

# Comparison of Grassland Type classification output (EUGW) with the Habitat Mapping Layer of Czechia

&

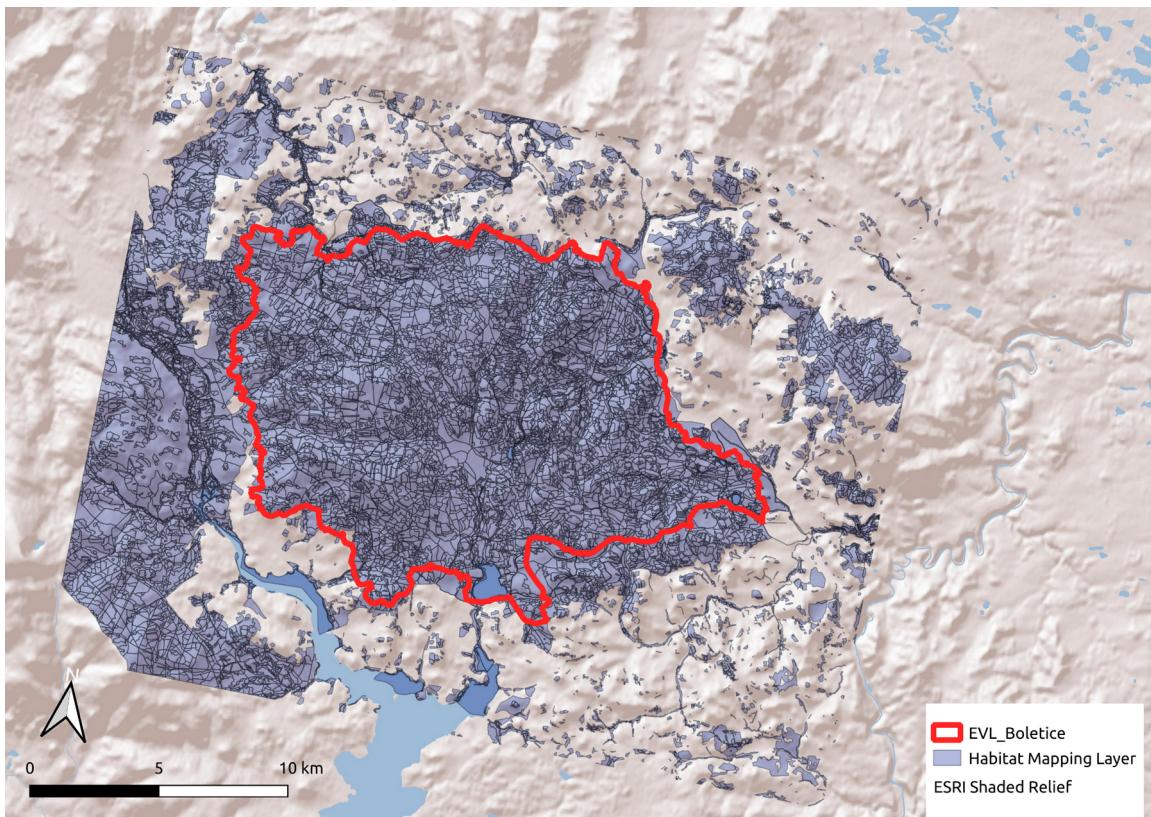
## Workflow for validation of other classification models outputs

