

# PRISM-DNN User Guide

## Introducing PRISM-DNN

**PRISM-DNN** (Pretrained pIxFel-based Spatial-Temporal Deep Neural Network) is a deep learning model designed to make accurate predictions for each satellite image pixel by learning from its **spatial neighbors** and **temporal history**. Whenever your task involves predicting something for every pixel over space and time, PRISM-DNN can help.

This ability is especially valuable when you have:

- **Few ground-truth labels**
- **Massive archives of satellite images**

## Why Is It Unique?

### 1. Trained to Understand Space and Time Together

Most models treat each pixel separately. PRISM-DNN is different: It learns how **what happens nearby (in space) and what happened before and after (in time)** affect the pixel you're looking at now.

### 2. Uses Large-Scale Unlabeled Data for Smarter Learning

When real-world labels are limited, PRISM-DNN first trains on **unlabeled satellite images** to learn general spatial-temporal patterns. This pretraining helps the model **make better use of limited labeled samples**, improving both **prediction accuracy** and **data efficiency** in low-resource scenarios.

### 3. Highly Transferable Across Tasks

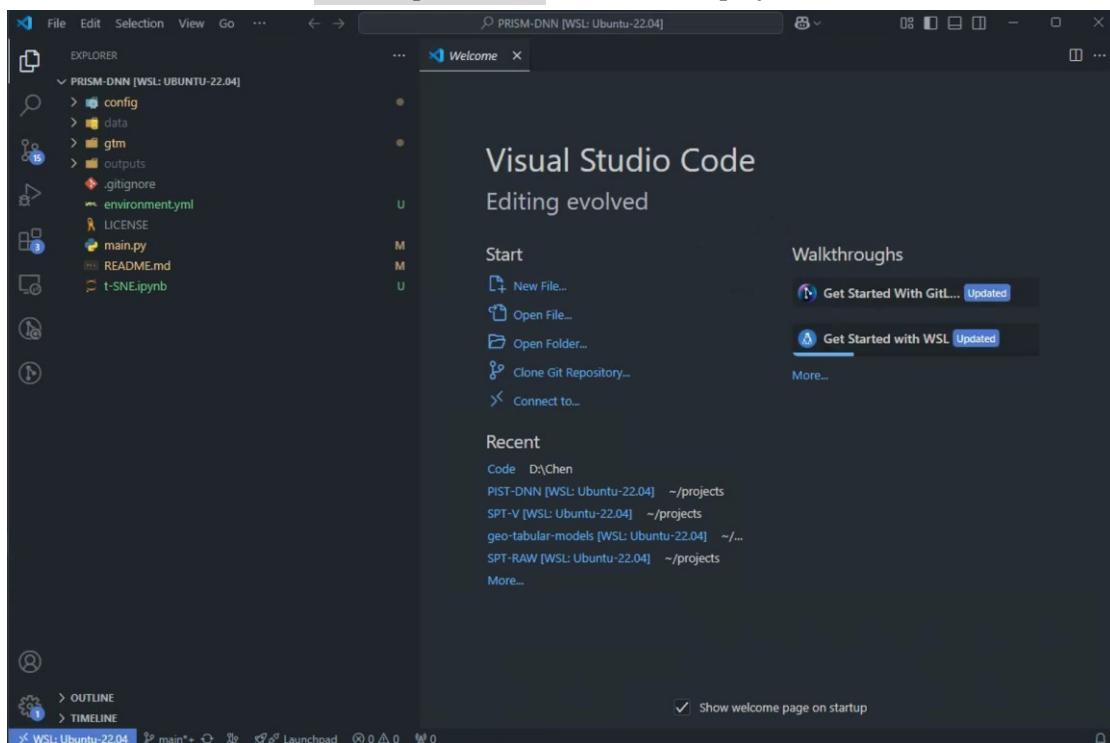
Once pretrained, the model can be adapted to many tasks: retrieving different pollutants, various aerosol components, or other related variables—you name it.

# Step-by-Step Tutorial for Using the Model

## Step 1: Model Setup

### 1.1 Open the Project in VSCode

- 1st. **Download the project source code** from GitHub, then unzip or clone it to a folder on your computer.
- 2nd. **Open VSCode**.
- 3rd. In VSCode, click on *File > Open Folder* and select the project folder.



### 1.2 Set Up Python Environment (Python $\geq$ 3.8)

- 1st. Make sure you have installed **Miniconda** or **Anaconda** on your system.
- 2nd. Open the VSCode terminal (*View → Terminal*), and navigate to your PRISM-DNN project folder:

```
cd ~/projects/PRISM-DNN
```

```
$ cd ~/projects/PRISM-DNN
(bnu)
# bnu @ DESKTOP-KG2BP7D in ~/projects/PRISM-DNN on git:main x [20:13:55]
```

- 3rd. Create the Conda environment using the provided configuration file:

```
conda env create -f environment.yml -n bnu
```

This will automatically create a new environment (You can replace *bnu* with any environment name you prefer) and install all necessary packages listed in the

