# Rab Protein Classification

Gaëtan, Luca, Sydney

Mentor: Deepak Yadav

Email: gaetan.ray@unil.ch, luca.paley@unil.ch, sydney.fleming@unil.ch

Classify the Rabs into their group

Rab-43 (Alpha fold)

**Group 1**  **Group 2**  **Group 3**  **Group 4**  **Group 5**  **Group 6**

Rab-43 (Alpha fold)

GTPase with functions in tethering in **vesicular trafficking**



David Goodsell

# What are Rab proteins?

GTPase →
exchange GDP against
GTP to activate effectors



Wang S. *et al*, 2017

Large **family** of proteins



Klöpper, T.H. *et al.* 2012

# … and why study them?

**Griscelli syndrome**
- partial albinism
- immunological defects
- primary neurological dysfunctions

**Choroideremia**
- Blindness

**Neurological disorders**

Rab-associated
**diseases**

**Pipeline and methods**

# What does our data look like?

```
>PoTr_Rab1C
MNPEYDYLFKLLLIGDSGVGKSCLLLRFADDSYLESYISTI
GVDFKIRTVEQDGKTIKLQIWDTAGQERFRTITSSYYRGAH
GIIVVYDVTDQESFNNVKQWLNEIDRYASENVNKLLVGNKC
DLTANKVVSYETAKAFADEIGIPFMETSAKNATNVEQAFMA
MAAEIKNRMASQPAMNNARPPTVQIRGQPVNQKSGCSS
>PaTe_Rab1_1
MSLQQEYDYLFKILLIGNSAVGKSSLLLRFADNVFNESFLP
TIGVDFKIRTFDLNGKTVKLQIWDTAGQERFKTITNSYYKG
AHGIILVYDVTDKQSFKDVENWLAEVEKYANENVVRVLVGN
KVDLESKREVTSEEGKELADSLNIRFIETSAKNSSNVEKAF
ITLANEIKAKVAKSSEAIPVKTGPRITPDQQQNTVKDTGCC
>BrMa_Rab1
MVSINPEYDYLFKLLLIGDSGVGKSCLLLRFADDTYTESYI
STIGVDFKIRTIDLNGKTIKLQIWDTAGQERFRTITSSYYR
GAHGIIVVYDITDQESFNNVKQWLQEIDRYACENVNKLLVG
NKCDLIIRRAVEHSAAKEYADQLGIPFLETSAKSSTNVEQA
FLTMASEIKNRMGPIQQVGTGPSVRIGGSQPVNEKKSGGCC
>CaAl_Ypt1
MNNEYDYLFKLLLIGDSGVGKSCLLLRFADDTYTPDYISTI
GVDFKIRTIELDGKTIKLQIWDTAGQERFRTITSSYYRGAH
GIIIVYDVTDQESFNNVKQWLQEIDRYATGGVMKLLVGNKA
DLSDKKIVEYTAAKEFADALDIPFLETSALSSTNVEQAFYT
MARQIKAQMTNNANAGNAANAKGKSNVNLRGESLTSNQSNS
CC
>PaTe_Rab1
MIKEYDYLFKLVIIGNSGVGKSSLLLRFADDQFSESYLTTI
GVDFRFRTLPIDGKNVKLQIWDTAGQERFRTITSAYYKGAD
GIVMVYDVTQGQSFDDIDKFWLHEVESYGEKNVQLLIIGNK
NDLDEQKQVETSKAEEYCKSHNMLFMECSAKTADHVNNAFL
ELSRKLMAKKDASQPPKTTNTTSNASQQSQSRGQTNTNTQQ
SKQLSAGNTNQKKQKDGGCC
```

## Number of Rab proteins per group



**Unequal** distribution

# Methods used


KNN


Decision tree


Random forest

# First feature: GAAC

- GAAC = grouped amino acid composition
- **Physico-chemical** properties

| | Aliphatic | Aromatic | Positively charged | Negatively charged | Uncharged |
|---|---|---|---|---|---|
| rab1 | 0,2 | 0,3 | 0,1 | 0,3 | 0,1 |
| rab2 | 0,3 | 0,1 | 0,3 | 0,1 | 0,2 |
| rab3 | 0,1 | 0,2 | 0,3 | 0,3 | 0,1 |

