

# Algorithm Design and Theoretical Basis Description

## Technical Note D2

Krasen Samardzhiev, Barbara Metzler, Martin Fleischmann, Dani Arribas-Bel

### Table of contents

<b>1</b>	<b>Exective summary</b>	<b>2</b>
<b>2</b>	<b>Theroetical basis</b>	<b>2</b>
2.1	Morphometric Classification Homogenisation . . . . .	2
2.2	AI Modelling using Satellite Imagery . . . . .	2
<b>3</b>	<b>Algorithm Design</b>	<b>2</b>
3.1	Morphometric Classification Homogenisation . . . . .	2
3.1.1	Model architecture . . . . .	2
3.1.2	Data preprocessing . . . . .	2
3.1.3	Morphometric characterisation . . . . .	2
3.1.4	Target labels . . . . .	2
3.1.5	Train/test/validation split . . . . .	2
3.1.6	Training and validation . . . . .	2
3.2	AI Modelling using Satellite Imagery . . . . .	2
3.2.1	Data preprocessing . . . . .	3
3.2.2	Train/test split . . . . .	4
3.2.3	Unbalanced dataset . . . . .	5
3.2.4	Model architectures . . . . .	5
3.2.5	Evaluation Metrics . . . . .	12
3.2.6	Preliminary results . . . . .	13
<b>4</b>	<b>Discussion</b>	<b>14</b>
<b>5</b>	<b>Next steps</b>	<b>14</b>

# **1 Exective summary**

Summary goes here.

## **2 Theroetical basis**

Kind of lit review I suppose? Probably could be adapted stuff we have in the proposal.

### **2.1 Morphometric Classification Homogenisation**

A brief theoretical background.

### **2.2 AI Modelling using Satellite Imagery**

A brief theoretical background.

## **3 Algorithm Design**

### **3.1 Morphometric Classification Homogenisation**

#### **3.1.1 Model architecture**

One paragraph summarising the whole thing.

#### **3.1.2 Data preprocessing**

#### **3.1.3 Morphometric characterisation**

#### **3.1.4 Target labels**

#### **3.1.5 Train/test/validation split**

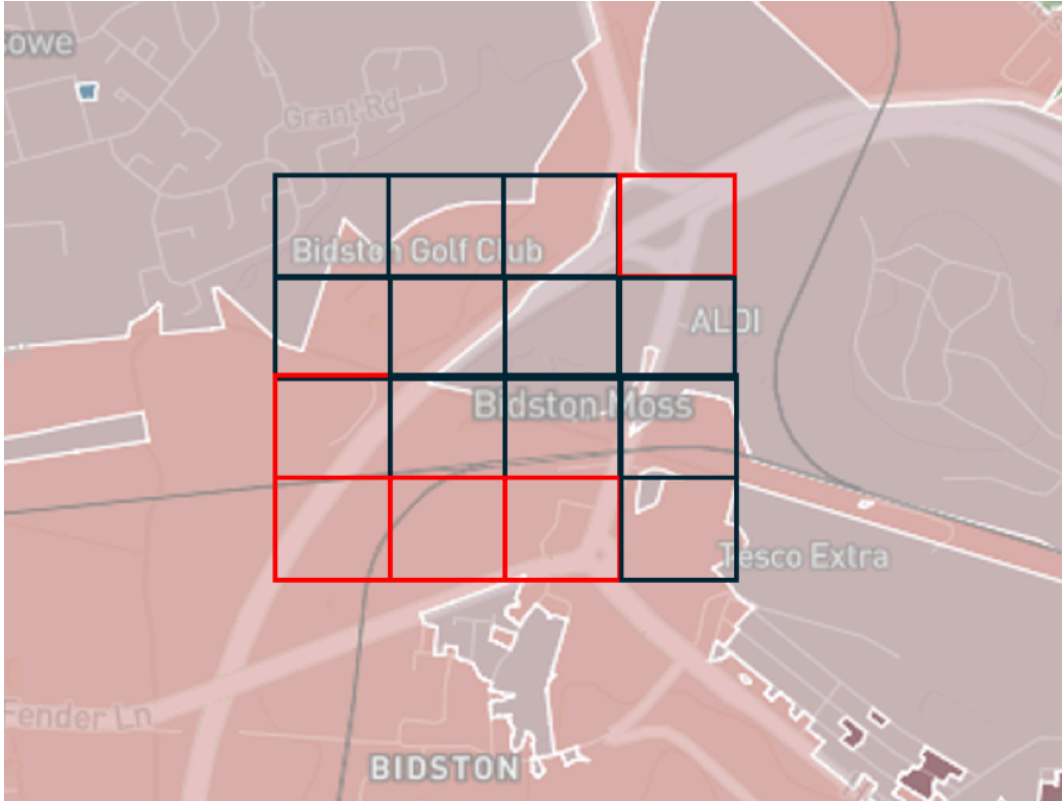
#### **3.1.6 Training and validation**

### **3.2 AI Modelling using Satellite Imagery**

Most of the stuff from `technical_part_turing.qmd` except the data description.