*environment.yml* file. You can see that the *environment.yml* file includes libraries from *pytorch*, *nvidia*, and *conda-forge* channels, ensuring GPU support and compatibility.

```

environment.yml
...
channels:
  - pytorch
  - nvidia
  - conda-forge
dependencies:
  - _libgcc_mutex=0.1=conda_forge
  - _openmp_mutex=4.5=2_kmp_llvm
  - _py-xgboost-mutex=2.0=cpu_0
  - absl-py=1.4.0=pyhd8ed1ab_0
  - aiofiles=22.1.0=pyhd8ed1ab_0
  - aiohttp=3.8.4=py39h72bdee0_0
  - aiosignal=1.3.1=pyhd8ed1ab_0
  - aiosqlite=0.18.0=pyhd8ed1ab_0
  - alembic=1.10.2=pyhd8ed1ab_0
  - anyio=3.6.2=pyhd8ed1ab_0
  - argon2-cffi=21.3.0=pyhd8ed1ab_0
  - argon2-cffi-bindings=21.2.0=py39hb9d737c_3
  - arrow-cpp=11.0.0=ha770c72_10_cpu
  - asttokens=2.2.1=pyhd8ed1ab_0

```

4th. After installation, activate the environment (bnu):

```

conda activate bnu
# bnu @ DESKTOP-KG2BP7D in ~/projects/PRISM-DNN on git:main x [20:13:55]
$ conda activate bnu
(bnu)

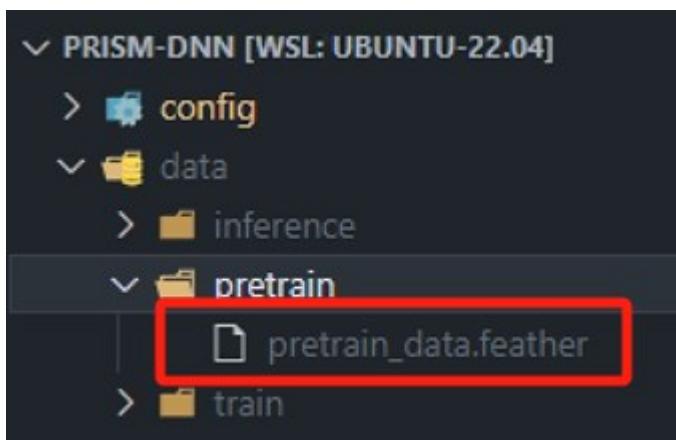
```

Now you're ready to run the code.

## Step 2: Pretraining the Model

### 2.1 Preparing Inputs for Pretraining

(e.g., *pretrain\_data.feather* in *data/pretrain/*)



1st. Two Types of Input Variables (**No Missing Values Allowed**):

- **Categorical** variables (e.g. *Month*, *Season*, *Land cover type*)
- **Continuous** variables (e.g. *Satellite reflectance*, *Brightness temperature*)

2nd. For each pixel (called a **target pixel**), PRISM-DNN needs:

- The **current value** of the target pixel

- The values of its **8 surrounding neighbors** (in space,  $3 \times 3$  window)
- The same set of variables from the **previous and next day** (in time)

3rd. Export your final data table to **CSV, PKL, or Feather** format (recommended: Feather for large-scale data)

## 2.2 Running Pretraining

Once your input file (e.g. `pretrain_data.feather`) is ready, you can start pretraining the PRISM-DNN model. The model reads all its settings from a configuration file: `config/pretrain.yaml`.

