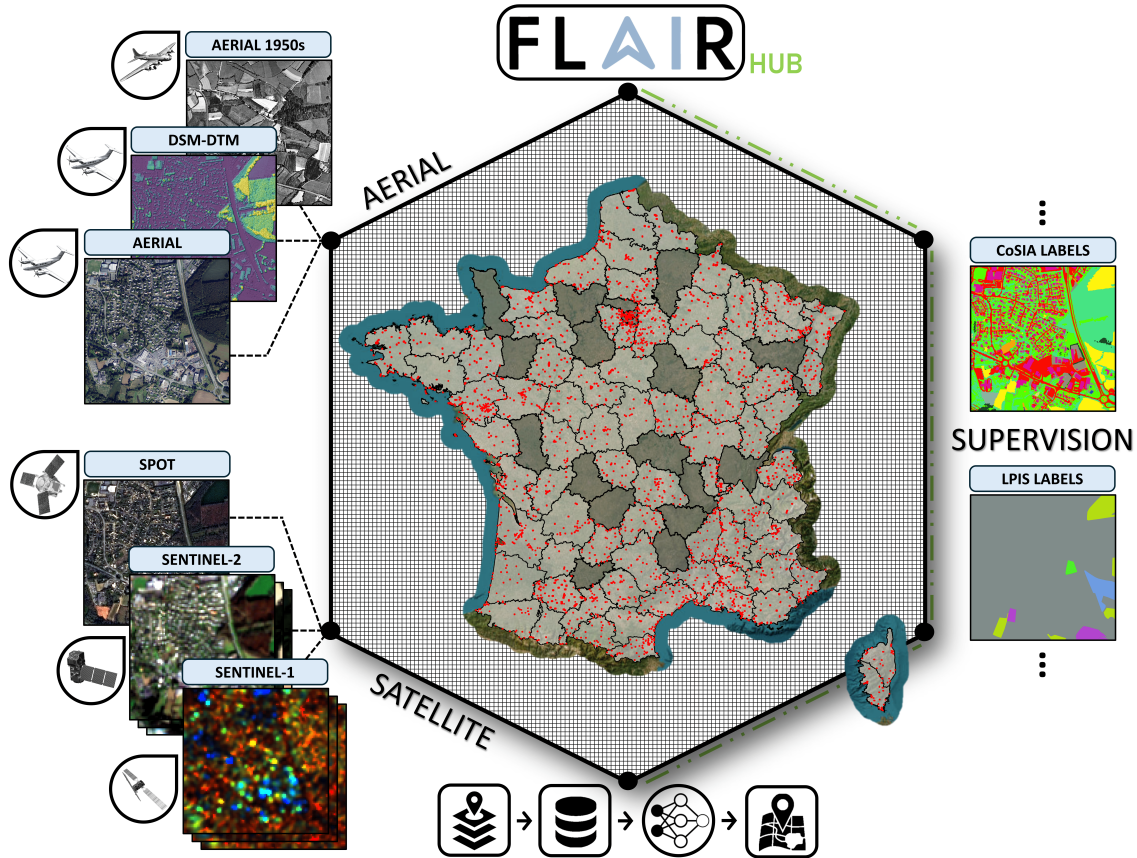


FLAIR-HUB: Large-scale Multimodal Dataset for Land Cover and Crop Mapping

Anatol Garioud, Sébastien Giordano, Nicolas David, Nicolas Gonthier
 Institut national de l'information géographique et forestière (IGN), France
flair@ign.fr



Abstract

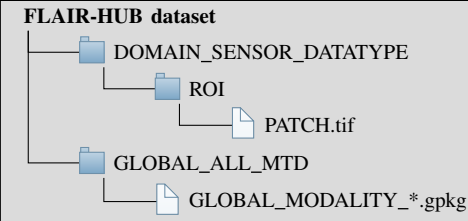
The growing availability of high-quality Earth Observation (EO) data enables accurate global land cover and crop type monitoring. However, the volume and heterogeneity of these datasets pose major processing and annotation challenges. To address this, the French National Institute of Geographical and Forest Information (IGN) is actively exploring innovative strategies to exploit diverse EO data, which require large annotated datasets. IGN introduces FLAIR-HUB, the largest multi-sensor land cover dataset with very-high-resolution (20 cm) annotations, covering 2528 km² of France. It combines six aligned modalities: aerial imagery, Sentinel-1/2 time series, SPOT imagery, topographic data, and historical aerial images. Extensive benchmarks evaluate multimodal fusion and deep learning models (CNNs, transformers) for land cover or crop mapping and also explore multi-task learning. Results underscore the complexity of multimodal fusion and fine-grained classification, with best land cover performance (78.2% accuracy, 65.8% mIoU) achieved using nearly all modalities. FLAIR-HUB supports supervised and multimodal pretraining, with data and code available at <https://ignf.github.io/FLAIR/FLAIR-HUB/flairhub>.

Dataset overview

Figures

- 63 202 918 400 pixels annotated at 0.20 m spatial resolution
- 241 100 patches (512×512)
- 74 spatio-temporal domains and 2 822 areas covering 2 528 km²
- 15 land cover semantic classes (+ 4 optional ones)
- 23/31/46 crop types in 3-level class hierarchy
- 256 221 Sentinel-2 acquisitions
- 532 696 Sentinel-1 acquisitions
- 1.6 m SPOT images aligned
- Aligned historical aerial images
- 20 cm resolution DSM and DTM aligned

Structure



I. Context

In recent years, remote sensing and Earth Observation (EO) had a growing impact on many scientific fields and economic sectors. Extracting information about the Earth’s surface from the sky or space is a key research area. This topic is involved in 11 of the 17 United Nations Sustainable Development Goals [1, 2]. In particular, the automatic analysis of EO images plays an important role in mapping human activities and their impact on the environment. For example, it is useful for applying the European regulation on products derived from deforestation [3], for achieving France’s no net land take target [4, 5], or to monitor soils degradation [6].

An increasing number of regional and national mapping agencies deployed image recognition models to monitor urbanisation, agricultural and forestry areas, risk prevention, and public policy [7, 8]. In this context the French National Institute of Geographical and Forest Information (IGN) [9], in response to the growing availability of high-quality EO data, is actively exploring innovative strategies to integrate these data with heterogeneous characteristics, especially to monitor land cover and crop type across the territory of France and provide reliable and up-to-date geographical reference datasets.

The FLAIR #1 dataset [10], which focused on aerial imagery for semantic segmentation, was released to facilitate research in the field. Building upon this dataset, the FLAIR #2 dataset [11, 12] extends the capabilities by incorporating a new input modality, namely Sentinel-2 satellite image time series, and introduces a new test dataset. Both FLAIR #1 and #2 datasets are part of the currently used by IGN to produce the

French national land cover map reference *Occupation du sol à grande échelle* (OCS-GE) [13].

In this paper, we introduce **FLAIR-HUB**, the increased version of the “French Land cover from Aerospace ImageRy” dataset, the largest multi-sensor land-cover dataset with very-high-resolution annotations. FLAIR-HUB combines very-high-resolution (VHR, 20 cm) images, photogrammetry-derived surface models, and optical Sentinel-2 and SAR Sentinel-1 multi-spectral satellite time series, high-resolution SPOT satellite images and historical analog aerial images from the 1950’s.

These acquisitions’ diverse spatial, spectral, and temporal resolutions offer valuable complementary perspectives for land cover and crop analysis. Over 63 billion pixels have been hand-annotated by geospatial experts, using a nomenclature of 19 land-cover classes and 23 crop type classes. The data spans 2 528 km² across French sub-regions featuring diverse bioclimatic attributes at various times of the year, thus displaying complex and challenging domain shifts.

In addition to this new dataset, we provided an extensive evaluation of the gain made with this dataset on different experiments of multimodal fusion with a recent computer vision backbone [14].

FLAIR-HUB combines heterogeneous and diverse data aiming to foster the development of new large-scale semantic segmentation methods. Given its scale and the complexity of the task it exhibits, it presents an exciting challenge for the machine learning communities. It is also an excellent dataset for multimodal self-supervised methods [15, 16, 17, 18] or data fusion methods [12, 19, 20] thanks to the spatial alignment between modalities. It is also a dataset that will evolve with the addition of new aligned modalities (e.g., hyperspectral, LIDAR) or new annotations (e.g., hedgerows, land use classes).

II. Related Work

A. Semantic Segmentation for remote sensing imagery

Land cover and crop type mapping can be translated into multi-classes pixel classification, also known as semantic segmentation. For almost a decade, deep learning is the de facto solution for semantic segmentation task [38], also in the case of remote sensing imagery [39]. The communities have created more efficient models based on artificial intelligence from FCN [40] to UNet [41], DeepLab [42], Vision Transformer [43] and Swin Transformer [14]. The objective was to find a way to train bigger and better models by adding more long-distance dependency in the inner features and more parameters. This has been done by adding attention mechanics and better optimization schemes.

B. Land cover datasets

Numerous land-cover datasets have been introduced to train semantic segmentation methods, see Table I. Existing datasets usually present a trade-off: they either offer high-resolution annotations but cover a small extent (like Vaihingen [21]), or

TABLE I: Land Cover Datasets. Publicly available datasets for semantic segmentation of land cover using optical remote sensing imagery. Topo refers to topographic information such as DSM, DEM, or slope. ‡ : The FLAIR dataset is included in the new FLAIR-HUB dataset.

| Dataset | Land Cover Annotation | | | | Acquisition | | |
|--------------------------|-----------------------|--------------|--------------|------------------------------|---------------------------|---------------------------|--|
| | Pixels $\times 10^6$ | Resolution | Classes | Source | Resolution | Extent (km ²) | Source |
| Vaihingen [21] | 82 | 8 cm | 6 | visual interpretation | 8 cm | 1 | aerial |
| EuroSAT [22] | 110 | 50 m | 10 | EU Urban Atlas [23] | 10 m | 11 059 | Sentinel-2 |
| MultiSenGE [24] | 534 | 10 m | 14 | visual interpretation | 10 m | 57 433 | Sentinel-1&-2 |
| Landcovernet [25] | 589 | 10 m | 7 | semi-automatic (MODIS [26]) | 10 m | 58 982 | Sentinel-2 |
| MiniFrance [27] | 1 510 | 50 m | 14 | EU Urban Atlas [23] | 50 cm | 53 000 | aerial |
| DynamicEarthNet [28] | 1 889 | 3 m | 7 | visual interpretation | 3 m | 16 986 | Sentinel-1&-2, PlanetFusion |
| LoopNet [29] | 3 133 | 30 m | 6 | visual interpretation | 30 m | 2 820k | Landsat 8 |
| OpenEarthMap [30] | 4 931 | 25–50 cm | 8 | visual interpretation | 25–50 cm | 799 | aerial, UAV., satellite |
| Five-Billion-Pixels [31] | 5 000 | 4 m | 24 | visual interpretation | 4 m | 50 000 | Gaofen-2 |
| LoveDA [32] | 6 000 | 30 cm | 7 | visual interpretation | 30 cm | 536 | aerial |
| Dynamic World [33] | 6 348 | 50m | 9 | semi-automatic / iteratif | 10 m | 634k | Sentinel-2 |
| DeepGlobe [34] | 6 867 | 50 cm | 7 | visual interpretation | 50 cm | 1 717 | Worldview-2/3, GeoEye-1 |
| BigEarthNet [35] | 8 500 | 100 m | 19 | semi-automatic (CLC [36]) | 10 m | 850 k | Sentinel-1&-2 |
| FLAIR [‡] [12] | 20 385 | 20 cm | 19 | visual interpretation | 20 cm / 10 m | 817 | aerial, Topo, Sentinel-2 |
| CatLC [7] | 25 600 | 1 m | 41 | visual interpretation | 1 m | 25 600 | aerial, Sentinel-1&-2, Topo |
| SeasoNet [37] | 63 353 | 10 m | 33 | semi-automatic (CLC [36]) | 10 m | 748k | Sentinel-2 |
| FLAIR-HUB (Ours) | 63 203 | 20 cm | 19/23 | visual interpretation | 20 cm 1.6/10 m | 2 528 | aerial, Topo, SPOT6-7 Sentinel-1&-2 |

provide large-extent coverage but with low-resolution annotations (such as BigEarthNet [35], SeasoNet [37] or SEN12MS [44]), some of them are even the output of fully automatic process [44, 45]. In contrast, FLAIR [12] and FLAIR-HUB offer very high-resolution annotations (20 cm) and covers a large portion of the French territory. FLAIR-HUB is equivalent to SEN12MS [44] in number of pixels but it provides high-quality very high resolution annotation compared to automatic annotation at 10m ground sampling distance.

FLAIR-HUB comprises over 63 billion manually annotated pixels, which is more than 3 times more than FLAIR [12] or CatLC [7], the closest counterparts to our dataset.

The spatial resolution of the annotation is crucial in land-cover analysis. Insufficient resolution prevents the precise measurements of surfaces and boundaries. Furthermore, small-scale features, such as individual houses, lone trees or roads, may not be captured accurately, limiting the potential applications of the derived segmentation. This new dataset tends to answer some of the current challenges about operational very high resolution land cover [46] that are spatial and semantic accuracies, upscaling and data fusion.

C. Crop Type datasets

Monitoring agricultural land cover using Earth Observation involves several key tasks and targets, including parcel delineation, crop type classification or segmentation.

Parcel delineation is similar to an image segmentation task, typically involving two or three classes (*e.g.*, boundary, interior, exterior). This task is especially critical in countries lacking a Land Parcel Information System (LPIS) [49]. Datasets for parcel delineation rely on two main sources of

labels: existing open-data parcel databases such as LPIS [61], and manually annotated parcel boundaries [60]. Most parcel delineation datasets are based on Sentinel-2 imagery at 10 m spatial resolution. Some also incorporate PlanetScope data at 3 m spatial resolution, while AI4Boundaries [61] additionally includes aerial imagery at 1 m spatial resolution, though this higher-resolution data is not available for all years. Since parcel delineation is particularly valuable in areas where labelled data is sparse, recent datasets aim to combine both manual and declarative sources to leverage their respective advantages. However, due to the limited availability of manually labelled data compared to LPIS, careful sampling is required to ensure geographic balance in the dataset [62].

Crop type classification is typically framed as a time series classification task and was one of the earliest approaches to crop mapping. The input time series can be constructed either at the pixel level or at the object level (*i.e.*, parcels), when parcel boundaries are available. The development of crop type classification datasets has been significantly accelerated by the open availability of Satellite Image Time Series (SITS), such as Sentinel-1 and Sentinel-2, as well as parcel vector data like the LPIS in the European Union [55, 56, 59]. The large sample sizes enabled by these data sources have made it possible to apply deep learning architectures effectively to the crop type classification task. Since time series classification does not require dense pixel-wise annotations, labels can also be derived from statistical surveys, such as the LUCAS Survey in the EU [72].

Crop type segmentation extends pixel-level time series classification by incorporating spatio-temporal elements into the

model. Although Sentinel-2 imagery has a spatial resolution of 10 m, frequent acquisitions and rich spectral bands produce datasets that are often one to two orders of magnitude larger than those used for standard time series classification. The availability of extensive label data and satellite imagery has enabled the creation of very large datasets. For instance, the dataset in [52] exceeds 10 TB in size, as it includes all available data with minimal curation [55, 56]. However, as suggested by Roscher et al. [73], using a curated subset of the data can improve usability, achieve better class balance, and maintain comparable classification performance.

Across all dataset types shown in Table II, a key challenge in their creation lies in homogenizing existing parcel and crop-type databases. Crop classification schemes are often country- or region-specific, making direct integration into a unified dataset difficult. To address this, classifications must be standardized before datasets from different sources can be used

together effectively. For smaller datasets, this harmonization is typically done manually by researchers. The resulting crop taxonomy is often a simplified version—either a subset of the original crop types or a grouping based on agronomic similarities. This issue has prompted specific efforts focused on data harmonization. For example, the FIBOA project (Field Boundaries for Agriculture) provides standardized parcel delineation and is used by the Fields of the World dataset [62]. Similarly, the EuroCrops project [74] has harmonized LPIS data across the European Union and has served as a source of standardized crop labels in recent datasets [57, 59]. EuroCrops introduced the HCAT nomenclature, developed in connection with the EAGLE matrix [75] by the European Environment Agency. Other datasets, such as [52] and [54], use crop classifications derived from the Indicative Crop Classification (ICC) developed by the FAO [76]. In addition to LPIS, some transnational harmonized sources like the LUCAS

TABLE II: Agricultural Land Datasets. Publicly available datasets for monitoring agricultural land. *Multi-year — Multi-temporal* : A multi-year dataset includes data from different years over different areas, whereas a multi-temporal dataset provides data from multiple years or times for each area. *TS — SITS* : SITS (Satellite Image Time Series) indicates instance or segmentation datasets that include one image per timestamp in the time series. *TS (Time Series)* refers to classification datasets using tabular time-series data, where each time series represents values of spectral indices (e.g., NDVI) aggregated at the object (e.g., parcel) level.

| Type / Dataset | Data | | | RoI | | | Annotation | | | Size |
|------------------------------------|-------------------------------------|--------------------|----------|--|-------|-----------------------------|--------------------|-------------|--------------------|------------------|
| | Source | Patch Size | #Samples | Extend | Areas | temporality | Parcels | Classes | Source | |
| SITS / Crop type | | | | | | | | | | |
| MunichCrops [47] | Sentinel-2 | 48×48 | - | 4 284 km ² Munich | 1 | multi-year 2016-2017 | 137 k | 17 | LPIS | 42 Go |
| Crop Type Mapping Ghana [48] | Sentinel-1 Sentinel-2 Planets | 32×32 | - | - Ghana | - | 2016 | 8 937 | 4-24 | Survey | 310 Go |
| CV4A Kenya [49] | Sentinel-2 | 2016×3035 | 4 | ~2 450 km ² Western Kenya | 4 | 2019 | >3 000 | 7 | Survey | 3.5 Go |
| ZueriCrops [50] | Sentinel-2 | 24×24 | 28 k | 2400 km ² Zurich - Swiss | 1 | 2019 | 116 k | 5-14-48 | LPIS FOAG | - |
| Pastis-R [51] | Sentinel-1 Sentinel-2 | 128×128 | 2433 | ~4 000 km ² France | 4 | 2019 | 124 k | 18 | LPIS FR | 54 Go |
| Sen4AgriNet [52] | Sentinel-2 | 366×366 | 225 k | All France Catalonia | 1 | multi-year 2016-2020 | 42 M | 9-158 | LPIS FR-Ca | 10 To |
| DENETHOR [53] | Sentinel-1 Sentinel-2 Planets | - | - | 1 152 km ² Germany | 2 | multi-temporal 2018-2019 | 4 500 | 9 | LPIS | 254 Go |
| AgriSen-COG [54] TS-SITS* | Sentinel-2 | 366×366 | 41 000 | - 5 Country | 5 | multi-temporal 2019-2020 | ~7 M | 11-* | LPIS | 28 Go (Label) |
| TS / Crop type | | | | | | | | | | |
| BreizhCrops [55] | Sentinel-2 L2A - L1C | - | 610k | 27200 km ² Britanny France | 1 | 2017 | 610 k | 9 | LPIS FR | 3.2 Go 8.5 Go |
| TimeSen2Crop [56] | Sentinel-2 | - | 1 100 k | 84k km ² Austria | 1 | multi-temporal 2018-2019 | 1 100 k | 16 | LPIS | 1.2 Go |
| CropDeepTrans [57] (early crop) | Sentinel-2 Crop Rot | - | ~7.6 M | 270k km2 FR-NL | 2 | multi-temporal 2016-2020 | 7.65 M | 24—32 | LPIS FR-NL | - |
| Sen4Map [58] | Sentinel-2 | 64×64 | 335 k | Europe | 335 k | 2018 | 20 LC ~35 Crops | Stat Survey | LUCA | 1.2 To |
| EuroCropsML[59] (few shot) | Sentinel-2 | - | 707 k | 3 EU Countries EE-LV-PT | 2 | 2021 | 707 k | 176 | LPIS EuroCrops | 4.7 Go |
| Parcel delineation | | | | | | | | | | |
| AI4SmallFarms [60] | Sentinel-2 | 1 000×1 000 | 62 | Vietnam Cambodia | 62 | t | 439 k | - | image labelling | 1.4 Go |
| AI4Boundaries [61] | Aerial Ortho(1 m) Sentinel-2 | 512×512 256×256 | 7 831 | Europe 7 countries 166 k km ² | 7 831 | 2019 S2-composite | 14.8 M | - | LPIS | 38 Go |
| Fields of the World [62] | Sentinel-2 multi-date | 256×256 | 70.5 k | 4 Continents 24 countries | ~80 | multi-Year S2 - 1 date | 1.63 M | - | FIBOA EuroCrops | - |

statistical survey [72] have been employed, notably in the Sen4Map dataset [58]. An alternative to manual or rule-based harmonization is the development of foundation models with few-shot learning capabilities, which could adapt to different national classification systems without requiring extensive retraining [59].

Another critical challenge stems from the temporal variability of crop types, which introduces significant transfer learning issues. A model trained on a specific year and region may perform well within its domain but generalize poorly across different years—even within the same area. This limitation has motivated the creation of adapted crop type segmentation datasets. Two strategies have emerged to address this issue. Multiyear datasets contain data from only one year per Region of Interest (ROI), but the year varies across ROIs [52]. In contrast, multi-temporal datasets include data from multiple years for each ROI, supporting more robust temporal generalization [54, 57]. Finally, a particularly relevant variant of the temporal generalization problem is early crop classification, where the objective is to identify crop types as early as possible in the growing season, rather than retrospectively at the end of the year. In this context, historical crop type information, *i.e.*, the crop grown on the parcel in the previous season—can provide valuable prior knowledge, as crop rotations are rarely random and often follow agronomic patterns.

D. Multimodal remote sensing imagery datasets

Multimodal datasets are valuable resources for developing models that effectively integrate diverse remote sensing modalities, each contributing complementary information. Designing architectures that can fully exploit the specific strengths of

different sensors remains an open research challenge [63, 77, 78, 79].

Beyond model design, multimodal datasets play a crucial role in self-supervised pretraining of deep learning models [16, 80, 81, 82, 83], as well as in thematic applications such as forest monitoring [65] and super-resolution [68]. An increasing number of datasets have become available for pretraining purposes, typically without annotations and featuring uniform global coverage. These are often based on observations from single-sensor satellite constellations like Landsat and Sentinel [18, 69]. However, truly multimodal datasets—involving multiple sensors—are still relatively rare [84, 85, 86, 87, 88, 89, 90], and only a few of them include ground-truth annotations [91, 92].

Some multimodal datasets provide data from different sensors over non-overlapping regions, resulting in no spatial alignment across modalities [91], while others rely on automatically generated labels [85]. Despite these limitations, multimodal datasets can also be leveraged in cross-modal supervision setups, where one modality is used to predict another. For example, Sentinel-2 and PALSAR-2 data can be used to estimate biomass, supervised by GEDI measurements [93], or to perform cloud removal tasks [94].

As argued by Roscher et al. [73], Earth observation models benefit more from high-quality, diverse, and well-curated datasets than from massive but uniform data acquisitions. The design of FLAIR-HUB is aligned with these findings, offering a dataset that combines extensive modal diversity and large-scale coverage with curated, high-quality annotations.