# GAAC predictor using KNN



Confusion matrix from KNN model with k = 5 on Rab extracted GAAC features

**Accuracy: 53%**
Precision: 0.523
Recall: 0.526
F1-Score: 0.519

# Second feature: CKSAAP

- CKSAAP = Composition of **k-spaced** amino acid pairs

| k | AA | LT |
|---|------|-------|
| k=1 | AxA | LxT |
| k=3 | AxxxA | LxxxT |
| k=5 | AxxxxxA | LxxxxxT |

Protein sequence : DSAWAELSGCIKT

# CKSAAP predictor using KNN



Confusion matrix from KNN model with k = 5 on Rab extracted CKSAAP features

**Accuracy: 94.7%**
Precision: 0.950
Recall: 0.947
F1-Score: 0.945

PCA

**Accuracy: 93.4%**
Precision: 0.936
Recall: 0.934
F1-Score: 0.934

15

Decision Tree



Confusion matrix from decision tree model with max depth = 22 on Rab extracted CKSAAP features after PCA (10PCs)

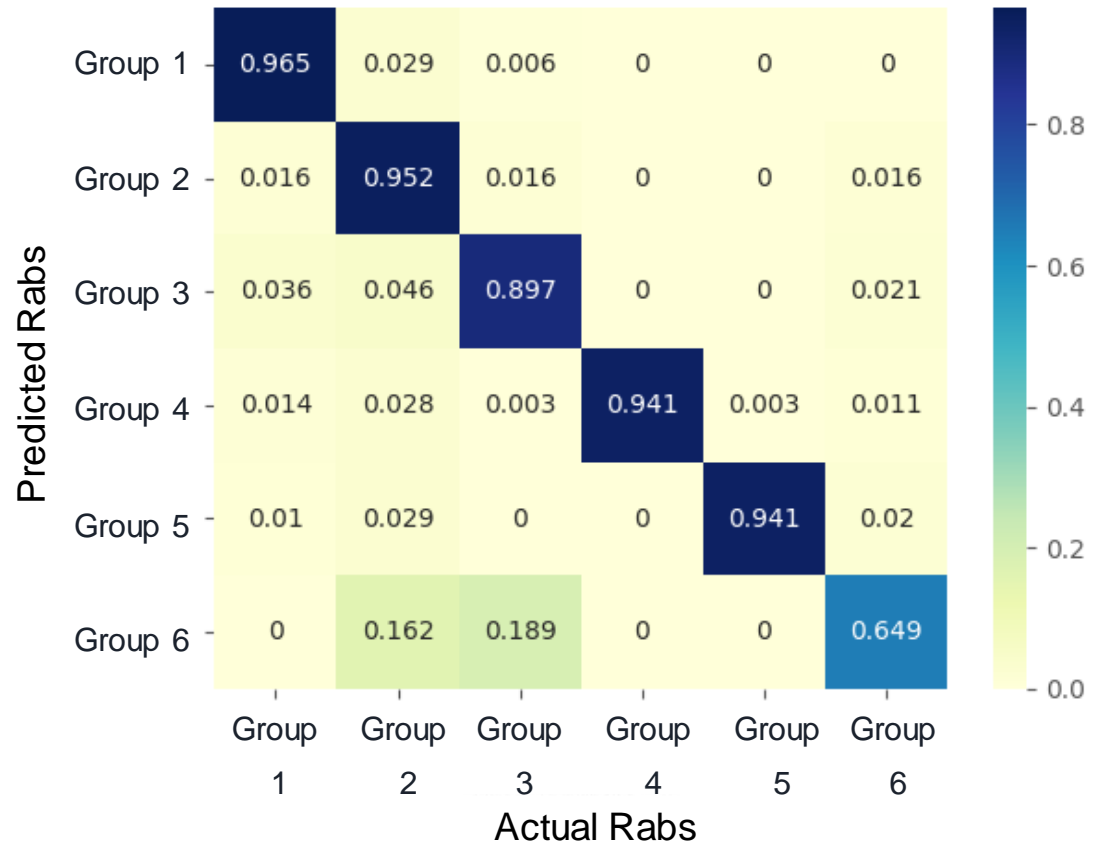|  | Group 1 | Group 2 | Group 3 | Group 4 | Group 5 | Group 6 |
|---|---|---|---|---|---|---|
| Group 1 | 0.92 | 0.02 | 0.02 | 0.03 | 0 | 0.01 |
| Group 2 | 0.05 | 0.81 | 0.07 | 0.02 | 0.01 | 0.04 |
| Group 3 | 0.02 | 0.06 | 0.89 | 0 | 0 | 0.04 |
| Group 4 | 0.02 | 0.03 | 0.01 | 0.93 | 0 | 0 |
| Group 5 | 0.01 | 0.06 | 0 | 0 | 0.92 | 0.01 |
| Group 6 | 0.04 | 0.27 | 0.1 | 0.02 | 0.02 | 0.55 |