### 1st. Locate the Configuration File

In your project folder, open: `config/pretrain.yaml`, You'll see a structure like this:

```
You, 1 second ago | 2 authors (You and one other)
seed: 42

pretrain_data_path: "data/pretrain/pretrain_data.feather"      You, 1 second ago * Uncommitted changes
output_folder_path: "outputs/pretrain"

dataset:
  input_cont_cols: ['MOD08_C_1', 'MOD08_C_2', 'MOD08_C_3', 'MOD08_C_4', 'MOD08_C_5', 'MOD08_C_6', 'MOD08_C_7', 'MOD08_C_8', 'MOD08_C_9', 'MOD08_C_10',
    'MOD08_C_11', 'MOD08_C_12', 'MOD08_C_13', 'MOD08_C_14', 'MOD08_C_15', 'blh_C', 'ps_C', 't2m_C', 'wind_C', 'rh_C',
    'bc_C', 'oc_C', 'dust25_C', 'so4_C', 'ssa_C', 'ssa25_C']
  input_cate_cols: ['Month']
  space_target_cols: ['MOD08_P1', 'MOD08_P2', 'MOD08_P3', 'MOD08_P4', 'MOD08_P6', 'MOD08_P7', 'MOD08_P8', 'MOD08_P9',
    'MOD08_P1_1', 'MOD08_P2_1', 'MOD08_P3_1', 'MOD08_P4_1', 'MOD08_P6_1', 'MOD08_P7_1', 'MOD08_P8_1', 'MOD08_P9_1',
    'MOD08_P1_2', 'MOD08_P2_2', 'MOD08_P3_2', 'MOD08_P4_2', 'MOD08_P6_2', 'MOD08_P7_2', 'MOD08_P8_2', 'MOD08_P9_2',
    'MOD08_P1_3', 'MOD08_P2_3', 'MOD08_P3_3', 'MOD08_P4_3', 'MOD08_P6_3', 'MOD08_P7_3', 'MOD08_P8_3', 'MOD08_P9_3',
    'MOD08_P1_4', 'MOD08_P2_4', 'MOD08_P3_4', 'MOD08_P4_4', 'MOD08_P6_4', 'MOD08_P7_4', 'MOD08_P8_4', 'MOD08_P9_4',
    'MOD08_P1_5', 'MOD08_P2_5', 'MOD08_P3_5', 'MOD08_P4_5', 'MOD08_P6_5', 'MOD08_P7_5', 'MOD08_P8_5', 'MOD08_P9_5',
    'MOD08_P1_6', 'MOD08_P2_6', 'MOD08_P3_6', 'MOD08_P4_6', 'MOD08_P6_6', 'MOD08_P7_6', 'MOD08_P8_6', 'MOD08_P9_6',
    'MOD08_P1_7', 'MOD08_P2_7', 'MOD08_P3_7', 'MOD08_P4_7', 'MOD08_P6_7', 'MOD08_P7_7', 'MOD08_P8_7', 'MOD08_P9_7',
    'MOD08_P1_8', 'MOD08_P2_8', 'MOD08_P3_8', 'MOD08_P4_8', 'MOD08_P6_8', 'MOD08_P7_8', 'MOD08_P8_8', 'MOD08_P9_8',
    'MOD08_P1_9', 'MOD08_P2_9', 'MOD08_P3_9', 'MOD08_P4_9', 'MOD08_P6_9', 'MOD08_P7_9', 'MOD08_P8_9', 'MOD08_P9_9',
    'MOD08_P1_10', 'MOD08_P2_10', 'MOD08_P3_10', 'MOD08_P4_10', 'MOD08_P6_10', 'MOD08_P7_10', 'MOD08_P8_10', 'MOD08_P9_10',
    'MOD08_P1_11', 'MOD08_P2_11', 'MOD08_P3_11', 'MOD08_P4_11', 'MOD08_P6_11', 'MOD08_P7_11', 'MOD08_P8_11', 'MOD08_P9_11',
    'MOD08_P1_12', 'MOD08_P2_12', 'MOD08_P3_12', 'MOD08_P4_12', 'MOD08_P6_12', 'MOD08_P7_12', 'MOD08_P8_12', 'MOD08_P9_12',
    'MOD08_P1_13', 'MOD08_P2_13', 'MOD08_P3_13', 'MOD08_P4_13', 'MOD08_P6_13', 'MOD08_P7_13', 'MOD08_P8_13', 'MOD08_P9_13',
    'MOD08_P1_14', 'MOD08_P2_14', 'MOD08_P3_14', 'MOD08_P4_14', 'MOD08_P6_14', 'MOD08_P7_14', 'MOD08_P8_14', 'MOD08_P9_14',
    'MOD08_P1_15', 'MOD08_P2_15', 'MOD08_P3_15', 'MOD08_P4_15', 'MOD08_P6_15', 'MOD08_P7_15', 'MOD08_P8_15', 'MOD08_P9_15',
    'blh_P1', 'blh_P2', 'blh_P3', 'blh_P4', 'blh_P6', 'blh_P7', 'blh_P8', 'blh_P9',
    'ps_P1', 'ps_P2', 'ps_P3', 'ps_P4', 'ps_P6', 'ps_P7', 'ps_P8', 'ps_P9',
    't2m_P1', 't2m_P2', 't2m_P3', 't2m_P4', 't2m_P6', 't2m_P7', 't2m_P8', 't2m_P9',
    'wind_P1', 'wind_P2', 'wind_P3', 'wind_P4', 'wind_P6', 'wind_P7', 'wind_P8', 'wind_P9',
    'rh_P1', 'rh_P2', 'rh_P3', 'rh_P4', 'rh_P6', 'rh_P7', 'rh_P8', 'rh_P9',
    'bc_P1', 'bc_P2', 'bc_P3', 'bc_P4', 'bc_P6', 'bc_P7', 'bc_P8', 'bc_P9']
```

### 2nd. Key Configuration Sections

\* These column names must exactly match the ones in your input file.

- `pretrain_data_path`:

Path to the input file for pretraining. Should point to your `feather/.csv/.pkl` file.

- `output_folder_path`:

Where to save model weights and logs after training. Suggested: leave as

`"outputs/pretrain"` unless you want a custom folder.

- `input_cont_cols`:

All **continuous input variables** from the central pixel at the current time step, used for learning spatial-temporal patterns (e.g., satellite bands, temperature, etc.).

- `input_cate_cols`:

List of **categorical inputs**, such as Month, Land cover, Sensor ID, etc.

- `space_target_cols`:

These columns contain values from the **8 neighboring pixels** around the target pixel (in a  $3 \times 3$  window).

- `space_target_cols`:

These columns contain values from the **8 neighboring pixels** around the target pixel.

- `time_target_cols`:

Columns used as temporal prediction targets, representing the values of the central pixel from the neighboring time steps (e.g. **previous and next day**).

- Dataloader:

Controls batch size and data loading speed:

```
batch_size: 1024          # Number of samples per batch (adjust based on memory)

pin_memory: True          # Speeds up GPU transfer

persistent_workers: True  # Recommended for large datasets
```

- Model:

Defines model architecture:

```
d_model: 128            # Embedding size

n_tf_head: 4             # Number of transformer heads

n_tf_layer: 4            # Number of transformer layers

p_tf_drop: 0.2           # Dropout rate in transformer

p_mlp_layer: 2           # MLP depth

p_mlp_drop: 0.2          # Dropout in MLP

lr: 0.0003               # Learning rate
```

- Callbacks:

Tells the model when and how to stop:

```
save_top_k: 1             # Save the best model only

monitor: "valid_loss"     # Monitor validation loss

mode: "min"                # Lower is better

patience: 10              # Stop early if no improvement after 10 epochs
```

- Trainer:

Sets how the model will run:

```
max_epochs: 200           # Max training epochs
```

```

accelerator: "gpu"           # Use GPU (recommended)

devices: 1                   # Number of GPUs (set to 1 unless using multi-GPU)

```

### 3rd. Run Pretraining in Terminal

Make sure your Conda environment is activated, then run:

```
python main.py -t pretrain -c config/pretrain.yaml
```

The model will start reading the input file, building the model, and saving output to *outputs/pretrain/*.

```

EXPLORER          config > pretrain.yaml M X
PRISM-DNN [WSL: UBUNTU-22.04]
  config
    finetune_inference.yaml
    finetune.yaml
    pretrain.yaml
  data
  gtm
  outputs
  .gitignore
  environment.yaml
  LICENSE
  main.py
  README.md
  t-SNE.ipynb

config > pretrain.yaml
You, 11 hours ago | 2 authors (You and one other)
1 seed: 42
2
3 pretrain_data_path: "data/pretrain/pretrain_data.feather"
4 output_folder_path: "outputs/pretrain"
5
6 dataset:
7   input_cont_cols: ['MOD08_C','MOD09_C_1','MOD09_C_2','MOD09_C_3','MOD09_C_4','MOD09_C_11','MOD09_C_12','MOD09_C_13','MOD09_C_14','bc_C','oc_C','dust_C','dust25_C','so4_C','ssa_C']
8   input_cate_cols: ['Month']
9   space_target_cols: ['MOD08_P1', 'MOD08_P2', 'MOD08_P3', 'MOD08_P4', 'MOD09_P1_1', 'MOD09_P2_1', 'MOD09_P3_1', 'MOD09_P4_1', 'MOD09_P1_2', 'MOD09_P2_2', 'MOD09_P3_2', 'MOD09_P4_2', 'MOD09_P1_3', 'MOD09_P2_3', 'MOD09_P3_3', 'MOD09_P4_3', 'MOD09_P1_4', 'MOD09_P2_4', 'MOD09_P3_4', 'MOD09_P4_4', 'MOD09_P1_5', 'MOD09_P2_5', 'MOD09_P3_5', 'MOD09_P4_5', 'MOD09_P1_6', 'MOD09_P2_6', 'MOD09_P3_6', 'MOD09_P4_6', 'MOD09_P1_7', 'MOD09_P2_7', 'MOD09_P3_7', 'MOD09_P4_7']
10
11
12
13
14
15
16
17
18

PROBLEMS   OUTPUT   DEBUG CONSOLE   TERMINAL   PORTS   GITLENS
(bnu)
# huu @ DESKTOP-KC9BPPD in ~/projects/PRISM-DNN on git:main x [20:17:16]
$ python main.py -t pretrain -c config/pretrain.yaml
Global seed set to 42
{'callbacks': {'early_stopping': {'min_delta': 0.0,
                                    'mode': 'min',
                                    'monitor': 'valid_loss',

```

Running screenshot:

```

You are using a CUDA device ('NVIDIA GeForce RTX 4090') that has Tensor Cores. To properly utilize them, you should set `torch.set_float32_matmul_precision('medium' | 'high')` which will trade-off precision for performance. For more details, read https://pytorch.org/docs/stable/generated/torch.set_float32_matmul_precision.html#torch.set_float32_matmul_precision
LOCAL_RANK: 0 - CUDA_VISIBLE_DEVICES: [0]

| Name      | Type       | Params | |
|---|---|---|---|
| geo_encoder | GeoTransformerEncoder | 538 K |
| 1 | sparse_decoder | MLPDecoder | 41.9 K |
| 2 | time_decoder | MLPDecoder | 46.2 K |
| 3 | criterion | MSELoss | 0 |

641 K Trainable params
0 Non-trainable params
641 K Total params
2.6M Total estimated model params size (MB)
Epoch 100: 1000 | 3878/3878 [01:12<00:00, 53.58it/s, loss=0.319, v_num=2, valid_loss=0.234]

```

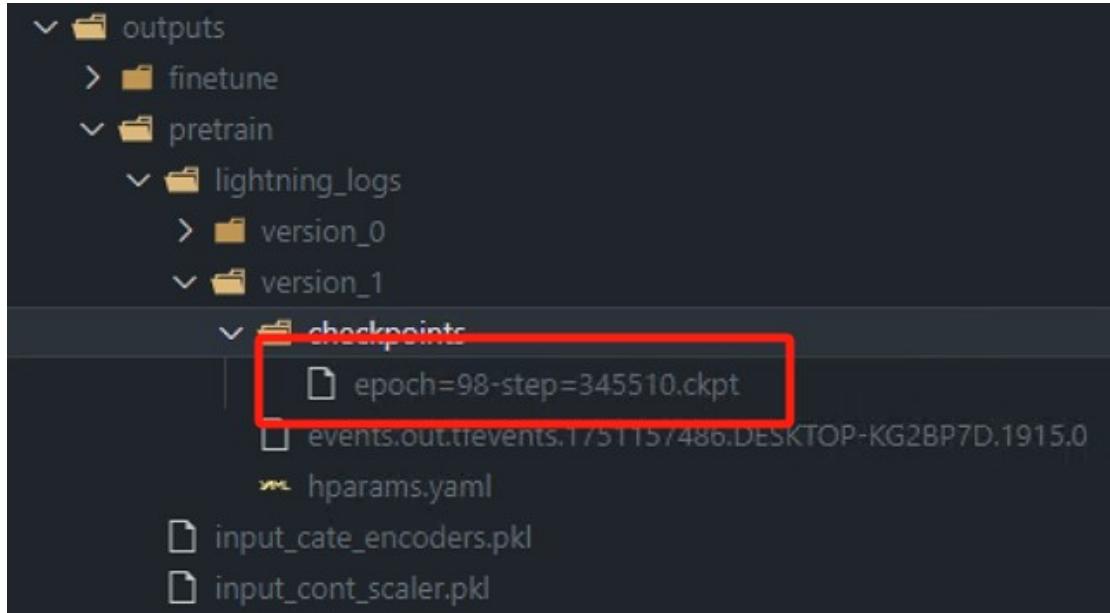
### 4th. Pretrain Finished

When the pretraining process completes, the following files will be generated automatically in the *outputs/pretrain/* directory:

- *input\_cate\_encoders.pkl* – encoders for categorical variables;
- *input\_cont\_scaler.pkl* – scalers for continuous variables;
- *lightning\_logs/* – training logs and versioned model runs.

\*Each time you run the pretraining process, a new subfolder (e.g., version\_0, version\_1, ...) will be created under outputs/pretrain/lightning\_logs/. For each version, the best-performing model checkpoint is saved inside:

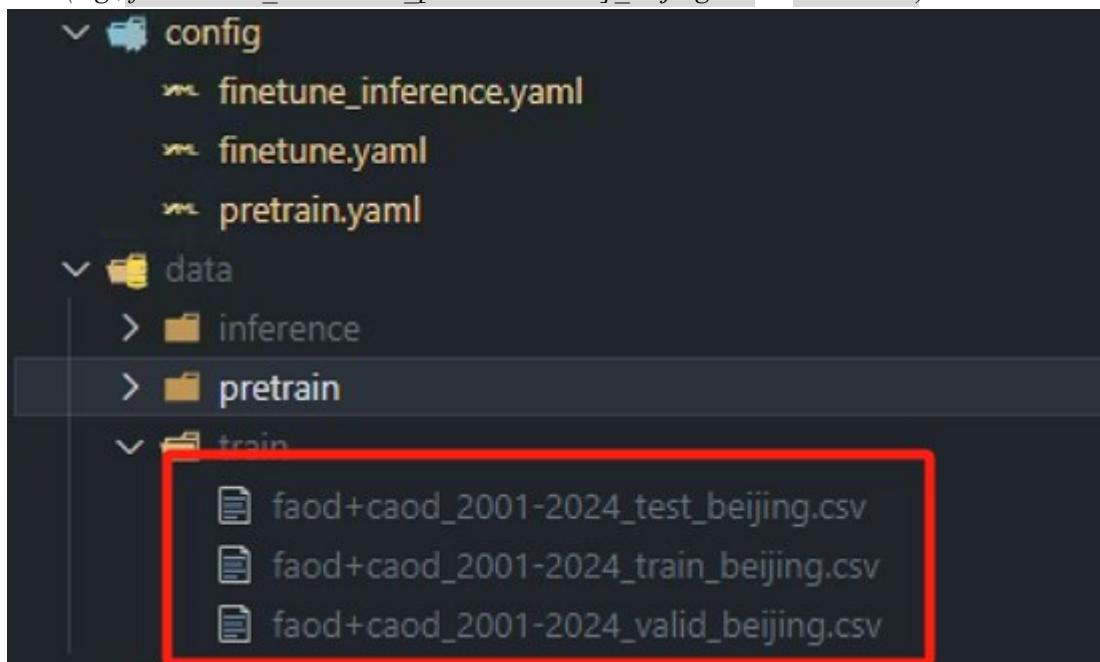
`outputs/pretrain/lightning_logs/version_X/checkpoints/`



### Step 3: Fine-tuning for a Target Task

#### 3.1 Preparing Inputs for Fine-tuning

(e.g., `faod+caod_2001-2024_[train/valid/test]_beijing.csv` in `data/train/`)



1st. Two Types of Input Variables (No Missing Values Allowed):

- **Categorical** variables (e.g. *Month, Season, Land cover type*)
- **Continuous** variables (e.g. *Satellite reflectance, Brightness temperature*)

2nd. Fine-tuning requires a labeled dataset, where each sample includes both:

- **Input features** (from the central pixel only)
- **Target value** (the real variable you want the model to learn to predict, e.g. fAOD, cAOD, FMF)

\*This step is different from pretraining. Here, we use a **standard supervised learning format**, similar to how models like LightGBM or Random Forest are trained.

