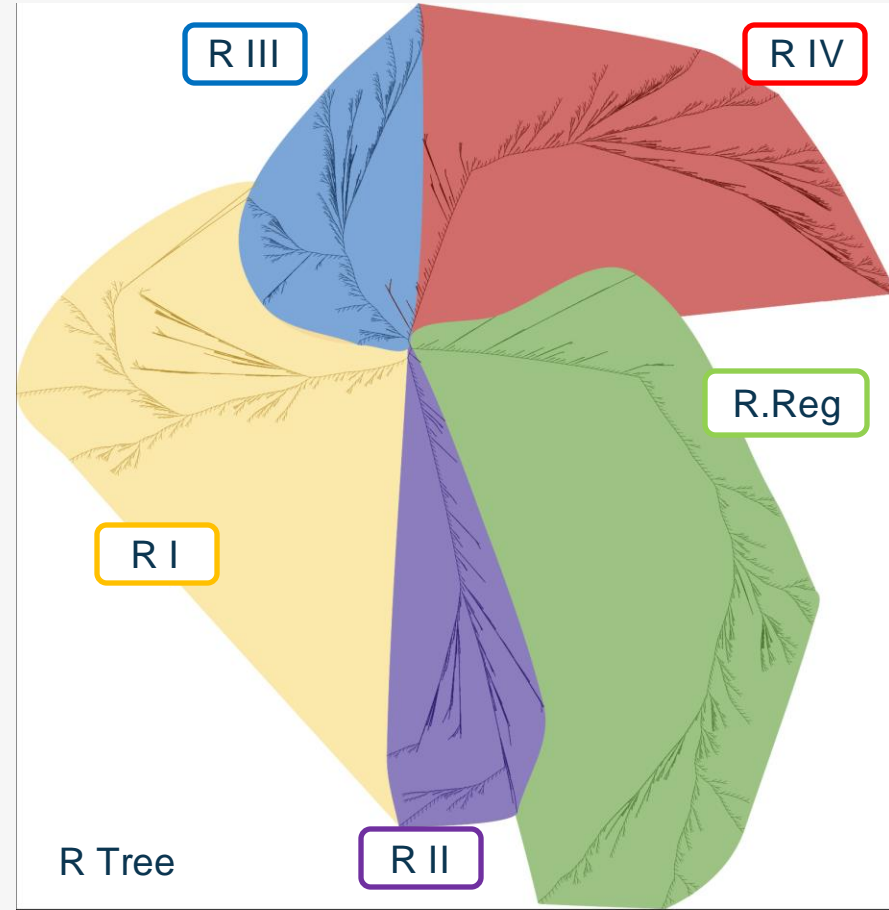  - Opisthokonta --> Metazoa and Fungi

# Alignment



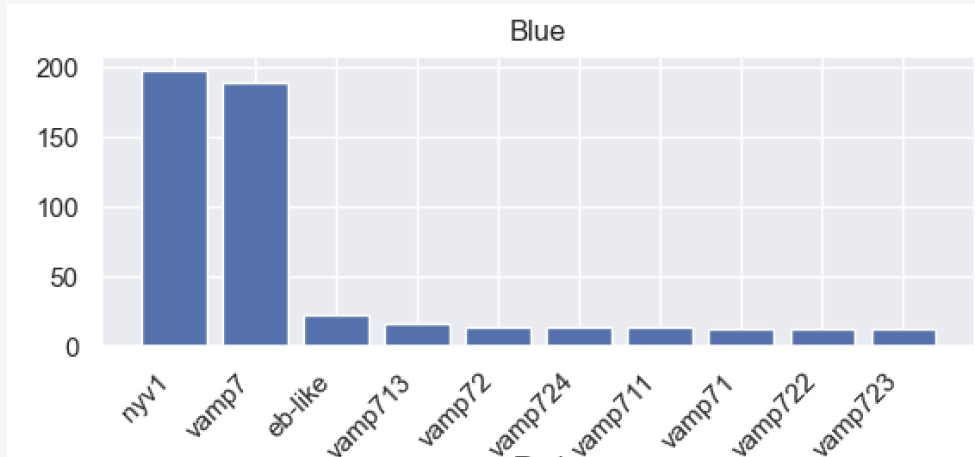HoSa_Cellub

Sequence with motif:

M S T G P T A A T G S N R R L Q Q T Q N Q V D E V V D I M R V N V D K V L E R D Q
K L S E L D D R A D A L Q A G A S Q F E T S A A K L K R K Y W W K N C K M W A I G I
T V L V I F I I I I I V W V V S S