Predicted Rabs / Actual Rabs

**Accuracy: 89,28%**
Precision: 0.839
Recall: 0.835
F-score: 0.837

16

# CKSAAP predictor using **random forest**



Confusion matrix from random forest model with max on Rab extracted CKSAAP features after PCA (10PCs)

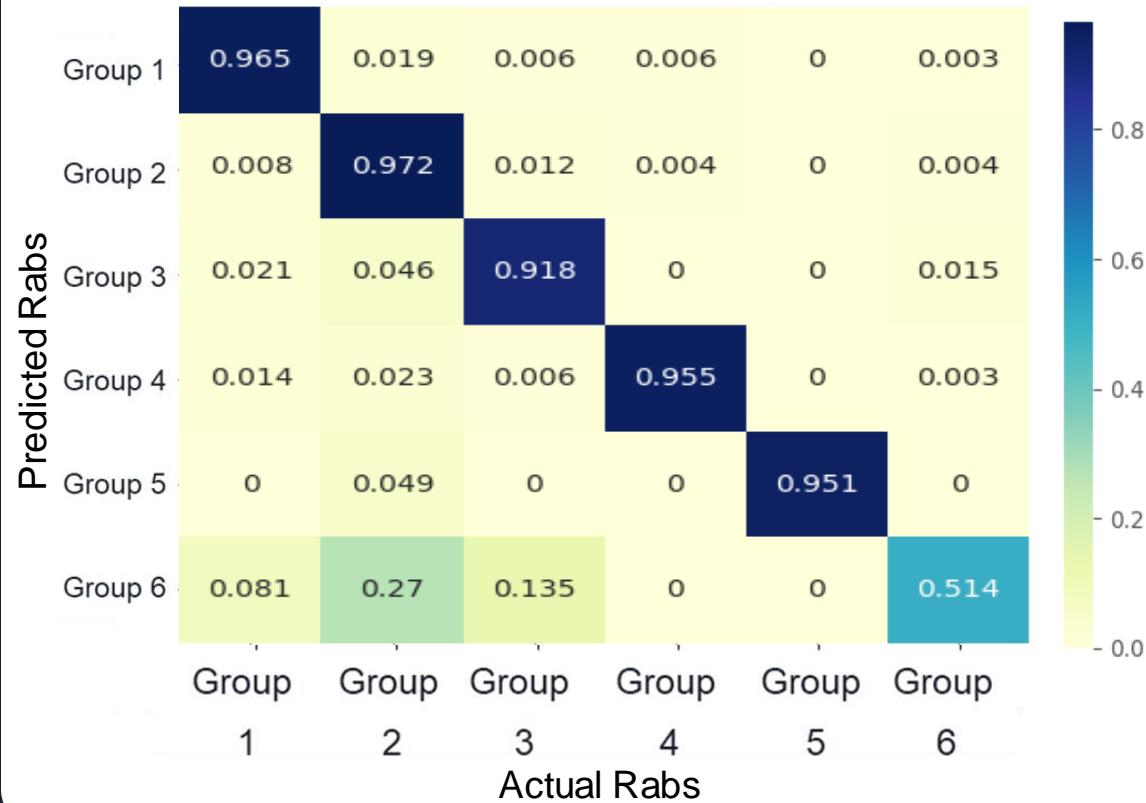| Predicted Rabs | Group 1 | Group 2 | Group 3 | Group 4 | Group 5 | Group 6 |
|---|---|---|---|---|---|---|
| Group 1 | 0.965 | 0.019 | 0.006 | 0.006 | 0 | 0.003 |
| Group 2 | 0.008 | 0.972 | 0.012 | 0.004 | 0 | 0.004 |
| Group 3 | 0.021 | 0.046 | 0.918 | 0 | 0 | 0.015 |
| Group 4 | 0.014 | 0.023 | 0.006 | 0.955 | 0 | 0.003 |
| Group 5 | 0 | 0.049 | 0 | 0 | 0.951 | 0 |
| Group 6 | 0.081 | 0.27 | 0.135 | 0 | 0 | 0.514 |

Actual Rabs

**Accuracy : 94.2%**
Precision : 0.942
Recall : 0.941
F-score : 0.941

Number of Rab proteins per group

Oversampling: Sample with replacement

## Number of Rab proteins per group



**Testing set**

**Training set**

**Oversampling:** Sample with replacement

Confusion matrix from KNN model with k = 5 on Rab extracted CKSAAP after PCA (10PCs) with oversampling