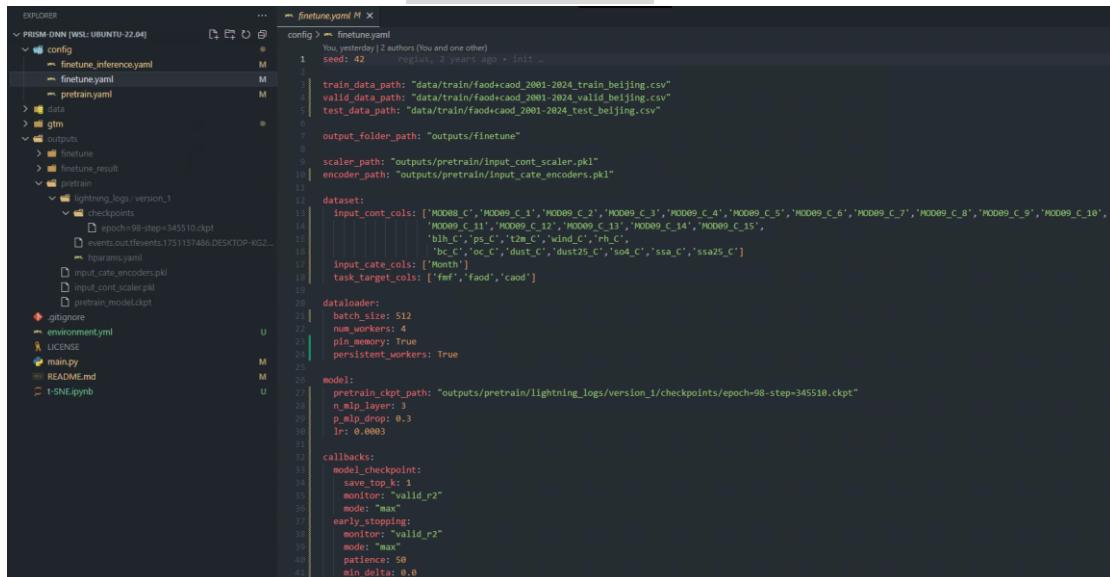
3rd. Export your final data table to **CSV** format.

### 3.2 Running Fine-tuning

Once you have prepared your labeled dataset and filled out the fine-tuning config file (*config/fine\_tune.yaml*), you can run supervised training to adapt the pretrained model to your specific prediction task—such as estimating fAOD, cAOD, FMF or other variables.

1st. Locate the Configuration File

In your project folder, open: *config/fine\_tune.yaml*. You'll see a structure like this:



```

EXP ORDER
└ PRISM-DNN (WSL: UBUNTU-22.04)
  └ config
    └ finetune.yaml
      You, yesterday | 2 authors (You and one other)
      1 seed: 42
      2 regday, 2 years ago + init -
      3 train_data_path: "data/train/faod+caod_2001-2024_train_beijing.csv"
      4 valid_data_path: "data/train/faod+caod_2001-2024.valid_beijing.csv"
      5 test_data_path: "data/train/faod+caod_2001-2024_test_beijing.csv"
      6
      7 output_folder_path: "outputs/finetune"
      8
      9 scalar_path: "outputs/pretrain/input_cont_scaler.pkl"
      10 encoder_path: "outputs/pretrain/input_cate_encoders.pkl"
      11
      12 dataset:
        13   input_cont_cols: ["MOD08_C_1", "MOD09_C_1", "MOD09_C_2", "MOD09_C_3", "MOD09_C_4", "MOD09_C_5", "MOD09_C_6", "MOD09_C_7", "MOD09_C_8", "MOD09_C_9", "MOD09_C_10",
        14     "MOD09_C_11", "MOD09_C_12", "MOD09_C_13", "MOD09_C_14", "MOD09_C_15",
        15     "tbl_C", "ps_C", "t2m_C", "wind_C", "dust25_C", "sod_C", "ssa_C", "ssa25_C"]
        16
        17   input_cate_cols: ["Month"]
        18   task_target_cols: ["fmf", "faod", "caod"]
        19
        20 dataloader:
        21   batch_size: 512
        22   num_workers: 4
        23   pin_memory: True
        24   persistent_workers: True
        25
        26 model:
        27   pretrain_ckpt_path: "outputs/pretrain/lightning_logs/version_1/checkpoints/epoch=98-step=345510.ckpt"
        28   n_mlp_layer: 3
        29   p_mlp_drop: 0.3
        30   lr: 0.0003
        31
        32 callbacks:
        33   mode_checkpoint:
        34     save_top_k: 1
        35     monitor: "valid_r2"
        36     mode: "max"
        37   early_stopping:
        38     monitor: "valid_r2"
        39     mode: "max"
        40     patience: 50
        41     min_delta: 0.0
  
```

2nd. Key Configuration Sections

\* These column names must exactly match the ones in your input file.

- ***train\_data\_path*, *valid\_data\_path*, *test\_data\_path*:**

Path to the input file for fine-tuning. Should point to your *.csv* file.

- ***scalar\_path*, *encoder\_path*:**

These files are generated during pretraining. By default, they are saved in the directory: *outputs/pretrain/*

- ***input\_cont\_cols*:**

All **continuous input variables** (e.g. MODIS bands, temperature).

- ***input\_cate\_cols*:**

All **categorical variables** (e.g. Month).

- ***task\_target\_cols*:**

Your **label columns** (e.g. faod, caod, fmf).

- Dataloader:

Controls batch size and data loading speed:

```
batch_size: 1024          # Number of samples per batch (adjust based on memory)

pin_memory: True          # Speeds up GPU transfer

persistent_workers: True  # Recommended for large datasets
```

- Model:

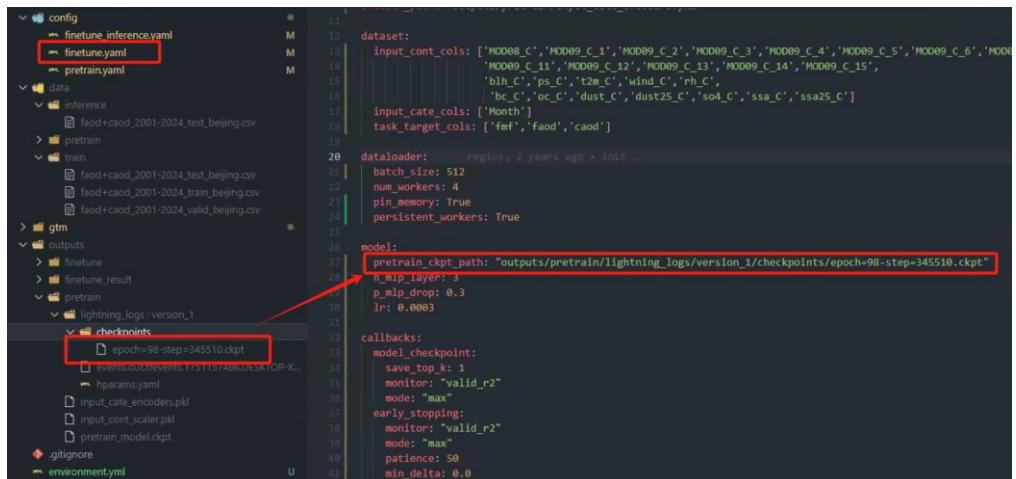
Model structure and learning rate:

```
pretrain_ckpt_path:
outputs/pretrain/lightning_logs/version_1/checkpoints/epoch-98-
step=345510.ckpt      # replace to the best checkpoint from your pretraining

n_mlp_layer: 3          # the number of fully connected (MLP) layer

p_mlp_drop: 0.3         # Dropout in MLP

lr: 0.0003              # Learning rate
```



- Callbacks:

Tells the trainer when to stop:

```
model_checkpoint:

    monitor: "valid_r2"    # Save model with best R2 score

early_stopping:

    patience: 50           # Stop early if no improvement after 50 epochs
```

- Trainer:

Runs the training loop:

```
max_epochs: 150          # Max training epochs

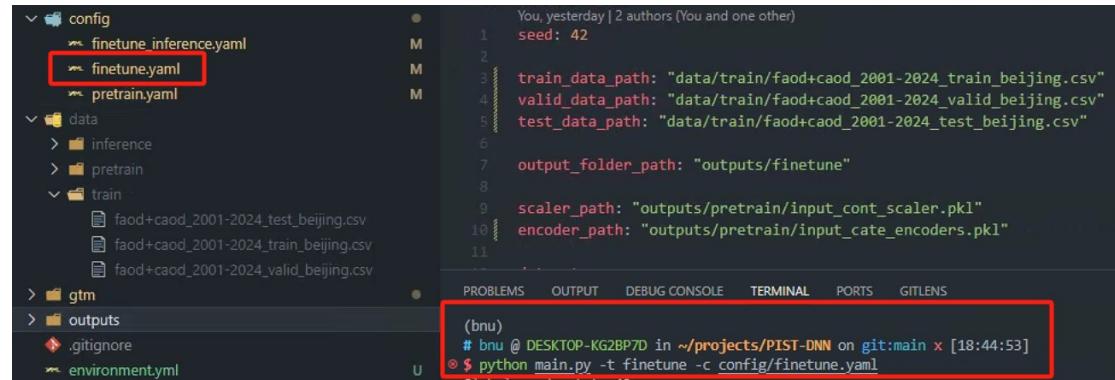
accelerator: "gpu"      # Use GPU (recommended)

devices: 1               # Number of GPUs (set to 1 unless using multi-GPU)
```

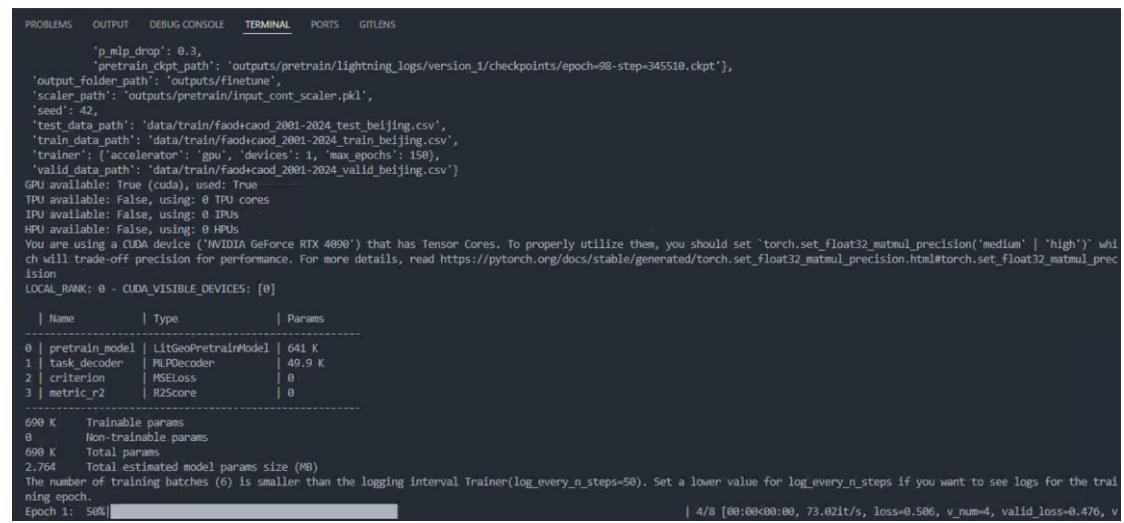
### 3rd. Run Fine-tuning in Terminal

Once your config is ready and your environment is active, run this command:

```
python main.py -t finetune -c config/finetune.yaml
```



Running screenshot:



### 4th. Fine-tune Finished

When the fine-tuning process completes, the terminal will display a summary of the prediction performance for each target variable. This includes metrics like  $R^2$  and RMSE, providing a quick overview of how well the model performed on the test set.

```

690 K Total params
2,764 Total estimated model params size (MB)
The number of training batches (6) is smaller than the logging interval Trainer(log_every_n_steps=50). Set a lower value for log_every_n_steps if you want to see logs for the training epoch.
Epoch 149: 100% [██████████] | 8/8 [00:00:00:00, 85.45it/s, loss=0.0196, v_num=3, valid_loss=0.0147, valid_r2=0.762]
"Trainer.fit" stopped: "max_epochs=150" reached.
Epoch 149: 100% [██████████] | 8/8 [00:00:00:00, 84.89it/s, loss=0.0196, v_num=3, valid_loss=0.0147, valid_r2=0.762]
You are using a CUDA device ('NVIDIA GeForce RTX 4090') that has Tensor Cores. To properly utilize them, you should set 'torch.set_float32_matmul_precision('medium' | 'high')' which will trade-off precision for performance. For more details, read https://pytorch.org/docs/stable/generated/torch.set\_float32\_matmul\_precision.html
Restoring states from the checkpoint at outputs/finetune/lightning_logs/version_3/checkpoints/epoch=143-step=864.ckpt
LOCAL_RANK: 0 - CUDA_VISIBLE_DEVICES: [0]
Loaded model weights from checkpoint at outputs/finetune/lightning_logs/version_3/checkpoints/epoch=143-step=864.ckpt
Predicting DataLoader 0: 100% [██████████] | 1/1 [00:00:00:00, 108.23it/s]
Task fm:
RMSE: 0.131
R2: 0.561
Task faod:
RMSE: 0.155
R2: 0.599
Task caod:
RMSE: 0.055
R2: 0.778
(bnu)
$ bnu @ DESKTOP-KG2BP7D in ~/projects/SPT-V on git:main ✘ [14:21:34]
$ [ ]

```

After that, the following files will be automatically saved in the `outputs/finetune/` directory:

- `faod+caod_2001-2024_test_beijing_pred.csv`

This file includes predictions for each target (e.g., faod\_PRED, caod\_PRED, fmf\_PRED); 128-dimensional pretrained feature embeddings for each sample (e.g., x\_feature\_0 to x\_feature\_127)

- `lightning_logs/`

Training logs, including validation scores, learning rate, epoch summaries, and checkpoints.



## That's it! The full training workflow is now complete.

### Step 4: Global-Scale Inference without Ground Truth

To perform inference in areas where no ground truth is available (e.g., pixels without station-based measurements), first place your input data in the `data/fine_tune_inference/` directory. This input file must have the **exact same column names and structure** as the one used during fine-tuning. Since there are no true target values, you should still include the target columns (e.g., faod, caod, fmf) but fill them with a constant value such as 0. These columns are required for compatibility with the model input format, but their values will not affect prediction results—as long as they are present and **non-empty**.

In your YAML configuration file (`fine_tune_inference.yaml`), set the `model_checkpoint_path` to the relative path of the best fine-tuned model checkpoint (e.g., `outputs/fine_tune/lightning_logs/version_3/checkpoints/epoch=143-step=864.ckpt`). The

settings under dataset: **must match exactly with those in the original fine-tuning YAML** to ensure consistent data preprocessing and feature encoding.

Then, execute the inference with the following command:

```
python main.py -t finetune_inference -c config/fine_tune_inference.yaml
```

```
finetune_inference.yaml M
config > finetune_inference.yaml
You, 1 hour ago | 1 author (You)
1 seed: 42
2
3 test_data_path: "data/inference"
4
5 output_folder_path: "outputs/fine_tune_result"
6
7 scaler_path: "outputs/pretrain/input_cont_scaler.pkl"
8 encoder_path: "outputs/pretrain/input_cate_encoders.pkl"
9 model_checkpoint_path: "outputs/fine_tune/lightning_logs/version_3/checkpoints/epoch=143-step=864.ckpt"
10
11 dataset:
12     input_cont_cols: ['MOD08_C_1', 'MOD09_C_1', 'MOD09_C_2', 'MOD09_C_3', 'MOD09_C_4', 'MOD09_C_5', 'MOD09_C_6', 'MOD09_C_7',
13         'MOD09_C_11', 'MOD09_C_12', 'MOD09_C_13', 'MOD09_C_14', 'MOD09_C_15',
14         'bil_C', 'ps_C', 't2m_C', 'wind_C', 'rh_C',
15         'bc_C', 'oc_C', 'dust25_C', 'so4_C', 'ssa_C', 'ssa25_C']
16     input_cate_cols: ['Month']
17     task_target_cols: ['fmf', 'faod', 'caod']
18
19 dataloader:
20     batch_size: 512
21     num_workers: 4
22     pin_memory: True
23     persistent_workers: True
24
25 trainer:
26     max_epochs: 162
27     accelerator: "gpu"
28     devices: 1
```

The model will output predictions for every pixel in the input file, enabling global-scale mapping without requiring any labeled data.

## Step 5: Visualizing Pretrained Features with t-SNE

To better understand what the model has learned during pretraining, you can **visualize the 128-dimensional feature embeddings** from the test predictions using t-SNE, a commonly used dimensionality reduction technique.

We provide a Jupyter Notebook (*t-SNE.ipynb*) that demonstrates how to load the prediction results, extract pretrained features, and visualize them in 3D space.

3D t-SNE of Sampled x\_features