# Identification of sub-groups

- Full sequences alignment to build the trees

- Average distance, visualisation with ITOL

- Labelling of sub-groups

- MSA for each of the sub-groups

# Sequences distribution across taxa

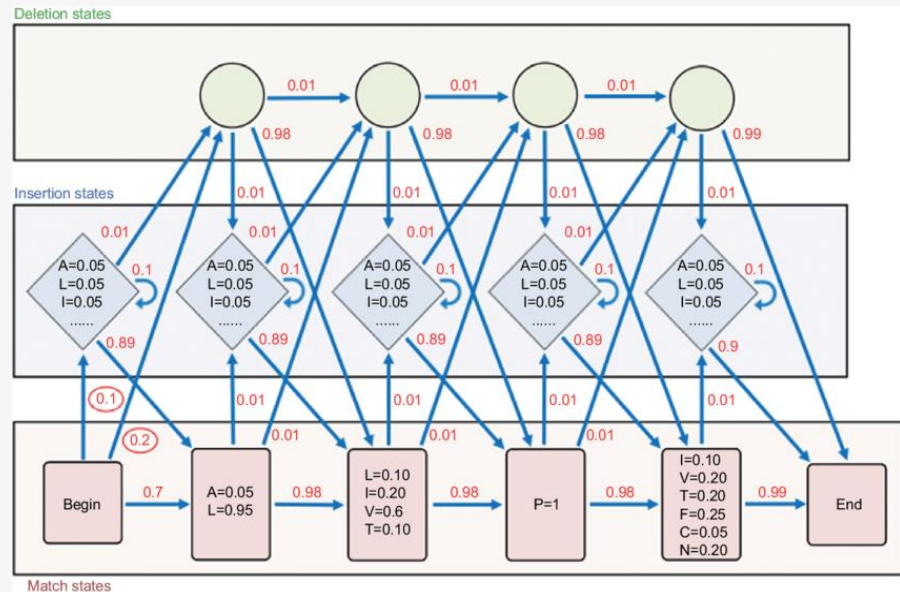|     | V   | M    | F    |
| --- | --- | ---- | ---- |
| Qa  | 471 | 3081 | 1481 |
| Qb  | 316 | 1042 | 972  |
| Qc  | 335 | 947  | 1327 |
| R   | 448 | 2294 | 1416 |
| SN  | 60  | 786  | 283  |

# 03

## HMM

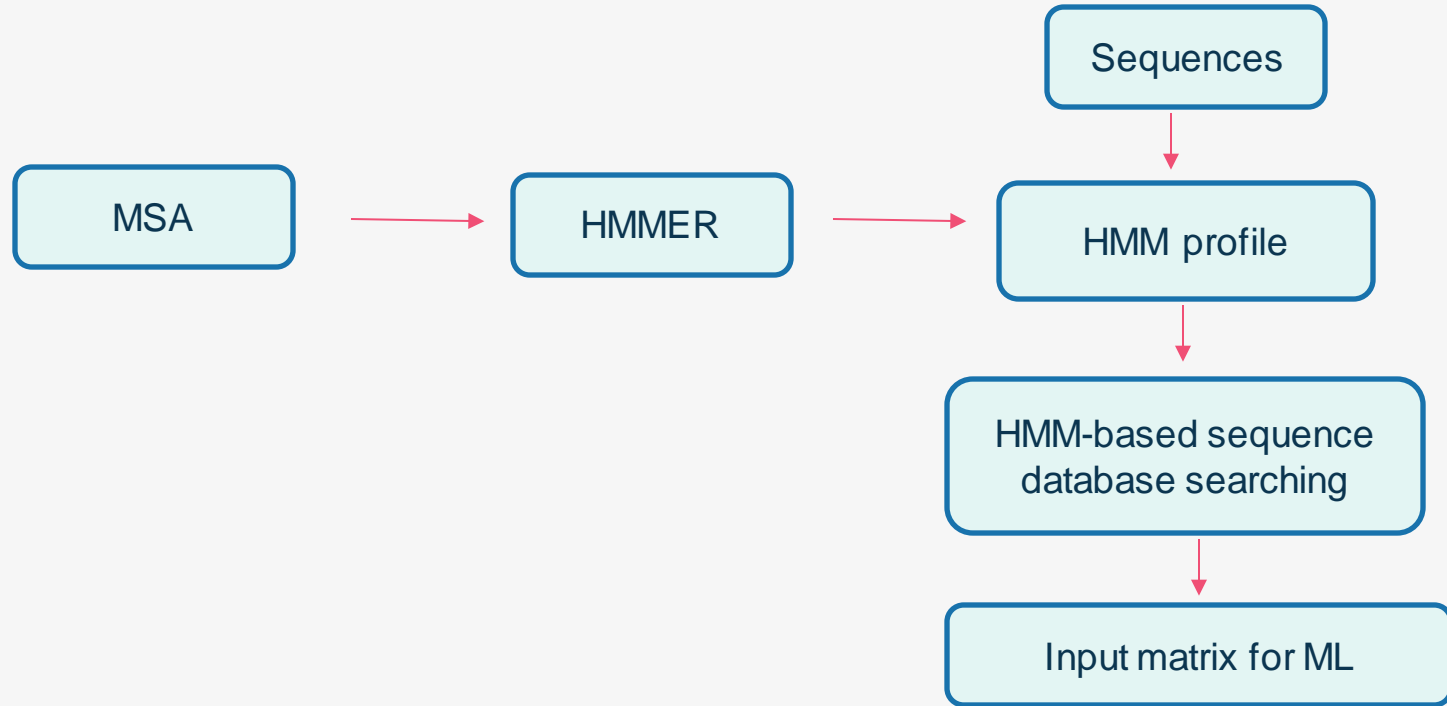Hidden Markov Models, Profiles and search

