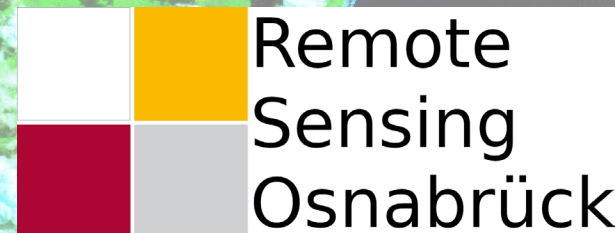# Investigating the Impact of Current Flood Map Validation Practices

**Antara Dasgupta, Björn Waske**

Metric sensitivity, sampling strategies, class imbalance
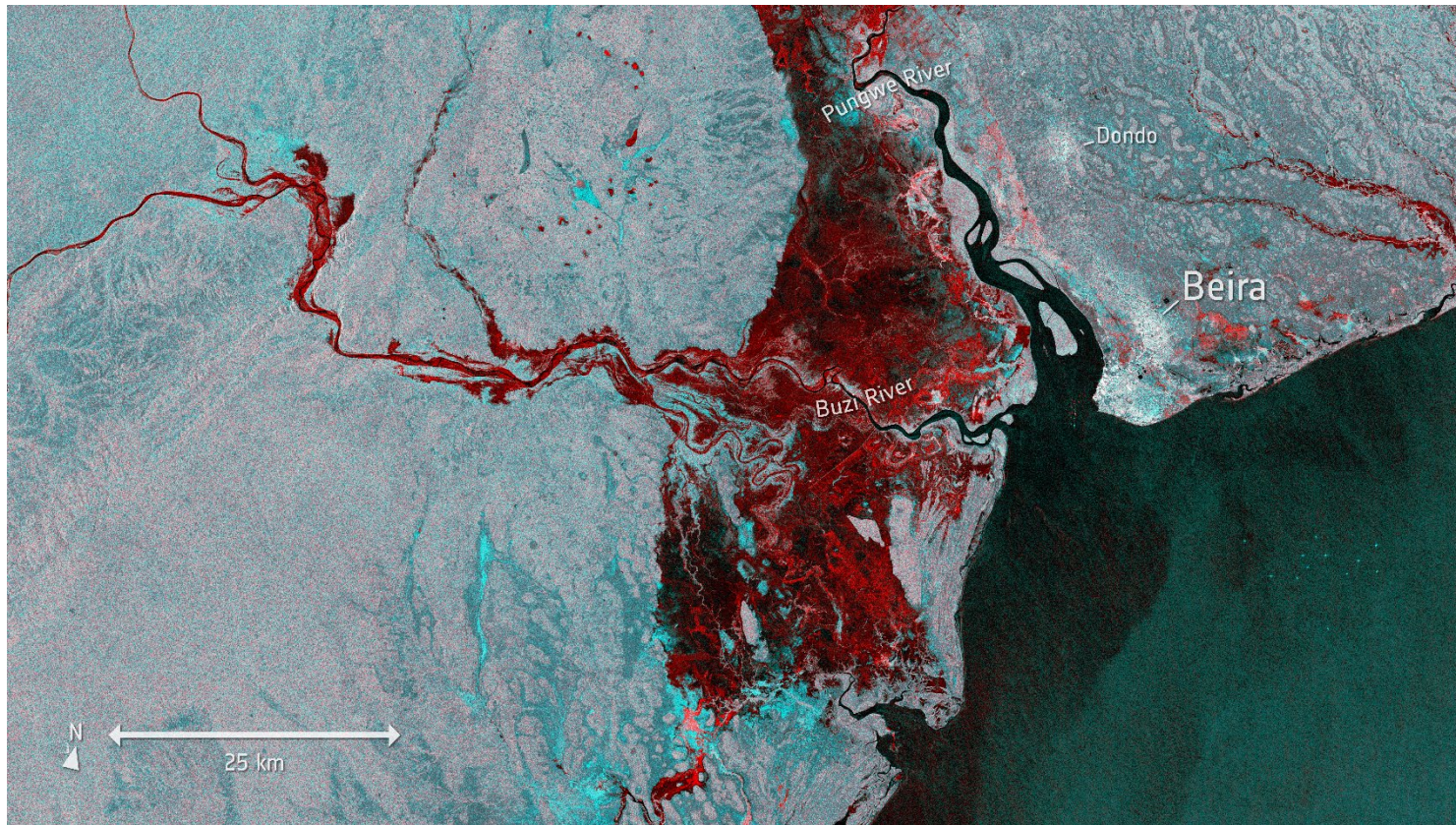
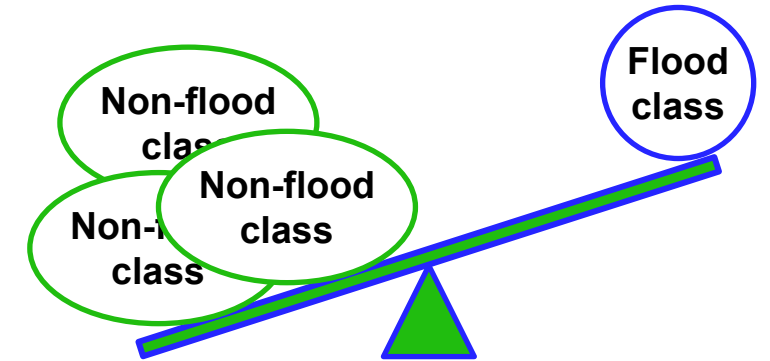# PROBLEMS WITH SATELLITE-BASED FLOOD MAP VALIDATION