Table III presents a comparison of multimodal datasets that include at least three spatially aligned modalities and

TABLE III: Multi modal Datasets. Publicly available datasets featuring aligned multi-modal data for Earth observation, with annotations and more than three modalities. Pixel counts are based on the highest-resolution modality. SITS (Satellite Image Time Series). S12: Sentinel-12. Topo: topographic information such as DSM, DEM, or slope. HS: hyperspectral imagery. VHR: aerial very high-resolution imagery. Historical: legacy VHR imagery. LU: land use. LC: land cover. ‡ : The FLAIR dataset is included in the new FLAIR-HUB dataset.

| Dataset | Number of Modalities | Modality | | | | | | | | | | | Task | Pixels $\times 10^9$ |
|-------------------------|----------------------|----------|----------|----------|----------|--------------|----------|------|------------|----|-------|------------|--|----------------------|
| | | VHR | S1 SITS | S2 SITS | SPOT | Landsat t.s. | Topo | ALOS | MODIS SITS | HS | LIDAR | Historical | | |
| DFC 2018 [63] | 3 | ✓ | | | | | | | | ✓ | ✓ | | LULC Segmentation | 2.0 |
| TSAT-TS [15, 64] | 3 | | ✓ | ✓ | ✓ | | | | | | | | Tree Species Classification | 4.7 |
| FLAIR‡ [12] | 3 | ✓ | | ✓ | | | ✓ | | | | | | LC Segmentation | 20 |
| TalloS [65] | 3 | | ✓ | ✓ | | | ✓ | | | | | | Tree Species Classification | 0.16 |
| PASTIS-HD [15, 51] | 3 | | ✓ | ✓ | ✓ | | | | | | | | Crop Type Segmentation | 7.5 |
| Neon Trees [66] | 3 | | ✓ | | | | | | | ✓ | ✓ | | Tree Detection | 363 |
| CatLC [7] | 4 | ✓ | ✓ | ✓ | | | ✓ | | | | | | LC Segmentation | 25.6 |
| S2NAIP [67, 68] | 4 | ✓ | ✓ | ✓ | | ✓ | | | | | | | "World Cover" Segmentation | 136 |
| SatlasPretrain [69] | 4 | ✓ | ✓ | ✓ | | ✓ | | | | | | | "World Cover" Segmentation | 3 087 |
| MADS [70] | 5 | | ✓ | ✓ | | | ✓ | | | ✓✓ | | | LC Segmentation | 1.9 |
| Planted [71] | 5 | ✓ | ✓ | | | ✓ | | ✓ | ✓ | | | | Tree Classification | 3.0 |
| FLAIR-HUB (Ours) | 6 | ✓ | ✓ | ✓ | ✓ | | ✓ | | | | | ✓ | LC & Crop Type Segmentation | 63 |

are accompanied by some form of ground truth annotation. Our dataset includes the highest number of aligned modalities, along with a large number of annotated pixels. This is especially notable given that several other datasets rely on automatic or uncurated label sources [67, 69]. Some datasets, such as FLAIR [12], PASTIS-HD [15], or Neon Trees [66], include only three modalities. Others, such as MADS [70] and Planted [71], do offer five aligned modalities but include fewer annotated pixels than FLAIR-HUB.

To our knowledge, FLAIR-HUB is the only dataset that combines historical aerial very high-resolution (VHR) imagery, SAR, and multispectral time series in a fully aligned setup. Though it does not include LiDAR and hyperspectral data, these modalities remain far less available in open-access, large-scale, well-annotated formats.

III. Dataset global description and naming conventions

The first level of the FLAIR-HUB dataset folder structure is defined by the *modalities per domain* information. This organization was chosen to facilitate the later addition of new domains, modalities or supervisions. More specifically the first level folder name is as follows : *DOMAIN_SENSOR_DATATYPE*.

- **DOMAIN:** The FLAIR-HUB dataset is composed of 74 spatio-temporal domains. We define a spatio-temporal domain as the conjunction of a geographical-extent and an acquisition date. The geographical extent corresponds to the French departments (≈ 100 departments), which are both an administrative subdivision of the territory and the management unit for very high spatial resolution aerial images acquisitions. When aerial surveys and the resulting orthoimages are managed together for two neighbouring departments, we assemble the two department identifiers in the domain's name (e.g., *D059062-2021*, corresponds to departments 59 and 62).

- **SENSOR_DATATYPE:** At the first level of the directory structure, the modality information is organized into two types: (i) image sources (7 folders) and (ii) supervisions (3 folders). Each domain therefore has on release 10 folders. For the image sources the *SENSOR* part relates either to the name of the sensors that acquired the images (SENTINEL1-ASC, SENTINEL1-DESC, SENTINEL2, SPOT) or products derived from aerial surveys (AERIAL-RLT, AERIAL, DEM). The *DATATYPE* gives a hint about the data format (TS: Times Series; PAN: Panchromatic; RGBI: Red, Green, Blue, Infrared channels; ELEV: Elevation). For supervisions, the value of the *SENSOR* field is specified to differentiate between supervisions that have a strong link with a specific sensor (for example where annotation was performed on the sensor's images) and those that can possibly be applied to all sensor images (in which case, the value *ALL* is used). The *DATATYPE* is used to exhibit the name of the supervision. In this name we distinguish between *LABEL* if it corresponds to

an external annotation (LABEL-COSIA, LABEL-LPIS) and MSK (masks) if it is a layer of information derived from a sensor data such as the SENTINEL2 snow and cloud mask (SENTINEL2_MSK-SC). A directory containing the metadata of all patches from all domains *GLOBAL_ALL_MTD* is also available. The format of each metadata file is a GeoPackage, easily readable with GeoPandas Python package containing both the geometric and attribute information of the patches. Metadata about acquisition dates, geometry or radiometry statistics are provided in the *GLOBAL_ALL_MTD* folder.

The second level of the directory structure represents the Regions Of Interest (ROI). A ROI is a set of contiguous patches. In FLAIR-HUB, all the patches are spatially aligned, meaning that regardless the modality and its spatial resolution, they represent the same coverage on the ground : $102.4\text{ m} \times 102.4\text{ m}$. This choice has been made to have patches of 512×512 pixels at 0.2 m spatial resolution. The concept of ROIs was introduced to facilitate the annotation process by pooling the labelling cost. FLAIR-HUB dataset is composed of 2822 ROIs. While patches have fixed sizes on the ground, ROIs have variable sizes. To ensure a perfect nesting of the patches in the ROIs, their width and length are multiples of 512 m . On average, a ROI covers 0.90 km^2 . Table IV provides information on the distribution of ROI sizes. The ROI names (name of the second level folder) is not unique. The ROI identifier corresponds to internal management information that is not important for the dataset, except for the first two characters, which describe the primary and secondary land use (A=Agricultural, F=Forest, N=Natural, U=Urban).

TABLE IV: Region of Interest (ROI) Size Distribution. The area are $512 \times 512\text{ m}$ multiples.

| ROI Size ($512 \times 512\text{ m}$) \times | 1 | 2 | 3 | 4 | 5 | 6 | 8 | 9 | 10 | 12 | 15 | 16 |
|--|-----|-----|-----|------|---|-----|----|----|----|----|----|----|
| Occurrence | 499 | 496 | 260 | 1279 | 6 | 180 | 25 | 18 | 10 | 19 | 23 | 7 |

The patches are finally available in each ROI folder at third level. Their file naming system was designed to be unique and to contain the maximum amount of information about their type. The file naming convention is structured as follows : *DOMAIN_SENSOR_DATATYPE_ROI_POSITION*. The *POSITION* information give the relative position of the patch in the ROI (*row-column*). This name is referred as the *patch_id* in the metadata files. The *GLOBAL_ALL_MTD_GEOM* metadata file provides : the spatial footprint of each patch (*patch_xmin*, *patch_ymin*, *patch_xmax*, *patch_ymax*) in the reference cartographic system of France RF93-Lambert93 (EPSG:2154), the identifier and spatial extent of the ROI to which the patch belongs (*ROI_id*, *ROI_xmin*, *ROI_ymin*, *ROI_xmax*, *ROI_ymax*), the position of the patch in the ROI (*patch_row*, *patch_col* and the size of the ROI in terms of number of patches (*ROI_nbpitch*, *ROI_nbpitchx*, *ROI_nbpitchy*). The *GLOBAL_ALL_MTD_MODALITIES* is a useful metadata for the patches. It provides, for each patch, the relative path to the patch file of each modality.

IV. Dataset modalities

A. Mono-temporal and multi-temporal modalities

In the following, we will detail the characteristics of each source image and supervision. We will focus on describing the spatial, spectral, temporal, and radiometric dimensions of the patches, as well as their origin. The metadata associated with each modality are explained. Table V provides a summary of these characteristics. The FLAIR #1 datapaper [10], had already described some of these modalities (AERIAL_RGBI, DEM_ELEV, AERIAL_LABEL-COSIA). Some of the following paragraphs is taken from this paper.

- **AERIAL_RGBI:** The AERIAL_RGBI modality was produced using the ORTHO HR[®] product; a mosaic of all individual images taken during an aerial survey and mapped onto a cartographic coordinate reference system. Different cameras are used for the aerial image acquisitions. This implies different sensors and consequently, different image characteristics. The final ORTHO HR[®] product has a spatial resolution of 0.20 m (R, G, B and NIR channels). By design, there are no clouds (and therefore no missing data) in the aerial images. Additionally, some radiometric processing methods are applied to obtain the final product. First, radiometric equalization methods are applied to the individual images. Then, a global radiometric correction is carried out on the merged images covering an entire spatial domain to provide a more satisfying colour balance between channels. We consider this radiometric equalization to be relative: within a domain, the radiometric properties are shared (even if the acquisition dates may differ), but there are shifts in radiometry observable between domains (due to both the date of acquisition and the specific radiometric corrections applied). Therefore, the radiometry of R, G, B and NIR images of the ORTHO HR[®] product cannot be considered as a physical measurement of channel reflectance. This radiometric information is encoded as an unsigned 8-bit integer. Each AERIAL_RGBI patch is structured with a

shape of $C \times H \times W$ where $C = 4$ channels and $H \times W = 512 \times 512$ pixels. To take this modality into account more precisely one can use the AERIAL_MTD_RADIO-STATS and the AERIAL_MTD_DATES metadata. AERIAL_MTD_RADIO-STATS provides the mean and standard deviation per patch for the Red, Green, Blue, and Infrared channels. AERIAL_MTD_DATES gives for each patch the date and time of acquisition, the name of the original individual image and the name of the camera used. The temporal distribution of the aerial images can be seen in Figure 1.

- **AERIAL-RLT_PAN:** Such as the AERIAL_RGBI modality, the AERIAL-RLT_PAN is an orthoimage produced with individual aerial surveys images. The image mosaic is stated to be from the 1950's, but in reality, the images intersecting the FLAIR-HUB patches date from 1947 to 1965 (see Figure 1). To keep things simple, we included the date 195X in the domain name of these modalities. 4 domains in the dataset correspond to the same department but with different acquisition dates. These double domains share the same associated AERIAL-RLT_PAN domain. As a result, AERIAL-RLT_PAN is the only modality with 70 domains instead of 74. During these years, aerial photography was not conducted at the department level but rather on a much more local scale, without any specified spatial resolution. Consequently, even though the data has been resampled to a pixel size of 0.4 m in this dataset, the actual spatial resolution can varies significantly from one area to another. Thanks to other unpublished metadata, we estimate that the actual spatial resolution can vary by a factor of 3, ranging from 0.4 m to 1.2 m. The AERIAL-RLT_MTD_DATES metadata provides the acquisition date of each patch and the name of the original historical aerial image. All of the images have been acquired in the Panchromatic interval, being sensitive from blue to red wavelengths. Then, intra-domain radiometric equalization is very challenging due to the nature of old film-based images (vignetting). Locally, there can be significant radiometric differences within a single domain. Additionally, statistical equalization techniques are applied, which result in

TABLE V: Overview of the different data modalities available across the dataset. We provide the details about spatial, temporal, spectral, and radiometric resolution for each modality.

| Modality | Spatial resolution | | Temporal resolution | Spectral Resolution | | | | Radiometric resolution | Volume |
|--------------------|--------------------|------------|---------------------|---------------------|---------|----------|----------------------------------|------------------------|-----------|
| | Patch size | Pixel size | | type | channel | per date | Calibration | | |
| AERIAL-RLT_PAN | 256×256 | 0.4 m | Mono-temporal | Panchromatic | 1 | Relative | Domain equalization (DN) | UInt8 | 11.15 Go |
| AERIAL_RGBI | 512×512 | 0.2 m | Mono-temporal | Multi-spectral | 4 | Relative | Domain equalization (DN) | UInt8 | 232.76 Go |
| DEM_ELEV | 512×512 | 0.2 m | Mono-temporal | Multi-channel | 2 | Absolute | Altitude (m) | Float32 | 365.18 Go |
| SENTINEL1-ASC_TS | 10×10 | 10.24 m | Time Series | Multi-channel | 2 | Absolute | σ_0 Backscatter (no unit) | Float32 | 16.68 Go |
| SENTINEL1-DESC_TS | 10×10 | 10.24 m | Time Series | Multi-channel | 2 | Absolute | σ_0 Backscatter (no unit) | Float32 | 17.96 Go |
| SENTINEL2_TS | 10×10 | 10.24 m | Time Series | Multi-spectral | 10 | Absolute | BOA Reflectance (%) | UInt16 | 41.61 Go |
| SPOT_RGBI | 64×64 | 1.60 m | Mono-temporal | Multi-spectral | 4 | Absolute | BOA Reflectance (%) | UInt16 | 6.24 Go |
| AERIAL_LABEL-COSIA | 512×512 | 0.2 m | Mono-temporal | Mono-channel | 1 | Absolute | Label | UInt8 | 2.20 Go |
| ALL_LABEL-LPIS | 512×512 | 0.2 m | Mono-temporal | Multi-channel | 3 | Absolute | Label | UInt8 | 6.35 Go |
| SENTINEL2_MSK-SC | 10×10 | 10.24 m | Time Series | Multi-channel | 2 | Absolute | Probability (%) | UInt16 | 8.39 Go |

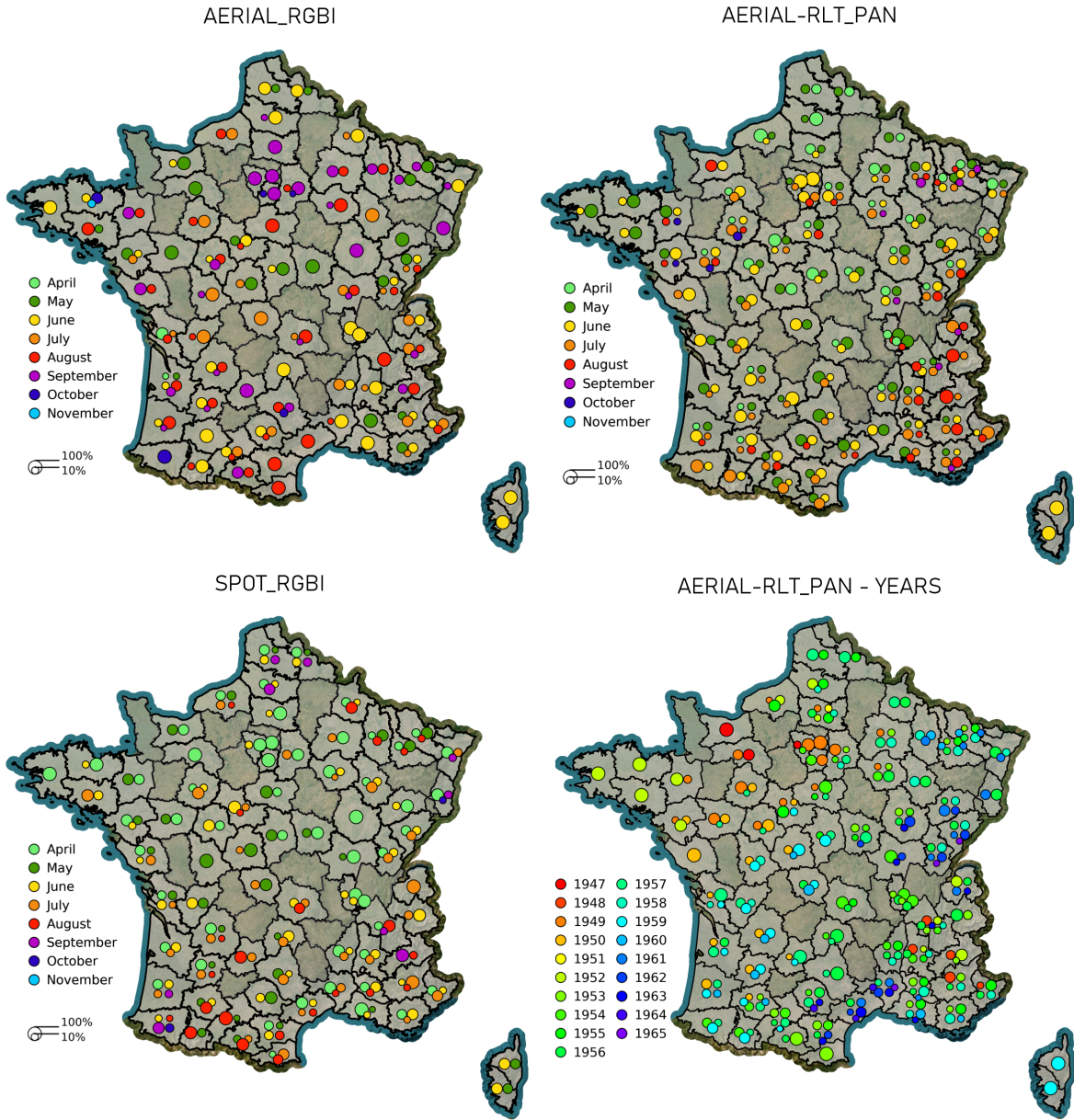


Fig. 1: Temporal Distribution of Mono-temporal Modality Acquisitions. From top left to bottom right, we provide the information about the month of acquisition for the Aerial VHR images, historical images, SPOT satellite images, and finally, the year of acquisition for the historical data. For each domain, the size of the circle is proportional to the amount of data for that month or year.

a significant reduction in global contrast. Therefore, handling the radiometry of this modality is particularly challenging. As with the AERIAL_RGBI modality, we consider the calibration to be relative and not corresponding to a physical measurement. The value is encoded as an 8-bit unsigned integer. The AERIAL-RLT_MTD_RADIO-STATS metadata provides mean and standard deviation information of the panchromatic channel of each patch. Each patch is structured with a shape of $H \times W = 256 \times 256$ pixels.

- **DEM_ELEV:** This modality has two channels : the Digital Surface Model (DSM) and the Digital Terrain Model (DTM).

The DSM gives the altitude, in meters (absolute calibration), for each pixel. Thanks to dense matching techniques, the DSM is derived from the same aerial survey that is used to produce the AERIAL_RGBI modality. This gives the DSM the same spatial resolution as the AERIAL_RGBI but more importantly prevents temporal shifts and ground cover changes between the two products making them temporally coherent. However, there still are small geometric differences between the two products, because orthoimages are projected on the DTM. Moreover, dense matching techniques are applied automatically, potentially introducing noise and artifacts. In particular, radiometrically homogeneous areas

(e.g., parts of the aerial image with little texture) tend to lead to locally false and varying 3D information. Apart from these artifacts, the vertical accuracy of the DSM can be considered to be twice the spatial resolution (0.4 m). The DTM information comes from the RGE ALTI Digital Terrain Model. This product is a national DTM available at a spatial resolution of 1 m. It is constructed from different sources such as dense matching of aerial images, airborne Lidar or, for mountainous areas, from airborne Synthetic Aperture Radar acquisitions. Depending on the source, the vertical accuracy of the RGE ALTI DTM can vary between 0.3 m and 7 m. DTM provides altitude from ground surface, removing buildings, trees, and other objects. The difference between the DSM and the DTM information is then related to the height of buildings or trees. Each patch is therefore sized $C \times H \times W = 2 \times 512 \times 512$. Each channel represents an altitude and is therefore encoded as a single-precision floating-point (Float32).

- **SPOT_RGBI**: This modality has been produced using images from satellite SPOT 6-7 that are acquired each year to produce the French SPOT annual mosaic. For each FLAIR-HUB patch, we chose to use images from the annual SPOT 6-7 mosaic of the same year as the AERIAL_RGBI images, ensuring temporal proximity between the two modalities (see Figure 1 for the temporal distribution of both modality). However, there can be a gap of several months between the two observations, leading to differences in object appearances in the images (e.g., agriculture, forest) or actual land cover changes. The acquisition date, time, and the name of the original SPOT 6-7 image are provided in the SPOT_MTD_DATES metadata table. Initially the images are distributed at 1.50 m resolution but were resampled in the FLAIR-HUB dataset to 1.60 m to be a multiple of 0.2 m. The images include four spectral channels: Red, Green, Blue, and Infrared. Such as Sentinel-2 images the radiometric information is calibrated to Level-2A bottom-of-the-atmosphere reflectance (absolute calibration, no unit). Mean and standard deviation of the 4 channels are given in the SPOT_MTD_RADIO-STATS metadata file. The reflectance percentage is encoded as a UInt 16 : 0 = 0% reflectance and to 10 000 = 100% reflectance. The shape of each patch is then $C \times H \times W = 4 \times 64 \times 64$.

- **SENTINEL-2_TS & _MSK-SC**: Yearly acquisitions from the Copernicus Sentinel-2A and Sentinel-2B satellites are provided for each area. The time-series (*_TS*) data correspond to Level-2A bottom-of-the-atmosphere reflectance. The dataset includes 10 spectral channels (B02, B03, B04, B05, B06, B07, B08, B8A, B11, B12), excluding the atmospheric bands with a 60 m spatial resolution. To ensure consistency with other modalities and fit the spatial extent of patches to a multiple of 0.2 m, the spatial resolution of Sentinel-2 images was resampled to 10.24 m. While the nominal revisit time at the equator is 5 days, the actual length of the time series varies significantly across different areas due to

orbits, ground segment data gaps, or acquisition failures. Consequently, the number of acquisitions in the dataset ranges from 20 to 146 as it can be seen on Figure 2. Level-2A data include geophysical masks for snow and cloud cover (*_MSK-SC*), which are associated with the time-series data to filter out unfavourable acquisition conditions. To reduce the number of files, each dataset patch is structured with a shape of $(T \times C) \times H \times W$, where T represents the acquisition dates stacked in the first dimension. $C = 10$ for *_TS* data and $C = 2$ for *_MSK-SC*. The SENTINEL2_MTD_DATES metadata is available for analysing time series. For each patch, it provides a dictionary containing the length of the time series and the acquisition date corresponding to each position in the series. This metadata is complemented by SENTINEL2_MTD_RADIO-STATS, which describes, in dictionary form, the means and standard deviations of the 10 bands per date. Pixels with a probability of being cloudy strictly greater than 0 are excluded from the mean and standard deviation calculation. When the number of samples for statistical computation is zero, the value *nan* is returned. On the contrary to FLAIR dataset [11, 12], we only consider aligned pixels and do not provide context information from Sentinel-2 time series. That means a reduction of about 93% of pixels from Sentinel-2 per VHR patches.

- **SENTINEL-1ASC & -1DESC_TS**: Yearly time series from the C-Band Copernicus Sentinel-1A and Sentinel-1B satellites are provided. Ground Range Detected (GRD) products are used in dual-polarization mode (VV and VH). Both ascending (*ASC_TS*) and descending orbits (*DESC_TS*) are included separately, as their incidence angles differ significantly. The Sigma nought (σ_0) backscattering coefficient, which represents the normalized radar backscatter intensity of the surface, is calculated for both polarization channels. It provides essential information on surface properties such as roughness, moisture content, and land cover type. No speckle filtering was applied and data averaging results from the GRD product's equivalent number of looks, which is approximately 4. Similar to Sentinel-2 time-series data, each Sentinel-1 patch is structured with a shape of $(T \times C) \times H \times W$, where T represents the acquisition dates stacked in the first dimension, and $C = 2$ corresponding to the two polarization channels (VV and VH). Figure 2 illustrates the number of acquisition available for these modalities. To align with the spatial extent requirements and ensure consistency across modalities, the spatial resolution of Sentinel-1 images was resampled to 10.24 m. Just like the metadata describing the Sentinel-2 series, the files SENTINEL1-ASC_MTD_DATES, SENTINEL1-DESC_MTD_DATES, SENTINEL1-ASC_MTD_RADIO-STATS, and SENTINEL1-DESC_MTD_RADIO-STATS describe the lengths, dates, and statistics of the Sentinel-1 series.