# What are HMM profiles?

- Based on multiple sequences alignements
- Probabilistic models used to represent a family of sequences

- They capture conserved and variable regions, as well as insertions and deletions

- Allow to determine how likely it is for a sequence to belong a specific group
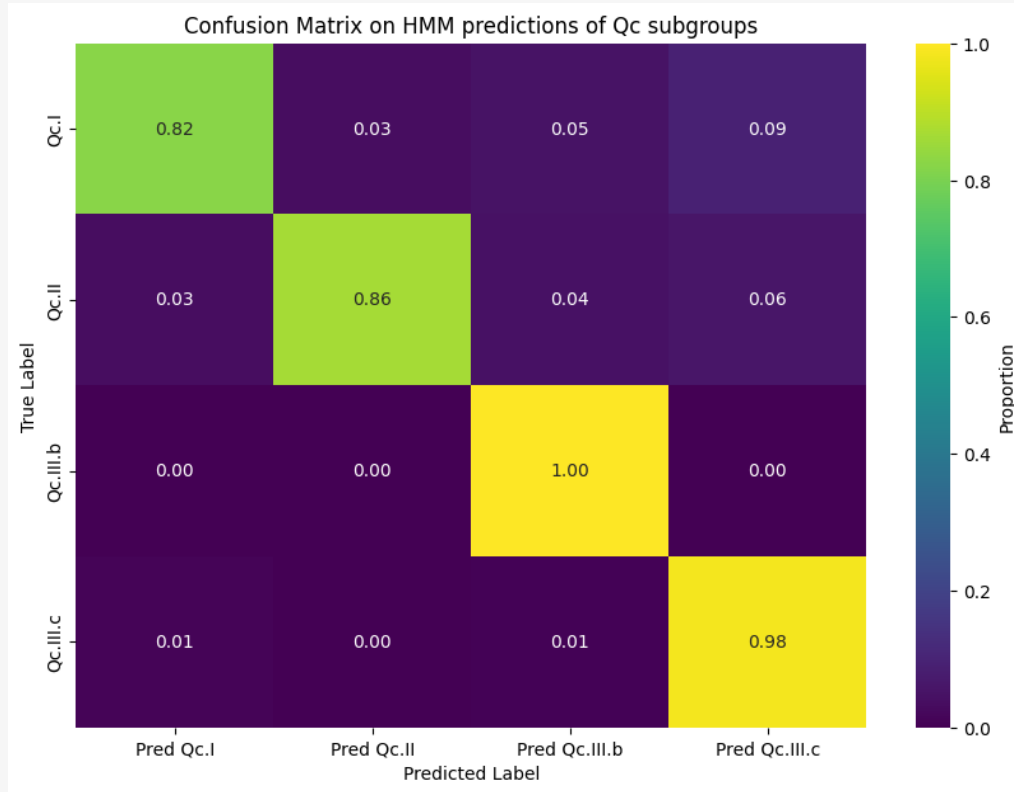
Arthur Gruber ResearchGate