Jakub & Dan, Habitat Pilot workshop in Oulu, June 2025

# Aims

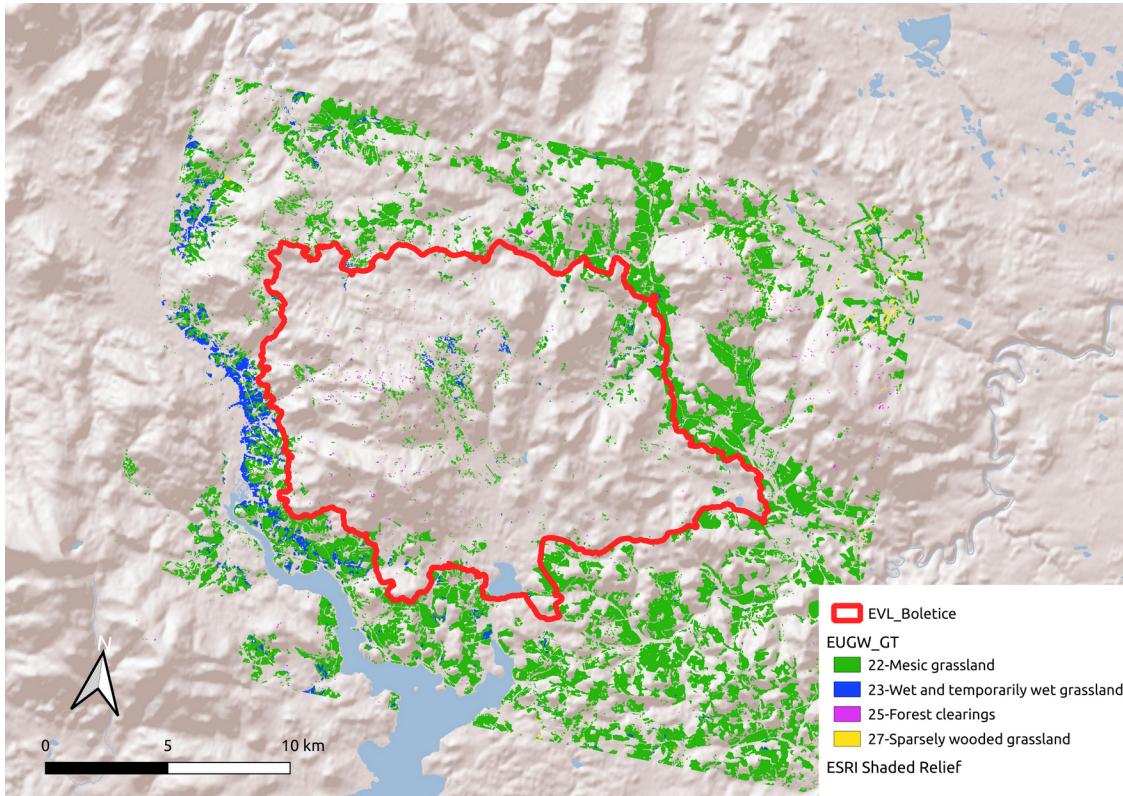
- **visual representation** of the comparison of the Grassland Type Characterization Component (GT) with the Habitat Mapping Layer of Czechia (VMB)
- Natura 2000 site Boletice and surrounding landscape
- proposal of validation process for classification models outputs based on vector polygon ground truth data

# Data – Habitat Mapping Layer of Czechia



- vector polygon data
- collected in field
- Czech national classification system
- Annex 1 also included