**Number of Rab proteins per group**
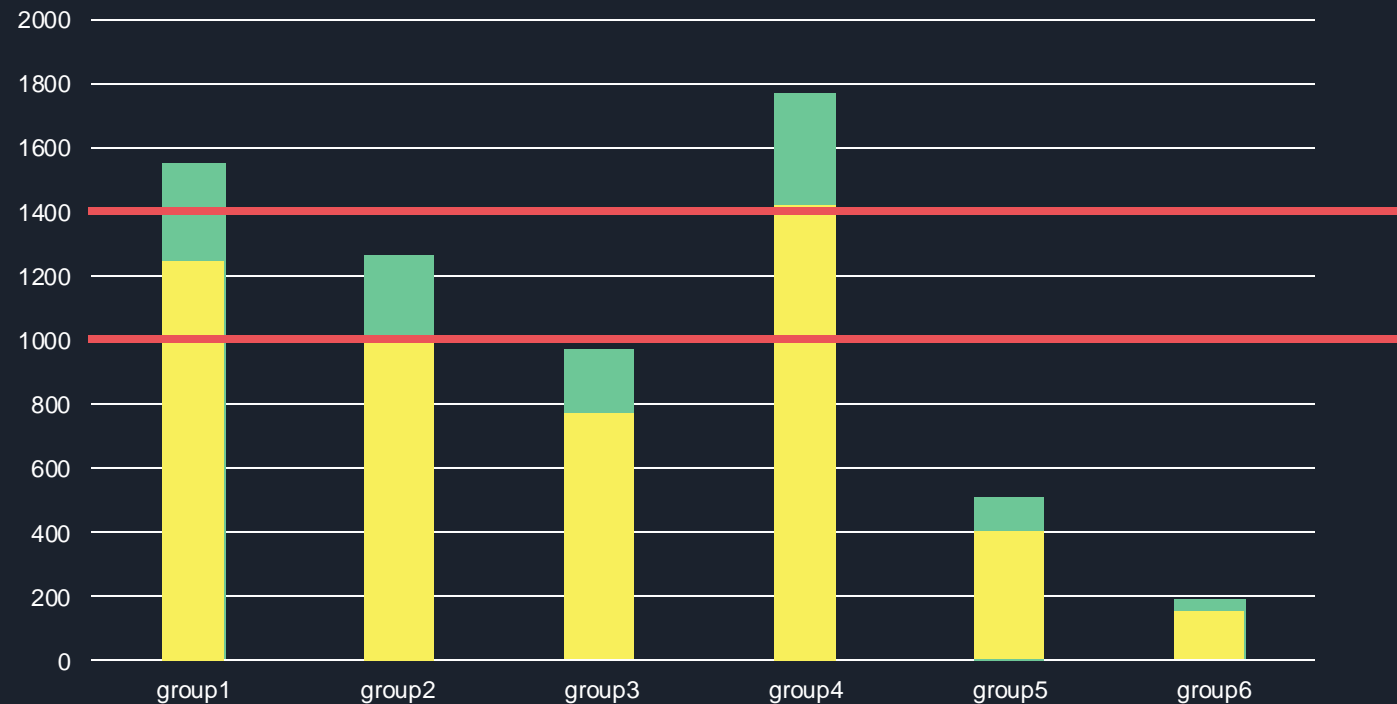
Testing set

Training set

## Number of Rab proteins per group



Testing set

Training set

Arf

SNARE

Ras

Group 0
Non-Rabs

k=11

# CKSAAP predictor with improved model

## KNN



Confusion matrix from KNN model with k = 11 on Rab extracted CKSAAP features with oversampling

**Accuracy: 90.2%**   Recall: 0.902
Precision: 0.912   F1-Score: 0.905

## Random forest



Confusion matrix by Random Forest model on Rab extracted CKSAAP feature with oversampling

**Accuracy: 88.2%**   Recall: 0.882
Precision: 0.891   F1-Score: 0.874

- **Physico-chemical** properties (CTD= Conjoint Triad Descriptors)

| | Hydro-phobicity | Polarity | Van der Waals volume | Charge | … |
|---|---|---|---|---|---|
| rab1 | 0,212 | 0,66454 | 0,117 | 0,5888 | |
| rab2 | 0,42588 | 0,1588 | 0,3665 | 0,05 | |
| rab3 | 0,61 | 0,258 | 0,557 | 0,255 | |