# Building and using the HMM profiles

```
Sequences
   ↓
MSA → HMMER → HMM profile
                   ↓
         HMM-based sequence
         database searching
                   ↓
         Input matrix for ML
```

# HMM profiles performance

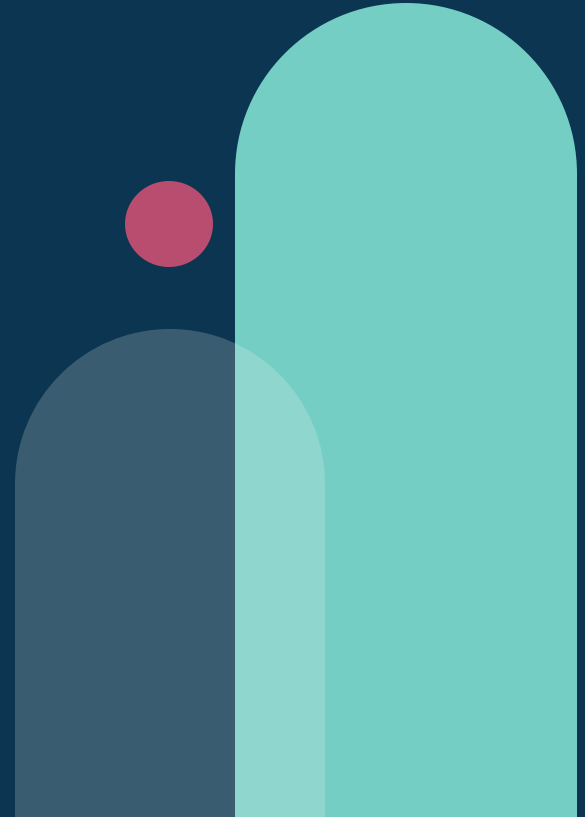

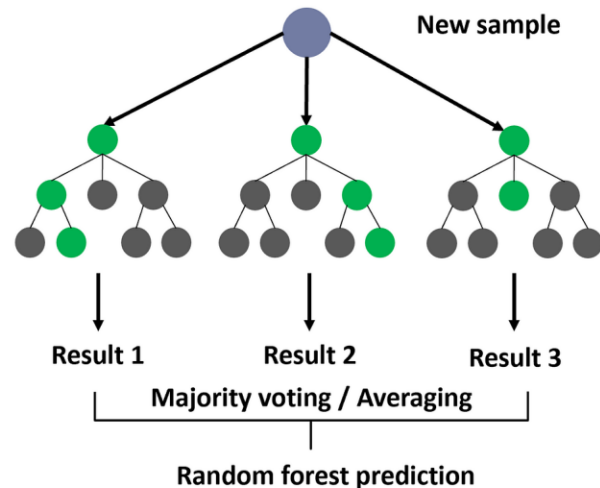Confusion Matrix on HMM predictions of Qc subgroups

Accuracy = 0.92

# 04

# Machine learning

Random forests

# Machine learning – Random forests

- Commonly used model for classification
- Ensemble learning technique
- Creation of multiple decision trees
- Aggregation of results
- High accuracy and reduced overfitting



Dr Roi Yeoshua, Medium

# Machine learning – input

| Sequence | Score/length Qa | EvalQa | Score/length Qb | Eval Qb | Score/length Qc | Eval Qc | Score/length R | Eval R | Label |
|---|---|---|---|---|---|---|---|---|---|
| DiOr_Syx1a | 0.86254 | 2.2e-79 | NaN | NaN | 0.076 | 6.6e-06 | NaN | NaN | Qa |
| TrVi_Bos1 | NaN | NaN | 0.47964 | 8.2e-40 | NaN | NaN | NaN | NaN | Qb |
| ChMy_Syx6 | 0.38248 | 9.4e-15 | 0.13983 | 0.0011 | 1.1024 | 6.2e-41 | NaN | NaN | Qc |
| … | … | … | … | … | … | … | … | … | … |