# Data – Grassland Type Characterization Component



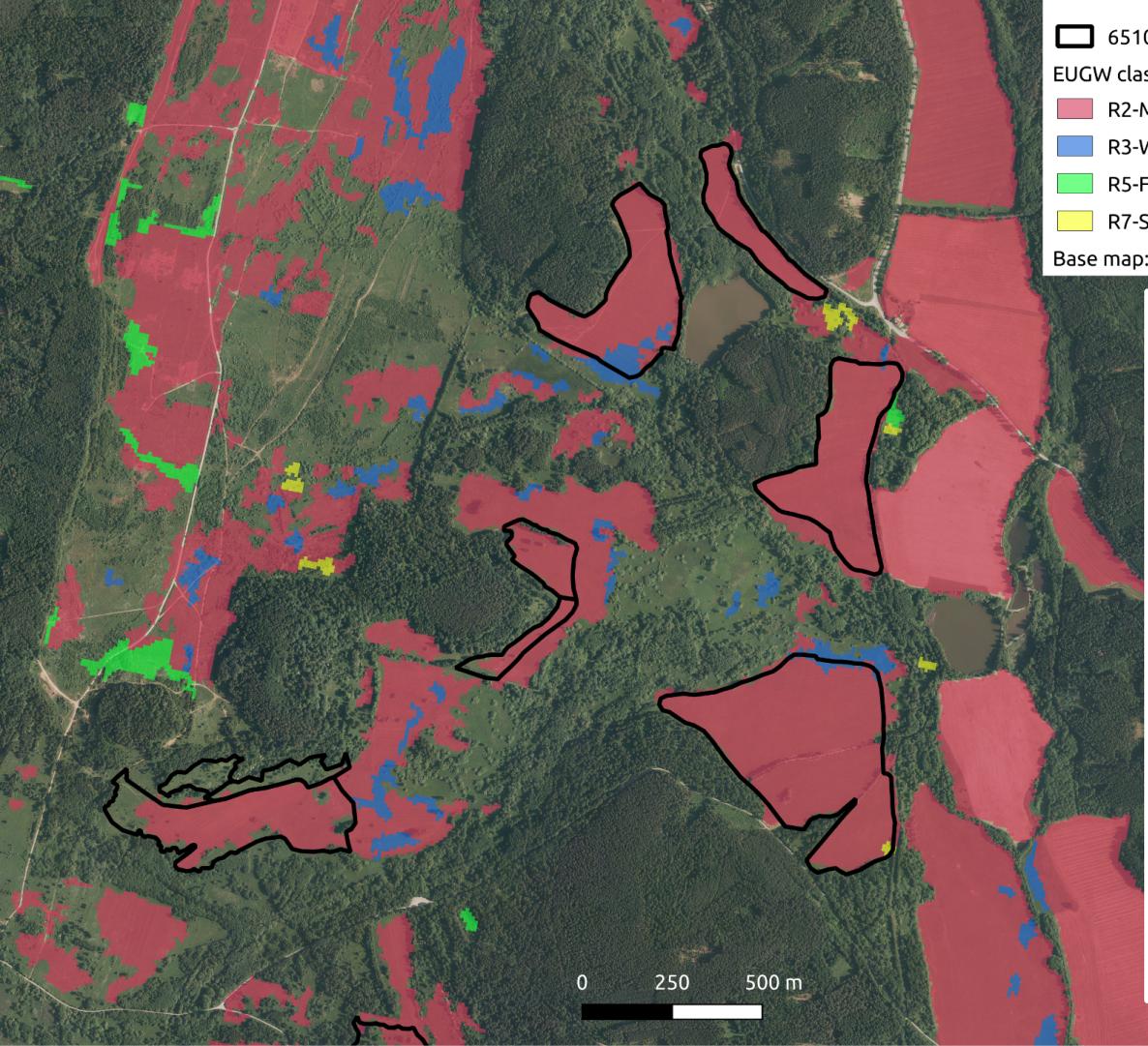
- raster layer
- grassland classification
- recognized 4 grassland types:
  - R2-Mesic grassland
  - R3-Wet and temporarily wet grassland
  - R5-Forest clearings
  - R7-Sparingly wooded grassland

# Habitat classification systems

EU-GW	GT 22 Mesic grasslands	GT 23 Wet and temporarily wet grassland	GT 25 Forest clearings	GT 27 Sparsely wooded grassland
EUNIS 2021	R2	R3	R5	R7
Annex 1	6270, 6510, 6520	6410, 6420, 6440, 6450, 6460, 6510	6430	6310, 6530, 9070
Czech	T1.1, T1.2, T1.3	T1.4, T1.5, T1.7, T1.9, T1.10	M5, M7, A4.1, A4.2, A4.3, T1.6, T1.8, T4.1, T4.2	NA

# Habitats present in study area

EU-GW	GT 22 Mesic grasslands	GT 23 Wet and temporarily wet grassland	GT 25 Forest clearings	GT 27 Sparsely wooded grassland
EUNIS 2021	R2	R3	R5	R7
Annex 1	6270, <b>6510</b> , <b>6520</b>	<b>6410</b> , 6420, 6440, 6450, 6460, <b>6510</b>	<b>6430</b>	6310, 6530, 9070
Czech	T1.1, T1.2, T1.3	T1.4, T1.5, T1.7, T1.9, T1.10	M5, M7, A4.1, A4.2, A4.3, <b>T1.6</b> , T1.8, T4.1, <b>T4.2</b>	NA



6510 habitats (belongs to R2)