Confusion matrix from KNN model with k = 11 on Rab extracted CTD after PCA (10PCs) feature with oversampling

**Accuracy: 69%**
Precision: 0.692
Recall: 0.686
F1-Score: 0.686

# CKSAAP and CTD predictors combined

## KNN



Confusion matrix from KNN model with k = 11 on Rab extracted CKSAAP and CTD features with oversampling

**Accuracy: 73%**　　Recall: 0.727
Precision: 0.731　　F1-Score: 0.727

## Random forest



Confusion matrix from Random Forest model on Rab extracted CKSAAP and CTD features with oversampling

**Accuracy: 73%**　　Recall: 0.730
Precision: 0.747　　F1-Score: 0.728

- **Physico-chemical** properties (KSCT= K spaced conjoint triad)

|  | Hydro-phobicity | Polarity | Van der Waals volume | Charge | … |
|---|---|---|---|---|---|
| rab1 | 0,212 | 0,66454 | 0,11 | 0,5888 |  |
| rab2 | 0,02588 | 0,0588 | 0,3665 | 0,005 |  |
| rab3 | 0,6 | 0,258 | 0,557 | 0,255 |  |

# KSCT predictor

## KNN



Confusion matrix from KNN model with k = 11 on Rab extracted KSCT feature with oversampling

**Accuracy: 87.9%**   Recall: 0.879
Precision: 0.888   F1-Score: 0.881

## Random forest



Confusion matrix from Random Forest model on Rab extracted KSCT feature with oversampling

**Accuracy: 89.2%**   Recall: 0.892
Precision: 0.900   F1-Score: 0.890
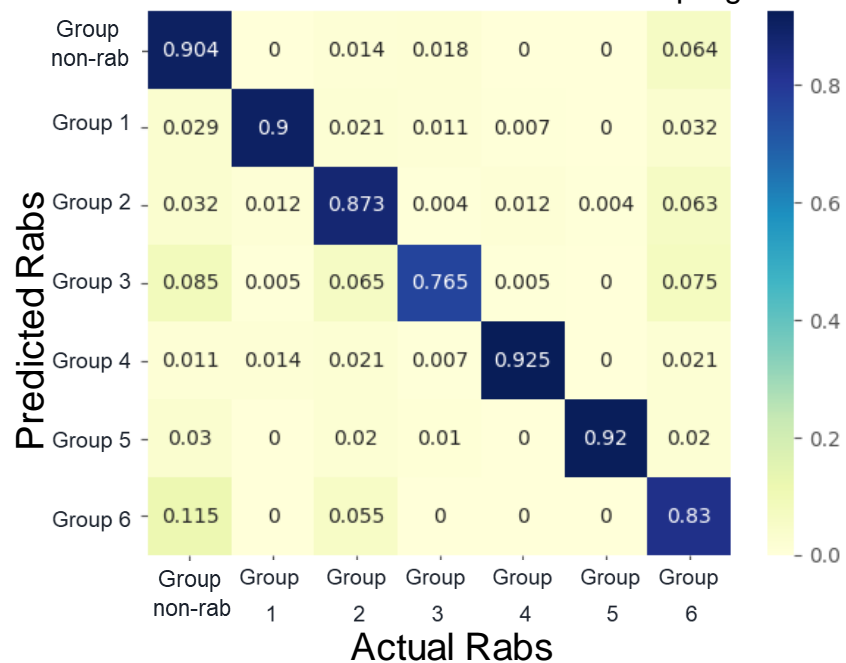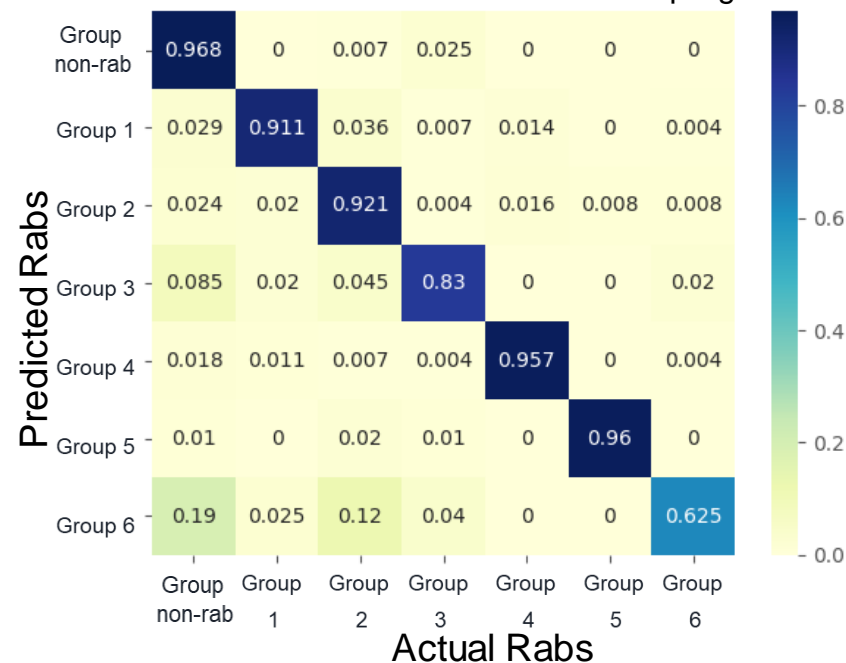
# KSCT and CKSAAP predictors combined

