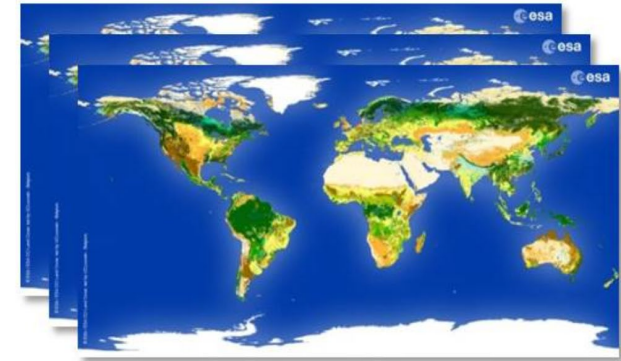# Lessons learned from building a training data set for land cover mapping at 10m

Myroslava Lesiv, Daniele Zanaga, Ruben Van De Kerchove, Nandika Tsendbazar, Martin Herold and Steffen Fritz
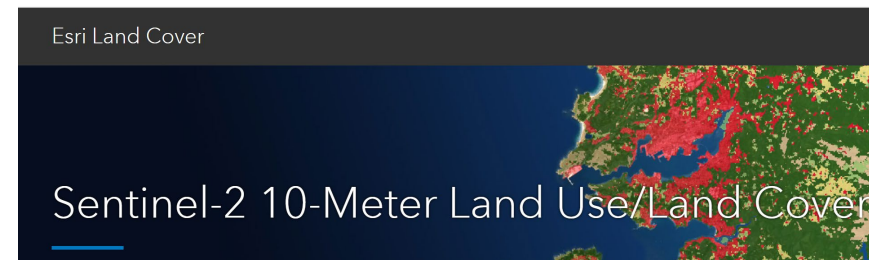
# Background

- High quality training data is a critical input for land cover/ land use mapping

- New requirements:
  - Very high resolution mapping
  - More thematic details
  - Change detections

- Different sources of training data available:
  - on-ground observations and
  - visually interpreted very high resolution images.
  - existing land cover/land use maps
  - automatic generation