EUGW classification raster

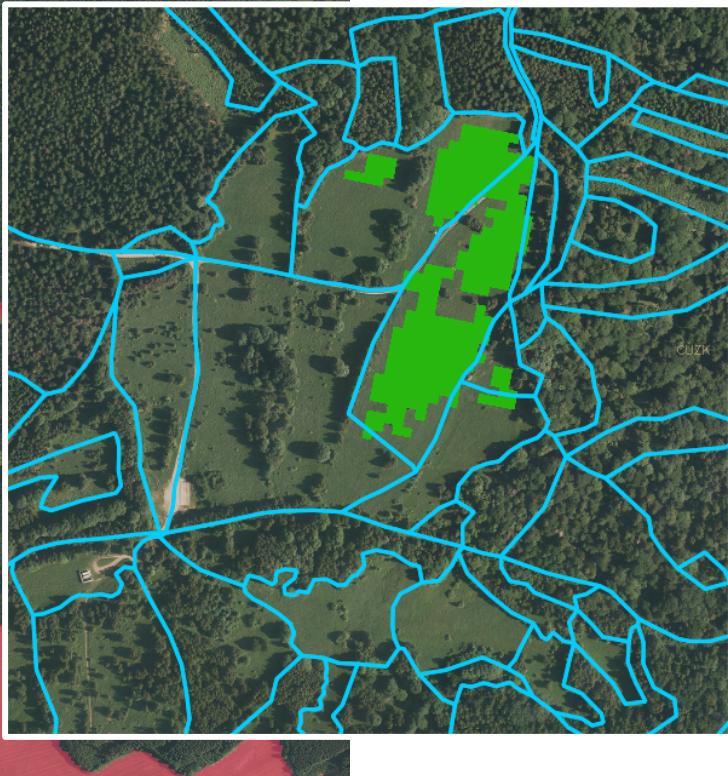
R2-Mesic grassland

R3-Wet and temporarily wet grassland

R5-Forest clearings

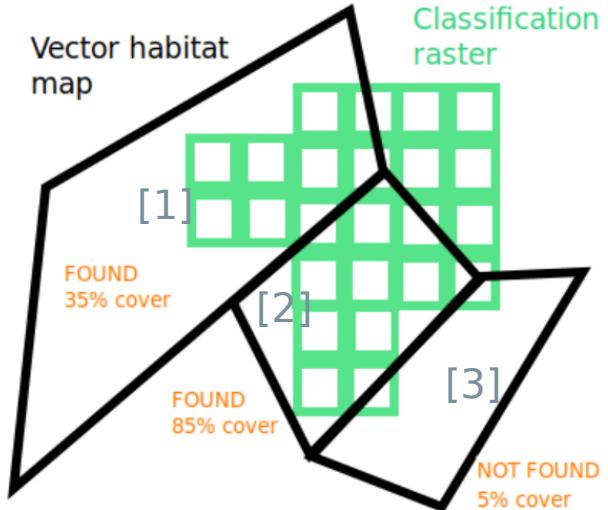
R7-Sparsely wooded grassland

Base map: CUZK, 2021



# Metrics

1. if the **polygon has been identified** as given category
  - binary output, threshold dependent
  - → → *proportion of polygons identified*
  
2. the **percentage of polygon area covered** by classified category
  - from 0 to 100%
  - → → *coverage distribution within habitat categories*