### 3.2.1 Data preprocessing

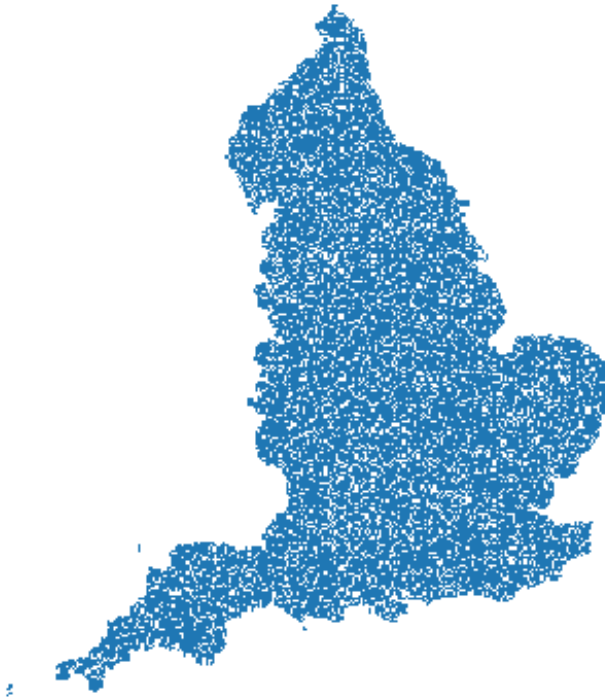
For our analysis, we use two datasets of image tiles at different scales: larger tiles (224 x 224 pixels, covering 2240 x 2240 meters) for segmentation tasks, and smaller tiles (56 x 56 pixels, covering 560 x 560 meters) for classification tasks. The segmentation dataset includes 26,753 tiles (21,402 for training and 5,351 for testing), while the classification dataset consists of 403,722 tiles (342,648 for training and 61,074 for testing). For consistent sampling and comparison, we only use tiles that fully overlap with the spatial signatures, ensuring that each tile aligns with the urban form and function typology in our labeling framework. This alignment supports robust comparison of classification and segmentation outcomes on a pixel-level.



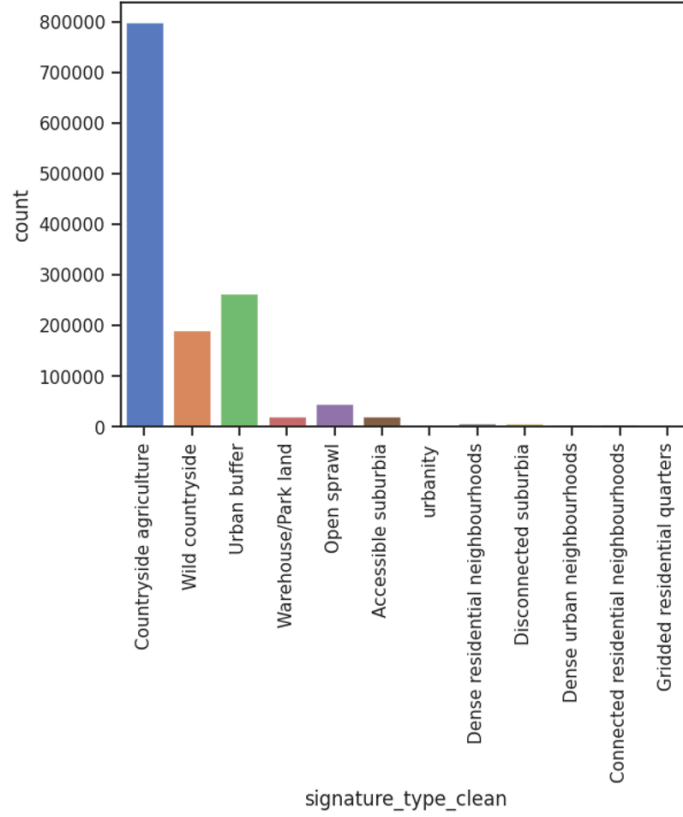
A significant challenge in our dataset is class imbalance, where certain urban fabric types are substantially more represented than others. This imbalance influenced our decisions regarding model architecture and loss function selection, leading us to explore specialized approaches for handling uneven class distributions.

### 3.2.2 Train/test split

We split the dataset into 80% train and 20% test data. The test datasets for segmentation and classification overlap.