- **AERIAL_LABEL-COSIA**: The AERIAL_LABEL-COSIA supervision consists in determining the land cover at the pixel-level. It is based on photo-interpretation of the AERIAL_RGBI

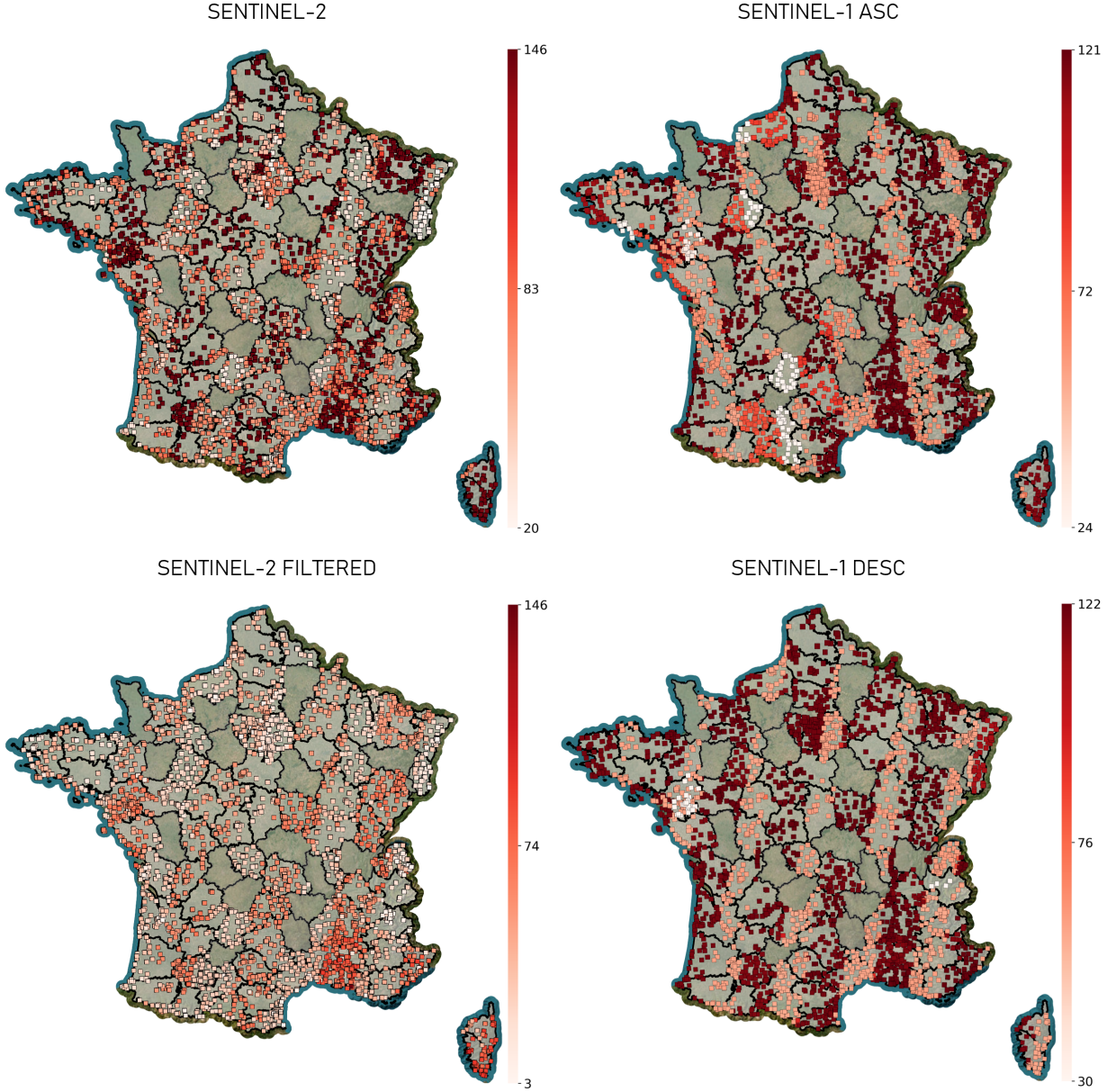





















Fig. 2: Spatio-temporal distribution of multi-temporal modality acquisitions. We plot the number of acquisitions per area for the different STIS; areas are buffered by 5 km for visualization purposes. The acquisition orbits can be distinguished.

images and has been manually produced by experts. An initial spatial multi-level image segmentation approach [95] was applied, simplifying the labelling at the cluster level. We note that this segmentation was not necessarily final, but was modified interactively when deemed appropriate. It was specified that movable objects (*e.g.*, cars, boats) are not to be annotated as such, but to be classified as the underlying cover. For example, a car on an asphalt road is labelled as an impervious surface [10]. We explicitly named this supervision AERIAL_LABEL-COSIA because it is both temporally and geometrically fully consistent with the AERIAL_RGBI images. The AERIAL_LABEL-COSIA patches are then shaped as $C \times H \times W = 512 \times 512$.

The land cover classification consists of 19 classes, ranging from 0 to 18. Table VI provides the list of classes and their respective label number and their frequency in terms of pixel count and percentage (for total, train, valid and test partitions). Please note that the class order has changed compared to previous versions of the dataset ([10, 11]) to ensure that labels start at 0, classes are thematically organized, and classes corresponding to weak labels appear at the end of the nomenclature. Only the first 15 classes are used in the FLAIR-HUB experiments. Classes 0 to 4 correspond to different types of anthropized areas. These are the categories of interest to be used, for example, to monitor soil land take. Classes 5 to 8 represent natural surfaces without agricultural

TABLE VI: Semantic classes and their frequency in the LABEL-COSIA nomenclature of the FLAIR-HUB dataset.

| | Class | LABEL-COSIA | Baseline | Pixels | % total | % train | % valid | % test |
|---|-----------------------|-------------|----------|----------------|---------|---------|---------|--------|
|  | Building | 0 | ✓ | 3 483 982 647 | 5.51 | 5.33 | 5.42 | 6.13 |
|  | Greenhouse | 1 | ✓ | 134 298 230 | 0.21 | 0.23 | 0.15 | 0.20 |
|  | Swimming pool | 2 | ✓ | 17 885 590 | 0.03 | 0.03 | 0.02 | 0.04 |
|  | Impervious surface | 3 | ✓ | 5 796 512 286 | 9.17 | 8.84 | 8.81 | 10.44 |
|  | Pervious surface | 4 | ✓ | 3 530 039 654 | 5.59 | 5.60 | 6.04 | 5.21 |
|  | Bare soil | 5 | ✓ | 2 539 309 904 | 4.02 | 4.21 | 3.96 | 3.49 |
|  | Water | 6 | ✓ | 3 308 863 698 | 5.24 | 5.27 | 5.60 | 4.86 |
|  | Snow | 7 | ✓ | 443 134 338 | 0.70 | 0.72 | 0.40 | 0.88 |
|  | Herbaceous vegetation | 8 | ✓ | 10 998 905 498 | 17.40 | 16.79 | 16.42 | 19.99 |
|  | Agricultural land | 9 | ✓ | 8 665 649 328 | 13.71 | 14.17 | 14.69 | 11.59 |
|  | Plowed land | 10 | ✓ | 1 733 051 984 | 2.74 | 3.04 | 3.03 | 1.63 |
|  | Vineyard | 11 | ✓ | 1 647 328 848 | 2.61 | 2.67 | 2.69 | 2.35 |
|  | Deciduous | 12 | ✓ | 12 731 200 586 | 20.14 | 20.23 | 19.95 | 20.04 |
|  | Coniferous | 13 | ✓ | 4 227 196 348 | 6.69 | 6.46 | 6.35 | 7.62 |
|  | Brushwood | 14 | ✓ | 3 451 432 094 | 5.46 | 5.63 | 5.66 | 4.80 |
|  | Clear cut | 15 | ✗ | 378 090 812 | 0.60 | 0.61 | 0.70 | 0.49 |
|  | Ligneous | 16 | ✗ | 2 809 408 | 0.00 | 0.00 | 0.00 | 0.00 |
|  | Mixed | 17 | ✗ | 36 603 366 | 0.06 | 0.05 | 0.01 | 0.12 |
|  | Undefined | 18 | ✗ | 74 348 348 | 0.12 | 0.12 | 0.09 | 0.13 |

intensive usage. Classes 9 to 11 relate to agricultural areas, while classes 12 to 17 correspond to forested areas. Class 16 (ligneous) and 17 (mixed) can be considered as weak labels. Due to the polygon-based annotation process, the photo-interpreter may sometimes be unable to distinguish between the *deciduous* and *coniferous* classes. In such cases, they use the *ligneous* label. Similarly, if both *deciduous* and *coniferous* are present within the same polygon, the *mixed* code is used. The geometric segmentation was fine enough to ensure that this weak labels remain very rare ($\approx 0.06\%$ of the annotated pixels). Finally, class 19 indicates cases where the photo-interpreter could not provide an annotation (e.g., shadows, uncertainty between certain classes). For a more in-depth analysis of the AERIAL_LABEL-COSIA labels, the AERIAL_MTD_LABEL-COSIA metadata file is available. This file provides, for each patch, the histogram of AERIAL_LABEL-COSIA labels (unit: pixels).

Even though annotations are made with photo-interpretation, some errors are unavoidable, especially for classes that are visually hard to distinguish, such as bare soil and pervious surfaces. Around 37 k randomly chosen polygons were manually annotated, remaining hidden from the annotating teams. This accounted for an area of 18.7 km² equivalent to approximately 468 million pixels. Annotation batches not achieving 95% accuracy were rejected and sent back for re-annotation. This iterative process fostered productive exchanges between the annotators and independent geography experts, ensuring a high-quality dataset.

• **ALL_LABEL-LPIS:** The AERIAL_LABEL-LPIS modality provides information on agricultural surfaces, structured into a semantic crop classification with three hierarchical levels. This annotation is derived from the

French LPIS (Land Parcel Identification System), which consists of declarative parcel data digitally submitted by farmers. These declarations are made within the framework of the European Common Agricultural Policy (CAP), which provides subsidies and support based on land use. However, farmers are not obligated to declare all their parcels but only those relevant to their subsidy claims. As a result, not all agricultural parcels present within a given ROI are included in the AERIAL_LABEL-LPIS modality, since only parcels declared under the CAP are represented. For instance, it is known that vineyard surfaces are frequently absent from the labels due to non-declaration.

LPIS data is digitized by farmers using an online platform, based on the most recent available aerial orthoimagery which is identical to that used in the AERIAL_RGBI modality. To ensure temporal consistency, LPIS data for each Region of Interest (ROI) was selected from the same year as the corresponding AERIAL_RGBI imagery. Consequently, the AERIAL_LABEL-LPIS modality offers temporally aligned supervision across multiple years.

The original LPIS classification includes over 230 crop types. To enhance consistency across years and simplify downstream tasks, these classes were harmonized and organized into a three-level hierarchical taxonomy. While this taxonomy diverges from HCAT [74] at the second and third levels, it remains broadly similar at the first level. The second and third level classes were determined based on crop availability within the ROIs. Crop types with limited spatial extent (e.g., present in < 10 ROIs or domains) were merged with similar crops based on phenological traits and growing seasons. Certain LPIS classes—those dominated by ligneous vegetation (e.g., pastured woodlands or oak/chestnut

TABLE VII: Semantic classes and their frequency in the LABEL-LPIS multilevel nomenclature of the FLAIR-HUB dataset.

| Class LV.1 | LV.1 | Baseline | Pixels | %total | %train | %valid | %test | Class LV.2 | LV.2 | Class LV.3 | LV.3 |
|-----------------------|------|----------|----------------|--------|--------|--------|-------|---------------------------|------|---------------------------|------|
| Grasses | 0 | ✓ | 6 318 583 131 | 10.0 | 9.36 | 8.66 | 12.91 | Grasses | 0 | Grasses monoculture | 0 |
| | | | | | | | | | | Grasses mixture | 1 |
| Wheat | 1 | ✓ | 1 639 060 670 | 2.59 | 2.84 | 3.18 | 1.41 | Wheat | 1 | Winter wheat | 2 |
| | | | | | | | | | | Spring wheat | 3 |
| Barley | 2 | ✓ | 658 266 930 | 1.04 | 1.09 | 0.94 | 0.96 | Barley | 2 | Winter barley | 4 |
| | | | | | | | | | | Spring barley | 5 |
| Maize | 3 | ✓ | 1 288 052 739 | 2.04 | 2.09 | 2.37 | 1.63 | Maize | 3 | Maize | 6 |
| Other cereals | 4 | ✓ | 383 132 669 | 0.61 | 0.62 | 0.63 | 0.55 | Sorghum/Millet | 4 | Sorghum | 7 |
| | | | | | | | | | | Millet / Foxtail millet | 8 |
| | | | | | | | | Other winter cereals | 5 | Winter durum wheat | 9 |
| | | | | | | | | | | Winter triticale | 10 |
| | | | | | | | | | | Winter oat | 11 |
| | | | | | | | | Other spring cereals | 6 | Winter rye | 12 |
| | | | | | | | | | | Spring oat | 13 |
| | | | | | | | | | | Other spring cereals | 14 |
| | | | | | | | | Other cereals | 7 | Other cereals | 15 |
| Rice | 5 | ✓ | 42 118 356 | 0.07 | 0.09 | 0.07 | 0.00 | Rice | 8 | Rice | 16 |
| Hemp/Flax/Tobacco | 6 | ✓ | 47 ,265 258 | 0.07 | 0.04 | 0.09 | 0.17 | Hemp/Tobacco | 9 | Hemp/Tobacco | 17 |
| | | | | | | | | Flax | 10 | Fiber flax | 18 |
| | | | | | | | | | | Other flax | 19 |
| Sunflower | 7 | ✓ | 473 323 562 | 0.75 | 0.90 | 0.98 | 0.12 | Sunflower | 11 | Sunflower | 20 |
| Rapeseed | 8 | ✓ | 366 394 829 | 0.58 | 0.57 | 0.91 | 0.35 | Rapeseed | 12 | Rapeseed | 21 |
| Other oilseed crops | 9 | ✓ | 7 823 820 | 0.01 | 0.02 | 0.00 | 0.00 | Other oilseed crops | 13 | Mustard | 22 |
| | | | | | | | | | | Other oilseed crops | 23 |
| Soy | 10 | ✓ | 122 451 006 | 0.19 | 0.27 | 0.10 | 0.04 | Soy | 14 | Soy | 24 |
| Other protein crops | 11 | ✓ | 147 499 283 | 0.23 | 0.25 | 0.39 | 0.06 | Other protein crops | 15 | Spring peas | 25 |
| | | | | | | | | | | Winter protein crops | 26 |
| | | | | | | | | | | Other protein crops | 27 |
| Fodder legumes | 12 | ✓ | 385 503 065 | 0.61 | 0.67 | 0.60 | 0.45 | Alfalfa | 16 | Alfalfa | 28 |
| | | | | | | | | Other fodder legumes | 17 | Clover | 29 |
| | | | | | | | | | | Other fodder legumes | 30 |
| Beetroots | 13 | ✓ | 117 492 248 | 0.19 | 0.20 | 0.26 | 0.09 | Beetroots | 18 | Beetroots | 31 |
| Potatoes | 14 | ✓ | 67 363 659 | 0.11 | 0.12 | 0.13 | 0.06 | Potatoes | 19 | Potatoes | 32 |
| Other arable crops | 15 | ✓ | 307 449 041 | 0.49 | 0.59 | 0.39 | 0.26 | Fruits and vegetables | 20 | Fruits and vegetables | 33 |
| | | | | | | | | Aromatic/Medicinal plants | 21 | Aromatic/Medicinal plants | 34 |
| | | | | | | | | Other arable crops | 22 | Buckwheat | 35 |
| | | | | | | | | | | Other arable crops | 36 |
| Vineyard | 16 | ✓ | 1 141 947 919 | 1.81 | 1.92 | 1.65 | 1.57 | Vineyard | 23 | Vineyard | 37 |
| Olive groves | 17 | ✓ | 51 694 187 | 0.08 | 0.09 | 0.04 | 0.09 | Olive groves | 24 | Olive groves | 38 |
| Fruit orchards | 18 | ✓ | 504 280 801 | 0.80 | 0.87 | 1.00 | 0.42 | Fruit orchards | 25 | Fruit orchards | 39 |
| Nut orchards | 19 | ✓ | 80 408 145 | 0.13 | 0.16 | 0.03 | 0.11 | Nut orchards | 26 | Nut orchards | 40 |
| Other permanent crops | 20 | ✓ | 97 546 687 | 0.15 | 0.03 | 0.50 | 0.27 | Lavandin | 27 | Lavandin | 41 |
| | | | | | | | | Other permanent crops | 28 | Berries | 42 |
| | | | | | | | | | | Other permanent crops | 43 |
| Mixed crops | 21 | ✓ | 230 723 268 | 0.37 | 0.37 | 0.39 | 0.32 | Mixed crops | 29 | Mixed crops | 44 |
| Background | 22 | ✓ | 48 724 537 127 | 77.09 | 76.84 | 76.68 | 78.15 | Background | 30 | Background | 45 |

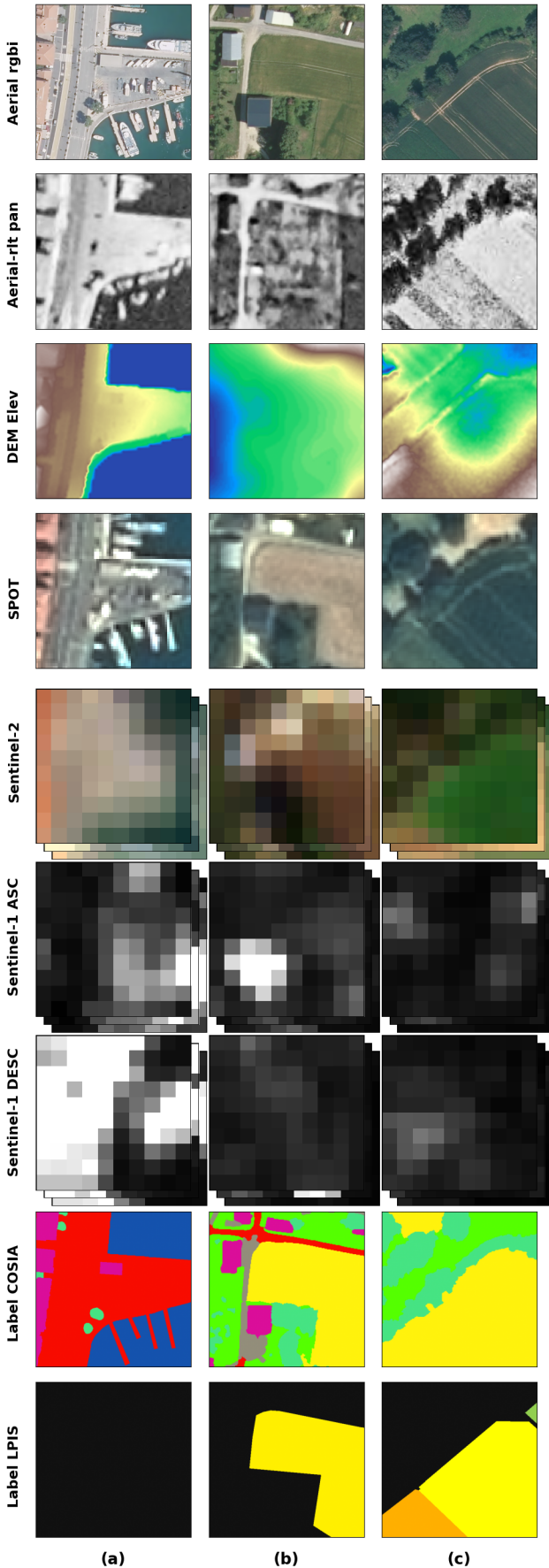


Fig. 3: Example patches from the dataset, illustrating all available modalities. We only plot one image per satellite time series. ASC stands for ascending and DESC for descending.

groves), aquatic areas (e.g., salt marshes or reed beds), or artificial surfaces (e.g., greenhouses) were grouped under the background class. In contrast, mixed-crop classes were deliberately left unmerged with homogeneous crop classes. These mixed classes are better interpreted as unknown and are therefore excluded from training loss computation. The final taxonomy contains 23, 31, and 46 classes at levels 1, 2, and 3 respectively. These have been rasterized as separate channels within a single TIFF file, with a spatial resolution of 20 cm—matching that of the AERIAL_RGBI modality

Spatial discrepancies between AERIAL_LABEL-LPIS and AERIAL_LABEL-COSIA are expected due to the declarative nature of LPIS data integrating land use alongside land cover. Since parcel boundaries in LPIS are manually digitized by farmers and may not align precisely with the actual land cover, pixel-level correspondence between the two labels is not guaranteed. As a result, pixels within a single LPIS parcel can belong to different AERIAL_LABEL-COSIA classes, particularly along the edges or in areas with internal heterogeneity.

As previously mentioned, LPIS data are declarative in nature, provided by farmers, and therefore lack an associated recall threshold at the object (parcel) level. Nonetheless, certain quality criteria must be met for declared parcels. First, the overall error rate of the declared crop type should be approximately 2%, and remain below 5%. Second, with regard to the geometric accuracy of parcel boundaries, the precision for parcel block boundaries separating a parcel from the background is expected to be around 1 m. In contrast, the precision for internal boundaries between adjacent parcels may be lower, with discrepancies of up to 5 m. This reduced accuracy stems from the fact that digitization is performed prior to crop sowing, using ortho-images from previous years, as current-year imagery is not yet available. While external parcel block boundaries tend to be temporally stable and can therefore be digitized more precisely, internal parcel boundaries may vary annually and may not correspond to visible features in historical imagery. Finally, as shown in Table VII, the selected ROIs in the current dataset were not chosen to ensure a minimum area for all crop types. Consequently, this may affect the performance of models trained using LPIS labels.

Figure 3 provides a visual comparison of the different data modalities available in the dataset, shown over three patches sampled from geographically distinct regions. This illustration highlights the diversity of sensor inputs and annotations in terms of spatial, temporal, and radiometric resolutions.

B. Official dataset partitions

Information regarding the usage of patches for training, validation, or testing is provided in the metadata file *GLOBAL_ALL_MTD_SPLIT*. This file establishes the correspondence between each *patch_id* and its

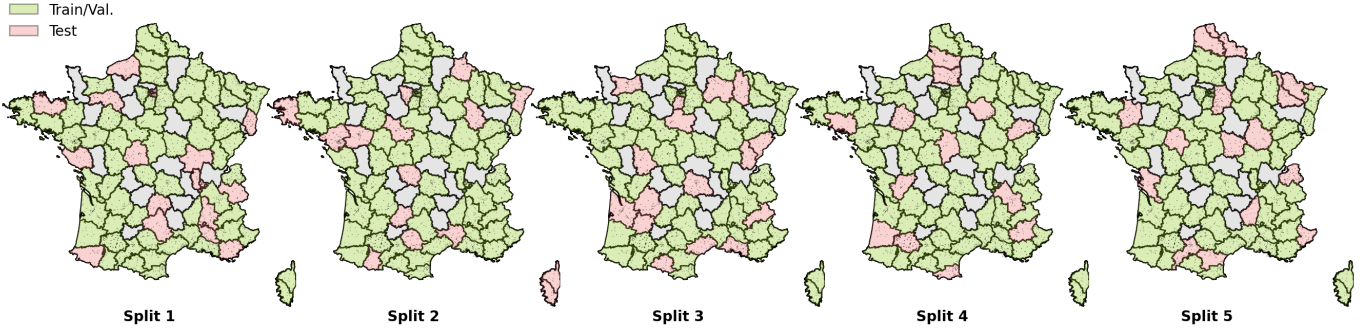


Fig. 4: Spatial distribution of splits in a k-fold configuration. Split 1 corresponds to the official FLAIR-HUB split (also named *split_flairhub*).

associated data split. Seven predefined splits are available: *split_flairhub*, *split_1*, *split_2*, *split_3*, *split_4*, *split_5*, and *split_flairchallenge*.

- **split_1, ..., split_5:** The general approach for creating these splits was to expertly define five clusters of domains. Each cluster is designed to represent all types of landscapes and climates. A separation between train+validation and test is then performed domain-wise. Split_1 used cluster 1 for testing and clusters 2 to 5 for training and validation. Similarly, for splits 2 to 5, the test partition was rotated accordingly. Next, the separation between training and validation was performed ROI-wise. This choice was made to ensure that both the training and validation sets contain a sufficient number of domains while maintaining spatial independence between patches. Indeed, contiguous patches belonging to the same ROI all share the same usage, either training or validation. The train/validation split was performed to obtain a 80% / 20% distribution between train and validation, respectively. Splits 1 to 5 are therefore the splits we will use when a model needs to be evaluated through cross-validation. Figure 4 illustrates the five splits and the rotation of the test set.

- **split_flairhub:** The split named *split_flairhub* is the one primarily used for the experiments presented in this paper. It is identical to *split_1* from the cross-validation setup. For convenience, this split has been duplicated under a separate name. In this split, the TRAIN set comprises 152 225 patches, the VALIDATION set contains 38 175 patches, and the TEST set includes 50 700 patches. Table VI shows the relative proportions of AERIAL_LABEL-COSIA across the training, validation, and test sets, while Table VII reports the corresponding distribution for ALL_LABEL-LPIS.

- **split_flairchallenge:** The entirety of the FLAIR #1 [10] and FLAIR #2 [11, 12] datasets is included within the FLAIR-HUB dataset. FLAIR #1 and FLAIR #2 datasets differ only in their test partitions. The *split_flairchallenge* metadata enables the reproduction of experiments from both publications by indicating, for each FLAIR-HUB patch, whether it was not used in either dataset (*none*), or whether it was used for training (*train*), validation (*valid*), the FLAIR #1 test set

(*test-1*), or the FLAIR #2 test set (*test-2*). The FLAIR #1 and FLAIR #2 test sets are included in the test set of the official split : *split_flairhub*.

V. Baseline architecture

The baseline architecture, namely **UPerFuse**, integrates multiple feature extraction and a fusion strategy which are shown in Figure 5. It is built upon four primary components: a Swin Transformer feature extractor, a UTAE spatio-temporal encoder, a fusion mechanism, and a UPerNet decoder for segmentation.