- No hits from the HMM profile do not return anything
- Need to set an appropriate value for the missing data for the model to work
- Different replacement values were tested, from 0.001 to 0.95
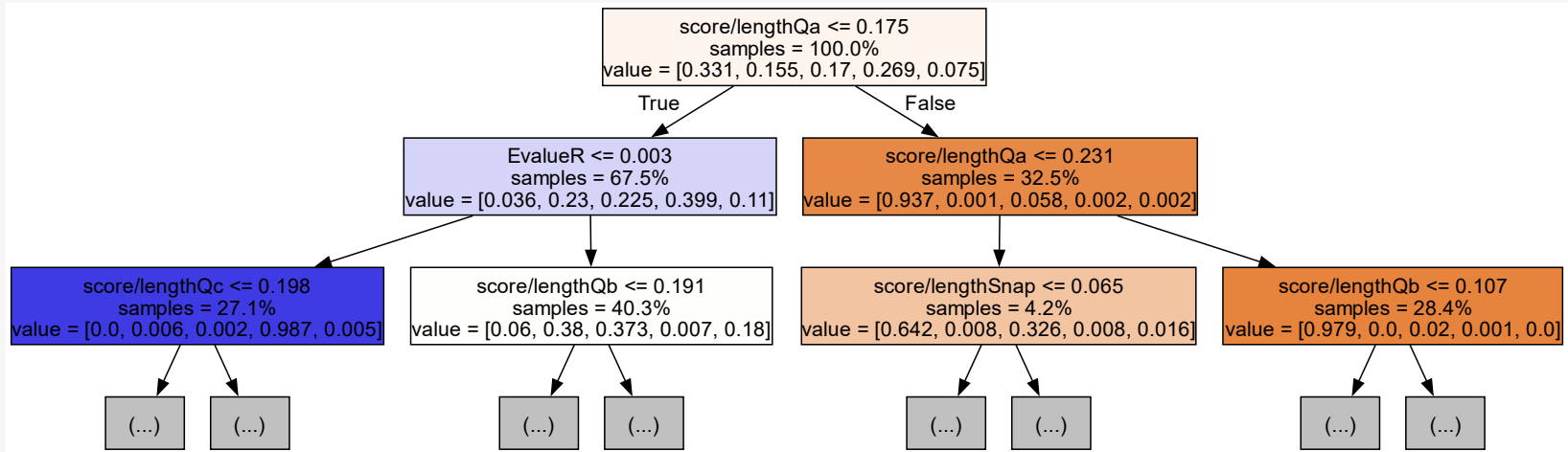- Specific to the score
- Specific to the E-value

# Machine learning – input

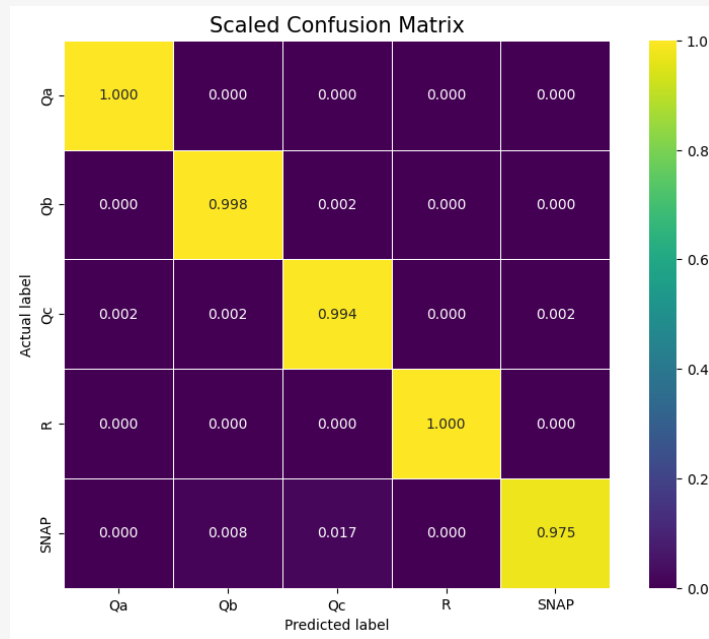| Sequence | Score/length Qa | EvalQa | Score/length Qb | Eval Qb | Score/length Qc | Eval Qc | Score/length R | Eval R | Label |
|---|---|---|---|---|---|---|---|---|---|
| DiOr_Syx1a | 0.86254 | 2.2e-79 | 0.01 | 0.9 | 0.076 | 6.6e-06 | 0.01 | 0.9 | Qa |
| TrVi_Bos1 | 0.01 | 0.9 | 0.47964 | 8.2e-40 | 0.01 | 0.9 | 0.01 | 0.9 | Qb |
| ChMy_Syx6 | 0.38248 | 9.4e-15 | 0.13983 | 0.0011 | 1.1024 | 6.2e-41 | 0.01 | 0.9 | Qc |
| … | … | … | … | … | … | … | … | … | … |