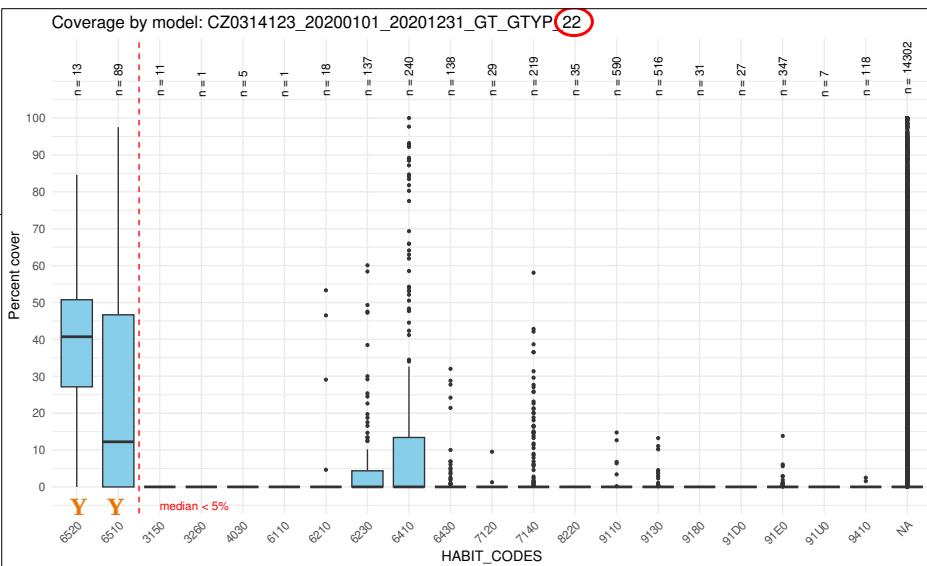
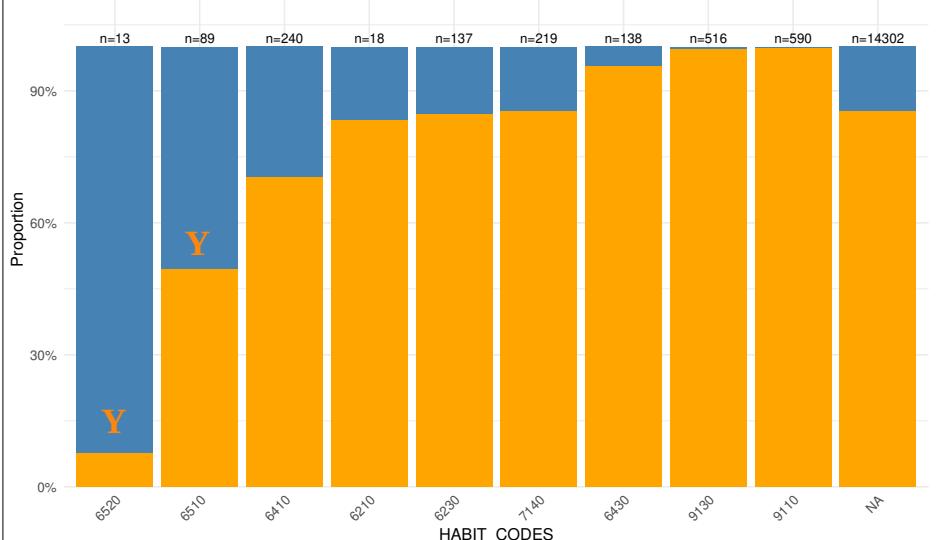
→ → computed for each polygon separately, than summarised in a graphs

# Results for Boletice, prototype data of EUGW grassland type characterization component (GT)

## GT 22 – Mesic grasslands

Top 10 HABIT\_CODES by found (>= 10% of cover) polygon proportion

Model: CZ0314123\_20200101\_20201231\_GT\_GTYP\_22



→ classification prefers correct habitats, but does not capture their full size

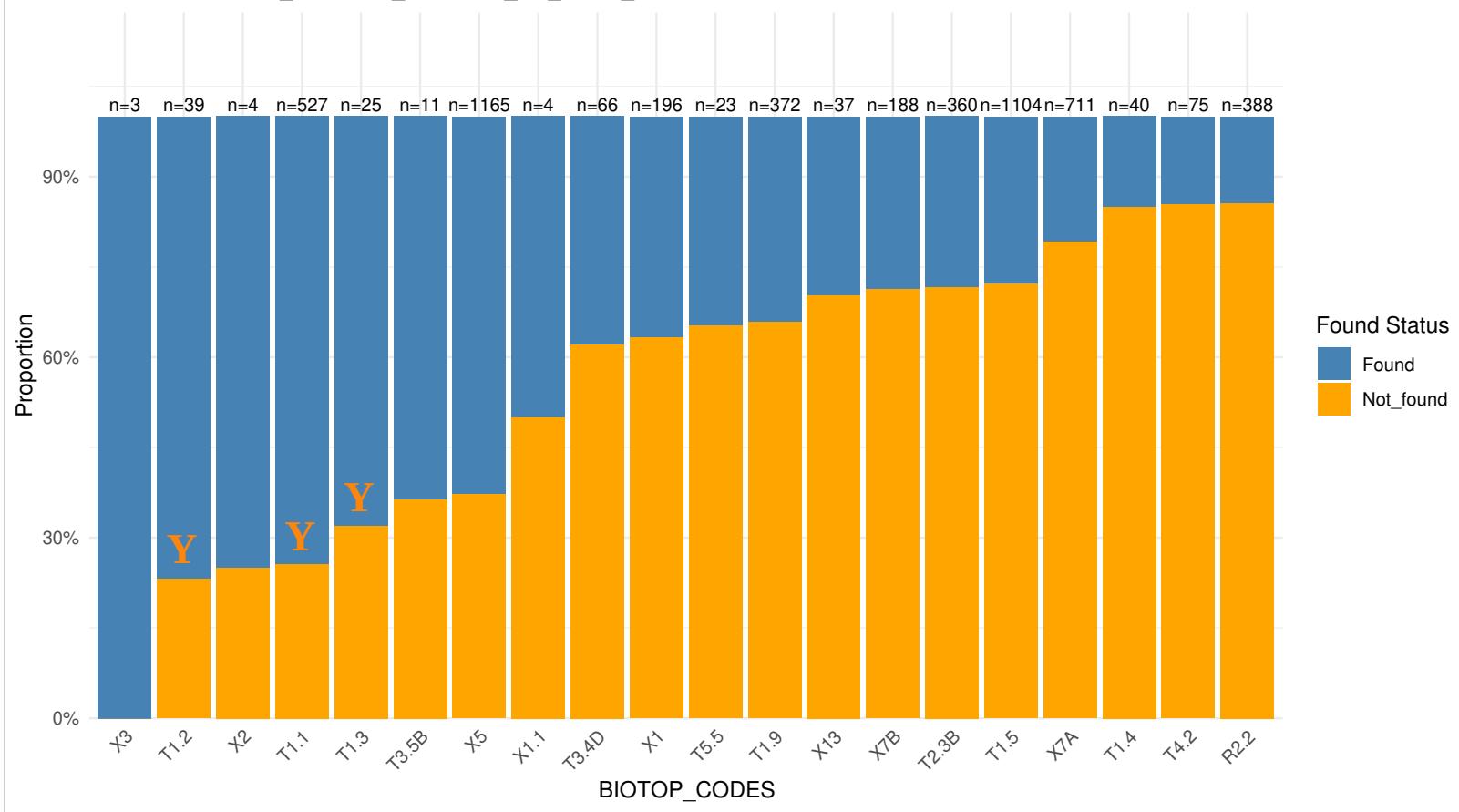
→ → NA here represents polygons, which are not Annex1 habitats, but they have czech cathegory assigned