Swin Transformer: the Swin Transformer [14] module processes mono-temporal data and is designed for hierarchical spatial feature extraction. The input imagery is first partitioned into patches, which undergo linear embedding to transform them into feature representations. These representations are processed through multiple Swin Transformer stages, interspersed with patch merging operations. The feature representations from the patches are processed through a series of Swin Transformer stages, each containing a specific number of Swin Transformer blocks (*i.e.*, 2, 2, 6 and 2). Each stage alternates between layers of regular and shifted window-based multi-head self-attention. These stages include patch merging operations that reduce spatial resolution while increasing the channel dimension. The Swin Transformer blocks within these stages leverage layer normalization, multi-layer perceptrons, and attention mechanisms to capture long-range dependencies and hierarchical feature representations. Skip connections are used to facilitate gradient flow and preserve information across layers, providing a comprehensive understanding of spatial structures.

UTAE: the UTAE (U-Net with Temporal Attention Encoder, [96]) module processes multi-temporal data. It consists of a series of Down-Convolution blocks that progressively reduce the spatial resolution while enhancing feature representation. A temporal attention mechanism is applied to capture dependencies across different time steps in multi-temporal imagery. The extracted features are then upsampled through multiple Up-Convolution blocks to restore spatial resolution

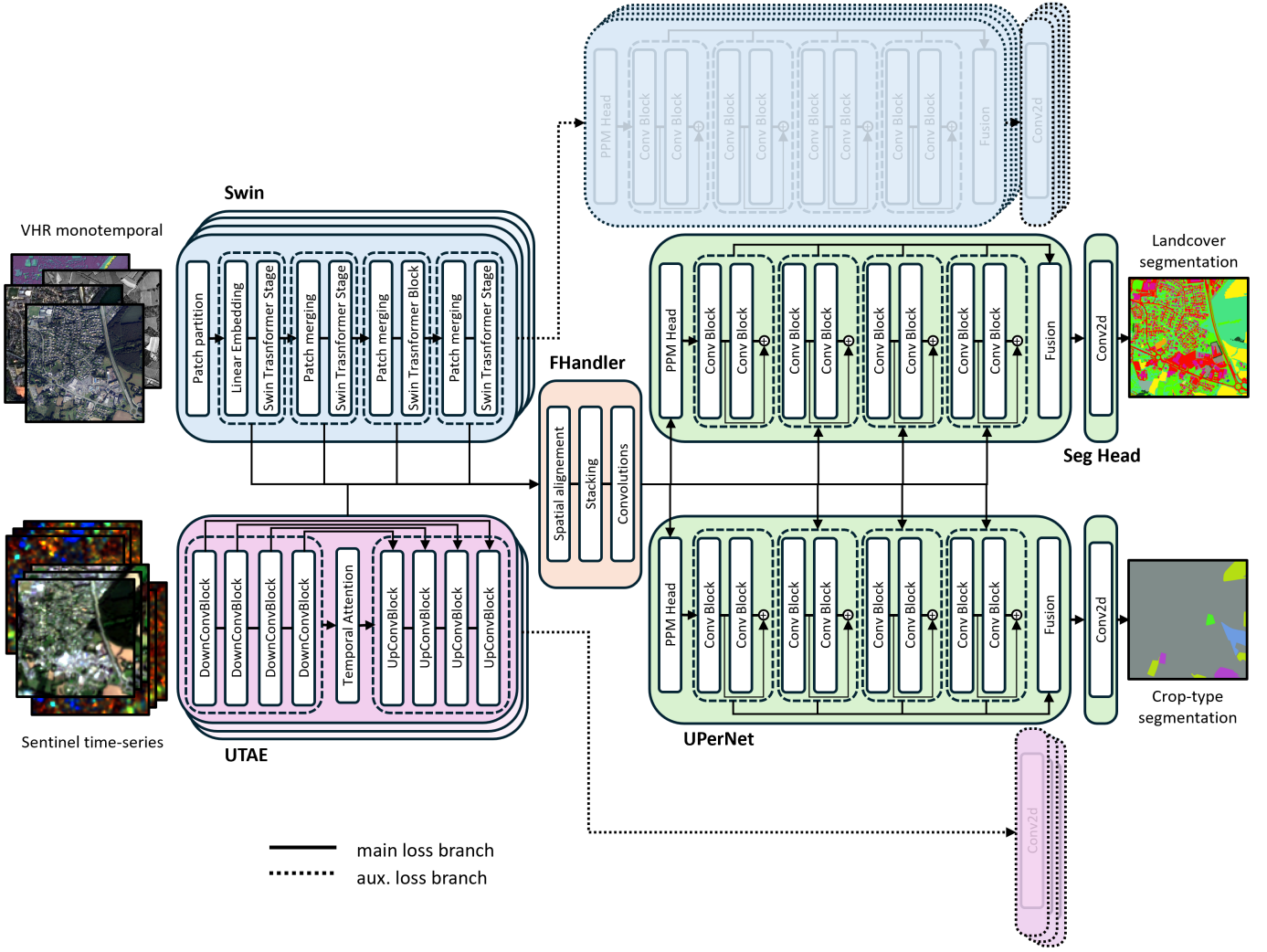


Fig. 5: Architecture of the baseline UPerFuse model designed for multimodal fusion and multi-task semantic segmentation. The transparent modules correspond to auxiliary loss branches.

while preserving temporal context.

Fusion Mechanism: the outputs from the Swin Transformer and UTAE modules are combined through a dedicated FusionHandler module. This fusion mechanism performs spatial alignment via interpolation, feature stacking, and convolutional refinement.

UPerNet Decoder: the fused features are processed using the UPerNet decoder [97], which consists of a Pyramid Pooling Module (PPM) and multiple convolutional blocks. The PPM aggregates multi-scale contextual information, which is subsequently refined through a series of convolutional blocks. Skip connections and hierarchical fusion ensure the preservation of fine-grained spatial details.

Segmentation Head: Final segmentation predictions are obtained through a series of convolutional layers applied to either the UPerNet or UTAE output, depending on the input

modality. If only multi-temporal data is used, the UTAE logits outputs are directly passed to the segmentation heads. The workflow dynamically adapts based on the input modalities (mono-temporal or multi-temporal). If any mono-temporal modality is active, all feature maps are passed through the FusionHandler, enabling per-stage alignment and merging of feature maps. The merged feature maps are then processed for each task (e.g., land cover or crop-type mapping) using a common UPerNet decoder followed by a segmentation head. When only multi-temporal modalities are utilized, the UTAE logits outputs are directly employed for segmentation without requiring additional fusion.

Auxiliary branches are integrated for each modality to improve gradient flow. These branches process encoder feature maps independently through separate decoders, bypassing the fusion step, and directly feeding into segmentation heads. This design ensures robust feature extraction while maintaining flexibility in handling diverse remote sensing data inputs.

Network supervision: The total loss function of the proposed model is designed to handle multiple tasks, multiple modalities, and auxiliary losses, with weighted contributions for each component.

Main Loss: for each task t , let \mathcal{L}_t be the primary loss computed using a task-specific criterion. Given the ground truth labels Y_t and the predicted logits Z_t , the main loss is defined as:

$$\mathcal{L}_t^{\text{main}} = \mathcal{L}_{\text{task}}(Z_t, Y_t) \quad (1)$$

where $\mathcal{L}_{\text{task}}$ represents the task-specific loss function (*e.g.*, cross-entropy loss with class weighting).

Auxiliary Loss: Auxiliary losses are introduced to enhance feature learning by deep supervision [98] and to prevent the multimodal model from focusing on one specific input data. Let \mathcal{M} be the set of auxiliary modalities. For each task t , the auxiliary loss is computed as:

$$\mathcal{L}_t^{\text{aux}} = \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} \mathcal{L}_{\text{aux}}(Z_t^m, Y_t) \quad (2)$$

where Z_t^m represents the auxiliary logits derived from modality m , and \mathcal{L}_{aux} follows the same loss formulation as the main loss.

Total Loss Function: the final loss function balances the main and auxiliary losses with task-specific weights w_t and a global auxiliary loss weight β :

$$\mathcal{L}_{\text{total}} = \sum_{t \in \mathcal{T}} w_t (\mathcal{L}_t^{\text{main}} + \beta \mathcal{L}_t^{\text{aux}}) \quad (3)$$

where \mathcal{T} denotes the set of all tasks, w_t is the weight associated with task t , and β is a global coefficient that balances auxiliary losses against the main task losses.

VI. Benchmark framework and metric

Framework and settings: our implementation is based on PyTorch Lightning [99], leveraging the segmentation-models-pytorch (SMP) [100] library to access pretrained Timm encoders [101]. Additionally, we incorporate the U-TAE network from its official repository [96], using the default architecture but with increased encoder and decoder widths.

The model is optimized using AdamW [102], which decouples weight decay from the optimization step to improve stability. We set the weight decay to 0.01 and use β parameters of (0.9, 0.999), which control the moving averages of gradient moments. Learning rate scheduling is managed by OneCycleLR, which dynamically adjusts the learning rate throughout training to improve convergence, incorporating an initial warm-up phase of 20%. The training procedure spans 150 epochs with a batch size of 5 and an initial learning rate of 0.00005. Experiments are conducted on a high-performance computing (HPC) cluster using 4 to 6 nodes, each equipped with 4 NVIDIA Tesla V100 (32GB memory), A100 or H100 (80GB memory) GPUs. The Distributed Data Parallel strategy in PyTorch Lightning is used to ensure efficient distributed training. Data augmentation techniques include cloud removal and temporal averaging.

For the AERIAL_RGBI and SPOT_RGBI imagery, only three channels were used in the experiments: Infrared, Red, and Green. Although an initial setup included the four channels, only these three were retained based on performance evaluations. For both modalities and the DEM_ELEV channels, as well as AERIAL_RLT-PAN, input data were normalized using the statistics (mean and standard deviation) reported in Table VIII, which were computed over the combined TRAIN and VAL partitions of the *split_fairhub*.

Metric: The performance of the semantic segmentation

TABLE VIII: Radiometric statistics of the mono-temporal modalities. Mean and standard deviation of the mono-temporal modalities across different dataset partitions of the *split_fairhub*.

| | | TRAIN | | VAL | | TRAIN+VAL | | TEST | | ALL | |
|----------------|-----|---------|--------|---------|--------|-----------|--------|---------|--------|---------|--------|
| | | mean | std | mean | std | mean | std | mean | std | mean | std |
| AERIAL_RGBI | R | 105.66 | 52.40 | 105.68 | 51.55 | 105.66 | 52.23 | 104.20 | 52.80 | 105.35 | 52.36 |
| AERIAL_RGBI | G | 111.36 | 45.81 | 111.30 | 44.87 | 111.35 | 45.62 | 111.16 | 45.65 | 111.31 | 45.63 |
| AERIAL_RGBI | B | 102.19 | 44.61 | 102.12 | 43.00 | 102.18 | 44.30 | 101.75 | 44.74 | 102.09 | 44.39 |
| AERIAL_RGBI | I | 106.83 | 39.71 | 105.64 | 40.06 | 106.59 | 39.78 | 104.36 | 40.82 | 106.12 | 40.01 |
| DEM_ELEV | DSM | 321.15 | 549.12 | 270.82 | 486.63 | 311.06 | 537.55 | 384.14 | 523.28 | 326.43 | 535.41 |
| DEM_ELEV | DTM | 317.16 | 549.41 | 266.78 | 486.82 | 307.06 | 537.82 | 380.10 | 523.90 | 322.42 | 535.75 |
| SPOT_RGBI | R | 433.63 | 314.09 | 431.79 | 307.38 | 433.26 | 312.76 | 427.38 | 366.04 | 432.03 | 324.70 |
| SPOT_RGBI | G | 509.12 | 286.21 | 507.29 | 278.16 | 508.75 | 284.61 | 502.06 | 352.90 | 507.34 | 300.28 |
| SPOT_RGBI | B | 468.3 | 228.49 | 465.67 | 215.86 | 467.77 | 226.02 | 462.90 | 283.88 | 466.75 | 239.36 |
| SPOT_RGBI | I | 1134.58 | 530.77 | 1146.81 | 589.68 | 1137.03 | 543.11 | 1085.83 | 502.72 | 1126.26 | 535.28 |
| AERIAL_RLT_PAN | PAN | 125.84 | 38.68 | 126.24 | 37.5 | 125.92 | 38.45 | 125.95 | 39.57 | 125.92 | 38.69 |

models is evaluated using the mean Intersection over Union (mIoU) and Overall Accuracy (O.A.) metrics. For the land-cover task, we exclude the ill-defined classes (see Table VI) and thus evaluate the results over the remaining 15 classes. For the LPIS crop-type task, two classes, rice and other oilseed crops, are absent from the test set, so mIoU is computed over the remaining classes.

VII. Benchmark results

A. CNN-based versus Transformer-based Model

In Table IX, we compare different architectures for the encoder and decoder part using aerial imagery only (3 bands: infrared, red, and green) as input. For this configuration, we evaluate several seminal encoder backbones, including convolutional neural networks (CNNs) such as ResNet [103], ResNeXt [104], ConvNeXtV2 [105] and HRNet [106], as well as transformer-based models such as the Swin Transformer [14] and the Mix Transformer (MiT-BX) encoders [107]. To ensure a fair comparison, we selected versions of these architectures with a similar number of parameters. For the Swin Transformer, we evaluated different model sizes, from "Tiny" to "Large", to explore the trade-off between model complexity and performance. The Tiny variant is optimized for lightweight and fast inference, while the Large variant prioritizes accuracy at the cost of increased computational load. On the decoder side, we compare four well-established architectures: UNet [41], SegFormer [107], DeepLabV3 [108] and UPerNet [97].

Decoder performance across fixed encoders: In the first three blocks of results in Table IX, we alternately fixed the encoder (ResNet50, ConvNeXtV2, and Swin-Base) while varying the decoder (UNet, SegFormer, DeepLabV3, UPerNet). The choice of decoder is particularly important when using the ResNet-50 encoder. Indeed, there is a 10.2% gap in mIoU and a 2.9% difference in OA between the best and worst-performing decoders. Spatial context is known to be very important in the case of semantic segmentation for land cover classification. As a consequence, decoders that incorporate mechanisms to enhance spatial context, such as DeepLabV3 and UPerNet, outperform U-Net. These results had already been observed in [12] for a ResNet-34 encoder and FPN, DeepLabV3 decoders. Although the SegFormer decoder is designed to capture spatial context, its relatively lower parameter count appears to limit its ability to produce strong results. The impact of decoder choice on performance is considerably reduced when using Swin-Base or ConvNeXtV2 encoders. With the ConvNeXtV2 encoder, the performance difference between decoders is only 1.1% in mIoU and 0.7% in OA; with the Swin-Base encoder, the difference is 3.5% in mIoU and 1.9% in OA. We hypothesize that when the encoder effectively captures spatial context and learns spatial relationships between objects, the specific decoding strategy becomes less critical. Notably, with these encoders, the U-Net decoder consistently achieves the best results.

Encoder performance with fixed decoder: In the fourth block of results in Table IX, the decoder is fixed to UPerNet, while several encoders are tested to evaluate their impact on

TABLE IX: Per-Class Evaluation for Land Cover Segmentation – Comparison of Architectures. Class-wise IoU scores for different encoders and decoders baselines using aerial imagery only.

| ENC. | DEC. | mIoU | O.A. | building | greenhouse | pool | imperv. | pervious | bare soil | water | snow | herbaceous | agriculture | plowed | vineyards | deciduous | coniferous | brushwood | PARA. | EP. |
|---------------|-----------|-------|------|----------|------------|------|---------|----------|-----------|-------|------|------------|-------------|--------|-----------|-----------|------------|-----------|-------|-----|
| ResNet50 | UNet | 51.5 | 72.6 | 80.3 | 32.5 | 56.2 | 70.8 | 44.3 | 46.6 | 82.8 | 0.6 | 48.1 | 52.4 | 30.2 | 72.7 | 69.2 | 58.7 | 26.7 | 32.5 | 48 |
| ResNet50 | SegFormer | 51.6 | 72.5 | 77.3 | 56.1 | 52.6 | 69.2 | 40.8 | 46.7 | 81.0 | 0.9 | 48.7 | 54.0 | 25.3 | 71.5 | 69.4 | 57.5 | 23.6 | 24.8 | 19 |
| ResNet50 | DeepLabV3 | 61.7 | 75.5 | 81.1 | 70.1 | 58.0 | 73.1 | 54.0 | 58.6 | 86.8 | 72.1 | 50.8 | 54.9 | 32.6 | 75.1 | 70.4 | 60.6 | 27.0 | 39.6 | 76 |
| ResNet50 | UPerNet | 58.9 | 74.2 | 80.0 | 71.0 | 53.6 | 72.0 | 52.4 | 57.8 | 85.7 | 47.8 | 48.0 | 53.4 | 33.2 | 74.4 | 69.2 | 58.4 | 27.1 | 30.0 | 49 |
| ConvNeXtV2 | UNet | 64.2 | 77.2 | 84.2 | 76.2 | 60.0 | 75.2 | 56.2 | 63.0 | 89.0 | 72.5 | 54.2 | 57.1 | 36.3 | 77.5 | 71.3 | 60.4 | 29.3 | 92.8 | 46 |
| ConvNeXtV2 | SegFormer | 63.1 | 76.5 | 83.6 | 74.0 | 59.8 | 74.5 | 56.5 | 62.6 | 88.8 | 67.1 | 51.9 | 55.2 | 32.9 | 77.1 | 71.3 | 61.3 | 29.9 | 88.5 | 39 |
| ConvNeXtV2 | DeepLabV3 | 63.25 | 76.8 | 83.7 | 75.4 | 59.2 | 75.2 | 56.3 | 60.6 | 89.2 | 65.4 | 52.8 | 56.9 | 34.6 | 77.5 | 71.2 | 60.5 | 30.3 | 96.2 | 35 |
| ConvNeXtV2 | UPerNet | 63.8 | 77.0 | 83.5 | 76.5 | 59.4 | 74.8 | 56.5 | 63.0 | 89.5 | 67.8 | 53.8 | 57.3 | 34.7 | 78.5 | 70.8 | 61.2 | 29.2 | 90.2 | 57 |
| Swin - Base | UNet | 64.8 | 77.9 | 84.7 | 79.0 | 62.2 | 76.2 | 57.5 | 64.2 | 90.6 | 63.8 | 54.9 | 58.3 | 37.6 | 78.3 | 72.0 | 62.5 | 30.1 | 92.0 | 100 |
| Swin - Base | SegFormer | 64.4 | 77.4 | 83.5 | 77.2 | 61.1 | 75.4 | 57.4 | 63.4 | 89.2 | 67.3 | 53.5 | 58.0 | 38.4 | 78.4 | 71.4 | 62.8 | 29.7 | 87.6 | 49 |
| Swin - Base | DeepLabV3 | 61.3 | 76.0 | 80.2 | 73.6 | 44.7 | 72.0 | 55.2 | 60.3 | 89.4 | 64.0 | 51.7 | 57.4 | 34.8 | 77.8 | 69.6 | 61.7 | 27.8 | 95.4 | 106 |
| Swin - Base | UPerNet | 64.1 | 77.5 | 83.9 | 78.4 | 61.6 | 75.7 | 57.2 | 62.9 | 90.3 | 63.4 | 54.3 | 57.1 | 34.8 | 77.7 | 71.7 | 62.6 | 30.2 | 89.4 | 79 |
| ResNet50 | UPerNet | 58.9 | 74.2 | 80.0 | 71.0 | 53.6 | 72.0 | 52.4 | 57.8 | 85.7 | 47.8 | 48.0 | 53.4 | 33.2 | 74.4 | 69.2 | 58.4 | 27.1 | 30.0 | 49 |
| ResNext50 | UPerNet | 58.5 | 74.5 | 81.0 | 73.0 | 54.6 | 72.2 | 53.1 | 57.8 | 86.2 | 32.5 | 49.1 | 55.5 | 31.0 | 77.0 | 68.5 | 57.7 | 27.7 | 29.4 | 87 |
| HRNet32 | UPerNet | 61.2 | 75.1 | 81.6 | 69.9 | 58.1 | 73.1 | 53.6 | 60.8 | 86.6 | 66.6 | 50.4 | 55.2 | 33.0 | 74.8 | 68.9 | 56.0 | 28.9 | 33.4 | 90 |
| ConvNextV2(t) | UPerNet | 62.7 | 76.4 | 82.6 | 75.3 | 59.1 | 73.8 | 55.1 | 60.2 | 88.6 | 64.8 | 53.2 | 55.8 | 35.4 | 76.1 | 70.9 | 60.6 | 29.5 | 29.8 | 43 |
| MiT-B2 | UPerNet | 62.7 | 76.2 | 83.2 | 77.9 | 57.8 | 74.8 | 56.3 | 62.2 | 88.4 | 57.3 | 51.3 | 56.5 | 37.4 | 79.1 | 69.7 | 58.4 | 30.1 | 25.6 | 81 |
| Swin - Tiny | UPerNet | 62.2 | 76.2 | 82.4 | 72.1 | 58.7 | 74.3 | 55.7 | 60.9 | 88.5 | 64.4 | 52.6 | 55.4 | 30.8 | 76.4 | 70.9 | 60.4 | 28.9 | 29.4 | 111 |
| Swin - Small | UPerNet | 63.2 | 76.9 | 83.5 | 77.0 | 60.8 | 75.0 | 56.4 | 61.4 | 89.4 | 56.8 | 53.5 | 57.1 | 37.7 | 77.6 | 70.9 | 61.9 | 28.7 | 50.7 | 114 |
| Swin - Base | UPerNet | 64.1 | 77.5 | 83.9 | 78.4 | 61.6 | 75.7 | 57.2 | 62.9 | 90.3 | 63.4 | 54.3 | 57.1 | 34.8 | 77.7 | 71.7 | 62.6 | 30.2 | 89.4 | 79 |
| Swin - Large | UPerNet | 64.8 | 77.7 | 84.1 | 77.4 | 61.5 | 75.9 | 57.6 | 64.1 | 90.4 | 68.5 | 54.4 | 58.2 | 36.1 | 79.0 | 71.7 | 63.0 | 30.2 | 199.4 | 106 |

performance. Among the encoders, the ConvNextV2 yields the best performance with approximately 30 M parameters, comparable in size to the Swin Tiny. As expected, transformer-based models, such as Swin-Tiny and MiT-B2, outperform traditional CNN-based architectures, such as ResNet50, ResNeXt50, and HRNet32, confirming recent trends in semantic segmentation literature. However the ConvNextV2 demonstrates competitive performance with Swin Transformers, despite being a convolution-based architecture. This highlights the significance of modern optimization paradigms compared to new module (such as attention), which appear to play a crucial role in achieving state-of-the-art results.

Despite this, we ultimately selected the Swin Transformer as our default encoder due to its favourable balance between accuracy and computational efficiency across different configurations, and its greater maturity in modular integration with existing frameworks. We opted for UPerNet as the decoder in our final pipeline due to its flexibility and compatibility with multi-scale feature fusion. Moreover, the Swin-UPerNet architecture won the FLAIR #1 challenge [10] and, we observed that the hierarchical nature of Swin encoders helps to mitigate stitching issues between overlapping patches during inference on very large areas and demonstrates greater robustness within our production lines.

Scaling Swin encoders: performance vs. complexity: Finally, in the fifth part of Table IX, we evaluate the impact of encoder size by testing variants of the Swin encoder architecture. The best performance is achieved with the Swin-Large encoder. This can be attributed to the large volume of annotations available in the FLAIR-HUB dataset, which supports training more complex models. However, the performance gain compared to the Swin-Base encoder is modest, with only +0.7% in mIoU and +0.2% in OA, despite Swin Large having more than twice the number of parameters. This trade-off between model complexity and performance led us to adopt the Swin-Base encoder as the baseline configuration.

B. Multimodality Fusion

B.1 Land Cover Mapping

In Table X, we report land cover segmentation performance across various combinations of input modalities, using a fixed UPerFuse architecture (see Section V) and identical hyperparameters. The best results are achieved when incorporating nearly all available modalities (denoted LC-L), reaching 78.2% OA and 65.8% mIoU. Notably, the inclusion of historical imagery appears to slightly reduce performance (LC-L vs. LC-ALL). This outcome is expected, since historical images were primarily included to support future transfer learning tasks. Their purpose is to enable models trained on recent very high-resolution (VHR) data to generalize to much older imagery. However, the domain gap and temporal shift of these images introduces noise that affects segmentation accuracy.

LC-A vs. LC-L: Adding all available modalities yields only marginal improvements compared to using the Aerial VHR modality alone (OA: +0.9%, mIoU: +1.7%). Users applying the models in production environments must carefully consider whether these modest gains justify the additional complexity and preprocessing effort required for multimodal data. However, the current limited improvement is likely explained by the annotation process, which was performed on the Aerial VHR images. As a result, this modality benefits from strong geometric and temporal consistency with the reference data, which likely facilitates model learning and performance.