- No hits from the HMM profile do not return anything
- Need to set an appropriate value for the missing data for the model to work
- Best working values :
  Low replacement value for the score : 0.01
  High replacement value for the E-values : 0.9

# Machine Learning – Example tree
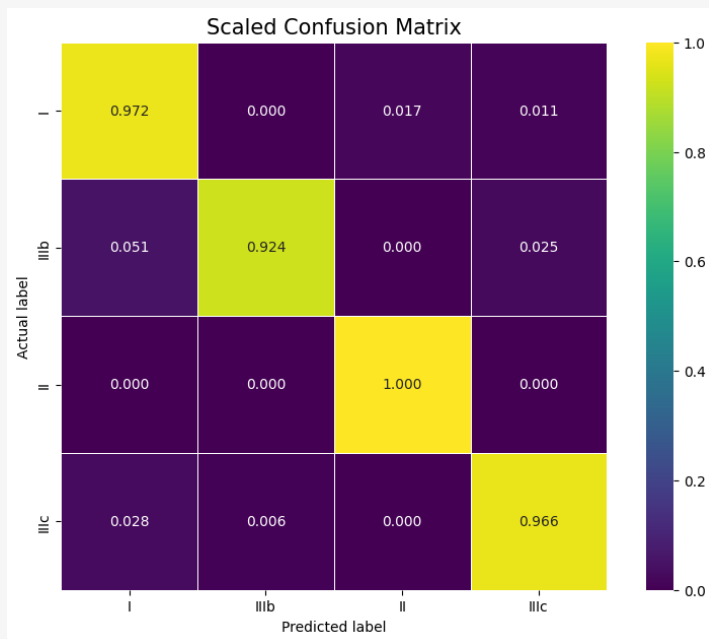
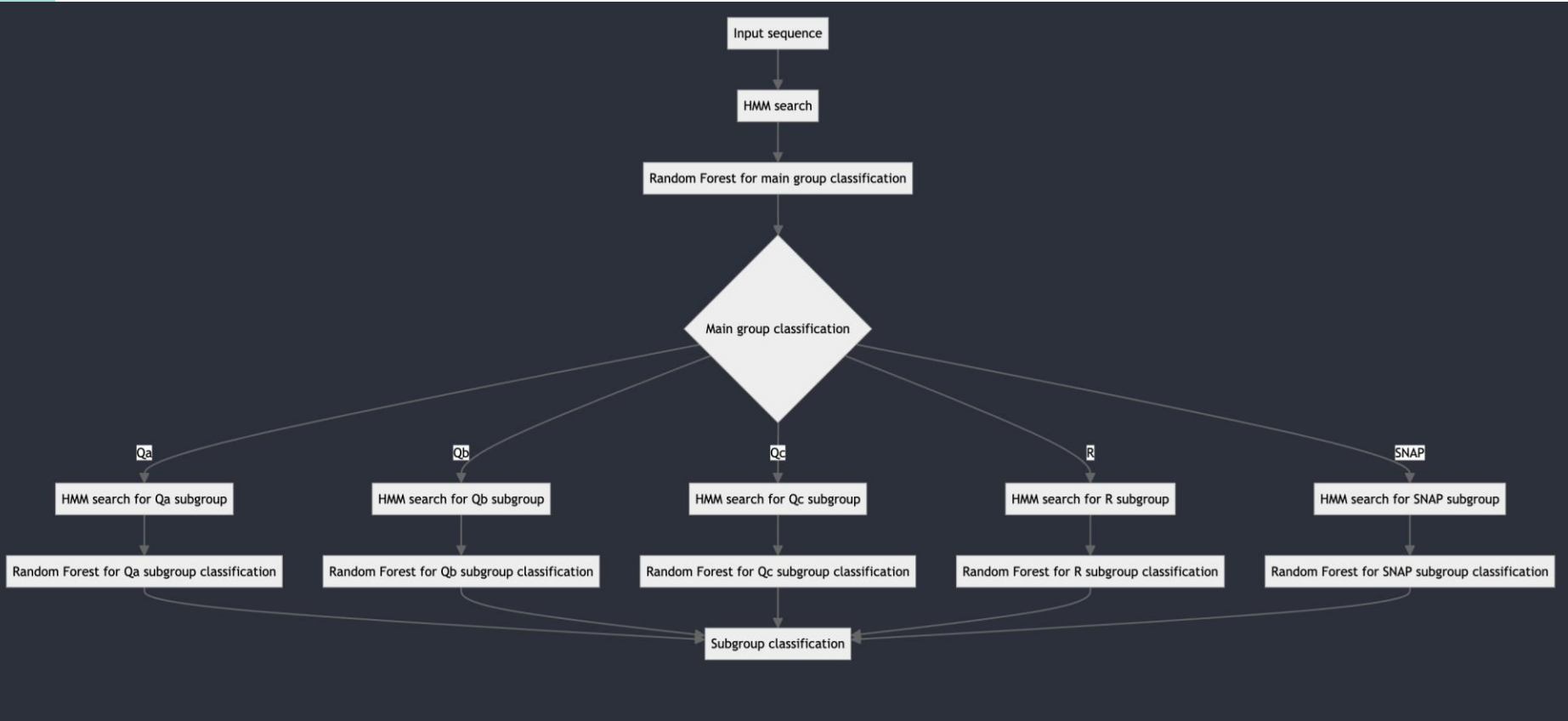# Machine learning – Main group results