### 3.2.3 Unbalanced dataset



### 3.2.4 Model architectures

As part of the AI model design, we tested three main experiments to analyze urban fabric classification and segmentation. First, we conduct a baseline experiment using image embeddings from the SatlasPretrain model <sup>1</sup>, which we fit to an XGBoost classifier to predict urban fabric classes (Approach A). Second, we fine-tune three different geospatial foundation models—SatlasPretrain, Clay, and IBM/NASA’s Prithvi model—to perform segmentation tasks (Approach B). Third, we take the best-performing geospatial foundation model from the segmentation experiments (Clay) and fine-tune it specifically for a classification task (Approach C). To evaluate and compare the results, we report weighted pixel-level accuracy, F1 score, and Intersection over Union (IoU) metrics across the experiments.

- Approach A: Image embeddings + XGBoost model
- Approach B: Fine-tuned geospatial foundation model (segmentation)

---

<sup>1</sup>Bastani, F. et al., 2023. Satlaspretrain: A large-scale dataset for remote sensing image understanding. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 16772-16782.

- Approach C: Fine-tuned geospatial foundation model (classification)

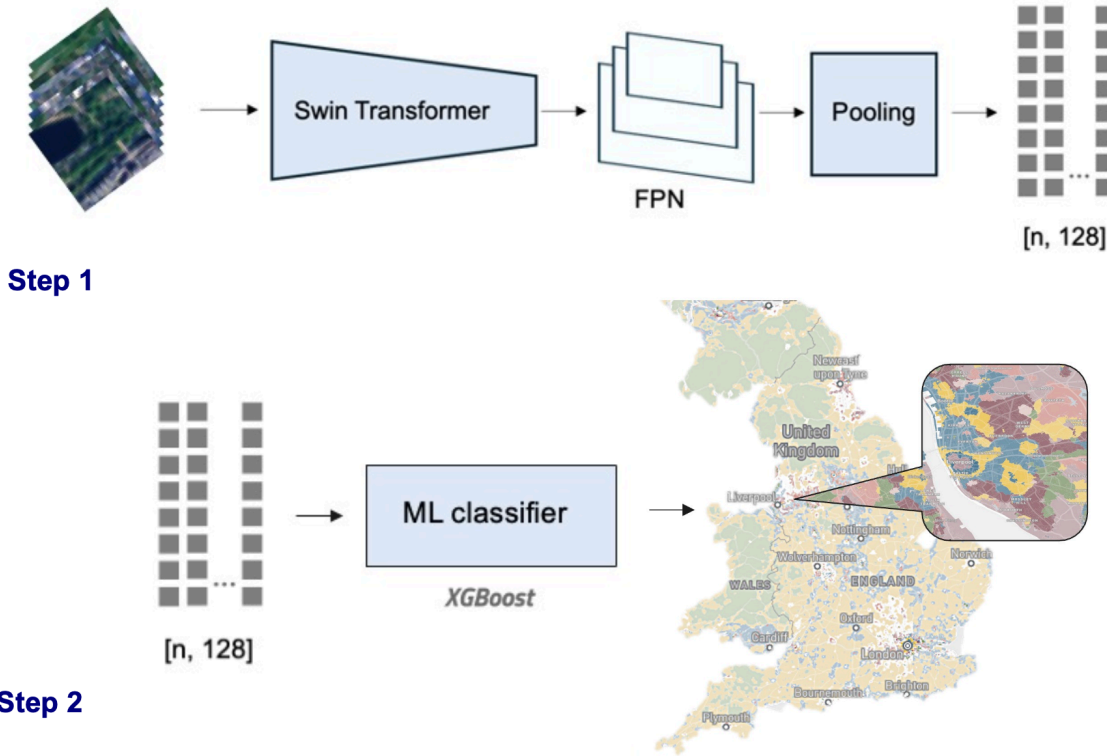
### 3.2.4.1 Baseline approach (Approach A)

The tiles are fed into the geospatial foundation model SatlasPretrain<sup>2</sup> that has been pretrained with more than 302 million labels on a range of remote sensing and computer vision tasks.

The model operates in two main steps:

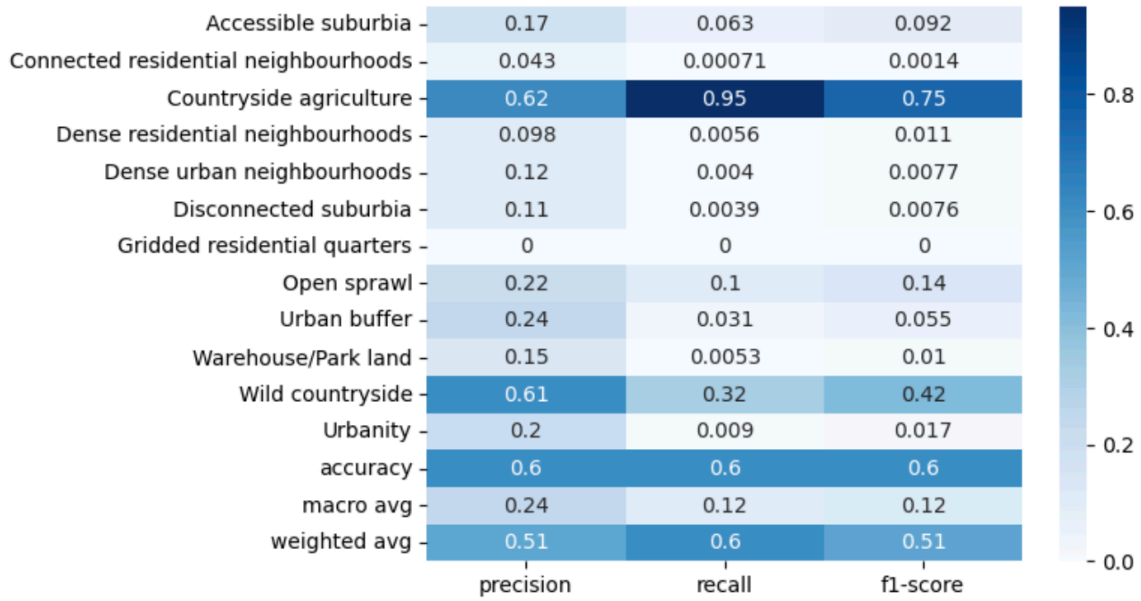
1. Foundation Model: A vision transformer model with a feature pyramid network (FPN) and a pooling layer is used to derive image embeddings—lower-dimensional representations of the images (Fig. 2).
2. Machine Learning Classifier: The image embeddings are then input into an XGBoost classifier to predict urban fabric classes across England.

Our baseline model achieved a moderate prediction accuracy of approximately 61% with varying accuracy across the various spatial signature classes as seen in the figure below.



<sup>2</sup>Bastani, F. et al., 2023. Satlaspretrain: A large-scale dataset for remote sensing image understanding. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 16772-16782.