# GT 22 – Mesic grasslands

# Czech classification system perspective

Top 20 BIOTOP\_CODES by found (>= 10% of cover) polygon proportion

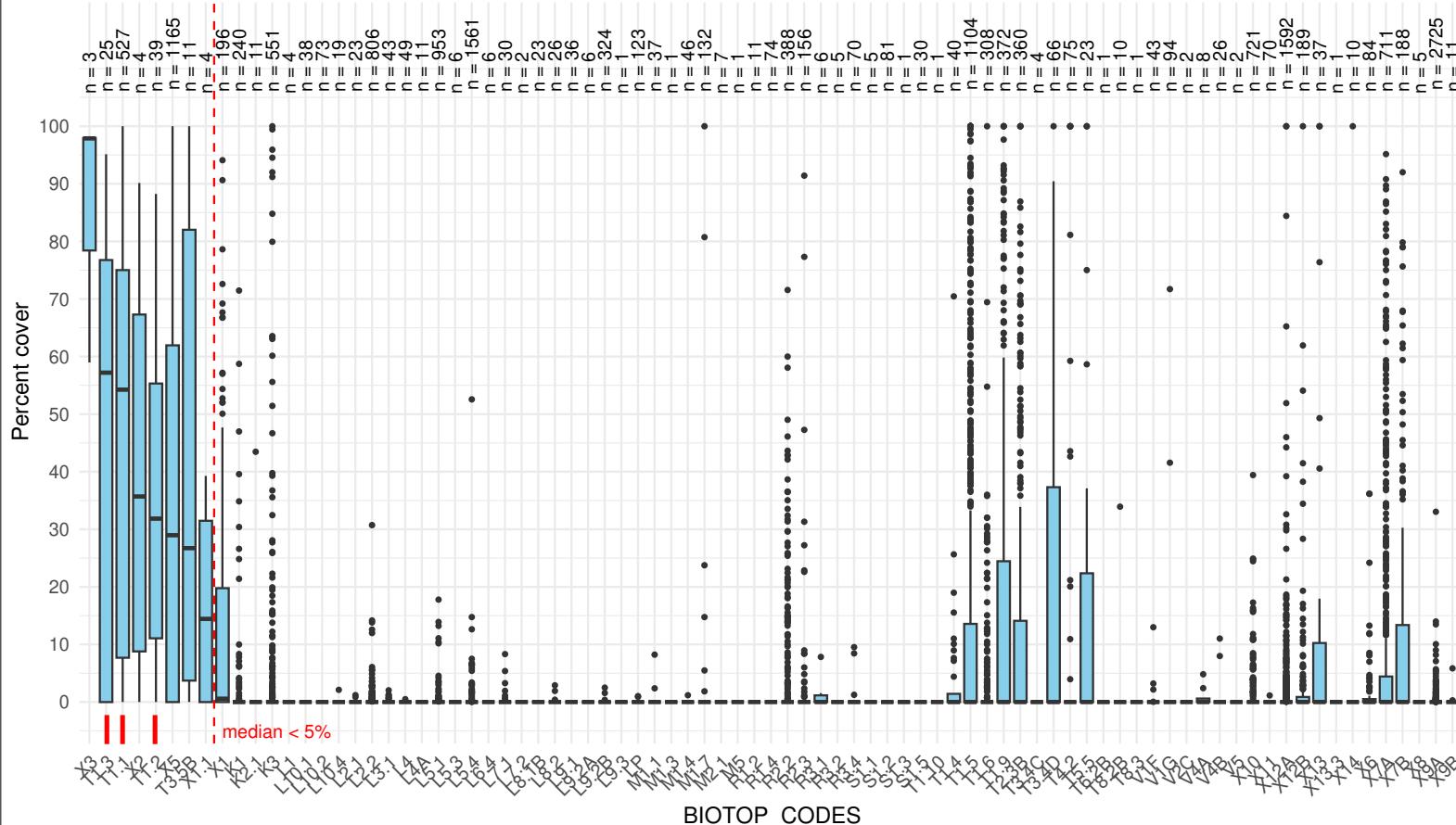
Model: CZ0314123\_20200101\_20201231\_GT\_GTYP\_22