- Confusion matrix



- Evaluation metrics:
- Accuracy: Correctly classified instances out of the total instances
- Precision: Ratio of TPs to the sum of TPs and FPs
- Recall: Ratio of TPs to the sum of TPs and FNs
- F1 score: Harmonic mean of precision and recall

- Accuracy: 0.9996
- Precision: 0.9996
- Recall: 0.9996
- F1: 0.9996

# Machine learning – Qc subgroups results

- Confusion matrix



- Evaluation metrics:
- Accuracy: Correctly classified instances out of the total instances
- Precision: Ratio of TPs to the sum of TPs and FPs
- Recall: Ratio of TPs to the sum of TPs and FNs
- F1 score: Harmonic mean of precision and recall

- Accuracy: 0.967
- Precision: 0.967
- Recall: 0.967
- F1: 0.967

# 05

# What's next ?

Future things to implement

# Challenges

- Using IqTree for better classification

- Using a tree to classify the main groups

- Investigating other ML models to upgrade the HMMs performance

# Feedback

- It was really interesting to work practically

- Projects had a real biological meaning and were based on real data
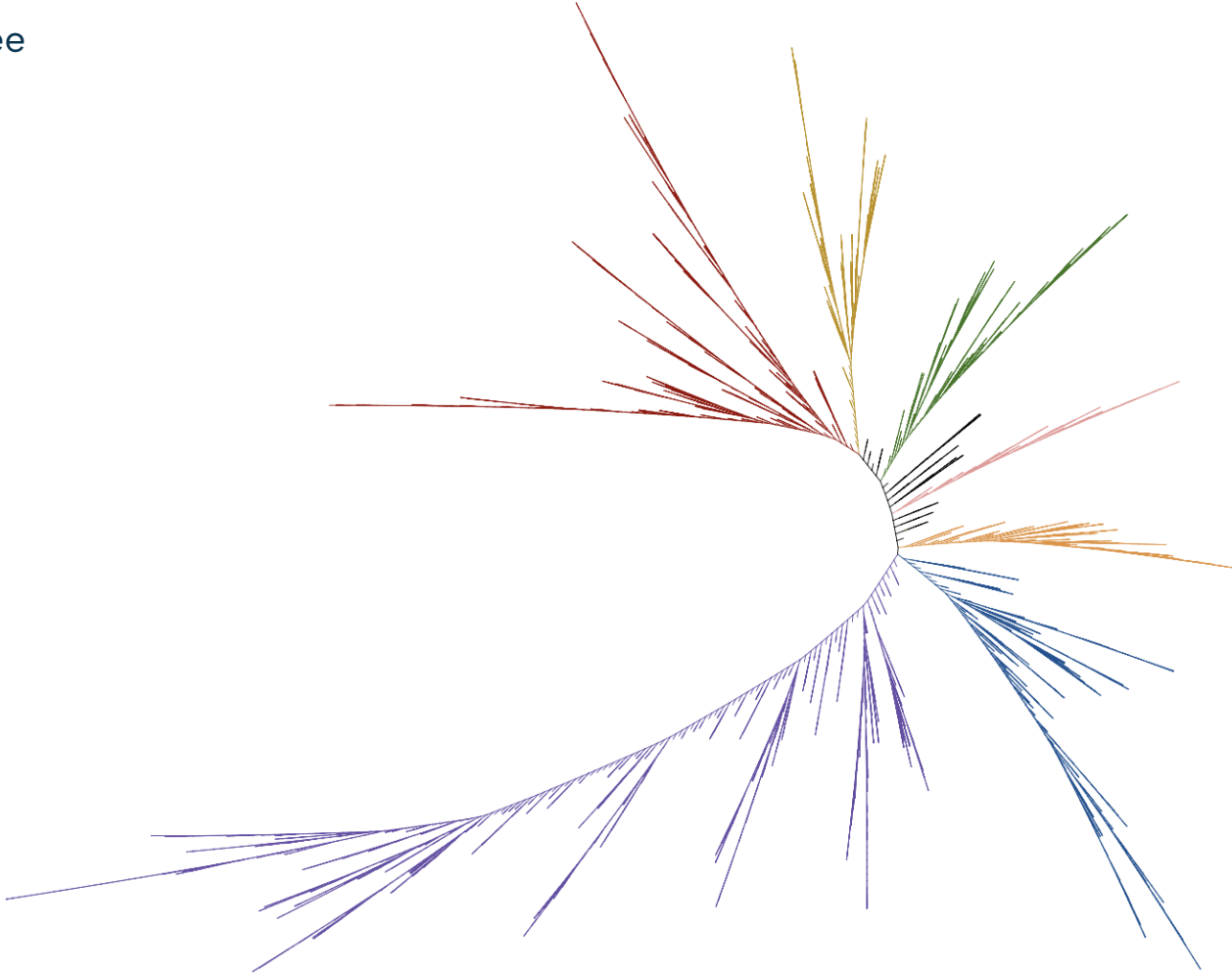
- More time for the presentation would be appreciated
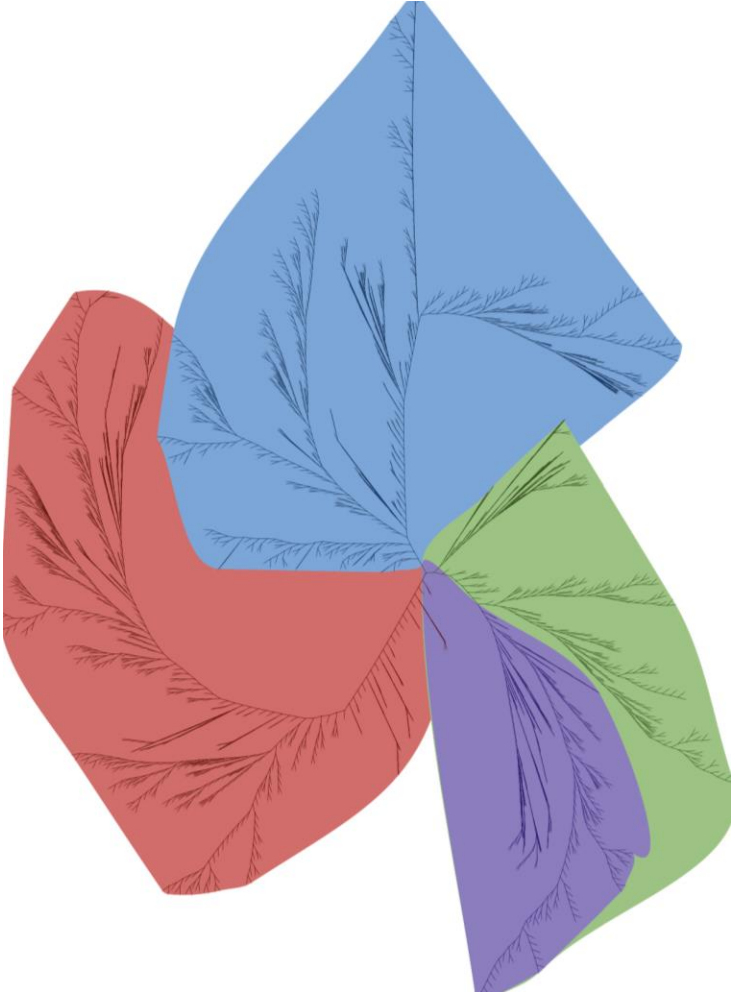
# Thanks!

**Do you have any questions?**

Gabriel Chiche, Marius Audenis, Leana Ortolani
Supervisor Carlos Pulido Quetglas

Qa Tree

# Qc Tree

SNAP Tree