# Problem 1: Metric Sensitivity



Validated Map / Classified Map

**Common area >**
**Map differences**
**(mostly!)**

# Problem 2: Class Imbalance



Copernicus Sentinel-1 Image showing the flooding from Cyclone Idai in red, around the port town of Beira in Mozambique on 19 March 2019, provided by the Copernicus Emergency Mapping Service (CEMS).
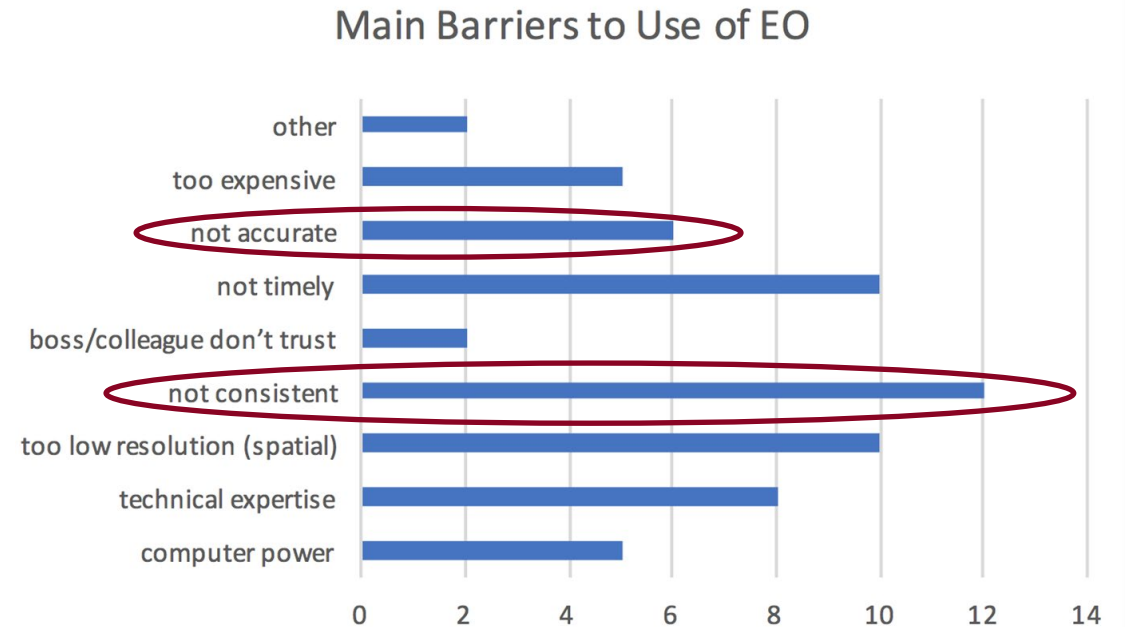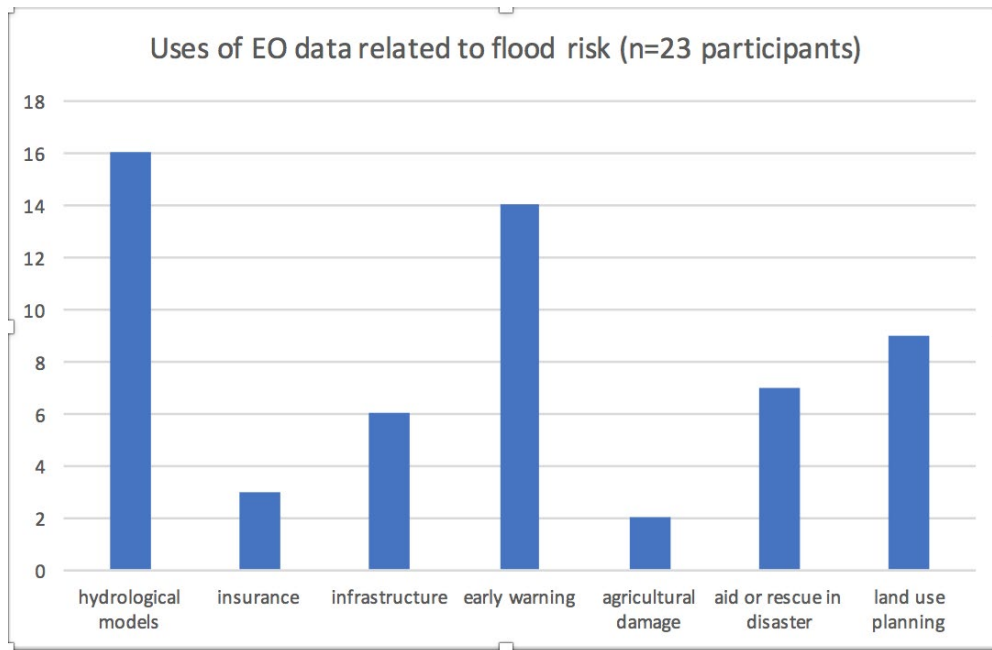
Common metrics designed for LULC assessment

Unable to deal with large class imbalance in binary classifications

# Problem 3: Over-confidence in maps confuses users!



Uses of EO data related to flood risk (n=23 participants)
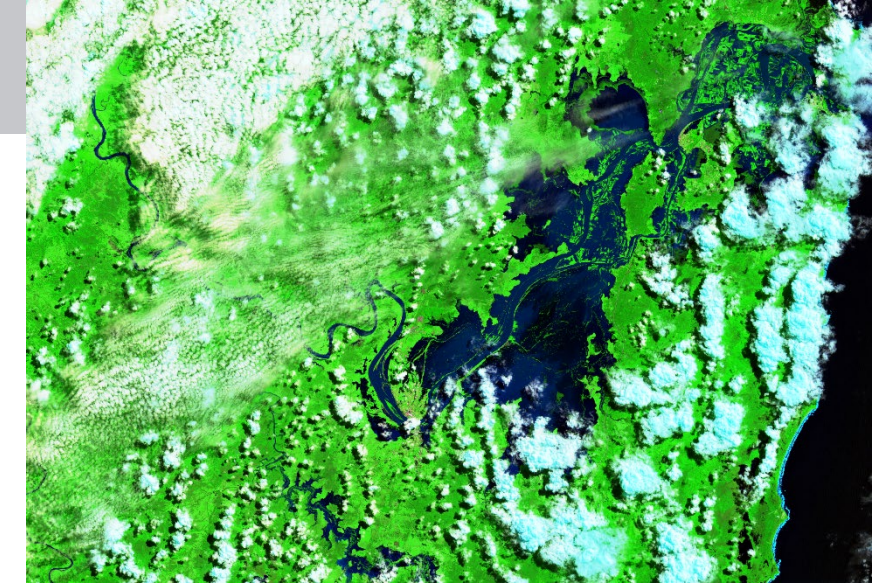


Main Barriers to Use of EO

Tellman, B. (2019). What flood event map accuracy is required to enable governments, aid agencies, and insurance companies to protect vulnerable lives and livelihoods? Global Flood Partnership 2019 Conference – 11 - 13 June 2019, Guangzhou (China).