LC-A vs. LC-B: Adding elevation information to the Aerial VHR images provides a consistent improvement in performance (OA: +0.6%, mIoU: +1.0%). As shown in Table XI, this benefit is also observed at the class level. Nearly all classes show increased IoU scores when 3D information

TABLE X: Quantitative Evaluation for Land Cover Segmentation. Performance of the UPerFuse architecture with different input modalities during training and testing. Auxiliary losses are used in configurations with more than one modality. PARA.: number of model parameters (in millions). EP.: epoch with best validation score. SITS: Satellite Image Time Series. S1/2: Sentinel-1/2.

| Model ID | Aerial VHR | Elevation | SPOT | S2 SITS | S1 SITS | Historical | PARA. | EP. | O.A. | mIoU |
|----------|------------|-----------|------|---------|---------|------------|-------|-----|-------------|-------------|
| LC-A | ✓ | | | | | | 89.4 | 79 | 77.5 | 64.1 |
| LC-B | ✓ | ✓ | | | | | 181.4 | 124 | 78.1 | 65.1 |
| LC-C | ✓ | ✓ | ✓ | | | | 270.6 | 131 | 78.2 | 65.2 |
| LC-D | ✓ | | | ✓ | | | 93.9 | 85 | 77.6 | 64.7 |
| LC-E | ✓ | | | | ✓ | | 95.8 | 98 | 77.7 | 64.5 |
| LC-F | ✓ | | | ✓ | ✓ | | 97.7 | 64 | 77.7 | 64.9 |
| LC-G | | | | ✓ | | | 0.9 | 89 | 57.8 | 34.2 |
| LC-H | | | | | ✓ | | 1.8 | 106 | 54.5 | 28.2 |
| LC-I | | | ✓ | | | | 89.2 | 94 | 64.1 | 43.5 |
| LC-J | | ✓ | | | | | 89.4 | 97 | 67.4 | 51.2 |
| LC-K | ✓ | | | | | ✓ | 181.4 | 45 | 77.6 | 64.3 |
| LC-L | ✓ | ✓ | ✓ | ✓ | ✓ | | 276.4 | 121 | 78.2 | 65.8 |
| LC-ALL | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 365.8 | 129 | 78.2 | 65.6 |

is included. The only exception is the snow class, for which IoU decreases. While altitude data is expected to reduce class confusion, its effectiveness appears limited in regions with low visual texture. In such cases, noise and artefacts in the elevation data may introduce additional sources of confusion, offsetting potential gains.

LC-B vs. LC-C: Incorporating SPOT imagery alongside VHR and elevation data is expected to provide additional value for two main reasons. First, the availability of a second temporal observation of ground objects, even at a coarser 1.6-meter resolution, may assist in distinguishing certain classes. Second, and more importantly, SPOT data are absolutely calibrated in ground surface reflectance. This radiometric consistency has the potential to reduce discrepancies between the radiometric properties of various VHR aerial images. Despite these theoretical advantages, the addition of SPOT data does not lead to any notable improvement in segmentation performance (OA: +0.1%, mIoU: +0.1%).

LC-A vs. LC-D,E,F: The addition of time series modalities results in only limited performance gains. Specifically, including Sentinel-2 time series (LC-D) yields an increase of 0.1% in overall accuracy and 0.6% in mIoU. Adding Sentinel-1 time series (LC-E) leads to a gain of 0.2% in overall accuracy and 0.4% in mIoU, while the use of Sentinel-1 and Sentinel-2 (LC-F) does not provide a change in overall accuracy and an increase of 0.4% in mIoU. In theory, temporal information could help reduce confusion between certain land cover classes, such as deciduous and coniferous forests or herbaceous and agricultural areas. However, the observed improvements are minimal and as shown in Table XI, the limited contributions are consistently small across classes and across the different temporal modalities

(S1, S2, and S1+S2). The only notable exception is the plowed class, for which mIoU increases by 2.1% with S2 images, 2.3% with S1 images, and 3.1% when both are used. This finding highlights a form of complementarity between optical and radar time series for this specific class. The results for LC-D can also be compared with those from FLAIR #2 [11, 12], where the addition of Sentinel-2 images led to larger improvements: an mIoU gain of 1.23% using fixed hyperparameters, and up to 3.90% when comparing best runs and optimized hyperparameter settings. A key difference between FLAIR-HUB and FLAIR #2 is the removal of the super-patch data, which involved the use of larger spatial footprints for Sentinel-2 patches. This suggests that the performance gains observed in FLAIR #2 were more likely due to an expanded spatial context rather than the model’s ability to learn temporal features.

LC-A vs. LC-G,H,I,J: When using only a single input modality, we can observe that the LC-A (Aerial VHR) configuration performs very well, coming close to the best results achieved with multiple sources of information. These results indicate that the shape and texture of land cover objects, well captured by very high spatial resolution, are key variables for their discrimination. The single-modality temporal configurations Sentinel-2 (LC-G) and Sentinel-1 (LC-H) are less effective, with respective mIoU scores of 34.2% and 28.2%. These results are disappointing but can be explained by the fact that the size of objects in several land cover classes from the nomenclature is on the same scale or smaller than the spatial resolution of Sentinel-2 and Sentinel-1 images. The class-wise metrics (Table XI) highlight, for example, the inability to effectively learn the *greenhouse* and *swimming pool* classes. Furthermore, it is worth noting that the chosen baseline for encoding temporal information, UTAE, has very few parameters (on the order

TABLE XI: Per-Class Evaluation for Land Cover Segmentation. Per-class performance of the UPerFuse architecture with different input modalities. Auxiliary losses are used in configurations with more than one modality.

| Model ID | mIoU | <i>building</i> | <i>greenhouse</i> | <i>pool</i> | <i>imperv.</i> | <i>pervious</i> | <i>bare soil</i> | <i>water</i> | <i>snow</i> | <i>herbaceous</i> | <i>agriculture</i> | <i>plowed</i> | <i>vineyards</i> | <i>deciduous</i> | <i>coniferous</i> | <i>brushwood</i> |
|----------|-------------|-----------------|-------------------|-------------|----------------|-----------------|------------------|--------------|-------------|-------------------|--------------------|---------------|------------------|------------------|-------------------|------------------|
| LC-A | 64.1 | 83.9 | 78.4 | 61.6 | 75.7 | 57.2 | 62.9 | 90.3 | 63.4 | 54.3 | 57.1 | 34.8 | 77.7 | 71.7 | 62.6 | 30.2 |
| LC-B | 65.1 | 85.1 | 80.4 | 62.5 | 76.2 | 58.1 | 64.1 | 90.8 | 62.6 | 55.2 | 57.8 | 37.8 | 79.6 | 72.3 | 63.1 | 31.3 |
| LC-C | 65.2 | 85.2 | 79.1 | 62.1 | 76.4 | 58.3 | 64.8 | 90.9 | 64.4 | 55.1 | 58.4 | 37.4 | 78.6 | 72.3 | 63.4 | 31.8 |
| LC-D | 64.7 | 84.0 | 78.9 | 61.2 | 75.8 | 57.5 | 63.0 | 90.5 | 68.3 | 54.4 | 57.5 | 36.9 | 78.1 | 71.9 | 62.9 | 29.4 |
| LC-E | 64.5 | 84.1 | 78.9 | 62.0 | 76.0 | 57.6 | 63.7 | 90.6 | 62.7 | 54.7 | 57.4 | 37.1 | 78.1 | 71.9 | 63.1 | 30.2 |
| LC-F | 64.9 | 84.0 | 79.3 | 61.1 | 75.6 | 57.7 | 63.8 | 90.5 | 68.1 | 54.9 | 56.9 | 37.9 | 78.1 | 71.7 | 63.7 | 29.6 |
| LC-G | 34.2 | 34.9 | 0.0 | 0.0 | 38.3 | 27.4 | 33.6 | 65.3 | 67.5 | 34.4 | 42.1 | 10.2 | 41.1 | 56.0 | 48.2 | 14.5 |
| LC-H | 28.2 | 42.4 | 1.3 | 0.0 | 35.7 | 23.0 | 36.8 | 57.5 | 10.7 | 29.4 | 42.3 | 5.5 | 25.2 | 53.1 | 46.5 | 13.7 |
| LC-I | 43.5 | 57.2 | 49.8 | 13.8 | 53.2 | 40.0 | 44.0 | 71.0 | 62.0 | 36.9 | 48.2 | 4.6 | 42.1 | 58.7 | 52.9 | 18.1 |
| LC-J | 51.2 | 76.1 | 70.3 | 27.0 | 58.5 | 38.8 | 49.6 | 82.1 | 81.7 | 37.5 | 50.9 | 11.9 | 50.4 | 63.3 | 46.0 | 23.8 |
| LC-K | 64.3 | 83.8 | 77.6 | 59.4 | 75.5 | 57.4 | 63.1 | 90.0 | 62.6 | 53.5 | 57.9 | 38.0 | 78.6 | 72.5 | 64.2 | 30.7 |
| LC-L | 65.8 | 85.3 | 79.1 | 62.0 | 76.6 | 58.2 | 64.7 | 90.5 | 73.4 | 55.1 | 58.6 | 37.5 | 78.6 | 72.3 | 63.5 | 31.1 |
| LC-ALL | 65.6 | 85.3 | 80.3 | 62.8 | 76.5 | 58.5 | 65.1 | 90.8 | 67.6 | 55.0 | 58.6 | 37.9 | 78.4 | 72.3 | 63.2 | 31.4 |

of one million), which is significantly smaller compared to the Swin-Base baseline, with around 90 million parameters. The results achieved using only the SPOT source (LC-I) are satisfying with 64.1% of OA and 43.5% mIoU. We still observe difficulties in learning classes with small objects, while the IoUs for the other classes remain close to the aerial configuration. Furthermore, the land cover annotation is temporally and spatially consistent with the Aerial VHR source. Therefore, it is not a true reference for computing metrics on the SPOT image. This is reflected in the poor results for the *plowed* class. The configuration using only the Elevation modality (LC-J) achieves good results, with 67.4% OA and 51.2% mIoU. The first channel of this modality, the DSM (Digital Surface Model), is derived from the same source images as those used to produce the Aerial VHR modality. However, we observed that combining Aerial VHR with Elevation yields results very similar to using Aerial VHR alone. This suggests that features such as object texture and shape can also be learned from the Elevation modality alone. The addition of color remains crucial for several classes, such as *swimming pools*, differences between types of non-vegetated soils (impervious, pervious, or plowed), and for distinguishing various types of vegetation (*vineyards*, *coniferous*).

In Table XII, we report the performance of various network architectures under the LC-F configuration, with and without enhancement strategies. These strategies include modality dropout (dropout), auxiliary losses (auxloss), and monthly temporal averaging of time series inputs (sentemp). To evaluate the variability of the training process, we also include the results of a 5-fold cross-validation, which are presented in the final rows of Table XII.

Enhancement strategies yield modest and inconsistent improvements across classes. The auxiliary loss configuration

achieves the highest overall mIoU (64.9%) and shows slight gains for classes such as greenhouse (+1.0%), snow (+6.3%), herbaceous (+1.0%), and coniferous (+0.6%). These results suggest that auxiliary losses may support better training dynamics by improving gradient flow and promoting more effective use of multimodal inputs. Temporal averaging yields the highest overall accuracy (77.8%) and small gains for classes like pool, impervious surfaces, agriculture, and plowed fields. While the latter improves slightly (from 38.0% to 38.2%), the limited magnitude of these changes makes it difficult to draw strong conclusions about the specific benefits of temporal averaging. The modality dropout configuration does not lead to significant overall improvement. Some class-level variations are observed, such as higher IoU for snow (+9.4%), but these may reflect training variance rather than a consistent effect. When combined with auxiliary loss, performance again reaches 64.9% mIoU, though with a different distribution of class-level gains. Overall, the enhancements show only limited impact, and their contributions appear to be context-dependent.

The 5-fold evaluation highlights the impact of dataset partitioning on model performance (see Section IV-B for details). While the average mIoU across folds is 66.8% with a standard deviation of ± 1.6 , split_1 stands out with notably lower scores (mIoU: 64.3%, OA: 77.5%) compared to the other folds, which exceed 66% mIoU and reach up to 69.1%. This drop is likely due to the larger validation set in split_1, which combines FLAIR #1 and FLAIR #2 test sets domains. Substantial variability is noticeable at the class level. For instance, snow IoU ranges from 61.8% in split_1 to 92.8% in split_2, largely due to differences in the presence of snow-covered areas in the training, validation and test sets. This uneven representation directly impacts the learning and evaluation and contributes to the observed fluctuation. Similarly, plowed class increases from 38.0% to 53.9%. Other classes such as greenhouse (± 7.5) and bare soil (± 4.6) also exhibit marked

TABLE XII: Per-Class Evaluation for Land Cover Segmentation – Ablation Study. Class-wise IoU scores for the Swin Base-UP baseline using aerial imagery and Sentinel-1/2 time series (denoted setting LC-F). Results include mean and standard deviation over a 5-fold training and evaluation procedure.

| Model ID | CONF. | mIoU | O.A. | building | greenhouse | pool | imperv. | pervious | bare soil | water | snow | herbaceous | agriculture | plowed | vineyards | deciduous | coniferous | brushwood | P.A.R. | EP |
|----------------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|--------|----|
| Evaluation on Fold 1 | | | | | | | | | | | | | | | | | | | | |
| LC-F | - | 64.3 | 77.5 | 84.1 | 78.3 | 61.5 | 75.8 | 57.1 | 63.7 | 90.5 | 61.8 | 53.9 | 56.9 | 38.0 | 78.1 | 71.8 | 63.1 | 29.4 | 95.2 | 72 |
| | + dropout | 64.3 | 77.3 | 83.1 | 77.4 | 57.5 | 75.0 | 57.5 | 63.6 | 89.3 | 71.2 | 54.7 | 56.6 | 36.4 | 78.8 | 71.6 | 62.2 | 29.8 | 95.2 | 61 |
| | + auxloss | 64.9 | 77.7 | 84.0 | 79.3 | 61.1 | 75.6 | 57.7 | 63.8 | 90.5 | 68.1 | 54.9 | 56.9 | 37.9 | 78.1 | 71.7 | 63.7 | 29.6 | 97.7 | 64 |
| | + sentemp | 64.7 | 77.8 | 84.1 | 78.3 | 62.0 | 75.9 | 57.6 | 63.9 | 91.0 | 64.9 | 54.7 | 58.4 | 38.2 | 78.3 | 71.7 | 63.0 | 29.3 | 95.2 | 86 |
| | + aux & drop | 64.9 | 77.3 | 83.4 | 78.3 | 58.4 | 75.1 | 57.1 | 65.5 | 90.6 | 76.3 | 53.3 | 57.2 | 36.5 | 78.3 | 71.8 | 62.3 | 29.2 | 97.7 | 79 |
| 5-fold Evaluation | | | | | | | | | | | | | | | | | | | | |
| LC-F | split_1 | 64.3 | 77.5 | 84.1 | 78.3 | 61.5 | 75.8 | 57.1 | 63.7 | 90.5 | 61.8 | 53.9 | 56.9 | 38.0 | 78.1 | 71.8 | 63.1 | 29.4 | 95.2 | 72 |
| LC-F | split_2 | 67.4 | 79.7 | 82.7 | 58.2 | 57.4 | 75.1 | 59.5 | 67.2 | 90.5 | 92.8 | 54.4 | 65.2 | 47.0 | 81.1 | 76.2 | 62.4 | 41.0 | 95.2 | 43 |
| LC-F | split_3 | 66.9 | 78.6 | 84.6 | 75.9 | 60.0 | 75.4 | 59.7 | 69.2 | 86.4 | 78.0 | 54.1 | 62.3 | 42.5 | 82.4 | 73.2 | 59.9 | 39.9 | 95.2 | 66 |
| LC-F | split_4 | 69.1 | 80.0 | 86.1 | 78.0 | 60.0 | 75.0 | 52.4 | 73.4 | 91.6 | 81.2 | 55.3 | 67.5 | 53.9 | 89.0 | 75.2 | 63.8 | 33.5 | 95.2 | 49 |
| LC-F | split_5 | 66.2 | 79.2 | 84.7 | 73.5 | 60.6 | 76.6 | 58.8 | 59.8 | 88.7 | 62.4 | 55.0 | 63.4 | 53.7 | 82.4 | 75.1 | 63.1 | 34.9 | 95.2 | 78 |
| | Average | 66.8 | 79.0 | 84.4 | 72.8 | 60.0 | 75.6 | 57.5 | 66.7 | 89.5 | 75.2 | 54.5 | 63.1 | 47.0 | 82.6 | 74.3 | 62.5 | 35.7 | - | 62 |
| | \pm | 1.6 | 0.9 | 1.1 | 7.5 | 1.4 | 0.6 | 2.7 | 4.6 | 1.8 | 11.8 | 0.5 | 3.5 | 6.2 | 3.6 | 1.6 | 1.4 | 4.3 | - | 13 |

TABLE XIII: Quantitative Evaluation for Land Cover Segmentation. Performance of the best configuration (UPerNet architecture) on the FLAIR #1 [10] and FLAIR #2 (also FLAIR) [11, 12] test sets. The metrics are computed only on 12 classes (and not 15). PARA.: number of model parameters (in millions).

| Model | Input Data | Supervision Pixel $\times 10^9$ | PARA. | FLAIR #1 mIoU | FLAIR #2 mIoU |
|-----------------|-----------------------------------|------------------------------------|-------|------------------|------------------|
| U-Net [10] | aerial, Topo | 20.3 | 24.4 | 55.7 | X |
| U-T&T [12] | aerial, Topo, Sentinel-2 | 20.3 | 27.3 | X | 58.6 |
| UPerFuse (LC-L) | aerial, Topo, SPOT, Sentinel-1&-2 | 63.2 | 276.3 | 64.1 | 65.0 |

differences, indicating that spatial and seasonal heterogeneity in both training and test splits significantly influences class-level generalization.

Table XIII reports the performance of our best configuration, UPerFuse (LC-L), on the FLAIR #1 [10] and FLAIR #2 [11, 12] test sets. The model achieves 64.1% mIoU on FLAIR #1, compared to 55.7% with U-Net, and 65.0% on FLAIR #2, compared to 58.6% with U-T&T. These improvements are notable but must be interpreted with care: UPerFuse was trained on a much larger dataset, with nearly three times more annotated pixels, and supervision was applied over 15 classes, although the metrics here are computed on 12. The input configuration also includes additional modalities such as SPOT and Sentinel-1, which were not used in the previous baselines. The gains observed are therefore the result of several combined factors, including model capacity, input diversity, and training data volume.

B.2 Crop Type Mapping

Table XIV presents the performance of various configurations for the crop mapping task. Some classes

are excluded from the mIoU computation due to their absence in the test set. Overall, the results underscore the difficulty of the task, particularly as the number of target classes increases. This is especially evident in the lower mIoU scores observed in the last rows of the table, which correspond to more detailed nomenclatures. The performance degradation is largely attributable to strong class imbalance and the presence of rare or sparsely represented classes. A more granular analysis is provided in Table XV, which reports per-class IoU scores for the Level-1 crop mapping task across different input modality configurations. These results are based on the Swin-Base UPerNet baseline, with auxiliary losses applied in all multimodal settings. This table highlights the impact of input modality combinations on class-wise performance and reveals the high variability in segmentation accuracy, particularly for rare crop types.

Severe Imbalance in Crop Type Distribution: It is important to note that the ROIs in the dataset were originally selected for the land cover mapping task. As a result, the class distribution is highly skewed for crop classification: the background class alone accounts for approximately 78%

TABLE XIV: Quantitative Evaluation for crop mapping. Performance of the UPerFuse architecture with different input modalities. PARA.: number of model parameters (in millions). EP.: epoch with best validation score. SITS: Satellite Image Time Series. S1/2: Sentinel-1/2. Classes with zero pixels in the test set are excluded from mIoU computation.

| Model ID | Aerial | VHR | SPOT | S2 | SITS | S1 | SITS | PARA. | EP. | O.A. | mIoU |
|---------------------------------------|--------|-----|------|----|------|----|------|-------|-----|------|------|
| LV.1 - 23 classes (2 classes removed) | | | | | | | | | | | |
| LPIS-A | ✓ | | | | | | | 89.4 | 91 | 86.6 | 24.4 |
| LPIS-B | ✓ | | ✓ | | | | | 181.2 | 99 | 87.1 | 26.1 |
| LPIS-C | ✓ | | | ✓ | | | | 93.9 | 100 | 87.5 | 29.8 |
| LPIS-D | ✓ | | | ✓ | ✓ | | | 97.7 | 80 | 88.0 | 36.1 |
| LPIS-E | ✓ | | ✓ | ✓ | | | | 183.1 | 46 | 87.6 | 30.3 |
| LPIS-F | | | | ✓ | | | | 0.9 | 62 | 85.3 | 23.8 |
| LPIS-G | | | | | ✓ | | | 1.8 | 77 | 84.5 | 18.1 |
| LPIS-H | | | | ✓ | ✓ | | | 2.8 | 61 | 84.9 | 23.8 |
| LPIS-I | | | ✓ | ✓ | ✓ | | | 97.5 | 49 | 87.2 | 39.2 |
| LPIS-J | ✓ | | ✓ | ✓ | ✓ | | | 186.9 | 53 | 88.0 | 35.4 |
| LPIS-K | | | ✓ | | | | | 89.2 | 14 | 84.5 | 15.1 |
| LV.2 - 31 classes (3 classes removed) | | | | | | | | | | | |
| LPIS-I | | | ✓ | ✓ | ✓ | | | 97.5 | 74 | 87.5 | 29.6 |
| LV.3 - 46 classes (8 classes removed) | | | | | | | | | | | |
| LPIS-I | | | ✓ | ✓ | ✓ | | | 97.5 | 111 | 87.3 | 21.4 |

TABLE XV: Per-Class Evaluation for 1st Level Crop Mapping. Class-wise IoU scores for the Base-UP baseline with different input modalities. Auxiliary losses are used in configurations with more than one modality. The *rice* and other *oilseed crops* classes are excluded from mIoU computation due to having zero pixels in the test set.

| Model ID | mIoU | grasses | wheat | barley | maize | o. cereals | flax/hemp tobacco | sunflower | rapeseed | soy | o. protein c. | fodder legumes | beetroots | potatoes | o. arable c. | vineyard | olive groves | fruits orchards | nut orchards | o. permanent c. | mixed c. | background |
|----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------------|-------------|-------------|-------------|---------------|----------------|-------------|-------------|--------------|-------------|--------------|-----------------|--------------|-----------------|------------|-------------|
| LPIS-A | 24.4 | 49.4 | 34.2 | 13.1 | 60.5 | 3.5 | 2.7 | 12.6 | 38.0 | 0.0 | 3.1 | 13.3 | 53.9 | 7.5 | 19.7 | 43.4 | 13.5 | 36.8 | 2.9 | 14.8 | 1.5 | 88.6 |
| LPIS-B | 26.1 | 50.1 | 41.8 | 16.3 | 62.0 | 3.1 | 2.3 | 16.6 | 45.1 | 0.1 | 3.0 | 22.6 | 57.9 | 11.6 | 14.7 | 42.2 | 14.4 | 37.5 | 1.9 | 13.8 | 3.2 | 88.8 |
| LPIS-C | 29.8 | 49.9 | 50.1 | 27.8 | 75.1 | 5.5 | 2.7 | 22.0 | 58.7 | 11.9 | 4.8 | 26.2 | 68.1 | 9.6 | 17.2 | 41.2 | 19.5 | 36.1 | 1.5 | 7.1 | 1.3 | 88.7 |
| LPIS-D | 36.1 | 51.3 | 59.9 | 40.9 | 77.5 | 7.0 | 9.3 | 50.2 | 77.8 | 28.7 | 10.1 | 24.3 | 79.6 | 7.4 | 20.2 | 42.3 | 13.0 | 36.7 | 15.8 | 12.8 | 4.0 | 88.7 |
| LPIS-E | 30.3 | 50.4 | 48.7 | 20.1 | 76.2 | 3.7 | 0.6 | 26.6 | 63.4 | 7.9 | 8.5 | 24.9 | 75.9 | 11.4 | 25.8 | 41.1 | 9.4 | 35.6 | 1.9 | 13.5 | 2.6 | 88.8 |
| LPIS-F | 23.8 | 43.1 | 59.6 | 48.8 | 68.8 | 2.6 | 0.0 | 28.0 | 70.9 | 12.4 | 20.9 | 22.8 | 1.5 | 0.0 | 10.1 | 24.5 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 86.3 |
| LPIS-G | 18.1 | 37.0 | 53.9 | 10.6 | 45.2 | 10.2 | 0.0 | 20.5 | 68.1 | 2.7 | 2.0 | 12.5 | 0.0 | 0.0 | 3.9 | 20.2 | 0.0 | 8.2 | 0.0 | 0.0 | 0.0 | 86.0 |
| LPIS-H | 23.8 | 44.6 | 57.5 | 48.7 | 61.7 | 6.7 | 0.0 | 42.7 | 69.6 | 5.1 | 15.4 | 16.7 | 0.0 | 0.0 | 7.9 | 21.4 | 0.0 | 16.6 | 0.0 | 0.0 | 0.0 | 86.2 |
| LPIS-I | 39.2 | 47.6 | 65.7 | 46.0 | 74.5 | 14.0 | 57.0 | 44.1 | 81.6 | 51.8 | 8.7 | 28.2 | 75.2 | 7.2 | 22.8 | 33.0 | 14.2 | 27.8 | 29.8 | 0.3 | 5.5 | 87.6 |
| LPIS-J | 35.4 | 52.0 | 57.4 | 31.0 | 78.3 | 8.2 | 10.6 | 45.8 | 71.9 | 33.7 | 8.9 | 27.2 | 75.3 | 14.4 | 22.1 | 44.6 | 16.4 | 36.6 | 6.6 | 12.1 | 2.3 | 88.7 |
| LPIS-K | 15.1 | 42.5 | 34.3 | 11.2 | 28.8 | 0.3 | 0.0 | 2.9 | 20.6 | 0.0 | 0.0 | 1.9 | 1.5 | 0.0 | 17.7 | 27.6 | 0.4 | 26.5 | 12.0 | 0.0 | 0.1 | 87.4 |

of the test set, grasses represent about 12%, and all other crop classes individually account for less than 2% (as it can be seen in Table VII). This imbalance significantly impacts learning and evaluation for most classes. Given the strong class imbalance in the dataset, OA is more reflective of the model’s training dynamics and performance on dominant classes, while mIoU offers complementary insight into class-wise behaviour. For the Level-1 nomenclature, the highest OA (88.0%) is achieved by both LPIS-D and LPIS-J, which include aerial VHR imagery combined with other modalities. The best mIoU (39.2%) is obtained by LPIS-I, which excludes aerial imagery and relies solely on SPOT and Sentinel-1/2 time series. Table XV illustrates the differing behaviour of crop type classes, which is closely correlated with their frequency at Level-1 in the dataset. The background class, which accounts for approximately 78% of the dataset, consistently achieves the highest IoU scores across all configurations (LPIS-A to LPIS-J), with a large margin exceeding 28% over the second-best class. Grasses, the second most frequent class (around 12%), also attain relatively high IoU values across all configurations, ranging from 37% to 52%. For the remaining classes, which are less frequent (between 2.5% and 0.5%), IoU scores exhibit greater variability across configurations. This variability appears to be somewhat correlated with the presence of specific modalities in the input configuration and therefore varying acquisition dates. However, for very rare classes, it becomes difficult to identify consistent patterns. A likely explanation is that the scarcity of training examples for these low-frequency classes results in high variability between training runs of the same configuration—variability that is comparable to that observed between different configurations thus limiting the ability to draw reliable conclusions.