### 3.2.4.2 Results across spatial signatures



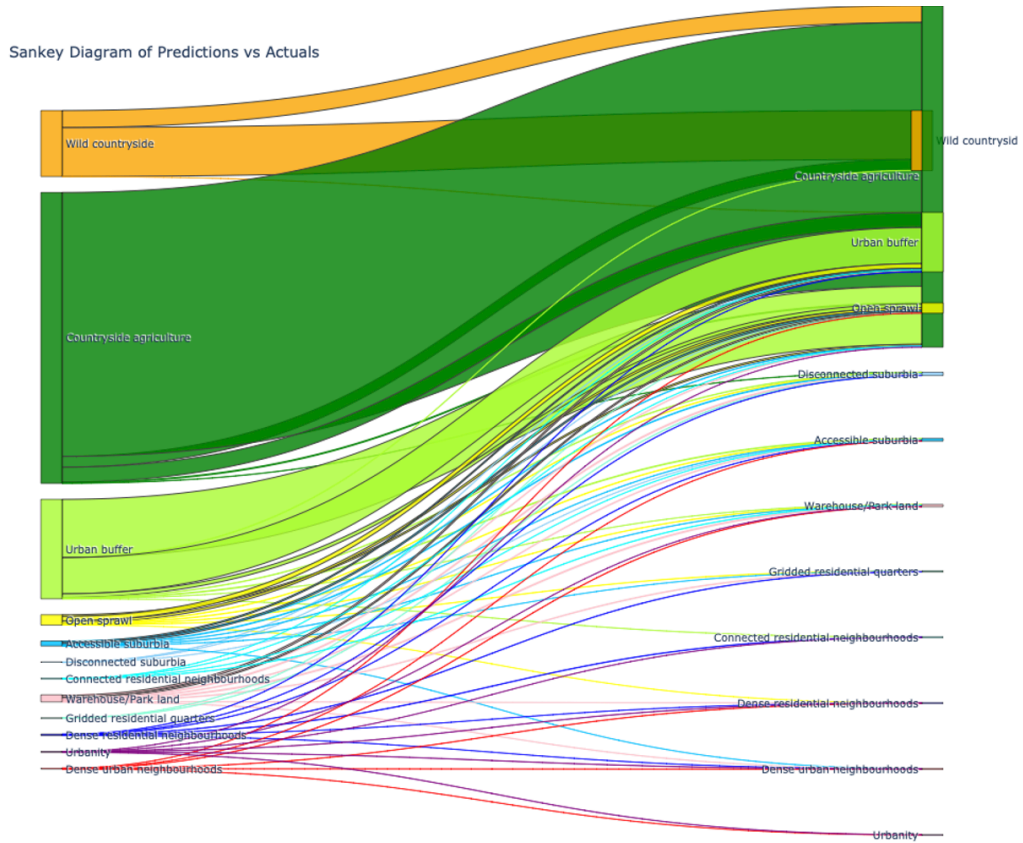
#### *Baseline approach (ordinal)*

In addition to the general classification task, we explored an ordinal regression task to account for the continuous nature of the spatial signatures, which are not strictly categorical. We applied the following ordinal mapping:

```
ordinal_mapping = {
    'Wild countryside': 0,      'Countryside agriculture': 1,
    'Urban buffer': 2,          'Open sprawl': 3,      'Disconnected suburbia': 4,
    'Accessible suburbia': 5,    'Warehouse/Park land': 6,    'Gridded residential quarters': 7,
    'Connected residential neighbourhoods': 8,
    'Dense residential neighbourhoods': 9,    'Dense urban neighbourhoods': 10,
    'Urbanity': 11, }

```

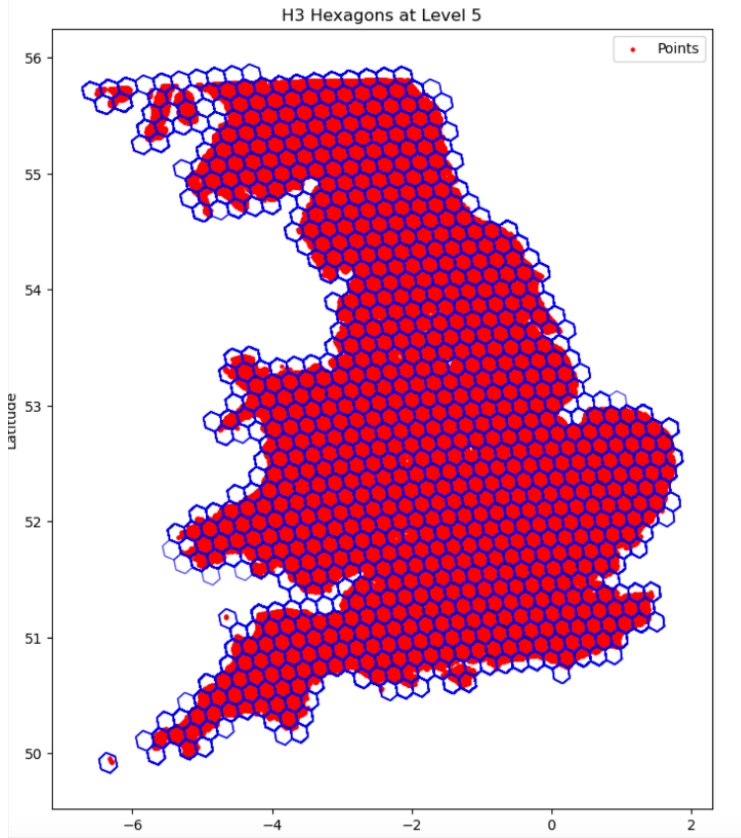
This ordinal approach produced a Mean Absolute Error (MAE) and Mean Squared Error (MSE) of 0.28, along with an  $R^2$  score of 0.62. The Sankey diagram below highlights the primary misclassifications, which tend to occur between similar classes.



### *Baseline approach + spatial context*

To enhance model performance, we incorporated spatial context, enabling the model to account for regional location. We included H3 Level 5 hexagon identifiers as a categorical variable, where each hexagon (approximately 560x560 meters) encompasses around 80 tiles.





### 3.2.4.3 Segmentation (Approach B)

In Approach B, we fine-tuned a state-of-the-art geospatial foundation model for a segmentation task. We used the 224x224x3 image tiles as input. We evaluated three state-of-the-art foundation models, each with unique characteristics as listed below:

Model	Architecture	Dataset Size	Image Sources
Satlas [1]	SwinT	302M labels	Sentinel-2
Clay <sup>3</sup>	MAE/ViT	70M labels	Multiple+
Prithvi <sup>4</sup>	MAE/ViT	250 PB	Sentinel-2/Landsat

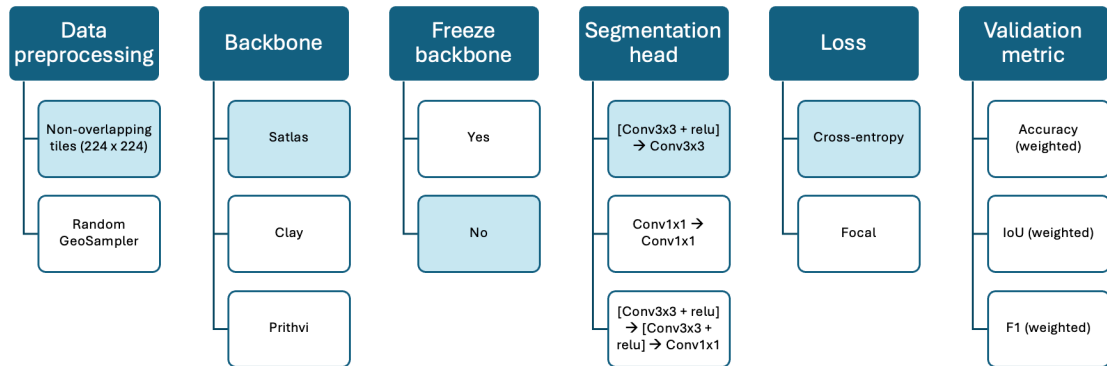
+Multiple sources include Sentinel-2, Landsat, NAIP, and LINZ