Kettner, A. J.,Schumann, G. J.-P., and Tellman, B. (2019), The push toward local flood risk assessment at a global scale, *Eos, 100*, https://doi.org/10.1029/2019EO113857. Published on 14 January 2019.
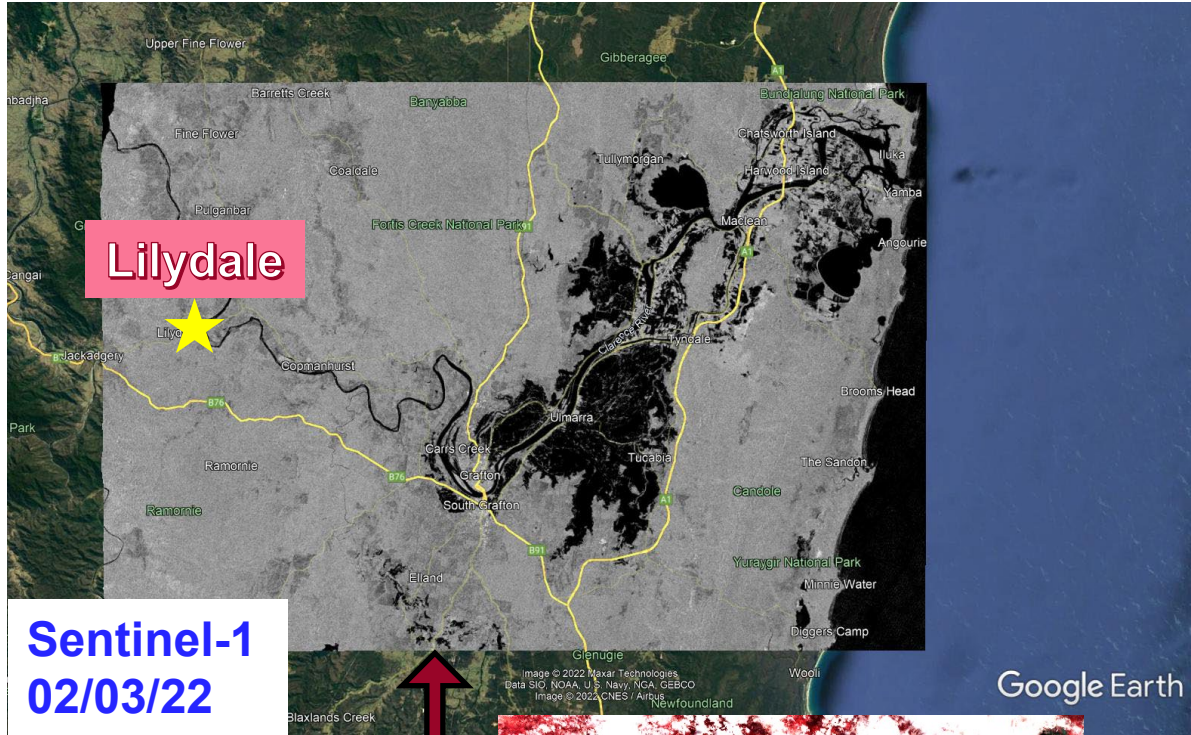
Test Site, Data, and Workflow

# STUDY AREA

# Test Case and Data: Clarence Valley Floods – 2022
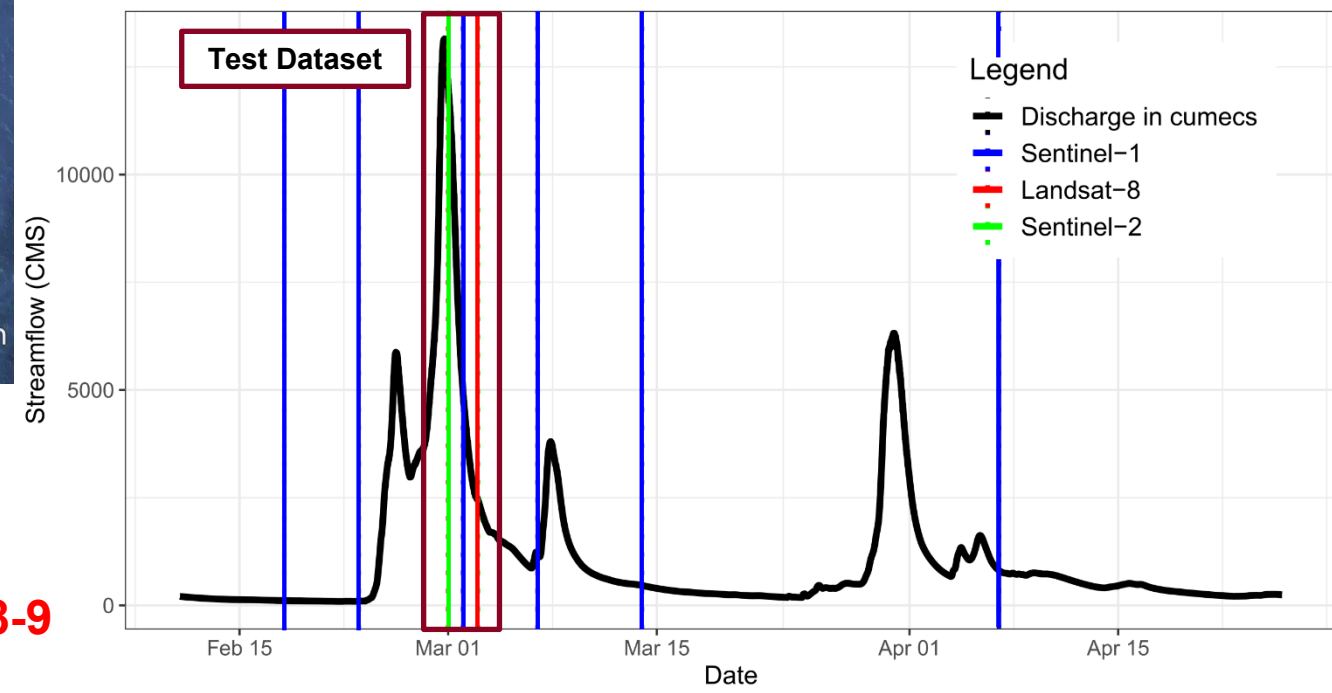


**Sentinel-2**
**01/03/22**

**Lilydale**

**Sentinel-1**
**02/03/22**

**Australia**

**Landsat 8-9**
**03/03/22**

### Discharge Hydrograph at Lilydale (2022)



Test Dataset

Legend
- Discharge in cumecs
- Sentinel-1
- Landsat-8
- Sentinel-2

Streamflow (CMS)