Comparing Crop Type and Land Cover Mapping Tasks: While OA increases from LC-A to LPIS-A, the mIoU drops significantly. This decrease is expected, as land cover classes

are specifically designed to be distinguishable using mono-temporal aerial imagery, whereas many LPIS classes require multi-temporal information for correct identification. As a result, class confusion is more pronounced in the LPIS labelling task. The increase in OA may seem counterintuitive. A likely explanation is that the LPIS classes dominating the pixel distribution (Background, Grasses, and Vineyards) can still be reliably segmented from aerial imagery alone. These three classes represent over 88% of the test set pixels, thus inflating OA despite lower class-wise performance. As expected, multi-temporal modalities prove more beneficial for the LPIS task than for land cover segmentation. The addition of Sentinel-1 or Sentinel-2 time series significantly improves mIoU, with an 11.9% gain from LPIS-A to LPIS-D. Moreover, configurations using only S1 or S2 (e.g., LPIS-F) perform comparably to the aerial-only LPIS baseline (just a 0.6% drop in mIoU), whereas they perform substantially worse on the land cover task (showing a 29.9% drop in mIoU from LC-A to LC-G).

Limitations of Modality Contribution: Configurations that combine multi-temporal inputs (S1, S2, or both) with high-resolution imagery (aerial VHR or SPOT) generally achieve the best performance for the LPIS task. This likely results from the complementarity of modalities: high-resolution imagery provides fine-grained textural features, while time series data capture phenological dynamics critical for distinguishing between crop types. However, performance differences across specific modality combinations are not always intuitive. For instance, the Aerial+S2 configuration (LPIS-C) does not consistently outperform either of the single-modality baselines such as Aerial (LPIS-A) or S2 (LPIS-F). In some cases, the differences are substantial: for example, barley and rapeseed exhibit IoU drops of 21% and 12.2%, respectively, in LPIS-C compared to LPIS-F. A similar pattern is observed when comparing LPIS-I and LPIS-J: adding aerial imagery to a configuration using SPOT and Sentinel-1/2 time series (LPIS-

I) results in large IoU decreases for Flax/Hemp (−40.6%) and Soy (−18.1%). These findings suggest that simply adding modalities does not guarantee improved performance. The observed inconsistencies may stem from limitations in both the dataset and the model architecture. Potential contributing factors include increased model complexity leading to overfitting, interactions between the main and auxiliary losses when adding modalities, and limitations of the current temporal encoder. These factors indicate a need for more effective fusion strategies and potentially stronger architectures to fully leverage multimodal data.

Aerial / SPOT modalities: As expected, the mono-temporal modalities and particularly the aerial imagery yield higher IoU scores for classes associated with textural patterns visible in very high-resolution (VHR) imagery. Vineyards and orchards are notable examples, with very low IoU when using only S1 and/or S2 modalities, but a significant IoU increase when aerial data is included in the configuration. More surprisingly, a similar trend is observed for beetroot and potato classes, which is more difficult to explain. One hypothesis is that these crops may exhibit unique spatial or structural characteristics during the aerial acquisition period that are detectable at 20 cm resolution. Some classes, such as wheat and maize, also show relatively high IoU scores using only aerial imagery. For maize, this may be because many aerial acquisitions occur when the crop is nearly fully grown, resulting in distinct textural patterns at VHR scale (see Figure 1). For wheat, the explanation is less straightforward. However, given the low IoU scores for other cereals (*e.g.*, barley) and the relatively higher frequency of wheat in the dataset, it is plausible that the model learns generalized cereal textures and defaults to classifying all cereals (except maize) as wheat. To better understand why LPIS-I (which excludes aerial imagery) outperforms LPIS-J (which includes all modalities), we introduced an additional configuration, LPIS-K, using only the SPOT modality. Unfortunately, this experiment did not provide clear answers. In LPIS-K, background IoU sits between that of LPIS-A (Aerial) and LPIS-F (S2), which aligns with expectations based on the relative spatial resolutions (20 cm, 1.5 m, and 10 m, respectively). For most other classes, LPIS-K produces lower or comparable IoU values relative to LPIS-A, with larger drops for classes characterized by fine spatial textures (*e.g.*, maize, vineyard) than for those with broader spatial features (*e.g.*, orchards).

Temporal modalities: As anticipated, S2 appears to be one

of the most informative sources for the LPIS task. When added to Aerial VHR, it leads to substantial performance gains: LPIS-C outperforms LPIS-A by +5.4% mIoU, and LPIS-E outperforms LPIS-B by +4.2%. S2 is also part of the best performing configuration, LPIS-I. Overall, the addition of S2 improves the IoU for most cereal, oleaginous, and proteinous crop classes, with many classes showing gains between +5% and +15% compared to the same configuration without S2. In contrast, Sentinel-1 (S1) performs poorly as a stand-alone modality, which aligns with expectations given its sensor characteristics. Its contribution becomes more nuanced when combined with other modalities. For instance, LPIS-D (Aerial VHR + S2 + S1) significantly improves over LPIS-C (Aerial VHR + S2), yielding a +6.5% mIoU gain. However, adding S1 to S2 in LPIS-H provides no benefit over LPIS-F, and the improvement from LPIS-E to LPIS-J is modest (+2.0%). These results suggest that S1 contributes most effectively when paired with a high-resolution modality, such as Aerial VHR or SPOT.

Table XVI presents class-wise IoU scores for the best configuration (LPIS-I) on both the validation and test sets, confirming previous observations. A clear drop in performance is observed on the test set, with the overall mIoU decreasing from 81.4% to 39.2%. Several classes, such as fodder legumes, beetroots, and mixed crops, show substantial declines, while others like soy and olive groves are no longer detected at all. This discrepancy highlights the greater difficulty of the test set and the model’s limited ability to generalize to unseen regions. This performance gap is primarily due to the severely under-represented classes and data split strategy, where the validation set shares zones with the training data, while the test set covers distinct domains. Furthermore, an additional challenge of the FLAIR-HUB crop mapping task is that the different *domains* span three years, introducing both spatial and temporal generalization issues. Crop type classes exhibit greater variability across years compared to land cover classes, making the task particularly sensitive to temporal shifts.

In Table XVII, the performance of the LPIS-I configuration is reported across the five folds of the KFold cross-validation. OA remains consistently high across splits, ranging from 85.2% to 88.7%, while mean IoU (mIoU) varies more substantially, from 34.2% to 42.8%, reflecting sensitivity to

TABLE XVI: Per-Class Evaluation for Crop Mapping. Class-wise IoU scores for the Base-UP baseline (setting LPIS-I) on the Validation and Test partition. The *rice* and other *oilseed crops* classes are excluded from mIoU computation due to having zero pixels in the test set.

| Class Set → ↓ | mIoU | grasses | wheat | barley | maize | o. cereals | rice | flax/hemp tobacco | sunflower | rapeseed | soy | o. protein c. | fodder legumes | beetroots | potatoes | o. arable c. | vineyard | olive groves | fruits orchards | nut orchards | o. permanent c. | mixed c. | background |
|------------------|------|---------|-------|--------|-------|------------|------|----------------------|-----------|----------|------|---------------|-------------------|-----------|----------|--------------|----------|--------------|--------------------|--------------|-----------------|----------|------------|
| Validation | 81.4 | 43.7 | 77.3 | 62.4 | 78.2 | 16.9 | 60.1 | 40.9 | 83.7 | 89.8 | 40.6 | 57.6 | 31.2 | 79.5 | 50.3 | 14.8 | 48.3 | 11.5 | 47.5 | 6.5 | 0.1 | 2.4 | 89.7 |
| Test | 39.2 | 47.6 | 65.7 | 46.0 | 74.5 | 14.0 | - | 57.0 | 44.1 | 81.6 | 51.8 | 8.7 | 28.2 | 75.2 | 7.2 | 22.8 | 33.0 | 14.2 | 27.8 | 29.8 | 0.3 | 5.5 | 87.6 |

TABLE XVII: Per-Class Evaluation for Crop Type Mapping – KFold evaluation. Class-wise IoU scores for the Swin Base-UP baseline using SPOT imagery and Sentinel-1/2 time series (denoted setting LPIS-I). mIoU is computed for each fold by excluding classes that have no corresponding pixels in the test set.

| Model ID | CONF. | mIoU | O.A. | grasses | wheat | barley | maize | o. cereals | rice | flax/hemp tobacco | sunflower | rapeseed | soy | o. protein c. | fodder legumes | beetroots | potatoes | o. arable c. | vineyard | olive groves | fruits orchards | nut orchards | o. permanent c. | mixed c. | background | PARAM. | EP. |
|----------|---------|------|------|---------|-------|--------|-------|------------|------|-------------------|-----------|----------|------|---------------|----------------|-----------|----------|--------------|----------|--------------|-----------------|--------------|-----------------|----------|------------|--------|-----|
| LPIS-I | split_1 | 39.2 | 87.2 | 47.6 | 65.7 | 46.0 | 74.5 | 14.0 | - | 57.0 | 44.1 | 81.6 | 51.8 | 8.7 | 28.2 | 75.2 | 7.2 | 22.8 | 33.0 | 14.2 | 27.8 | 29.8 | 0.3 | 5.5 | 87.6 | 97.5 | 49 |
| LPIS-I | split_2 | 34.2 | 88.7 | 44.8 | 69.8 | 56.4 | 74.4 | 14.8 | 2.7 | 1.9 | 61.4 | 85.1 | 32.5 | 17.1 | 31.6 | 69.7 | 16.0 | 24.8 | 45.4 | 6.2 | 36.5 | 7.2 | 1.2 | 3.5 | 89.5 | 97.5 | 108 |
| LPIS-I | split_3 | 39.2 | 88.0 | 42.3 | 75.6 | 68.0 | 67.2 | 10.8 | 30.2 | 22.7 | 60.2 | 91.4 | 34.3 | 39.1 | 27.8 | 82.8 | 4.3 | 9.9 | 46.4 | 14.3 | 37.3 | 6.6 | 0.8 | 1.3 | 89.0 | 97.5 | 119 |
| LPIS-I | split_4 | 40.7 | 85.2 | 37.9 | 73.3 | 54.9 | 76.8 | 6.5 | - | 17.9 | 74.8 | 81.3 | 48.8 | 47.2 | 20.4 | 74.2 | 38.7 | 13.4 | 37.8 | 20.8 | 34.0 | 0.3 | 2.9 | 5.4 | 86.6 | 97.5 | 57 |
| LPIS-I | split_5 | 42.8 | 87.4 | 45.0 | 72.4 | 56.0 | 73.7 | 9.5 | - | 8.3 | 77.3 | 81.6 | 49.5 | 36.8 | 20.8 | 68.7 | 42.5 | 11.1 | 46.5 | 17.3 | 47.4 | - | 2.0 | 2.2 | 88.0 | 97.5 | 43 |

specific domain characteristics, supervision availability and acquisition dates regarding phenologies. It is important to note that mIoU is computed only over the classes present in the test set for each split, excluding those marked with dashes. This introduces some complexity when comparing mIoU values across folds, as the class composition can differ. Notably, major crop classes such as wheat, maize, and rapeseed show relatively stable and high IoU scores, whereas rare or under-represented classes (*e.g.*, nut orchards, tobacco, and other permanent crops) exhibit large variability or near-zero performance. These results highlight the strong impact of class imbalance and domain-specific variation on class-wise segmentation performance.

C. Multitask training

Table XVIII presents a quantitative evaluation of the UPerFuse architecture in a multitask setting, comparing its performance on land cover and crop type mapping when trained either separately or jointly. Both tasks has been assigned the same weight for the experiments. Overall, the results indicate that multitask learning does not yield performance improvements for either task. In fact, crop mapping performance slightly decreases in the multitask setting, with LPIS mIoU dropping from 39.2% (single-task) to 36.1%, while land cover results remain relatively stable (65.8% to 64.7% mIoU). This suggests that land cover segmentation has a more robust learning, likely due to more balanced class distributions, greater visual separability in the data and sufficient learning data. The degradation in LPIS performance aligns with earlier observations about the complexity of the crop mapping task: the strong class

imbalance, scarcity of rare crop types, and reliance on subtle temporal dynamics make it more sensitive to architectural or training changes. Additionally, the moderate increase in model parameters in the multitask setting does not appear sufficient to offset this trade-off. These findings highlight the need for more tailored multitask architectures, as well as improved strategies to handle data imbalance in crop type segmentation.

D. Qualitative results

Figure 6 presents patch-level inferences on the test split across various models for both the land-cover (top) and crop-type (bottom) classification tasks. For the land-cover task, the results indicate that using aerial imagery alone already yields highly accurate predictions. As such, models incorporating additional modalities that include this source exhibit only marginal improvements. Nonetheless, certain confusions—such as between tree types or between agricultural and herbaceous covers—are slightly mitigated when temporal information from complementary modalities is introduced. In contrast, the LPIS task remains significantly more challenging, as previously discussed. Quantitative performance is notably lower, and several rare classes are often not retrieved, with predictions defaulting to one of the dominant categories in the supervision dataset. Despite this limitation, some parcels are accurately classified. The LPIS-J model, which integrates aerial, SPOT, Sentinel-1, and Sentinel-2 imagery, appears to yield the most visually coherent and accurate results.

TABLE XVIII: Quantitative Evaluation in the Multi-task Setting. Performance of the UPerFuse architecture for both land cover and crop mapping tasks. The evaluation is performed using the best input modality configuration for each task. PARA.: number of model parameters (in millions). EP.: epoch with best validation score. SITS: Satellite Image Time Series. S1/2: Sentinel-1/2.

| Training | Model ID | Aerial VHR | Elevation | SPOT | S2 SITS | S1 SITS | PARA. | EP. | LC mIoU | LC O.A. | LPIS mIoU | LPIS O.A. |
|-------------------|----------|------------|-----------|------|---------|---------|-------|-----|---------|---------|-----------|-----------|
| Only Land Cover | LC-L | ✓ | ✓ | ✓ | ✓ | ✓ | 276.4 | 121 | 65.8 | 78.4 | X | X |
| Only Crop Mapping | LPIS-I | | | ✓ | ✓ | ✓ | 97.5 | 53 | X | X | 39.2 | 87.2 |
| Multi-task | LC-L | ✓ | ✓ | ✓ | ✓ | ✓ | 286.6 | 81 | 64.7 | 77.9 | 34.7 | 88.1 |
| Multi-task | LPIS-I | | | ✓ | ✓ | ✓ | 102.6 | 87 | 47.8 | 66.9 | 36.1 | 87.6 |

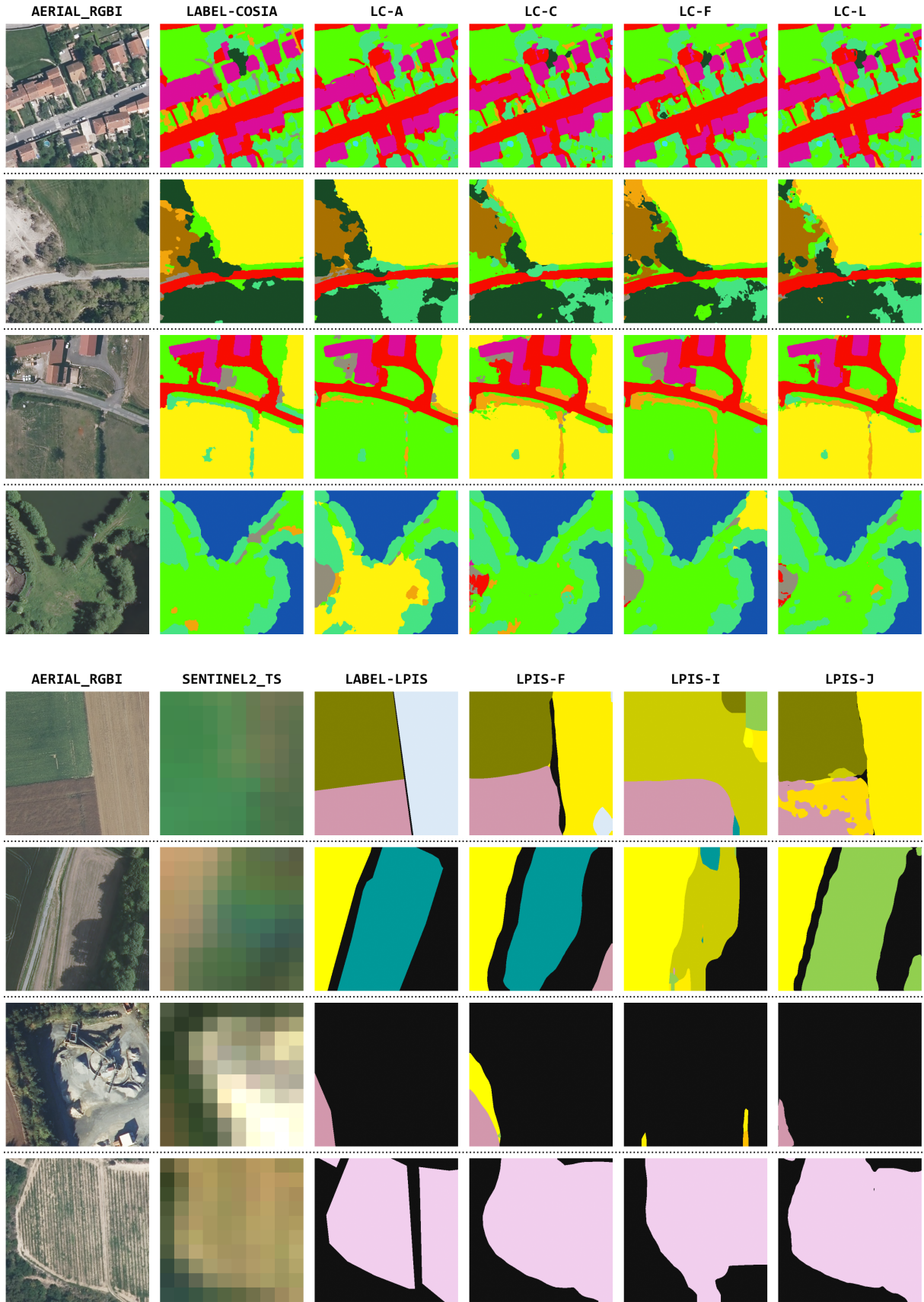


Fig. 6: Comparison of patch-level inference of LC and LPIS models. Top: LC results. Bottom: LPIS results. One can compare the ground truth (labels COSIA or LPIS) with the predictions of different unimodal or multimodal models. Details of the input data for each model can be found in the Tables [X](#) and [XIV](#).

To further evaluate model performance, predictions were extended to larger geographic regions. This qualitative assessment is crucial due to the spatial continuity inherent in geospatial data. The input ROI is divided into patches, and overlapping inferences are integrated to reduce edge effects. Results are shown over two ROI of the test-set in Figure 7 and Figure 8.

For the land-cover task, four models are visualized: aerial-

only (LC-A), Sentinel-2 (LC-G), SPOT-only (LC-I), and the best-performing configuration (LC-L). Region-level predictions are consistent with patch-level observations, affirming strong performance in the land-cover task. The Sentinel-2-only model (LC-G) delivers less precise results, as expected given its coarser spatial resolution of 10.24 m. Moreover, this model employs only the UTAE architecture and appears to generalize less effectively over broad areas compared to

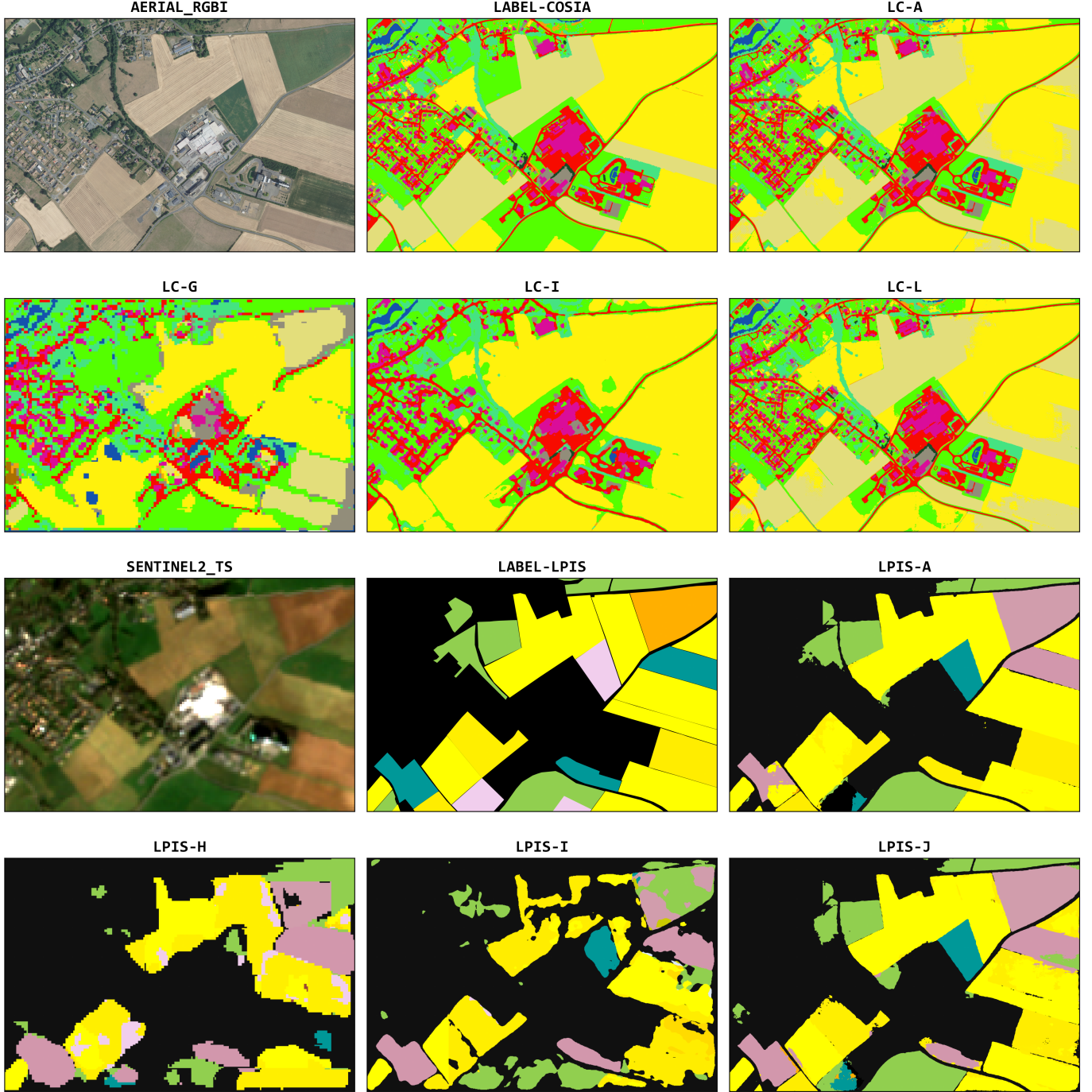


Fig. 7: ROI detections. We provide inference on large zones to illustrate the capacity of the different monomodal and multimodal models. Details of the input data for each model can be found in the Tables X and XIV.