**Copernicus Global Land Service**
*Providing bio-geophysical products of global land surface*

**RELEASE OF THE 10 M WORLDCOVER MAP**

Esri Land Cover

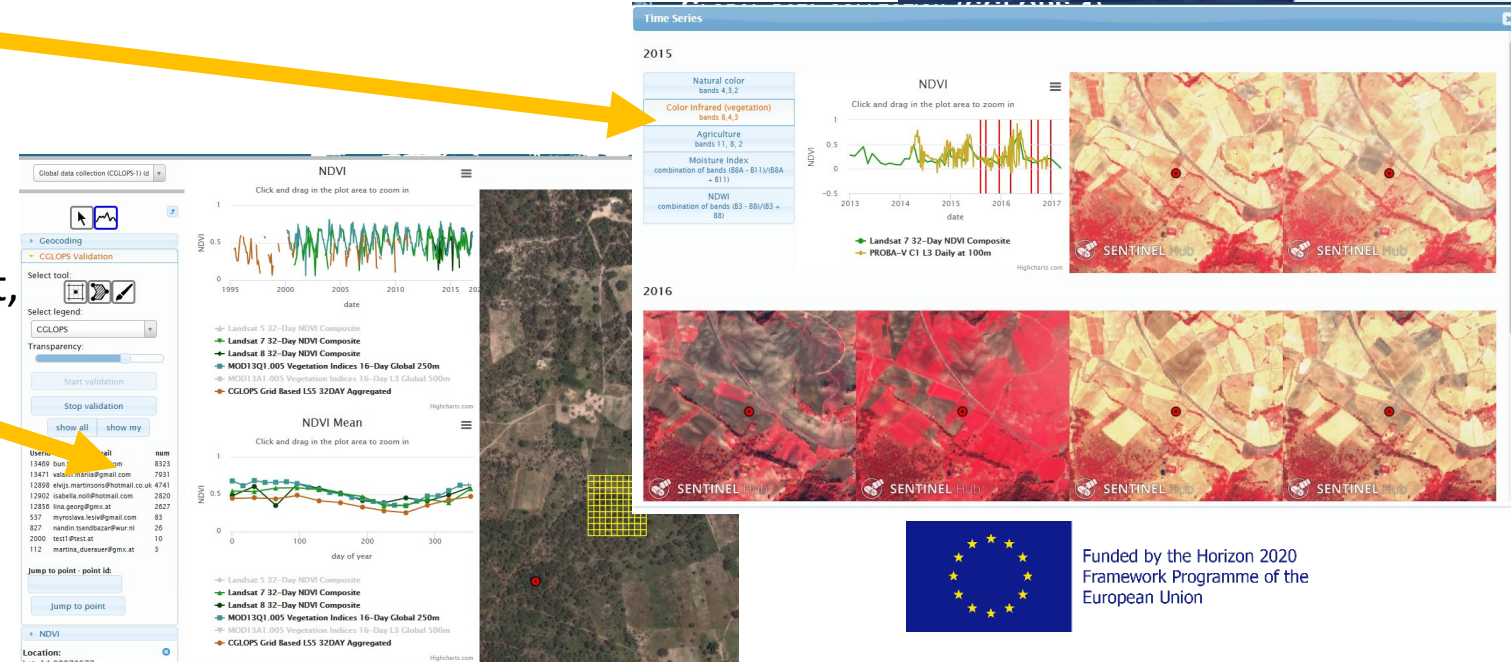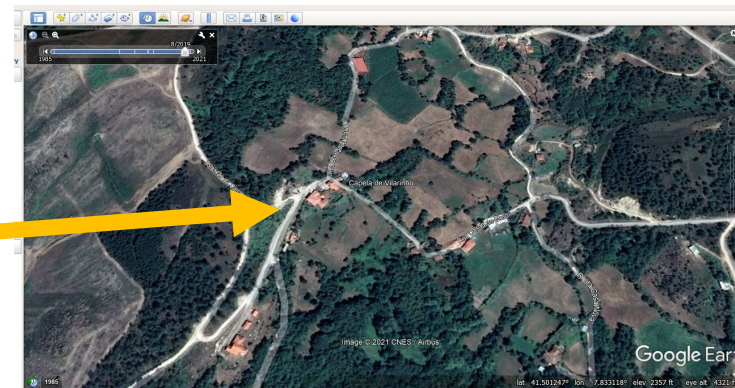Sentinel-2 10-Meter Land Use/Land Cover

# Challenges

- Unknown quality
  - Geolocation errors
  - Thematic errors
  - Timestamp
  - Not clear definitions