10000

5000

0

Feb 15    Mar 01    Mar 15    Apr 01    Apr 15

Date

# Validation Metrics Used

| Confusion Matrix | | Predicted condition | |
|---|---|---|---|
| | Total population = P + N | Positive (PP) | Negative (PN) |
| **Actual condition** Positive (P) | | True positive (TP), hit | False negative (FN), type II error, miss, underestimation |
| Negative (N) | | False positive (FP), type I error, false alarm, overestimation | True negative (TN), correct rejection |

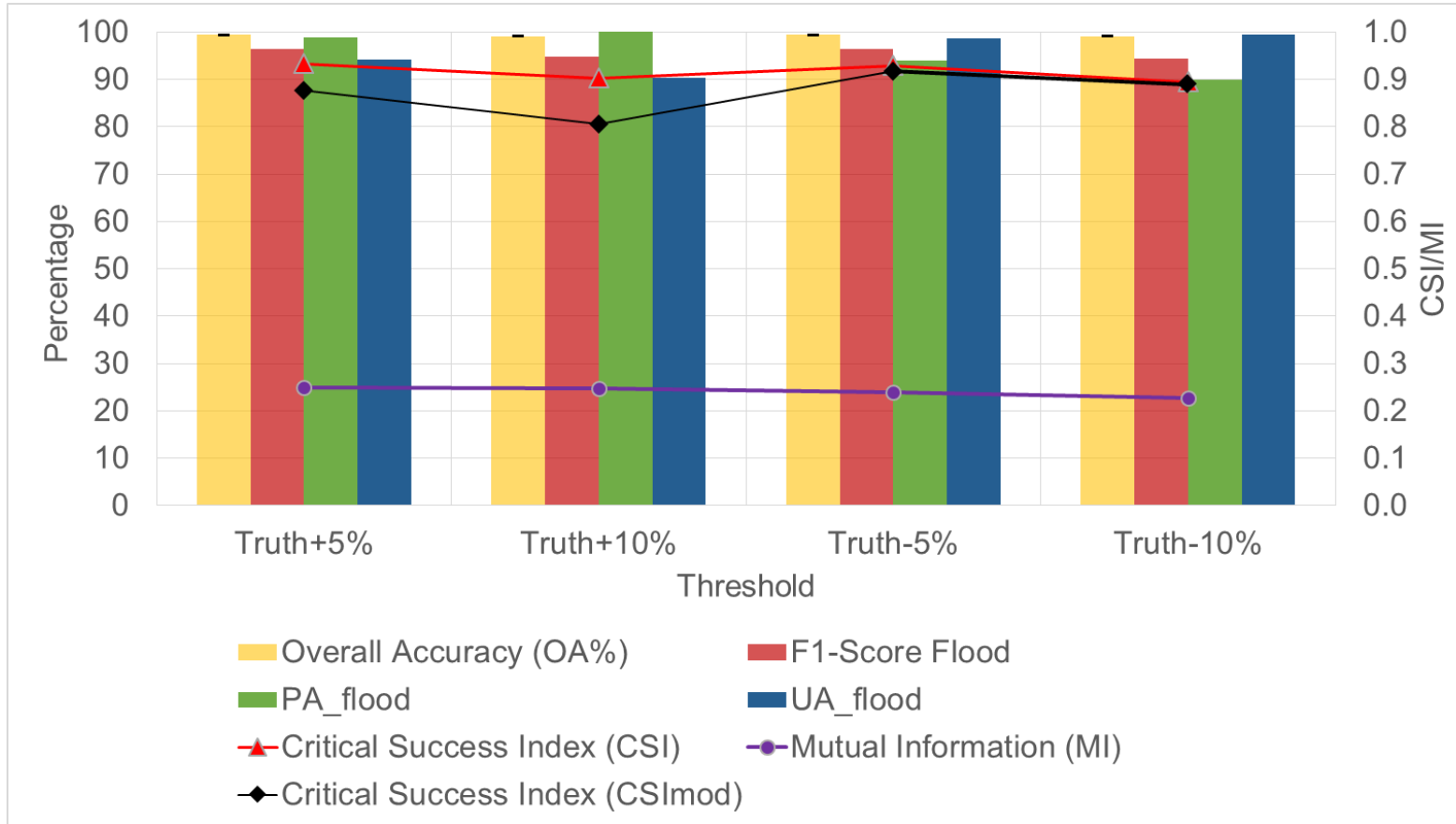| Metric Name | Formula |
|---|---|
| Prevalence | $\dfrac{P}{P+N}$ |
| Overall Accuracy | $\dfrac{TP+TN}{P+N}$ |
| User's Accuracy Flood/Precision | $\dfrac{TP}{TP+FP}$ |
| Producer's Accuracy Flood/Recall/True Positive or Hit Rate | $\dfrac{TP}{TP+FN}$ |
| F1-Score Flood | $\dfrac{2\times TP}{2\times TP+FP+FN}$ OR $\dfrac{2\times Precision \times Recall}{Precision+Recall}$ |
| Critical Success Index/Intersection over Union | $\dfrac{TP}{TP+FP+FN}$ |
| Critical Success Index modified | $\dfrac{TP-FP}{TP+FP+FN}$ |
| Mutual Information | $I(y;x) = H(x) - H(x|y)$ |
| *Rarely used metrics marked in red | where $I(Y;X) = \sum_{x\in\mathcal{X}}\sum_{y\in\mathcal{Y}} p(x,y)\log(\dfrac{p(x,y)}{p(x)p(y)})$ |

Research Question 1: How do metric choice, sampling design and sample size, influence accuracy assessment?