The following visualisations show the varying model configurations for the three different approaches tested for the segmentation task. The main difference is the varying backbone.

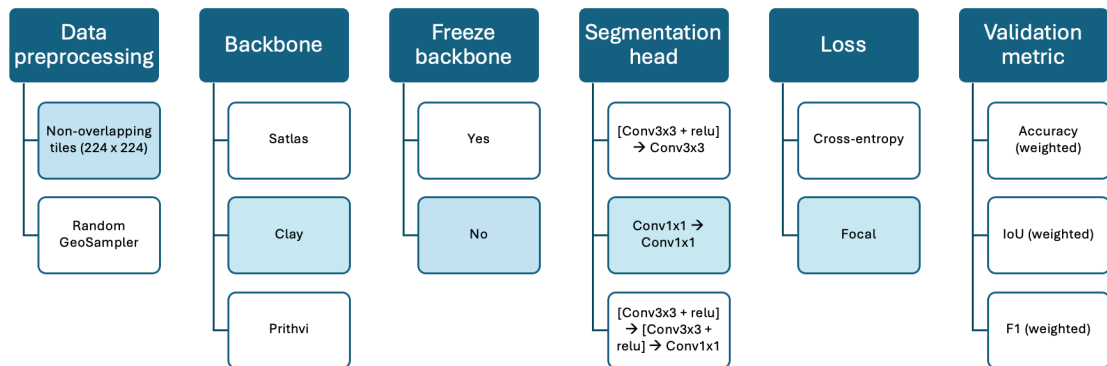
<sup>3</sup><https://huggingface.co/made-with-clay/Clay>

<sup>4</sup><https://huggingface.co/ibm-nasa-geospatial/Prithvi-100M>

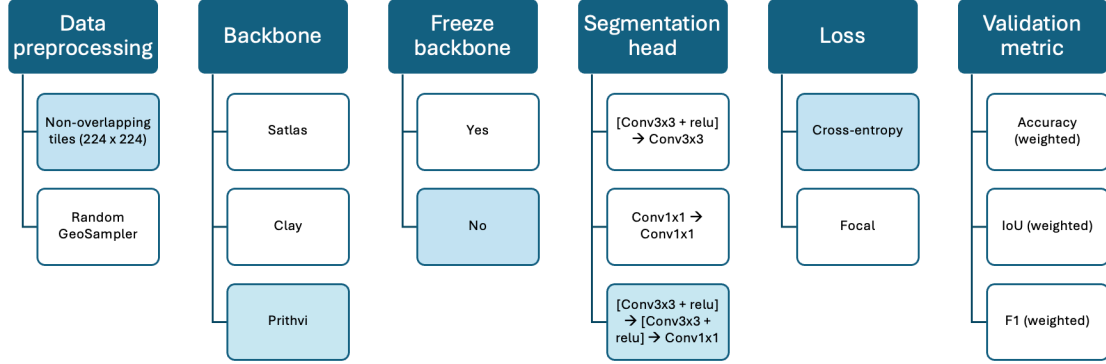
### 3.2.4.4 Model A: Satlas



### 3.2.4.5 Model B: Clay



### 3.2.4.6 Model C: Prithvi



After fine-tuning each foundation model for 10 epochs, we observed the following performance metrics:

Metric	Satlas	Clay	Prithvi
Weighted Accuracy	0.57	<b>0.72</b>	0.62
Weighted IoU	0.33	<b>0.58</b>	0.41
Weighted F1	0.41	<b>0.69</b>	0.58
Training Time/Epoch	9 mins	8 mins	20 mins
Parameters	90M	86M	120M
Implementation Score	5/10	6/10	7/10

The Clay model consistently outperformed other foundation models across all metrics, while also maintaining reasonable training times and computational requirements.