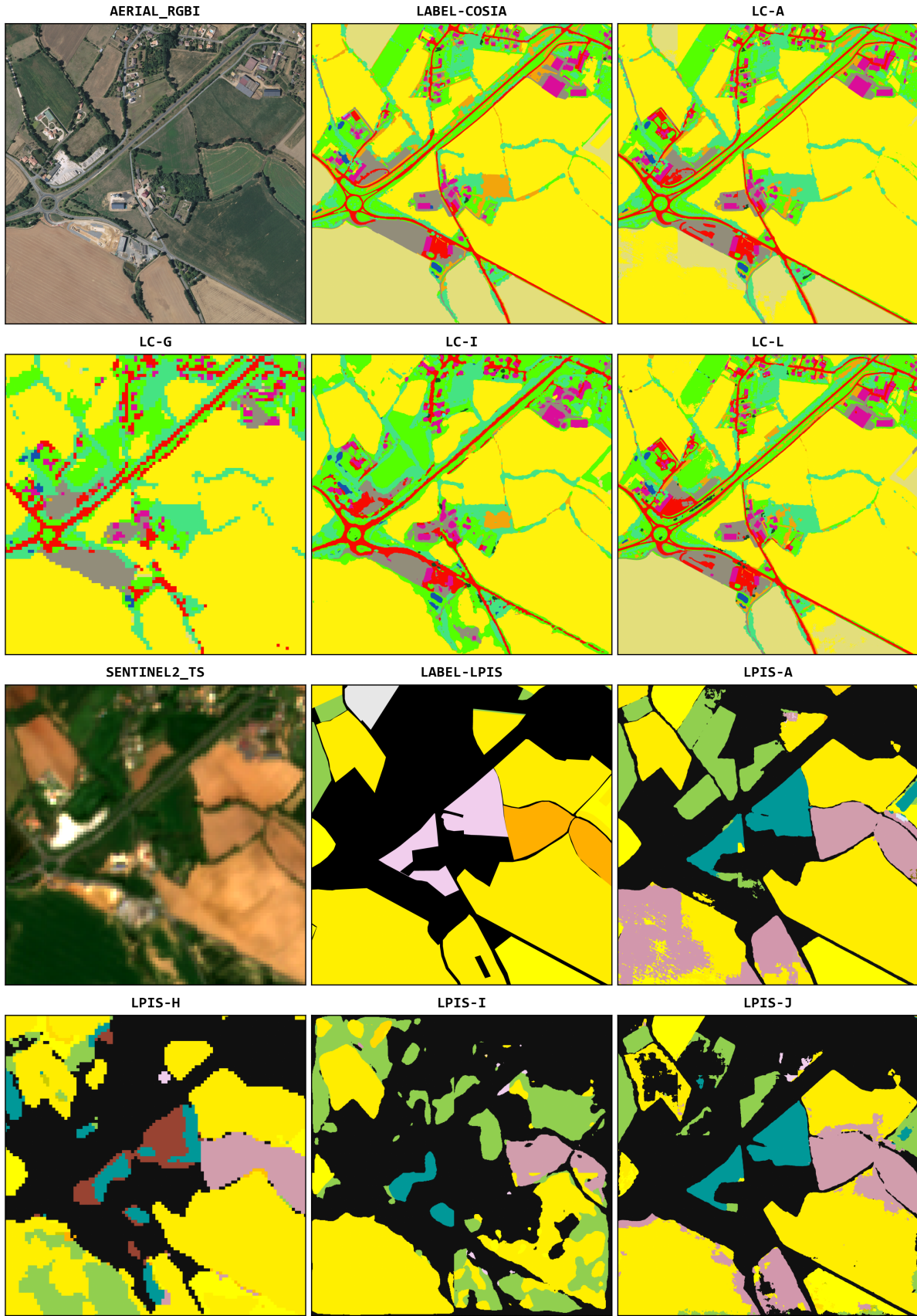


Fig. 8: ROI detections. We provide inference on large zones to illustrate the capacity of the different monomodal and multimodal models. Details of the input data for each model can be found in the Tables X and XIV.

attention-based alternatives. This limitation is particularly visible in agricultural zones, where class mixing is observed. In comparison, the SPOT-only model (LC-I), benefiting from a finer 1.6 m resolution, yields accurate predictions, though with less defined parcel boundaries than those seen in the 0.2 m aerial imagery (LC-A). It is worth noting that, due to differing acquisition dates relative to the aerial imagery used for ground-truth labelling, some classes—such as plowed land—are absent in SPOT-based predictions.

Results for the crop-type classification are markedly more variable. Models relying solely on aerial imagery (LPIS-A), Sentinel-1 and -2 (LPIS-H), SPOT combined with Sentinels (LPIS-I), and the best-performing configuration (LPIS-J) demonstrate differing behaviours. Notably, the aerial-only model performs well in delineating parcel boundaries and accurately classifies dominant crop types. In contrast, models using only Sentinel or SPOT data suffer from reduced spatial resolution, frequently omitting agricultural areas and exhibiting significant class confusions. However, adding very high-resolution aerial imagery to the SPOT and Sentinel modalities in the LPIS-J configuration substantially enhances visual coherence and parcel delineation, despite slightly lower overall quantitative scores compared to LPIS-I.

Overall, land-cover classification performs well due to a balanced dataset and clear annotations, while crop-type classification remains challenging. Temporal information is helpful but not fully exploited, and high spatial resolution appears key in both tasks.

VIII. Perspectives

A. Improving the deep learning model performances

The encoders used in the mono-temporal architectures of this work were pretrained on ImageNet. However, foundation models trained with self-supervised learning techniques are now widely recognized as more effective for pretraining and transfer [15, 16, 17, 18, 80, 82, 83]. Therefore, it would be relevant to assess the benefit of using off-the-shelf foundation models on the FLAIR-HUB dataset. Furthermore, given the large volume of FLAIR-HUB, training dedicated foundation models directly on this dataset appears feasible. This could pave the way for designing foundation models tailored to temporal or multi-modal tasks.

Several metadata provided with the FLAIR-HUB dataset were not used in the experiments presented. However, properly integrating this information into architectures could lead to performance gains [109, 110]. In particular, certain land cover classes, as well as all classes in the LPIS nomenclature, exhibit significant variations in appearance throughout the year, depending on the image sources. Encoding the acquisition date using the MTD_DATES metadata could therefore improve the temporal robustness of the models. Additionally, it could be beneficial to encode the geographical location, either through the MTD_GEOM metadata or by using information from the GeoTIFF file headers. Indeed, class appearances may vary

greatly across different French regions, and incorporating spatial context could help better model this geographic variability.

Both supervision nomenclatures exhibit significant variability in class frequency. This aspect was not addressed in the current study, although it is illustrated in Tables VI and VII. The classes were assigned binary weights of 1 or 0. Investigating weighted loss functions, exploring alternative loss functions [111], or applying different dataset sampling strategies [112] could lead to significant improvements. Hierarchical loss functions [113, 114] could also be considered for the LPIS task to take advantage of the nested structure of the three levels of the nomenclature.

In our experiments, we focus exclusively on a mid-level fusion strategy. However, other approaches, such as early fusion (input-level fusion) or late fusion (prediction-level fusion) could also be explored [78]. In addition, late fusion strategies may be particularly well suited to self-supervised pretraining performed independently for each modality.

In the selection of baselines for mono-temporal (Swin UPerNet) and multi-temporal (UTAE) image sources, we observe a significant imbalance in the number of parameters, to the disadvantage of multi-temporal architectures. For instance, as shown in Table X, UTAE models have around 1 million parameters, while mono-temporal setups can reach several hundred million. It would be valuable to conduct studies aimed at finding a better trade-off in parameter allocation between mono-temporal and multi-temporal branches.

B. Other Potential Uses of the FLAIR-HUB Dataset

The FLAIR-HUB dataset could also be of interest to researchers working on transfer learning and unsupervised domain adaptation. First, the pretrained weights can be used to fine-tune models on new visual categories or other types of sensors (*e.g.*, UAV, thermal imaging) or new tasks (*e.g.*, panoptic segmentation).

A thematically relevant application of transfer learning would be the development of AI models capable of generating land cover maps of the past [115]. To support this, we have included the AERIAL_195X modality in FLAIR-HUB. A promising direction would be to train models using annotations from recent aerial imagery and apply them to older aerial images. This type of experiment is also methodologically challenging due to the significant changes between the two modalities, including strong radiometric and domain shifts.

Furthermore, novel transfer methods could be evaluated based on their ability to train a model in a given spatial and/or temporal domain and transfer it to another in a supervised or unsupervised manner [116, 117]. FLAIR-HUB includes 74 distinct spatio-temporal domains along with the necessary metadata to support such experiments.

In the remote sensing community, many studies focus on super-resolution methods using multiple image sources. Since the FLAIR-HUB dataset provides image patches with spatially aligned modalities, it is particularly well suited for super-resolution methods using single or multiple images [68, 118, 119]. Specifically, models that aim to enhance the

spatial resolution of Sentinel-2 images using SPOT or aerial images could be effectively evaluated using FLAIR-HUB. The dataset is also tailored for the cloud removal task [94] on optical images thanks to SAR ones.

The FLAIR-HUB dataset’s Land Cover and LPIS annotations provide an interesting opportunity for remote sensing image synthesis [120, 121]. These annotations and metadata can serve as prompts for controlling generative models. These models could learn to capture the relationship link between land cover classes and the visual characteristics of the various sensors enabling the creation of augmented datasets. Such synthetic data could be used for tasks where annotations are scarce or expensive to obtain, like change detection [122].

C. Possible Future Extensions of the FLAIR-HUB Dataset

First, the FLAIR-HUB dataset is limited to metropolitan France. Although France’s territory is quite diverse, featuring oceanic, continental, Mediterranean, and mountainous bioclimatic regions, it does not contain tropical or desert areas. The aerial images were captured under favorable weather conditions between April and November, leading to a bias in the acquisition dates (see Figure 1). It could be interesting to expand the dataset to other countries or sensors (such as UAV) with a large variation of acquisition conditions (*e.g.*, angle, weather) but with an interoperable nomenclature to learn more generic models.

In the coming years, we plan to enhance FLAIR-HUB by adding new modalities, metadata, or tasks. For instance, we are waiting for the completion of the national coverage of France with high-definition LIDAR [123] (at a density of 10 to 20 points per square meter) to incorporate this point cloud modality as in [63, 66]. Hyperspectral images, such as those in [63, 70], could also be included. Additional auxiliary data such as weather, transportation, and socioeconomic indicators [124] could also be integrated into FLAIR-HUB for multimodal learning studies.

For approximately 200 ROIs in the FLAIR-HUB dataset, we are currently generating orthoimages from historical aerial imagery, spanning from 1960 to 2015. These images will be accompanied by temporally and spatially consistent land cover annotations for the semantic segmentation task. This new supervision dataset will be released upon completion. Historical aerial imagery presents significant variability in spatial resolution and radiometric characteristics.

Beyond its thematic interest such as enabling models to produce temporal series of LABEL-COSIA, this dataset also offers a valuable opportunity to study the generalization capabilities of AI models when faced with highly heterogeneous spatial and spectral inputs. Finally, we also have other types of labels that could be made available in FLAIR-HUB, including object detection tasks (*e.g.*, wind turbines, solar farms) or semantic segmentation tasks (*e.g.*, roads, hedgerows, isolated trees).

One of the most likely annotation extensions involves expanding the LPIS-labelled areas. As mentioned previously, a current limitation of the FLAIR-HUB dataset lies in the

significant imbalance of LPIS classes, which likely results in poor model training and high variability across training runs. Therefore, a priority for improving crop mapping performance lies more in enhancing the dataset than refining the model architecture, particularly by increasing the representation of low-frequency crop-type classes.

However, adding new ROIs or tasks raises several important questions. First, could the current best-performing land cover model be leveraged to assist in annotating new ROIs? Second, in the context of multitask learning, is it possible to integrate both land cover and crop-type annotations into a unified labelling scheme, and how would performance compare between this joint task and two separate tasks? We have initiated preliminary investigations to explore these directions, including strategies for selecting new ROIs to improve LPIS class balance (across all three annotation levels), and methods for merging land cover and crop-type labels into a comprehensive, multi-class label set. To reduce the need for manual annotation in new ROIs, we plan to evaluate semi-automated and soft-labelling approaches, using the current best land cover models. Once this improved dataset is available, it will facilitate the exploration of hierarchical classification strategies (particularly for the second and third levels of crop-type annotations) and support the creation of a unified benchmark dataset for evaluating multimodal fusion methods and architectures.

Beyond the crop mapping task, additional annotation types could further enrich the FLAIR-HUB dataset. For example, land cover predictions from FLAIR-HUB models are already employed by IGN to support the development of the OCSGE product [13]. This land use/land cover product features generalized geometries (with a minimum mapping unit of 200 m² for buildings and 500 m² for other classes) and results from conflation with existing databases such as LPIS. Initially, the decision was made not to include such generalized labels due to poor results in generalization. However, we now have access to OCSGE annotations for existing FLAIR-HUB ROIs, and they represent a valuable opportunity to assess model performance under varying degrees of label generalization.

In parallel, it may be worthwhile to introduce instance-level annotations for specific detection tasks (sometimes called panoptic segmentation). Notably, parcel boundary detection could be derived from raw LPIS vector data, while building detection could leverage existing resources such as the INRIA aerial labelling dataset [125] or the WHU building dataset [126]. Building detection is especially relevant, as current land cover labels prioritize the topmost visible cover (*e.g.*, trees over buildings), and do not preserve the high-quality building geometries found in dedicated building databases.

Finally, similarly to other work [127], we are also considering the generation of textual descriptions for ROIs to enable patch-level text annotations. This modality could serve to train CLIP-like models [128] and enhance few- and zero-shot capabilities, as demonstrated in recent multimodal frameworks [129].

IX. Conclusion

We presented FLAIR-HUB, the largest high-resolution multimodal dataset to date for land cover and crop type mapping. It contains over 63 billion annotated pixels across 2,528 km² of metropolitan France, with six spatially aligned modalities: aerial imagery, SPOT, Sentinel-1 and -2 time series, digital elevation models, and historical aerial photographs. These diverse sources capture a wide range of spatial, spectral, and temporal characteristics.

Through extensive benchmarks using state-of-the-art deep learning models, we highlighted both the challenges and opportunities of multimodal fusion. Our experiments show that combining complementary data sources significantly improves land cover and crop classification. At the same time, they underscore the difficulty of fine-grained crop mapping, multimodal integration, and multitask training in remote sensing.

FLAIR-HUB supports various learning settings, including supervised and self-supervised training, transfer learning, and domain adaptation.

By releasing this large-scale, extensively labelled dataset along with standardized benchmarks, we aim to support reproducible research and foster progress in the remote sensing, geospatial, and machine learning communities. FLAIR-HUB offers value for both methodological development and real-world applications.

Footprint of computations

The experiments presented in this article required computational resources equivalent to 27 311 hours on a single NVIDIA Tesla V100 GPU, producing 528.58 kg CO₂e. Based in France, this corresponds to a carbon footprint of 10.31 MWh, which is equivalent to 48.05 tree-years (calculated using green-algorithms.org v3.0 [130]).

Acknowledgment

The experiments conducted in this study were performed using HPC/AI resources provided by GENCI-IDRIS (Grant 2024-A0161013803, 2024-AD011014286R2 and 2025-A0181013803).

Data access

The dataset, pretrained models and codes are available at the following website: <https://ignf.github.io/FLAIR/flairhub>.

References

- [1] United Nations General Assembly. Transforming our world: the 2030 Agenda for Sustainable Development. *United Nations: New York, NY, USA*, 2015. 2
- [2] Ricardo Vinuesa, Hossein Azizpour, Iolanda Leite, Madeline Balaam, Virginia Dignum, Sami Domisch, Anna Felländer, Simone Daniela Langhans, Max Tegmark, and Francesco Fusco Nerini. The role of artificial intelligence in achieving the Sustainable Development Goals. *Nature Communications*, 11(1):233, 2020. 2
- [3] European Parliament. Regulation (EU) 2023/1115 of the European Parliament and of the Council of 31 May 2023 on the Making Available on the Union Market and the Export from the Union of Certain Commodities and Products Associated with Deforestation and Forest Degradation and Repealing Regulation (EU) No 995/2010 (Text with EEA Relevance). Technical report, European Parliament, 2023. Doc ID: 32023R1115. 2
- [4] R.F. Loi num. 2021-1104 du 22 août 2021 portant lutte contre le dérèglement climatique et renforcement de la résilience face à ses effets, 2021. [Online; accessed 16-May-2023]. 2
- [5] European Environment Agency. Soil degradation - Environment in EU at the turn of the century (Chapter 3.6). <https://www.eea.europa.eu/publications/92-9157-202-0/3.6.pdf>, 2020. [Online; accessed 4-June-2023]. 2
- [6] Food and FAO Agriculture Organization. Status of the World's Soil Resources: Main Report. <https://www.fao.org/3/i5199e/I5199E.pdf>, 2015. [Online; accessed 14-April-2023]. 2
- [7] Carlos García, Oscar Mora, Fernando Pérez-Aragüés, and Jordi Vitrià. Catlc: Catalonia multiresolution land cover dataset. *Scientific Data*, 9(1):554, 2022. 2, 3, 5
- [8] DLR. Data management system for large data volumes and AI applications for the BKG. 2
- [9] Institut national de l'information géographique et forestière. <https://www.ign.fr>, 2025. [Online; Accessed: 19-May-2025]. 2
- [10] Anatol Garioud, Stéphane Peillet, Eva Bookjans, Sébastien Giordano, and Boris Wattrélos. Flair #1: semantic segmentation and domain adaptation dataset. *arXiv*, 2022. 2, 7, 10, 14, 18, 21
- [11] Anatol Garioud, Apolline De Wit, Marc Poupée, Marion Valette, Sébastien Giordano, and Boris Wattrélos. Flair #2: textural and temporal information for semantic segmentation from multi-source optical imagery. *arXiv*, 2023. 2, 9, 10, 14, 19, 21
- [12] Anatol Garioud, Nicolas Gonthier, Loïc Landrieu, Apolline De Wit, Marion Valette, Marc Poupée, Sébastien Giordano, and Boris Wattrélos. FLAIR : a Country-Scale Land Cover Semantic Segmentation Dataset From Multi-Source Optical Imagery. *Advances in Neural Information Processing Systems*, 36:16456–16482, 2023. 2, 3, 5, 6, 9, 14, 17, 19, 21
- [13] IGN. OCS-GE. Un référentiel national utilisable aux différents échelons territoriaux pour la mise en place des politiques publiques d'aménagement du territoire et l'élaboration des documents d'urbanisme. <https://geoservices.ign.fr/ocsge>, 2023. [Online; accessed 10-December-2024]. 2, 29
- [14] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan

- Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision (CVPR)*, pages 10012–10022, 2021. 2, 14, 17
- [15] Guillaume Astruc, Nicolas Gonthier, Clement Mallet, and Loic Landrieu. Omnisat: Self-supervised modality fusion for earth observation. In *European Conference on Computer Vision*, pages 409–427. Springer, 2024. 2, 5, 6, 28
- [16] Guillaume Astruc, Nicolas Gonthier, Clement Mallet, and Loic Landrieu. Anysat: An earth observation model for any resolutions, scales, and modalities. In *Proceedings of the IEEE/CVF international conference on computer vision (CVPR)*, 2025. 2, 5, 28
- [17] Zhitong Xiong, Yi Wang, Fahong Zhang, Adam J Stewart, Joëlle Hanna, Damian Borth, Ioannis Papoutsis, Bertrand Le Saux, Gustau Camps-Valls, and Xiao Xiang Zhu. Neural plasticity-inspired foundation model for observing the earth crossing modalities. *arXiv e-prints*, pages arXiv–2403, 2024. 2, 28
- [18] Adam Stewart, Nils Lehmann, Isaac Corley, Yi Wang, Yi-Chia Chang, Nassim Ait Ali Braham, Shradha Sehgal, Caleb Robinson, and Arindam Banerjee. Ssl4eo-l: Datasets and foundation models for landsat imagery. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 5, 28
- [19] Vivien Sainte Fare Garnot, Loic Landrieu, and Nesrine Chehata. Multi-modal temporal attention models for crop mapping from satellite time series. *ISPRS Journal of Photogrammetry and Remote Sensing*, 187:294–305, 2022. 2
- [20] Nicolas Audebert, Bertrand Le Saux, and Sébastien Lefèvre. Semantic segmentation of earth observation data using multimodal and multi-scale deep networks. In *Asian conference on computer vision*, pages 180–196. Springer, 2016. 2
- [21] Franz Rottensteiner, Gunho Sohn, Markus Gerke, Jan Wegner, Uwe Breitkopf, and Jaewook Jung. Results of the ISPRS benchmark on urban object detection and 3D building reconstruction. *ISPRS Journal of Photogrammetry and Remote Sensing*, 2014. 2, 3
- [22] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2019. 3
- [23] Copernicus Land Monitoring Service. Urban Atlas. <https://land.copernicus.eu/local/urban-atlas>, 2018. [Online; accessed 01-June-2023]. 3
- [24] Romain Wenger, Anne Puissant, Jonathan Weber, Lhasane Idoumghar, and Germain Forestier. Multisenge: A multimodal and multitemporal benchmark dataset for land use/land cover remote sensing applications. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, 2022. 3
- [25] Hamed Alemohammad and Kevin Booth. LandCoverNet: A global benchmark land cover classification training dataset. In *Thirty-fifth Conference on Neural Information Processing Systems AI for Earth Sciences Workshop*, 2020. 3
- [26] Damien Sulla-Menashe and Mark A Friedl. User guide to collection 6 MODIS land cover (MCD12Q1 and MCD12C1) product. *United States Geological Survey*, 2018. 3
- [27] Javiera Castillo-Navarro, Bertrand Le Saux, Alexandre Boulch, Nicolas Audebert, and Sébastien Lefèvre. Semi-supervised semantic segmentation in earth observation: The MiniFrance suite, dataset analysis and multi-task network study. *Machine Learning*, 2022. 3
- [28] Aysim Toker, Lukas Kondmann, Mark Weber, Marvin Eisenberger, Andrés Camero, Jingliang Hu, Ariadna Pregel Hoderlein, Çağlar Şenaras, Timothy Davis, Daniel Cremers, Giovanni Marchisio, Xiao Zhu, and Laura Leal-Taixé. DynamicEarthNet: Daily multi-spectral satellite dataset for semantic change segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 3
- [29] Wuttichai Boonpook, Yumin Tan, Attawut Nardkulpat, Kritanai Torsri, Peerapong Torteeka, Patcharin Kamsing, Utane Sawangwit, Jose Pena, and Montri Jainan. Deep learning semantic segmentation for land use and land cover types using landsat 8 imagery. *ISPRS International Journal of Geo-Information*, 12(1), 2023. 3
- [30] Junshi Xia, Naoto Yokoya, Bruno Adriano, and Clifford Broni-Bediako. Openearthmap: A benchmark dataset for global high-resolution land cover mapping. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023. 3
- [31] Xin-Yi Tong, Gui-Song Xia, and Xiao Xiang Zhu. Enabling country-scale land cover mapping with meter-resolution satellite imagery. *ISPRS Journal of Photogrammetry and Remote Sensing*, 2023. 3
- [32] Junjue Wang, Zhuo Zheng, Ailong Ma, Xiaoyan Lu, and Yanfei Zhong. LoveDA: A remote sensing land-cover dataset for domain adaptive semantic segmentation. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. 3
- [33] Christopher F Brown, Steven P Brumby, Brookie Guzder-Williams, Tanya Birch, Samantha Brooks Hyde, Joseph Mazzariello, Wanda Czerwinski, Valerie J Pasquarella, Robert Haertel, Simon Ilyushchenko, et al. Dynamic world, near real-time global 10 m land use land cover mapping. *Scientific Data*, 9(1):251, 2022. 3
- [34] Ilke Demir, Krzysztof Koperski, David Lindenbaum, Guan Pang, Jing Huang, Saikat Basu, Forest Hughes, Devis Tuia, and Ramesh Raskar. DeepGlobe 2018: A challenge to parse the earth through satellite images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018. 3
- [35] Gencer Sumbul, Arne De Wall, Tristan Kreuziger, Filipe