# DIFFERENT THRESHOLD FLOOD MAPS VS. SYNTHETIC "TRUTH"

# Current Practice: Validation over the entire domain
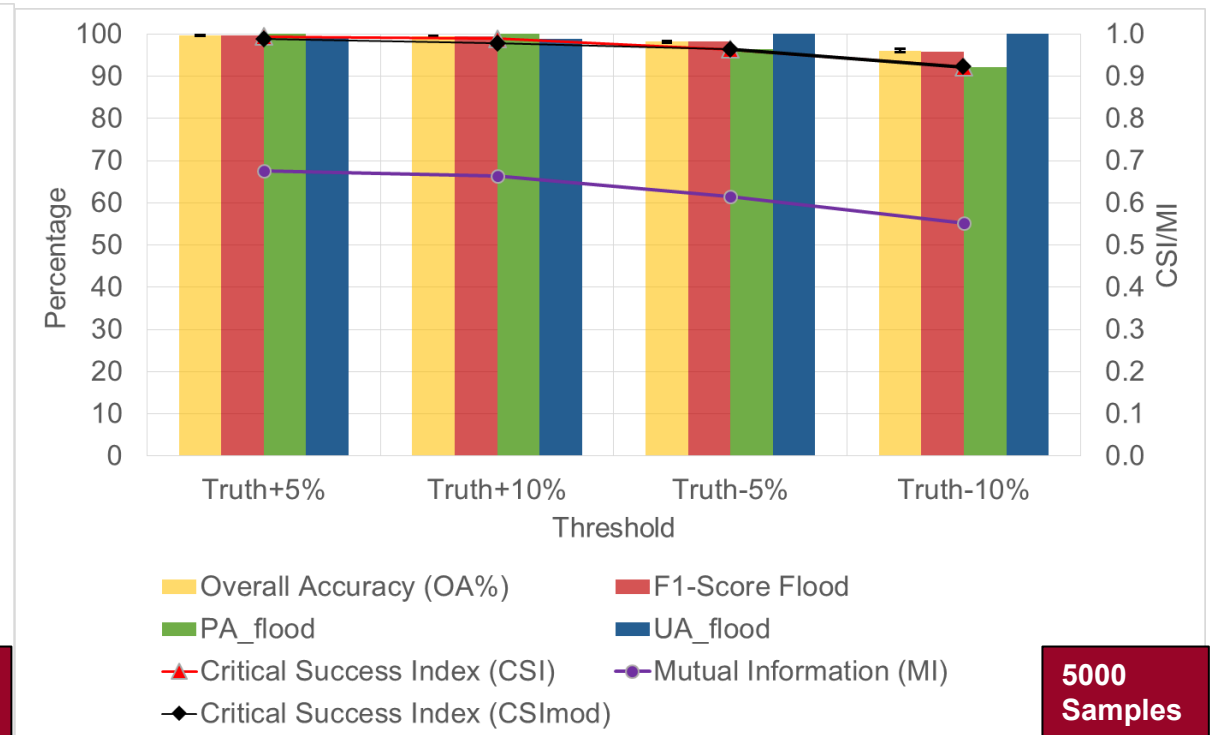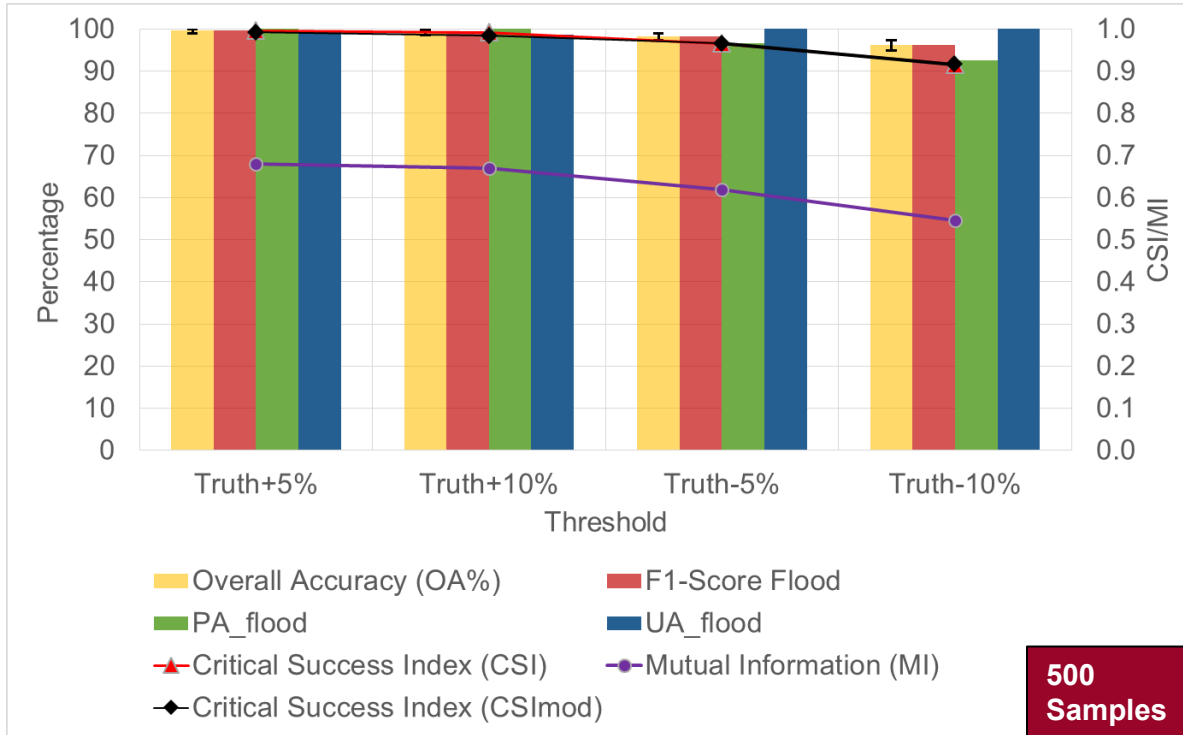