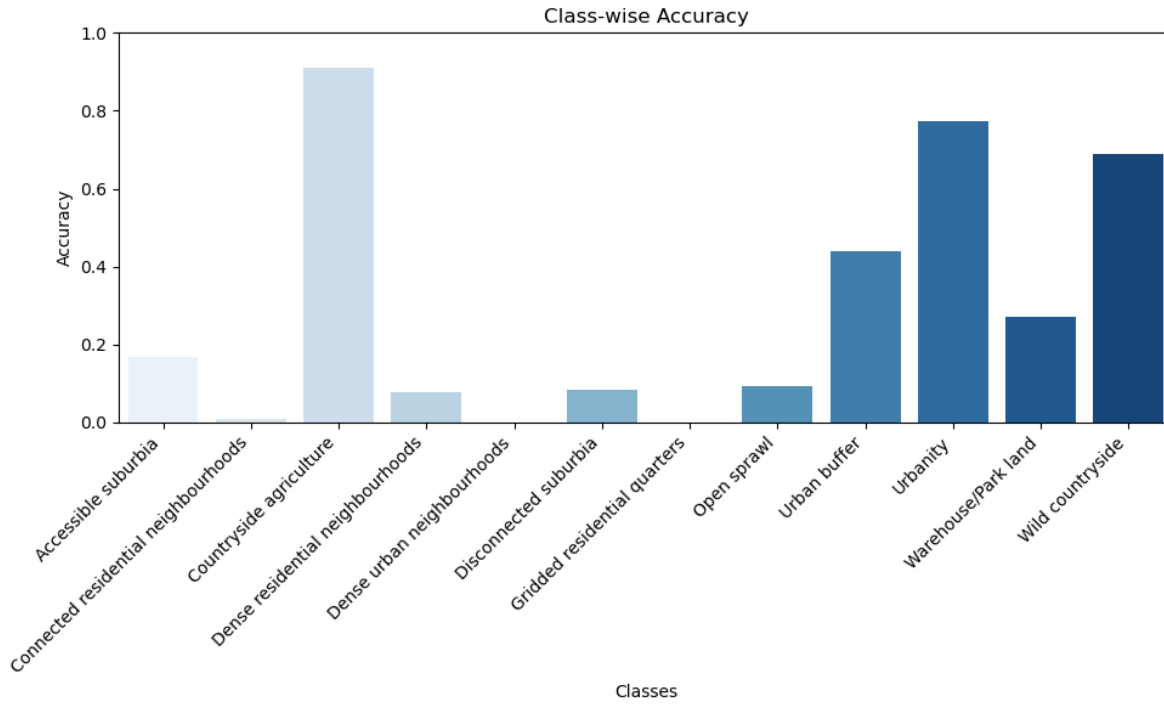
**Loss Function Impact:** The choice of loss function significantly influenced model performance. Focal loss proved particularly effective in handling class imbalance, especially when combined with the Clay model architecture.

### 3.2.4.7 Classification (Approach C)

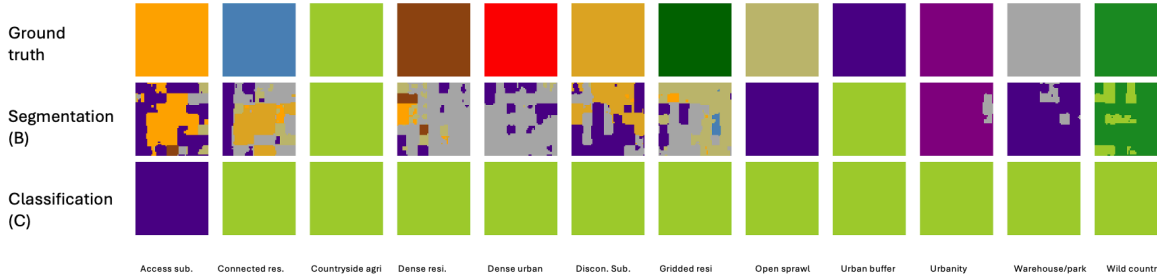
Finally, in Approach C, we fine-tuned a geospatial foundation model for a classification task. We used the 56x56x3 image tiles as input.

For this approach we only used the Clay model as backbone, since it performed the best in the previous experiments.

The accuracy varied across the different classes as seen in the figure below:



This figure shows a comparison between the predicted classes for the fine-tuned geospatial foundation model with segmentation approach (B) and the classification (C).



### 3.2.5 Evaluation Metrics

We employed multiple complementary metrics to evaluate model performance, as follows:

1. **Intersection over Union (IoU):** This metric measures the overlap between predicted and ground truth segmentations, ranging from 0 (no overlap) to 1 (perfect overlap). IoU

is calculated as the area of intersection divided by the area of union between the predicted and actual segmentation masks.

2. **Weighted F1 Score:** This metric provides a balanced measure of precision and recall, particularly important for imbalanced datasets. It is calculated as the harmonic mean of precision (how many of the predicted positives are correct) and recall (how many of the actual positives were identified), weighted by class frequencies.
3. **Weighted Accuracy:** This metric calculates the proportion of correct predictions, weighted by class frequencies to account for class imbalance. It provides a more representative measure of model performance across all classes, regardless of their frequency in the dataset.