## KNN

Confusion matrix from KNN model with k = 11 on Rab extracted KSCT and CKSAAP features with oversampling



**Accuracy: 89.3%**     Recall: 0.893
Precision: 0.901     F1-Score: 0.895

## Random forest

Confusion matrix from Random Forest model on Rab extracted KSCT and CKSAAP features with oversampling
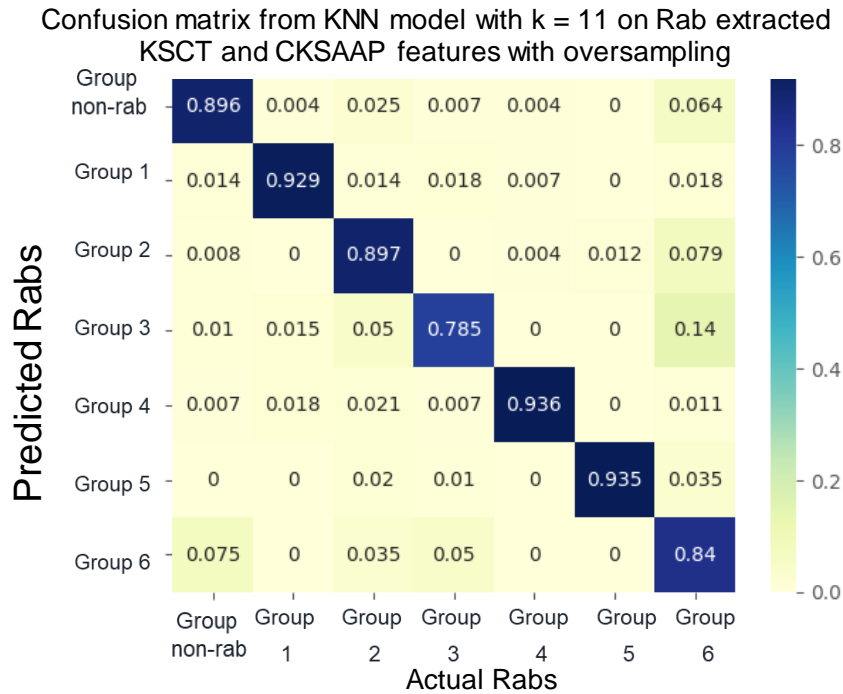


**Accuracy: 88.2%**     Recall: 0.882
Precision: 0.888     F1-Score: 0.876

Confusion matrix from KNN model with k = 11 on Rab extracted CKSAAP features with oversampling

**Accuracy: 90.2%**
Precision: 0.912
Recall: 0.902
F1-Score: 0.905

Best model: KNN with CKSAAP feature

Combos of features do not add to the model

# Challenges

- **Unequal distribution** in classes

- **Unintuitive results**: best model is not a combination (?)

- **Lots of testing** and probing, limited time

- **Google colab**… not so easy to colab

- **Different kind of coding** but lots of help from our mentor and forums/Chat GPT

# Feedback

- **Highly concentrated working period** blur lines between work and personal time

- **Grade disproportionally based on presentation,** rather than work done

- **Little time** between intermediate and final presentation but **helpful and inspiring**

- Very **helpful mentor**

- **Apply theory** seen in class on an interesting project

- Learned a lot about **coding and problem solving**

34

# Thanks!

Questions?

**CREDITS:** This presentation template was created by **Slidesgo,** including icons by **Flaticon,** and infographics & images by **Freepik**

Email: gaetan.ray@unil.ch,  luca.paley@unil.ch,  sydney.fleming@unil.ch