Low sensitivity of all metrics towards flooded area variations

All pixels in image used for the metric computation

Changes in maps hard to capture – overall accuracy remains almost the same (even for large drops in user's and producer's accuracy)
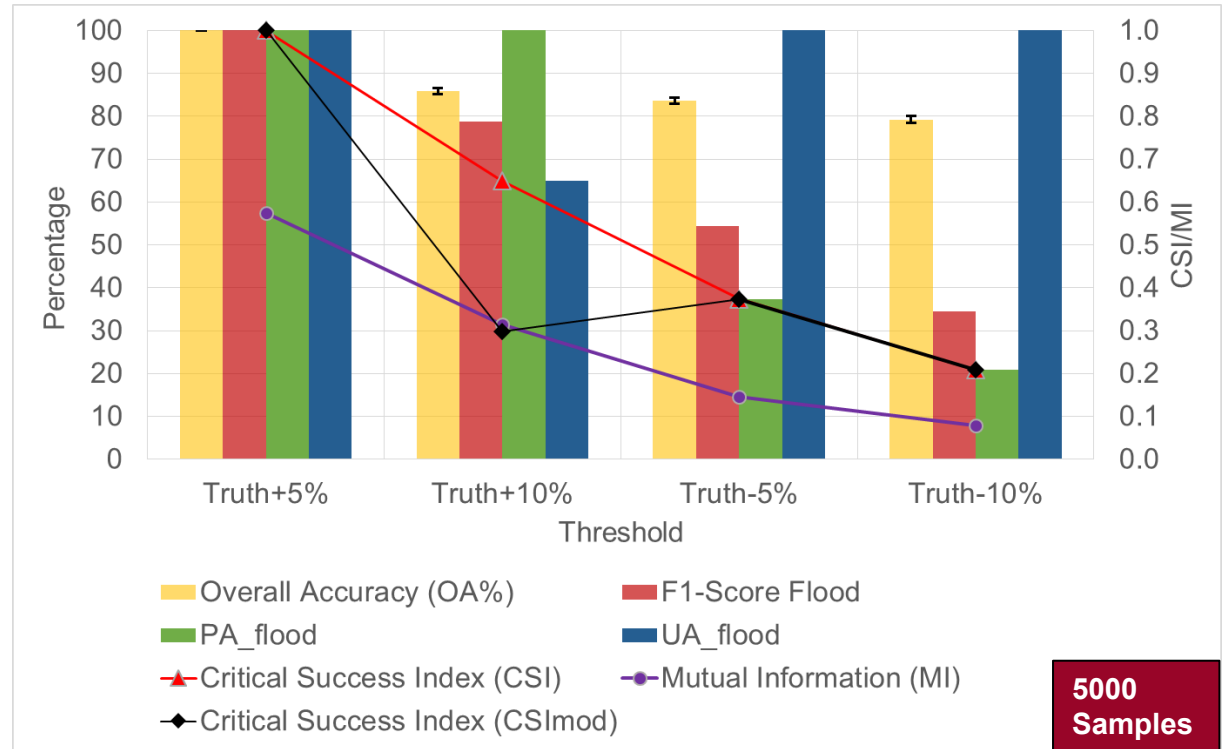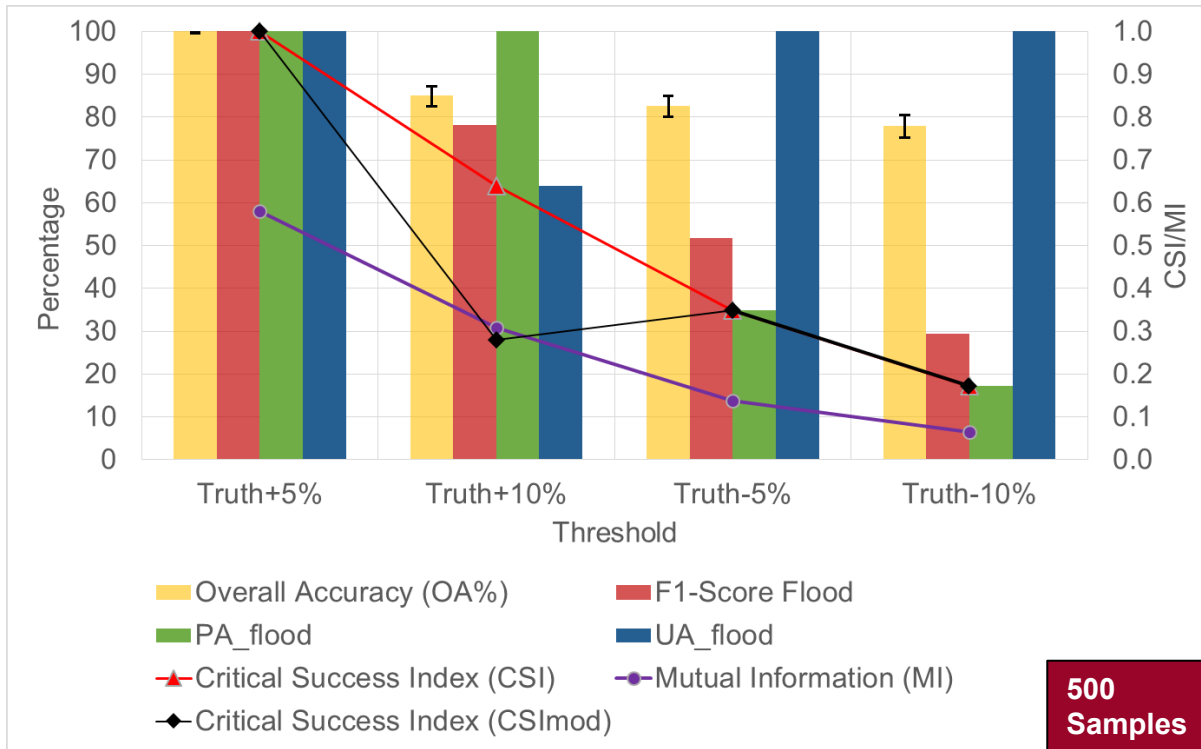
# Stratified Random Sampling*



All metrics biased towards over prediction

Low sensitivity demonstrated by low variation in metric values as a function of threshold variations