- Spatial distribution of data – overfitting issue

- Translation from one legend to another

- Translation of point observations into pixels

- Lack of data

- Access to data

# Geo-Wiki tool box

- Very high resolution (VHR) imagery from Google maps, Microsoft Bing, and ESRI

- Google Earth VHR historical images

- Planet time series of images

- Sentinel-2 time-series in False color

- Street level images from Google and Mappilary

- NDVI time series derived from Landsat, Proba-V and MODIS data

# Concept of a multipurpose data set

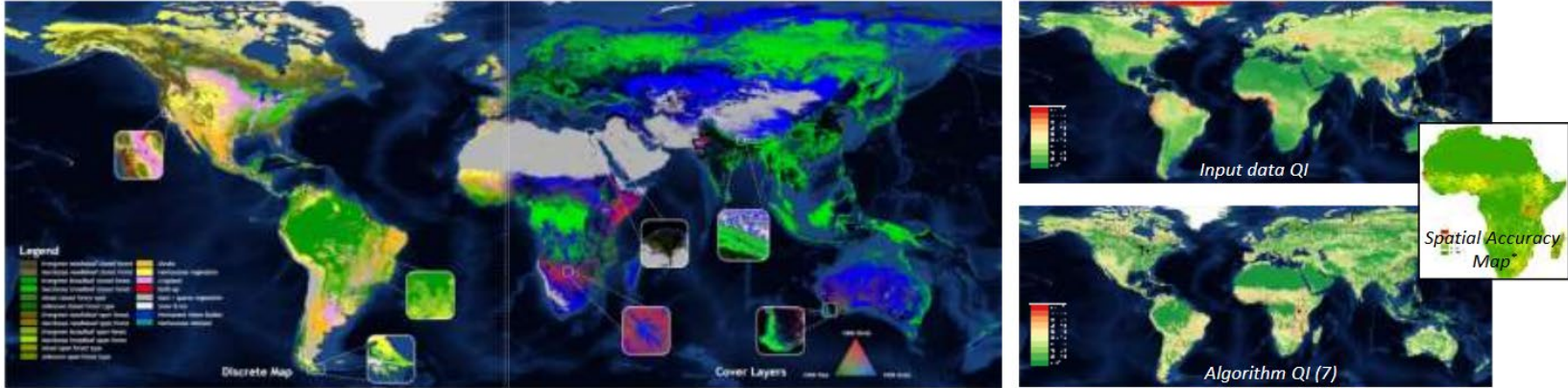Global Land Cover at 100m (JRC)

World Cover at 10m (ESA)