- Marcelino, Hugo Costa, Pedro Benevides, Mario Caetano, Begüm Demir, and Volker Markl. BigEarthNet-MM: A large-scale, multimodal, multilabel benchmark archive for remote sensing image classification and retrieval. *IEEE Geoscience and Remote Sensing Magazine*, 2021. 3
- [36] Corine Land Cover. Corine land cover. *European Environment Agency, Copenhagen*, 2000. 3
- [37] Dominik Koßmann, Viktor Brack, and Thorsten Wilhelm. Seasonet: A seasonal scene classification, segmentation and retrieval dataset for satellite imagery over germany. In *IGARSS 2022-2022 IEEE International Geoscience and Remote Sensing Symposium*, pages 243–246. IEEE, 2022. 3
- [38] Yujian Mo, Yan Wu, Xinneng Yang, Feilin Liu, and Yujun Liao. Review the state-of-the-art technologies of semantic segmentation based on deep learning. *Neurocomputing*, 493:626–646, 2022. 2
- [39] Xiaohui Yuan, Jianfang Shi, and Lichuan Gu. A review of deep learning methods for semantic segmentation of remote sensing imagery. *Expert Systems with Applications*, 169:114417, 2021. 2
- [40] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 2
- [41] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, pages 234–241. Springer International Publishing, 2015. 2, 17
- [42] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018. 2
- [43] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 2
- [44] M. Schmitt, L. H. Hughes, C. Qiu, and X. X. Zhu. SEN12MS – a curated dataset of georeferenced multi-spectral Sentinel-1/2 imagery for deep learning and data fusion. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, 2019. 3
- [45] Yindan Zhang, Gang Chen, Soe W Myint, Yuyu Zhou, Geoffrey J Hay, Jelena Vukomanovic, and Ross K Meentemeyer. Urbanwatch: A 1-meter resolution land cover and land use database for 22 major cities in the united states. *Remote Sensing of Environment*, 278:113106, 2022. 3
- [46] C. Mallet and A. Le Bris. Current Challenges in Operational Very High Resolution Land-cover Mapping. *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, XLIII-B2-2020:703–710, 2020. 3
- [47] Marc Rußwurm and Marco Körner. Multi-temporal land cover classification with sequential recurrent encoders. *ISPRS International Journal of Geo-Information*, 7(4):129, 2018. 4
- [48] Rose M Rustowicz, Robin Cheong, Lijing Wang, Stefano Ermon, Marshall Burke, and David Lobell. Semantic segmentation of crop type in africa: A novel dataset and analysis of deep learning methods. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition workshops*, pages 75–82, 2019. 4
- [49] Hannah Kerner, Catherine Nakalembe, and Inbal Becker-Reshef. Field-level crop type classification with k nearest neighbors: A baseline for a new kenya small-holder dataset. *arXiv preprint arXiv:2004.03023*, 2020. 3, 4
- [50] Mehmet Ozgur Turkoglu, Stefano D’Aronco, Gregor Perich, Frank Liebisch, Constantin Streit, Konrad Schindler, and Jan Dirk Wegner. Crop mapping from image time series: Deep learning with multi-scale label hierarchies. *Remote Sensing of Environment*, 264:112603, 2021. 4
- [51] Vivien Sainte Fare Garnot and Loic Landrieu. Panoptic segmentation of satellite image time series with convolutional temporal attention networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4872–4881, 2021. 4, 5
- [52] Dimitris Sykas, Ioannis Papoutsis, and Dimitrios Zografakis. Sen4agrinet: A harmonized multi-country, multi-temporal benchmark dataset for agricultural earth observation machine learning applications. In *2021 IEEE International Geoscience and Remote Sensing Symposium IGARSS*, pages 5830–5833. IEEE, 2021. 4, 5
- [53] Lukas Kondmann, Aysim Toker, Marc Rußwurm, Andres Camero Unzueta, Devis Peressuti, Grega Milcinski, Nicolas Longépé, Pierre-Philippe Mathieu, Timothy Davis, Giovanni Marchisio, et al. Denethor: The dynamicearthnet dataset for harmonized, inter-operable, analysis-ready, daily crop monitoring from space. In *35th Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, pages 1–13, 2021. 4
- [54] Teodora Selea. Agrisen-cog, a multicountry, multi-temporal large-scale sentinel-2 benchmark dataset for crop mapping using deep learning. *Remote Sensing*, 15(12):2980, 2023. 4, 5
- [55] Marc Rußwurm, Charlotte Pelletier, Maximilian Zollner, Sébastien Lefèvre, and Marco Körner. Breizhcrops: A time series dataset for crop type mapping. *The International Archives of the Photogrammetry, Remote*

- Sensing and Spatial Information Sciences*, XLIII-B2-2020:1545–1551, 2020. 3, 4
- [56] Giulio Weikmann, Claudia Paris, and Lorenzo Bruzzone. Timesen2crop: A million labeled samples dataset of sentinel 2 image time series for crop-type classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 14:4699–4708, 2021. 3, 4
 - [57] Valentin Barriere, Martin Claverie, Maja Schneider, Guido Lemoine, and Raphaël d’Andrimont. Boosting crop classification by hierarchically fusing satellite, rotational, and contextual data. *Remote Sensing of Environment*, 305:114110, 2024. 4, 5
 - [58] Surbhi Sharma, Rocco Sedona, Morris Riedel, Gabriele Cavallaro, and Claudia Paris. Sen4map: Advancing mapping with sentinel-2 by providing detailed semantic descriptions and customizable land-use and land-cover data. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2024. 4, 5
 - [59] Joana Reuss, Jan Macdonald, Simon Becker, Lorenz Richter, and Marco Körner. Eurocropsml: A time series benchmark dataset for few-shot crop type classification. *arXiv preprint arXiv:2407.17458*, 2024. 3, 4, 5
 - [60] Claudio Persello, Jeroen Grift, Xinyan Fan, Claudia Paris, Ronny Hänsch, Mila Koeva, and Andrew Nelson. Ai4smallfarms: A dataset for crop field delineation in southeast asian smallholder farms. *IEEE Geoscience and Remote Sensing Letters*, 20:1–5, 2023. 3, 4
 - [61] Raphaël d’Andrimont, Martin Claverie, Pieter Kempeneers, Davide Muraro, Momchil Yordanov, Devis Peressutti, Matej Batič, and François Waldner. Ai4boundaries: an open ai-ready dataset to map field boundaries with sentinel-2 and aerial photography. *Earth System Science Data*, 15(1):317–329, 2023. 3, 4
 - [62] Hannah Kerner, Snehal Chaudhari, Aninda Ghosh, Caleb Robinson, Adeel Ahmad, Eddie Choi, Nathan Jacobs, Chris Holmes, Matthias Mohr, Rahul Dodhia, et al. Fields of the world: A machine learning benchmark dataset for global agricultural field boundary segmentation. *arXiv preprint arXiv:2409.16252*, 2024. 3, 4
 - [63] Yonghao Xu, Bo Du, Liangpei Zhang, Daniele Cerra, Miguel Pato, Emiliano Carmona, Saurabh Prasad, Naoto Yokoya, Ronny Hänsch, and Bertrand Le Saux. Advanced multi-sensor optical remote sensing for urban land use and land cover classification: Outcome of the 2018 ieee grss data fusion contest. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(6):1709–1724, 2019. 5, 29
 - [64] S. Ahlswede, C. Schulz, C. Gava, P. Helber, B. Bischke, M. Förster, F. Arias, J. Hees, B. Demir, and B. Kleinschmit. *TreeSatAI Benchmark Archive*: a multi-sensor, multi-label dataset for tree species classification in remote sensing. *Earth System Science Data*, 15(2):681–695, 2023. 5
 - [65] Nikolaos Ioannis Bountos, Arthur Ouaknine, and David Rolnick. Fomo-bench: a multi-modal, multi-scale and multi-task forest monitoring benchmark for remote sensing foundation models. In *39th Annual AAAI Conference on Artificial Intelligence, AI for Social Impact track*, 2025. 5
 - [66] Ben G. Weinstein, Sarah J. Graves, Sergio Marconi, Aditya Singh, Alina Zare, Dylan Stewart, Stephanie A. Bohlman, and Ethan P. White. A benchmark dataset for canopy crown detection and delineation in co-registered airborne rgb, lidar and hyperspectral imagery from the national ecological observation network. *PLOS Computational Biology*, 17(7):1–18, 07 2021. 5, 6, 29
 - [67] Allen Institute for Artificial Intelligence (AI2). AI2-S2-NAIP. <https://huggingface.co/datasets/allenai/s2-naip>, 2024. [Online; accessed 01-Sept-2024]. 5, 6
 - [68] Piper Wolters, Favyen Bastani, and Aniruddha Kembhavi. Zooming out on zooming in: Advancing super-resolution for remote sensing, 2023. 5, 28
 - [69] Favyen Bastani, Piper Wolters, Ritwik Gupta, Joe Ferdinando, and Aniruddha Kembhavi. Satlaspretrain: A large-scale dataset for remote sensing image understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16772–16782, 2023. 5, 6
 - [70] Jingliang Hu, Rong Liu, Danfeng Hong, Andrés Camero, Jing Yao, Mathias Schneider, Franz Kurz, Karl Segl, and Xiao Xiang Zhu. Mdas: A new multimodal benchmark dataset for remote sensing. *Earth System Science Data*, 15(1):113–131, 2023. 5, 6, 29
 - [71] Luis Miguel Pazos-Outón, Cristina Nader Vasconcelos, Anton Raichuk, Anurag Arnab, Dan Morris, and Maxim Neumann. Planted: a dataset for planted forest identification from multi-satellite time series. In *IGARSS 2024-2024 IEEE International Geoscience and Remote Sensing Symposium*, pages 7066–7070. IEEE, 2024. 5, 6
 - [72] Raphaël d’Andrimont, Momchil Yordanov, Laura Martinez-Sanchez, Beatrice Eiselt, Alessandra Palmieri, Paolo Dominici, Javier Gallego, Hannes Isaak Reuter, Christian Joebges, Guido Lemoine, et al. Harmonised lucas in-situ land cover and use database for field surveys from 2006 to 2018 in the european union. *Scientific data*, 7(1):352, 2020. 3, 5
 - [73] Ribana Roscher, Marc Russwurm, Caroline Gevaert, Michael Kampffmeyer, Jefersson A. Dos Santos, Maria Vakalopoulou, Ronny Hänsch, Stine Hansen, Keiller Nogueira, Jonathan Prexl, and Devis Tuia. Better, not just more: Data-centric machine learning for Earth observation. *IEEE Geoscience and Remote Sensing Magazine*, 2024. 4, 5
 - [74] Maja Schneider, Tobias Schelte, Felix Schmitz, and Marco Körner. Eurocrops: The largest harmonized open crop dataset across the european union. *Scientific Data*, 10(1):612, 2023. 4, 11
 - [75] Stephan Arnold, Geoffrey Smith, Gerard Hazeu, Bar-

- bara Kosztra, Christoph Perger, Gebhard Banko, Tomas Soukup, Geir-Harald Strand, Nuria Valcarcel Sanz, and Michael Bock. The eagle concept: A paradigm shift in land monitoring. In *Land use and land cover semantics*, pages 107–144. CRC Press, 2018. 4
- [76] ICC indication crop classification. https://www.fao.org/fileadmin/templates/ess/documents/world_census_of_agriculture/appendix3_r7.pdf. Accessed: 2025-03-31. 4
- [77] Tong Wang, Guanzhou Chen, Xiaodong Zhang, Chenxi Liu, Xiaoliang Tan, Jiaqi Wang, Chanjuan He, and Wenlin Zhou. Lmfnet: An efficient multimodal fusion approach for semantic segmentation in high-resolution remote sensing. *arXiv preprint arXiv:2404.13659*, 2024. 5
- [78] Jiaxin Li, Danfeng Hong, Lianru Gao, Jing Yao, Ke Zheng, Bing Zhang, and Jocelyn Chanussot. Deep learning in multimodal remote sensing data fusion: A comprehensive review. *International Journal of Applied Earth Observation and Geoinformation*, 112:102926, 2022. 5, 28
- [79] Caleb Robinson, Kolya Malkin, Nebojsa Jojic, Huijun Chen, Rongjun Qin, Changlin Xiao, Michael Schmitt, Pedram Ghamisi, Ronny Hänsch, and Naoto Yokoya. Global land-cover mapping with weak supervision: Outcome of the 2020 ieeegrss data fusion contest. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 14:3185–3199, 2021. 5
- [80] Johannes Jakubik, Sujit Roy, CE Phillips, Paolo Fraccaro, Denys Godwin, Bianca Zadrozny, Daniela Szwarcman, Carlos Gomes, Gabby Nyirjesy, Blair Edwards, et al. Foundation models for generalist geospatial artificial intelligence. *arXiv preprint arXiv:2310.18660*, 2023. 5, 28
- [81] Guillaume Astruc, Nicolas Gonthier, Clement Mallet, and Loic Landrieu. Omnisat: Self-supervised modality fusion for earth observation. In *European Conference on Computer Vision (ECCV)*, pages 409–427, 2024. 5
- [82] Anthony Fuller, Koreen Millard, and James Green. Croma: Remote sensing representations with contrastive radar-optical masked autoencoders. *Advances in Neural Information Processing Systems*, 36:5506–5538, 2023. 5, 28
- [83] Gabriel Tseng, Ruben Cartuyvels, Ivan Zvonkov, Mirali Purohit, David Rolnick, and Hannah Kerner. Lightweight, pre-trained transformers for remote sensing timeseries. *arXiv preprint arXiv:2304.14065*, 2023. 5, 28
- [84] Alistair Francis and Mikolaj Czerkawski. Major tom: Expandable datasets for earth observation. In *IGARSS 2024-2024 IEEE International Geoscience and Remote Sensing Symposium*, pages 2935–2940. IEEE, 2024. 5
- [85] Giovanni Marchisio, Patrick Helber, Benjamin Bischke, Timothy Davis, Caglar Senaras, Daniele Zanaga, Ruben Van De Kerchove, and Annett Wania. Rapidai4eo: A corpus for higher spatial and temporal reasoning. In *2021 IEEE International Geoscience and Remote Sensing Symposium IGARSS*, pages 1161–1164. IEEE, 2021. 5
- [86] Utkarsh Mall, Cheng Perng Phoo, Meilin Kelsey Liu, Carl Vondrick, Bharath Hariharan, and Kavita Bala. Remote sensing vision-language foundation models without annotations via ground remote alignment. In *The Twelfth International Conference on Learning Representations*, 2024. 5
- [87] Diego Velazquez, Pau Rodriguez, Sergio Alonso, Josep M Gonfaus, Jordi Gonzalez, Gerardo Richarte, Javier Marin, Yoshua Bengio, and Alexandre Lacoste. Earthview: A large scale remote sensing dataset for self-supervision. In *Proceedings of the Winter Conference on Applications of Computer Vision*, pages 1228–1237, 2025. 5
- [88] Yi Wang, Nassim Ait Ali Braham, Zhitong Xiong, Chenying Liu, Conrad M Albrecht, and Xiao Xiang Zhu. Ssl4eo-s12: A large-scale multimodal, multi-temporal dataset for self-supervised learning in earth observation [software and data sets]. *IEEE Geoscience and Remote Sensing Magazine*, 11(3):98–106, 2023. 5
- [89] Xuyang Li, Danfeng Hong, Chenyu Li, and Jocelyn Chanussot. Seamo: A multi-seasonal and multimodal remote sensing foundation model. *arXiv preprint arXiv:2412.19237*, 2024. 5
- [90] Xin Guo, Jiangwei Lao, Bo Dang, Yingying Zhang, Lei Yu, Lixiang Ru, Liheng Zhong, Ziyuan Huang, Kang Wu, Dingxiang Hu, et al. Skysense: A multimodal remote sensing foundation model towards universal interpretation for earth observation imagery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27672–27683, 2024. 5
- [91] Matías Mendieta, Boran Han, Xingjian Shi, Yi Zhu, and Chen Chen. Towards geospatial foundation models via continual pretraining. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16806–16816, 2023. 5
- [92] Junshi Xia, Hongruixuan Chen, Clifford Broni-Bediako, Yimin Wei, Jian Song, and Naoto Yokoya. Openeearthmap-sar: A benchmark synthetic aperture radar dataset for global high-resolution land cover mapping. *arXiv preprint arXiv:2501.10891*, 2025. 5
- [93] Ghjulia Sialelli, Torben Peters, Jan D Wegner, and Konrad Schindler. Agbd: A global-scale biomass dataset. *arXiv preprint arXiv:2406.04928*, 2024. 5
- [94] Patrick Ebel, Vivien Sainte Fare Garnot, Michael Schmitt, Jan Wegner, and Xiao Xiang Zhu. UnCR-tainTS: Uncertainty Quantification for Cloud Removal in Optical Satellite Time Series. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. IEEE, 2023. 5, 29
- [95] Laurent Guigues, Jean Pierre Cocquerez, and Hervé Le Men. Scale-sets Image Analysis. *International Journal of Computer Vision*, 68(3):289–317, 2006. 10

- [96] Vivien Sainte Fare Garnot and Loic Landrieu. Panoptic segmentation of satellite image time series with convolutional temporal attention networks. *ICCV*, 2021. 14, 16
- [97] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding, 2018. 15, 17
- [98] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015. 16
- [99] PyTorch Lightning. <https://www.lightning.ai>, 2025. [Online; Accessed: 10-February-2025]. 16
- [100] Pavel Iakubovskii. Segmentation models Pytorch. https://github.com/qubvel/segmentation_models.pytorch, 2025. 16
- [101] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2025. 16
- [102] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *ICLR*, 2019. 16
- [103] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016. 17
- [104] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. 17
- [105] Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and Saining Xie. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16133–16142, 2023. 17
- [106] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3349–3364, 2021. 17
- [107] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in neural information processing systems*, 34:12077–12090, 2021. 17
- [108] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017. 17
- [109] Johannes Jakubik, Felix Yang, Benedikt Blumenstiel, Erik Scheurer, Rocco Sedona, Stefano Maurogiovanni, Jente Bosmans, Nikolaos Dionelis, Valerio Marsocci, Niklas Kopp, et al. Terramind: Large-scale generative multimodality for earth observation. *arXiv preprint arXiv:2504.11171*, 2025. 28
- [110] Valerio Marsocci, Nicolas Gonthier, Anatol Garioud, Simone Scardapane, and Clément Mallet. Geomulti-tasknet: remote sensing unsupervised domain adaptation using geographical coordinates. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2075–2085, 2023. 28
- [111] Reza Azad, Moein Heidary, Kadir Yilmaz, Michael Hüttemann, Sanaz Karimijafarbigloo, Yuli Wu, Anke Schmeink, and Dorit Merhof. Loss functions in the era of semantic segmentation: A survey and outlook. *arXiv preprint arXiv:2312.05391*, 2023. 28
- [112] Keiller Nogueira, Mayara Maezano Faita-Pinheiro, Ana Paula Marques Ramos, Wesley Nunes Gonçalves, José Marcato Junior, and Jefersson A Dos Santos. Prototypical contrastive network for imbalanced aerial image segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 8366–8376, 2024. 28
- [113] Liulei Li, Tianfei Zhou, Wenguan Wang, Jianwu Li, and Yi Yang. Deep hierarchical semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1246–1257, 2022. 28
- [114] Loic Landrieu and Vivien Sainte Fare Garnot. Leveraging class hierarchies with metric-guided prototype learning. In *British Machine Vision Conference (BMVC)*, 2021. 28
- [115] Arnaud Le Bris, Sébastien Giordano, and Clément Mallet. Cnn semantic segmentation to retrieve past land cover out of historical orthoimages and dsm: first experiments. In *ISPRS Annals of Photogrammetry, Remote Sensing and Spatial Information Sciences*, volume 2, pages pp–1013, 2020. 28
- [116] Damian Ibañez, Junshi Xia, Naoto Yokoya, Filiberto Pla, and Ruben Fernandez-Beltran. Inter-sensor high-resolution and multi-temporal image fusion for unsupervised domain adaptation in remote sensing. *IEEE Transactions on Geoscience and Remote Sensing*, 2025. 28
- [117] Benjamin Lucas, Charlotte Pelletier, Daniel Schmidt, Geoffrey I Webb, and François Petitjean. A bayesian-inspired, deep learning-based, semi-supervised domain adaptation technique for land cover mapping. *Machine Learning*, 112(6):1941–1973, 2023. 28
- [118] Ngoc Long Nguyen, Jérémy Anger, Axel Davy, Pablo Arias, and Gabriele Facciolo. Self-supervised multi-image super-resolution for push-frame satellite images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1121–1131, 2021. 28
- [119] Pawel Kowaleczko, Tomasz Tarasiewicz, Maciej Ziaja, Daniel Kostrzewa, Jakub Nalepa, Przemysław Rokita, and Michał Kawulok. A real-world benchmark for sentinel-2 multi-image super-resolution. *Scientific Data*,

10(1):644, 2023. 28

- [120] Samar Khanna, Patrick Liu, Linqi Zhou, Chenlin Meng, Robin Rombach, Marshall Burke, David B. Lobell, and Stefano Ermon. Diffusionsat: A generative foundation model for satellite imagery. In *The Twelfth International Conference on Learning Representations*, 2024. 29
- [121] Ankan Dash, Junyi Ye, and Guiling Wang. A review of generative adversarial networks (gans) and its applications in a wide variety of disciplines: from medical to remote sensing. *IEEE Access*, 12:18330–18357, 2023. 29
- [122] Yanis Benidir, Nicolas Gonthier, and Clément Mallet. The change you want to detect: Semantic change detection in earth observation with hybrid data generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. 29
- [123] IGN. LiDAR HD Une cartographie 3D du sol et du sursol de la France. <https://geoservices.ign.fr/lidarhd>, 2020. [Online; accessed 1-May-2025]. 29
- [124] Xiaoqin Yan, Zhangwei Jiang, Peng Luo, Hao Wu, Anning Dong, Fengling Mao, Ziyin Wang, Hong Liu, and Yao Yao. A multimodal data fusion model for accurate and interpretable urban land use mapping with uncertainty analysis. *International Journal of Applied Earth Observation and Geoinformation*, 129:103805, 2024. 29
- [125] Emmanuel Maggiori, Yuliya Tarabalka, Guillaume Charpiat, and Pierre Alliez. Can semantic labeling methods generalize to any city? the inria aerial image labeling benchmark. In *2017 IEEE International geoscience and remote sensing symposium (IGARSS)*, pages 3226–3229. IEEE, 2017. 29
- [126] Shunping Ji, Shiqing Wei, and Meng Lu. Fully convolutional networks for multisource building extraction from an open aerial and satellite imagery data set. *IEEE Transactions on geoscience and remote sensing*, 57(1):574–586, 2018. 29
- [127] Xiaoqiang Lu, Binqiang Wang, Xiangtao Zheng, and Xuelong Li. Exploring models and data for remote sensing image caption generation. *IEEE Transactions on Geoscience and Remote Sensing*, 56(4):2183–2195. 29
- [128] João Daniel Silva, João Magalhães, Devis Tuia, and Bruno Martins. Multilingual vision-language pre-training for the remote sensing domain. In *Proceedings of the 32nd ACM International Conference on Advances in Geographic Information Systems*, pages 220–232, 2024. 29
- [129] Jeremy Andrew Irvin, Emily Ruoyu Liu, Joyce Chuyi Chen, Ines Dormoy, Jinyoung Kim, Samar Khanna, Zhuo Zheng, and Stefano Ermon. TeoChat: A large vision-language assistant for temporal earth observation data. *ICLR*, 2025. 29
- [130] Loïc Lannelongue, Jason Grealey, and Michael Inouye. Green algorithms: Quantifying the carbon footprint of computation, 2020. 30