*Stratified sampling over entire domain
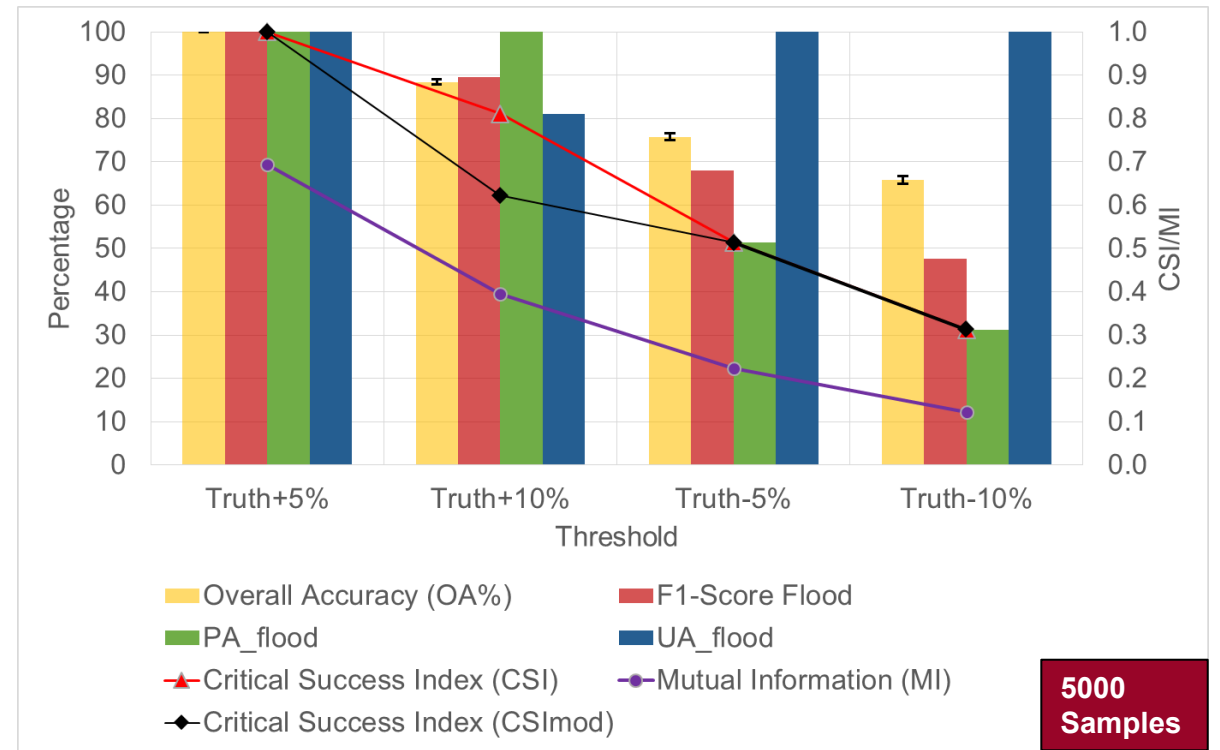
# Targeted Random Sampling*

CSImod alters the CSI metric effectively such that similar over- and under- predictions yield similar values

Targeted sampling results in higher variation in metric values i.e. greater metric sensitivity

*Sampling in areas with RF classification confidence < 0.9

# Targeted Stratified Random Sampling*



More flooded samples result in higher F1 and CSI scores for the same maps

Stratified sampling reduces variation in F1 and CSI and increases variation in MI and OA

*Stratified sampling in areas with RF classification confidence < 0.9

Research Question 2: How does the validation data error influence accuracy assessment and flooded area calculations?
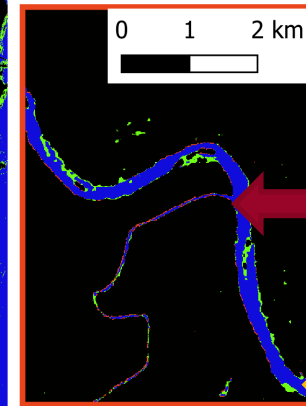
# SYNTHETIC VV-BINARIZED "TRUTH" VS. RANDOM FOREST BENCHMARKS

# Benchmark Sentinel-1 Classification



Random Forest Classification using **S1 VV, VH, elevation (CopDEM30), Otimized GLCM Texture PCs 1 and 2**

**False positives** at the edge of ephemeral streams

**False negatives** at channel banks

RF Benchmark vs. "Truth" Band 1 (Gray)
- True Negatives
- False Negatives
- False Positives
- True Positives

0    7.5    15 km

# Comparison against S-1 RF Benchmark

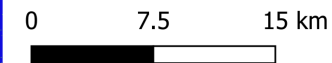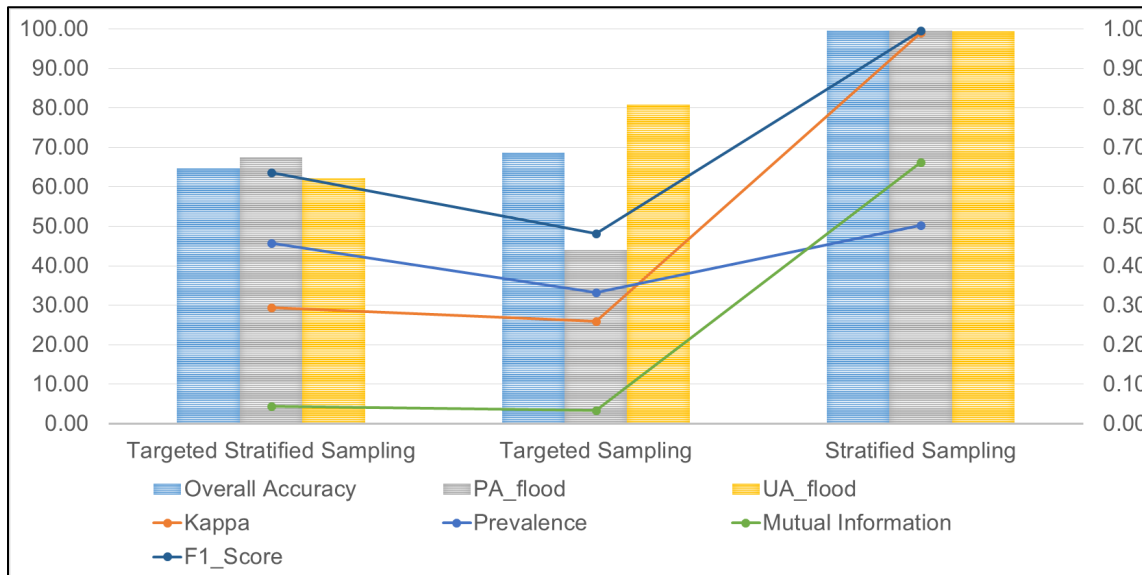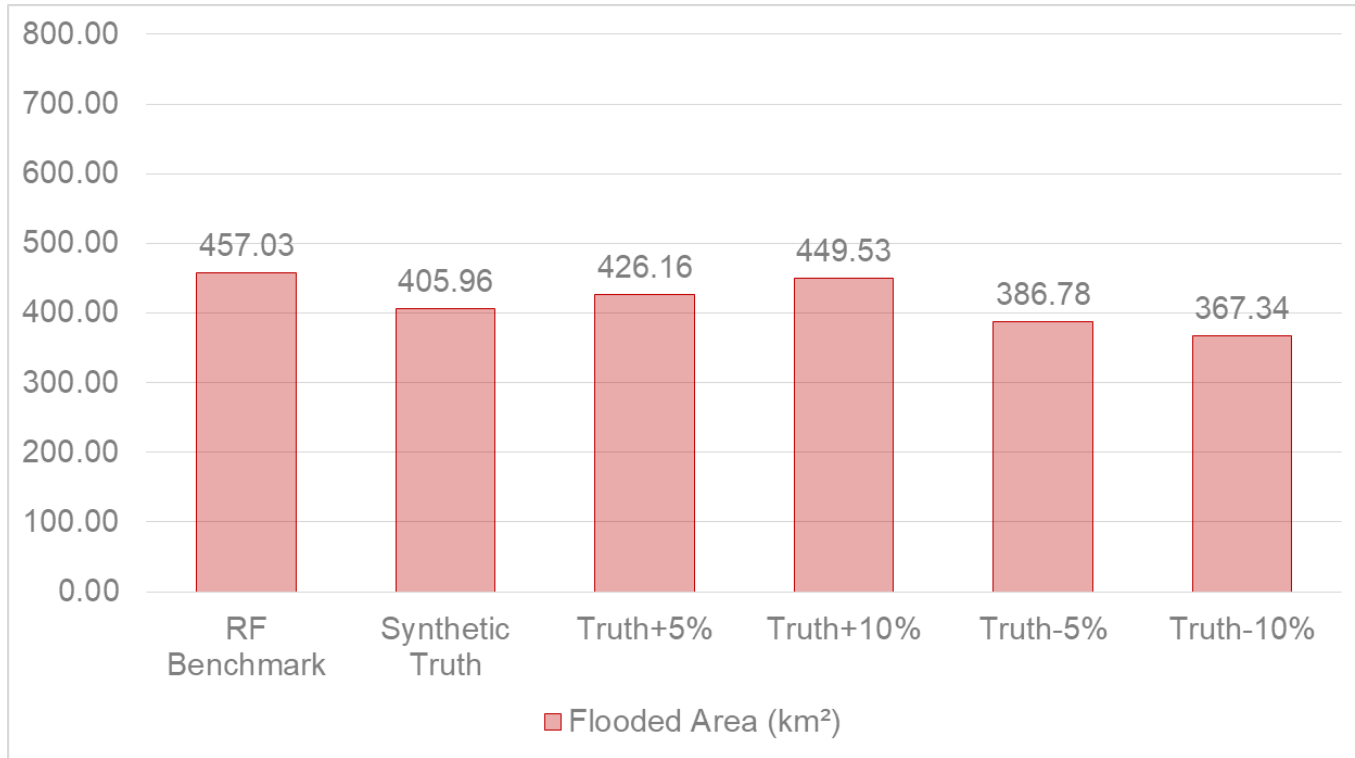| Sampling Strategy (500-1000 pts) | | Targeted Stratified Sampling | | Targeted Sampling | | Stratified Sampling | |
|---|---|---|---|---|---|---|---|
| **Confusion Matrices** | | | | **Reference** | | | |
| | | Flood | Non-flood | Flood | Non-flood | Flood | Non-flood |
| **"Truth"** | Flood | 308 | 204 | 73 | 64 | 501 | 3 |
| | Non-flood | 149 | 337 | 93 | 270 | 2 | 494 |