# GT 22 – Mesic grasslands

# Czech classification system perspective

Coverage by model: CZ0314123\_20200101\_20201231\_GT\_GTYP\_22



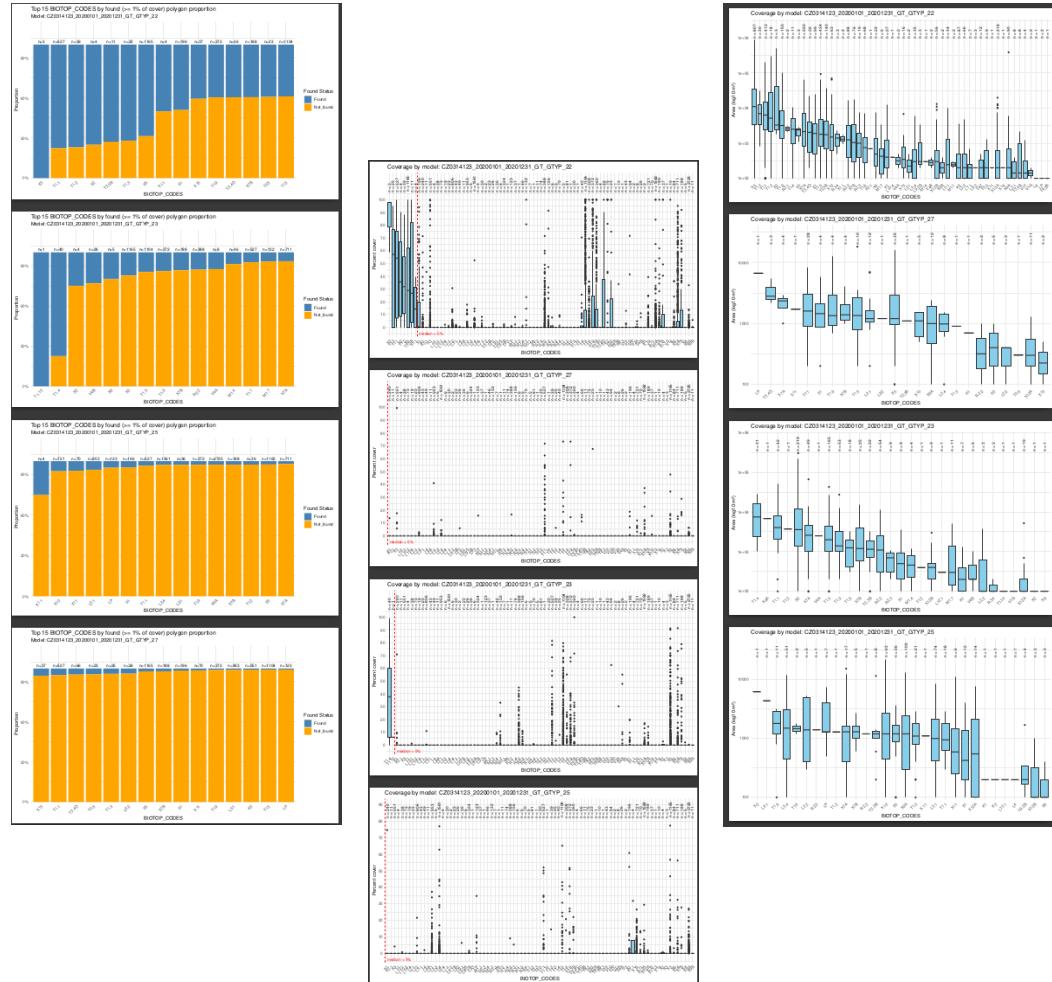
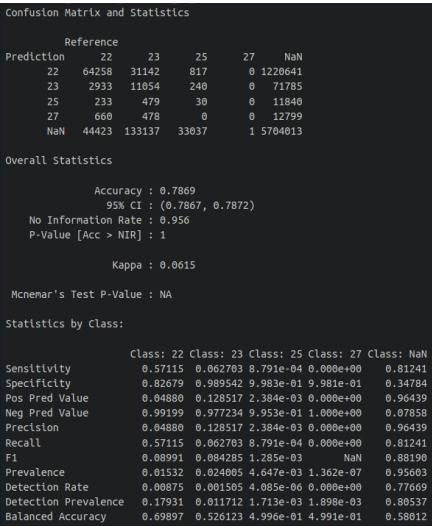
When plotted on **graphs**, it is possible to better assess how well the model performs on different habitat types, while **revealing trends in classification successes and failures**.