### 3.2.6 Preliminary results

Comparing the results is a non-trivial task because the image tiles do not correspond to each other and do not perfectly overlap (42px vs 224 px). To make a fair comparison we thus calculate the pixel-level accuracy scores across the approaches. For this purpose, we predict the full map of the test set and compare the overlapping tiles (as described in sampling). We then calculate the following metrics on a per-pixel level.

#### 3.2.6.1 Overall model performance comparison (Pixel-level)

Our comprehensive evaluation revealed varying levels of performance across the three different approaches:

Approach	Global Accuracy	Macro Accuracy	F1 Score	IoU
Classification (embeddings)	0.76 (0.66)	0.22 (0.13)	0.23	0.63
Classification + H3 level 5	<b>0.87</b> (0.82)	<b>0.42</b> (0.35)	<b>0.45</b>	<b>0.79</b>
Classification + H3 ordinal	0.80 (0.80)	0.26 (0.26)	0.26	0.69
Classification (Clay)	0.59 (0.68)	0.09	0.12	0.38
Segmentation (Clay)	0.73	0.31	0.30	0.58

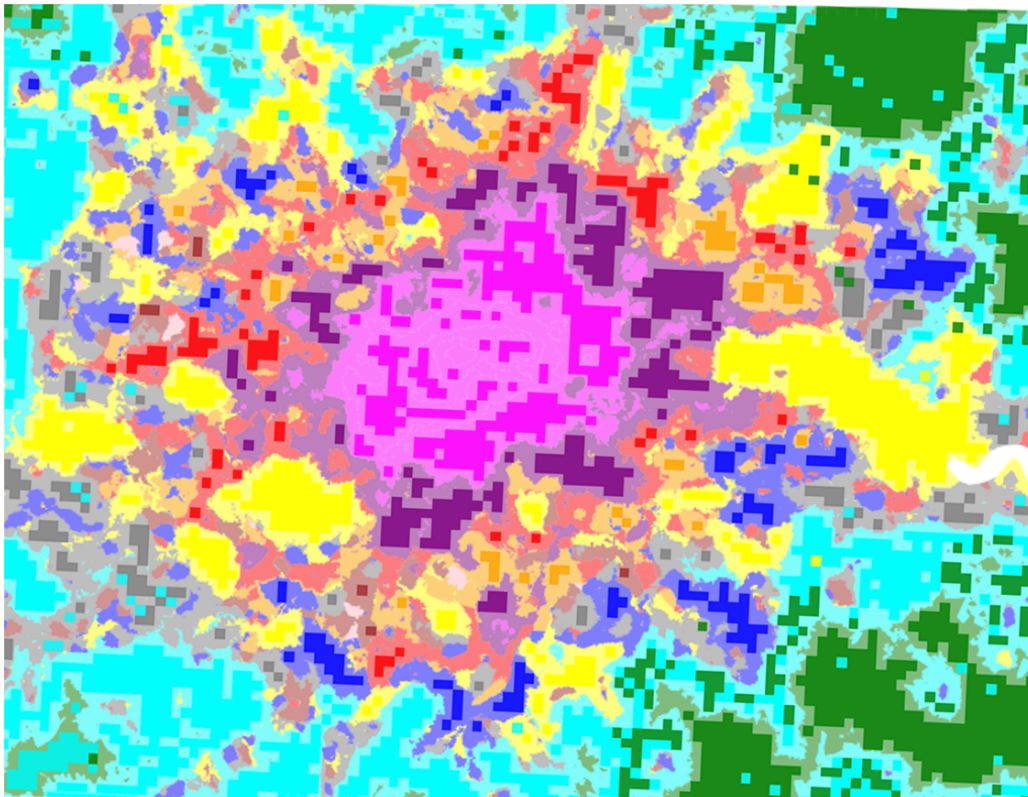
The baseline classification approaches demonstrated varying levels of success: - Basic embedding classification achieved 76% global accuracy (66% balanced) - Integration with H3 level 5 spatial indexing significantly improved performance to 87% global accuracy (42% balanced) - H3 level 5 ordinal classification reached 80% accuracy (26% balanced)

The fine-tuned geospatial foundation model performed better than the fine-tuned classification, with an accuracy score of 0.56 and 0.73 respectively.

Overall, the baseline approach with regional information performed best. This approach is not only the best performing but also relatively efficient to implement. Once the image embeddings are created, the downstream classification can be done in minutes.

#### **3.2.6.2 Prediction example: London**

The following figure shows an example of a prediction for the London area using the whole dataset. Each colour represents a different signature. The background colour represents the ground truth.



## **4 Discussion**

## **5 Next steps**