Stratified sampling provides spuriously high accuracy values
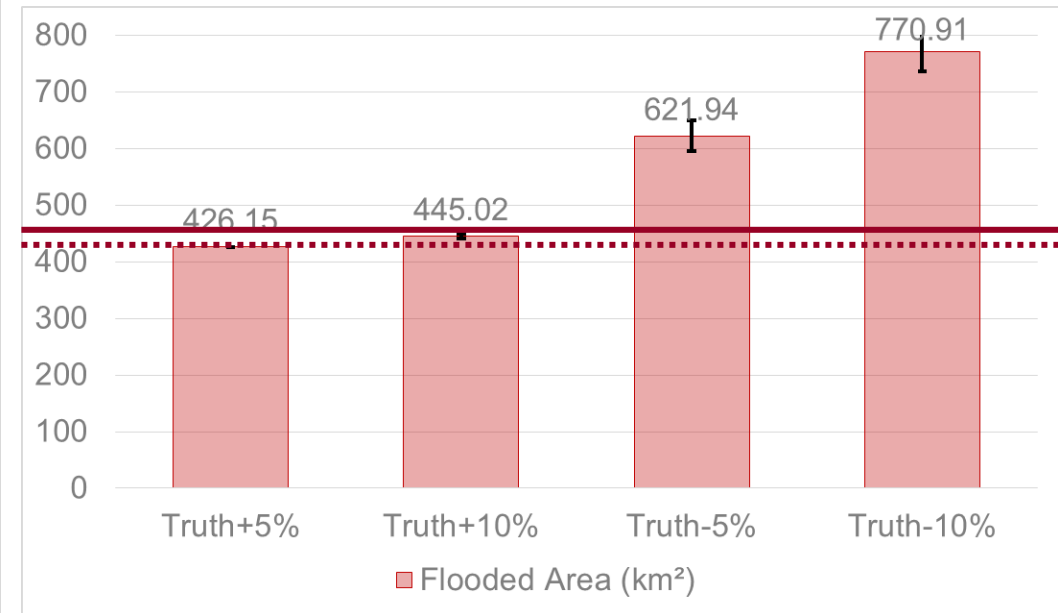
Targeted sampling proves better than other sampling designs

*Validation data errors "correlated" in this case as same input data source i.e. Sentinel-1 SAR

*In the absence of RF confidence or map uncertainty, a buffer along the flood boundary might be considered as the area to "target"

# Flooded Area Estimation



Direct map-based flooded area estimates

Bias-corrected area estimates after Olofsson et al. (2014)

—— RF Benchmark Flooded Area

········· Synthetic Truth Flooded Area

**\*Standard bias-correction techniques for change area estimation NOT directly applicable to flood mapping**

**\*Bias correction only applicable to random, systematic and stratified sampling designs**

Olofsson, P., Foody, G. M., Herold, M., Stehman, S. V., Woodcock, C. E., & Wulder, M. A. (2014). Good practices for estimating area and assessing accuracy of land change. *Remote Sensing of Environment*, *148*, 42–57. https://doi.org/10.1016/j.rse.2014.02.015

# Conclusions and Outlook

- Flood map accuracy assessments **strongly depend** on the choice of metrics, sampling strategy, and validation data quality.

- Increasing sample size **reduces the sensitivity** of accuracy estimation metrics.

- Confidence intervals could provide a clearer overview for decision-makers.

- Future work will focus on developing bias correction methods for flooded area calculations from satellite data.

- An assessment of impacts of land-use and elevation categories on accuracy assessments will also be considered.