

# The YAML: an overview

João Lucas de Sousa Almeida, Daniela Szwarcman

April 13, 2025

# The YAML format

If you are using the command-line interface (CLI) to run jobs using TerraTorch, so you must become familiar with YAML, the format used to configure all the workflow within the toolkit.

- Writing a YAML file is very similar to coding, because even if you are not directly handling the classes and others structures defined inside a codebase, you need to know how they work, their input arguments and their position in the pipeline.
- In this way, we could call it a "low-code" task.
- The YAML file used for TerraTorch has an almost closed format, since there are a few fixed fields that must be filled with limited sets of classes, which makes easier for new users to get a pre-existing YAML file and adapt it to their own purposes.
- In the next sections, we describe each field of a YAML file used for Earth Observation Foundation Models (EOFM) and try to make it clearer for a new user. However, we will not go into detail, since the complementary documentation (Lightning, PyTorch, ...) must fill this gap.

## Trainer: global arguments

In the section called **trainer** are defined the arguments that must be sent to the Lightning Trainer object. If you need a deeper explanation about this, check the [Lightning's documentation](#). In the first lines we have:

- **accelerator**: refers to the kind of device is being used to run the experiment. Usually **cpu** and **gpu**, but if you set **auto**, it will automatically select allocate the GPU is that is available or otherwise run on CPU.
- **strategy**: is related to the kind of parallelism is available. As we have usually ran the experiments using a single device for fine-tuning or inference, we do not care about it and choose the option **auto** by default.
- **devices**: indicates the list of available devices to use for the experiment. Leave it as **auto** if you are running with a single device.
- **num\_nodes**: is self-explanatory. We have mostly tested TerraTorch for single-node jobs, so, it is better to set it as 1 for now.
- **precision**: is the kind of precision used for your model. **16-mixed** have been an usual choice.

## Trainer: logging

Just below this initial stage, we have **logger**:

```
1 logger:
2     class_path: TensorBoardLogger
3     init_args:
4         save_dir: tests/
5         name: all\_ecos\_random
6
```

In this field we define the configuration for logging the model state. In this example we are using [Tensorboard](#), and saving all the logs in a directory *tests/all\_ecos\_random*. Others frameworks, as [MLFlow](#) are also supported. Check the [Lightning documentation about logging](#) for a more complete description.

## Trainer: callbacks

The **callbacks** field:

```
7  callbacks:
8    - class_path: RichProgressBar
9    - class_path: LearningRateMonitor
10      init_args:
11        logging_interval: epoch
12    - class_path: EarlyStopping
13      init_args:
14        monitor: val/loss
15        patience: 100
```

Represents a list of operations that can be invoked with determined frequency. The user is free to add others operations from Lightning or custom ones. In the current config we are basically defining: a progress bar to be printed during the model training/validation and a learning rate monitor, determined to call early-stopping when the model shows signals of overfitting.

## Trainer: final part

The rest of the arguments are:

```
16 max_epochs: 1
17 check_val_every_n_epoch: 1
18 log_every_n_steps: 20
19 enable_checkpointing: true
20 default_root_dir: tests/
```

- **max\_epochs**: the maximum number of epochs to train the model. Notice that, if you are using early-stopping, maybe the training will finish before achieving this number.
- **check\_val\_every\_n\_epoch**: the frequency to evaluate the model using the validation dataset. The validation is important to verify if the model is tending to overfit and can be used, for example, to define when update the learning rate, or to invoke the early-stopping.
- **enable\_checkpointing**: it enables the checkpointing, the action of periodically saving the state of the model to a file.
- **default\_root\_dir**: the directory used to save the model checkpoints.

# Datamodule

In this section, we start directly handling TerraTorch's built-in structures. The field **data** is expected to receive a **generic datamodule** or any other datamodule compatible with **Lightning Datamodules**, as those defined in our **collection of datamodules**.

In the beginning of the field we have:

```
21 data:
22     class_path: GenericNonGeoPixelwiseRegressionDataModule
23     init_args:
```

It means that we have chosen the generic regression datamodule and we will pass all its required arguments below **init\_args** and with one new level of indentation. The best practice here is to check the documentation of the datamodule class you are using (in our case, [here](#)) and verify all the arguments it expects to receive and then to fill the lines with **<argument\_name>: <argument\_value>**. As the TerraTorch and Lightning modules were already imported in the CLI script (**terratorch/cli\_tools.py**), you do not need to provide the complete paths for them. Otherwise, if you are using a datamodule defined in an external package, indicate the path to import the model, as **package.datamodules.SomeDatamodule**.

# Model

The field **model** is, in fact, the configuration for **task + model**:

```
24 model:
25   class_path: terratorch.tasks.PixelwiseRegressionTask
26   init_args:
27     model_args:
28       decoder: UperNetDecoder
29       pretrained: false
30       backbone: prithvi_eo_v2_600
31       backbone_drop_path_rate: 0.3
32       backbone_window_size: 8
33       decoder_channels: 64
34       num_frames: 1
35       in_channels: 6
36       bands:
37         - BLUE
38         - GREEN
39         - RED
40         - NIR_NARROW
41         - SWIR_1
42         - SWIR_2
```



## Model: continuation

Notice that there is a field **model\_args**, which it is intended to receive all the necessary configuration to instantiate the model itself, that means, the structure **backbone + decoder + head**. Inside **model\_args**, it is possible to define which arguments will be sent to each component by including a prefix to the argument names, as **backbone\_<argument>** or **decoder\_<other\_argument>**. Alternatively, it is possible to pass the arguments using dictionaries **backbone\_kwargs**, **decoder\_kwargs** and **head\_kwargs**.

```
43     model_args:
44         decoder: UperNetDecoder
45         pretrained: false
46         backbone: prithvi_eo_v2_600
47         backbone_kwargs:
48             drop_path_rate: 0.3
49             window_size: 8
```

The same recommendation made for the **data** field is repeated here, check the documentation of the **task** and model classes **backbones**, **decoders** and **heads** you are using in order to define which arguments to write for each subfield of **model**.

## Model: final part

The last fields are related to the others main parameters (as loss configuration):

```
1   loss: rmse
2   ignore_index: -1
3   freeze_backbone: true
4   freeze_decoder: false
5   model_factory: PrithviModelFactory
```

And the tiled inference. A kind of inference in which the model is applied to small pieces of the image instead of the full one.

```
1   tiled_inference_parameters:
2       h_crop: 224
3       h_stride: 192
4       w_crop: 224
5       w_stride: 192
6       average_patches: true
```

# Optimization and Learning Rate Scheduler

The last two fields of our example are the configuration of the optimizer and the lr scheduler. Those fields are mostly self-explanatory for users already familiar with machine learning:

```
1 optimizer:  
2   class_path: torch.optim.AdamW  
3   init_args:  
4     lr: 0.00013524680528283027  
5     weight_decay: 0.047782217873995426  
6 lr_scheduler:  
7   class_path: ReduceLROnPlateau  
8   init_args:  
9     monitor: val/loss
```

Check the [PyTorch documentation about optimization](#) to understand them more deeply.