# Copernicus Global Land Cover 2015-2019 PROBA-V 100m
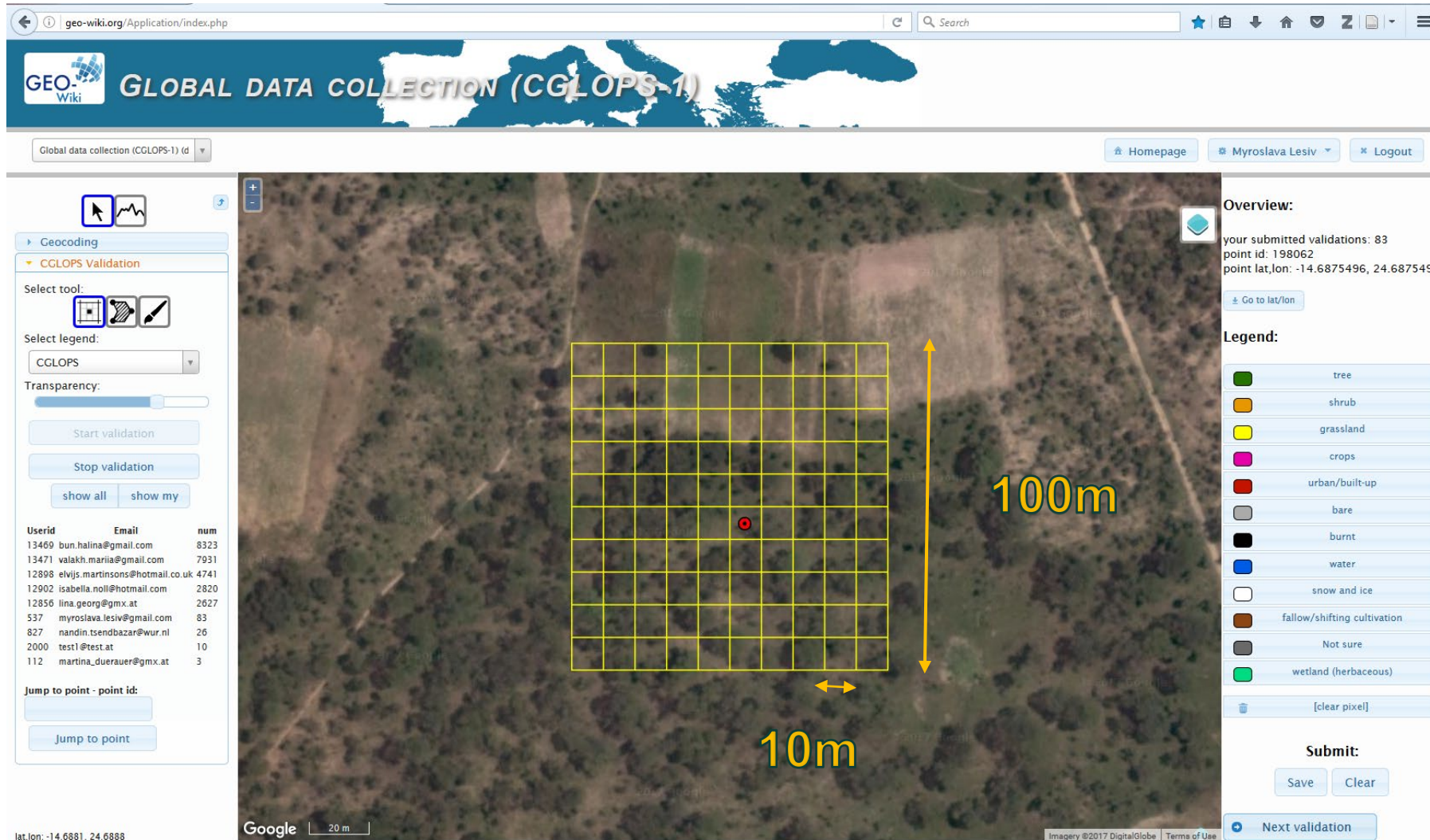


Discrete Map (23 classes)

10 Continuous Covers (0-100%)
*Permanent water is derived from GSWE (Pekel et al.)*
*Built-up is derived from WSF (Marconcini et al.)*

Quality Indicators
*(*) example over Africa, global maps under release test*

| Continuous Covers | |
|---|---|
| Bare | Snow |
| Crops | Tree |
| Grass | Urban |
| Moss | Permanent water |
| shrub | Seasonal water |

A systematic **SERVICE** providing a **DYNAMIC**, **YEARLY**, **USER- ORIENTED** product at **GLOBAL** scale @ **100m resolution** from 2015 onwards

land.copernicus.eu/global/lcviewer

remotesensing.vito.be

# Geo-wiki app



- Fractions at 100m
- Easy translation to discrete land cover classes
- Training data at 10m

# Data collection workflow

- Initial Geo-wiki training
    - Interface, tools, per class examples

- Weekly online seminars to check quality
    - discuss difficult locations
    - randomly revisit some classifications
        - Target <5 % of mistakes