On the other hand, this general assessment **does not allow for a comparison between different models**.

# Work in progress



- universal script
- vector polygon + raster layer(s)
- according to the user input produces the presented graphs
- pixel based confusion matrix



# Confusion Matrix

- function `caret::confusionMatrix`
- prediction and reference **must share classification system** and all levels must be present in data

Confusion Matrix and Statistics

Prediction	Reference				
	22	23	25	27	NaN
22	64258	31142	817	0	1220641
23	2933	11054	240	0	71785
25	233	479	30	0	11840
27	660	478	0	0	12799
NaN	44423	133137	33037	1	5704013

Overall Statistics

Accuracy : 0.7869  
95% CI : (0.7867, 0.7872)  
No Information Rate : 0.956  
P-Value [Acc > NIR] : 1

Kappa : 0.0615

McNemar's Test P-Value : NA

Statistics by Class:

	Class: 22	Class: 23	Class: 25	Class: 27	Class: NaN
Sensitivity	0.57115	0.062703	8.791e-04	0.000e+00	0.81241
Specificity	0.82679	0.989542	9.983e-01	9.981e-01	0.34784
Pos Pred Value	0.04880	0.128517	2.384e-03	0.000e+00	0.96439
Neg Pred Value	0.99199	0.977234	9.953e-01	1.000e+00	0.07858
Precision	0.04880	0.128517	2.384e-03	0.000e+00	0.96439
Recall	0.57115	0.062703	8.791e-04	0.000e+00	0.81241
F1	0.08991	0.084285	1.285e-03	NaN	0.88190
Prevalence	0.01532	0.024005	4.647e-03	1.362e-07	0.95603
Detection Rate	0.00875	0.001505	4.085e-06	0.000e+00	0.77669
Detection Prevalence	0.17931	0.011712	1.713e-03	1.898e-03	0.80537
Balanced Accuracy	0.69897	0.526123	4.996e-01	4.991e-01	0.58012

# Transferability, what data is needed?

- **Vector polygon layer** representing ‘ground truth’ (e.g. field mapped habitats)
  - unique identifier of polygon
  - habitat category represented by polygon
- **Classification raster(s)**
  - for ‘**confusion matrix** analysis’ there should be **identical levels in both** raster and vector layers (use of same classification system and all classified categories by model should be present in focal area – if not, its shouting errors, but I will improve it to stop make it happening)
  - for ‘**graphical** analysis’ **any classification system** could be used (since it does not really compare anything)
- code is already hanging on Github ([github.com/zbubster/vmb\\_validator](https://github.com/zbubster/vmb_validator)), but some essential parts still missing