- Comparison with regional products
    - Revise disagreeing locations

- Removing land cover class outliers based on spectral information
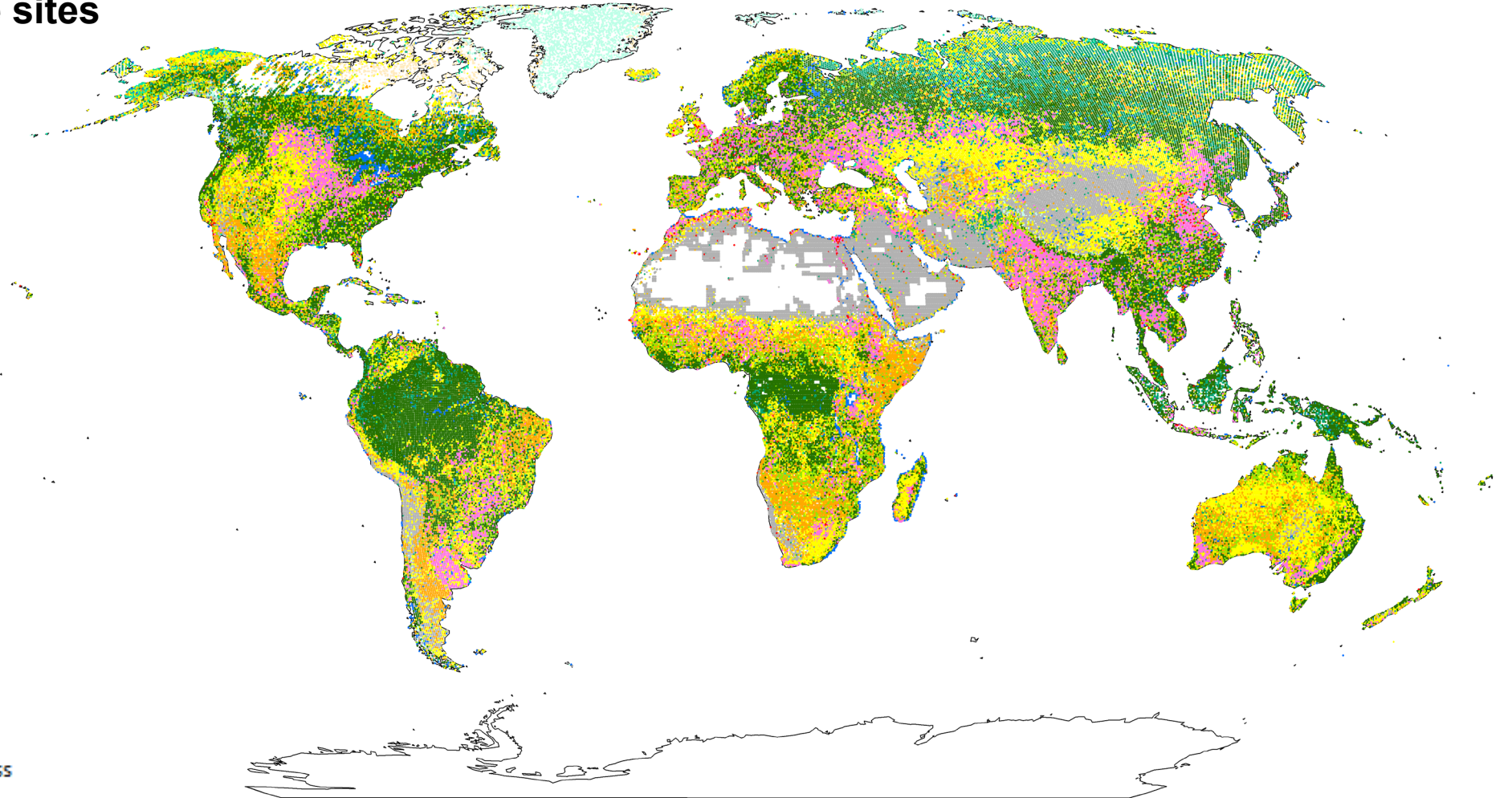    - homogeneous pixels

# Distribution of reference data 2015

**~180 000 sample sites**
by 20 experts

Sampling design:
- systematic
- uncertainty hotspots



- closed forest
- open forest
- shrubs
- herbs
- crops
- built-up
- bare
- snow and ice
- water
- wetland
- lichen and moss

# ESA World Cover 2020/2021 at 10m Sentinel 1 and Sentinel 2

**180 K  (at 100m) pixels  ~ 18 Millions (at 10m)**

Issues:

- geolocation errors of the underlying images used for visual interpretations

- land cover/land use changes that happened after 2015

- Correct fractions at 100m ~misclassifications at 10m

Landscape in Australia ➡️
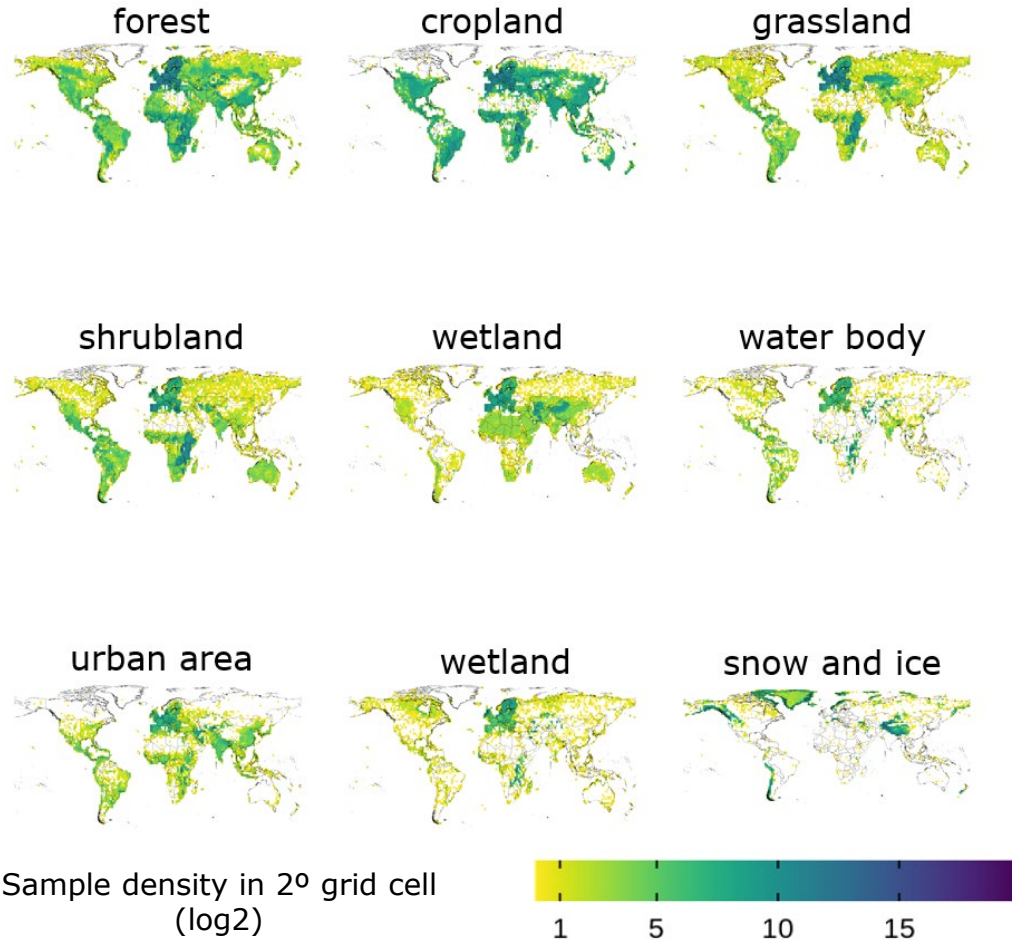
# Training data optimization

- Geolocation errors of labels?
  - subset only those pixels that are surrounded by pixels with the same label

- Land cover changes ? –
  - subset sample sites of potential changes by running BFAST model and revising these sites
  - Set of rules based on spectral information

- Misclassifications at 10m resolution?
  - Set of rules based on spectral information

# Lessons learned

- Having subpixel information about land cover is important for better defining classes at pixel level

- Homogeneous areas - it would be more convenient to label patches/segments rather than pixels

- There are always uncertainties associated with human labeling therefore additional data filtering is needed
    - E.g. taking into account spectral information

- The presented data will be made of open access

# Collection of existing reference data sets



forest | cropland | grassland

shrubland | wetland | water body

urban area | wetland | snow and ice

Sample density in 2º grid cell (log2)

1  5  10  15

**~ 7 million samples**

**years: 1951-2020**

**spatial units: 10-5000 m**

Sources:
LUCAS – Land use and land cover survey
GLIMS
Ramsar
GHS Urban Center Database
Global Croplands
JECAM
PRdataGO
...

caterina.barrasso@idiv.de | carsten.meyer@idiv.de

Macroecology and Society
A research group of iDiv

GEOKUR

# Thank you!

Myroslava Lesiv

International Institute for
Applied Systems Analysis

lesiv@iiasa.ac.at