

# Satellite NDVI Forecasting with Growth Curve Regression

A Deep Learning Framework for Long-Horizon Vegetation Dynamics  
Prediction

Antonio Henrique Xavier da Silva

January 2026

## **Abstract**

This thesis presents a novel deep learning framework for forecasting vegetation dynamics using multi-spectral Sentinel-2 satellite imagery, meteorological data, and land cover classification. Unlike conventional approaches that predict vegetation states iteratively step-by-step, the proposed method learns complete growth curve trajectories to enable interpretable 100-day forecasts. The architecture combines ConvLSTM-based spatiotemporal encoding with a parametric growth curve decoder, predicting saturation growth parameters (amplitude, rate, and offset) that model vegetation phenology explicitly. Trained on the GreenEarthNet benchmark dataset, the model demonstrates competitive performance while maintaining significantly fewer parameters than transformer-based alternatives. The growth curve formulation provides inherent interpretability, enabling analysis of predicted vegetation dynamics in terms of biophysically meaningful parameters.

# Contents

<b>1</b>	<b>Introduction</b>	<b>6</b>
1.1	Context and Motivation . . . . .	6
1.2	Problem Statement . . . . .	7
1.3	Research Objectives . . . . .	7
1.4	Contributions . . . . .	8
1.5	Thesis Structure . . . . .	9
<b>2</b>	<b>Literature Review</b>	<b>10</b>
2.1	Vision Transformers in Remote Sensing . . . . .	10
2.1.1	Adaptation to Satellite Imagery . . . . .	10
2.1.2	Hierarchical and Efficient Transformers . . . . .	11
2.2	Satellite Image Time Series Analysis . . . . .	11
2.2.1	Temporal-Spatial Factorization . . . . .	11
2.2.2	Lightweight Architectures . . . . .	12
2.3	Vegetation Forecasting . . . . .	12
2.3.1	Multi-Modal Learning for Geospatial Forecasting . . . . .	12
2.3.2	Explainable Earth Surface Forecasting . . . . .	13
2.4	Research Gaps and Opportunities . . . . .	14
2.4.1	Trajectory vs. Step-by-Step Prediction . . . . .	14
2.4.2	Regression-Focused Architectures . . . . .	14
2.4.3	Efficiency-Accuracy Trade-offs . . . . .	15
2.5	Summary . . . . .	15
<b>3</b>	<b>Materials and Methods</b>	<b>16</b>
3.1	Dataset: GreenEarthNet . . . . .	17

3.1.1	Data Sources . . . . .	17
3.1.2	Temporal Structure . . . . .	18
3.1.3	Dataset Splits . . . . .	18
3.2	Data Preprocessing Pipeline . . . . .	19
3.2.1	Cloud Masking and NaN Handling . . . . .	19
3.2.2	Best Available Pixel (BAP) Compositing . . . . .	19
3.2.3	Weather Feature Engineering . . . . .	20
3.2.4	Temporal Metadata . . . . .	21
3.2.5	Weather-Time Adjustment Mechanism . . . . .	21
3.3	Model Architecture . . . . .	22
3.3.1	Architecture Overview . . . . .	22
3.3.2	Latent Space Encoder . . . . .	22
3.3.3	Regression Parameter Head . . . . .	23
3.3.4	Weather-Time Adjustment MLP . . . . .	24
3.3.5	Growth Curve Layer . . . . .	24
3.3.6	Spatial Smoothing Layer . . . . .	25
3.4	Loss Function . . . . .	25
3.4.1	Masked Huber Regression Loss . . . . .	25
3.4.2	Variance Penalty . . . . .	25
3.4.3	kNDVI Loss . . . . .	25
3.4.4	Combined Loss . . . . .	26
3.5	Training Configuration . . . . .	26
3.5.1	Optimization . . . . .	26
3.5.2	Model Size . . . . .	27
3.5.3	Hardware . . . . .	27
<b>4</b>	<b>Results</b>	<b>28</b>
4.1	Error Analysis . . . . .	28
4.1.1	Per-Band Error Metrics . . . . .	28
4.1.2	Error by Prediction Step . . . . .	29
4.1.3	Error by Land Cover Class . . . . .	31
4.2	Comparison with Contextformer . . . . .	32
4.2.1	Nash-Sutcliffe Efficiency . . . . .	32

4.2.2	Efficiency Analysis . . . . .	33
4.3	Interpretability Analysis . . . . .	33
4.3.1	Growth Curve Parameter Interpretation . . . . .	33
4.3.2	NIR vs. Red Band Dynamics . . . . .	34
4.3.3	Regional Vegetation Dynamics . . . . .	34
4.4	Summary . . . . .	35

# List of Figures

3.1	Model architecture overview showing the encoder-decoder structure with growth curve parameterization. . . . .	22
4.1	kNDVI RMSE evolution over the 100-day forecast horizon (20 steps at 5-day intervals). Lower values indicate better prediction accuracy. . . .	29
4.2	Nash-Sutcliffe Efficiency evolution over the 100-day forecast horizon. Higher values indicate better predictive skill (1.0 is perfect, 0.0 equals predicting the mean). . . . .	30
4.3	Prediction error by ESA WorldCover land cover class. Vegetated classes (Tree cover, Shrubland, Grassland) show varying prediction difficulty based on their phenological complexity. . . . .	31
4.4	Spatial distribution of growth curve parameters for a sample validation scene. Top row: Amplitude ( $A$ ) for NIR and Red bands. Bottom row: Rate ( $\lambda$ ) and Offset ( $B$ ). Vegetated areas show distinct parameter patterns corresponding to their phenological state. . . . .	34

# List of Tables

3.1	Sentinel-2 spectral bands used in this study . . . . .	17
3.2	E-OBS climate variables . . . . .	18
3.3	GreenEarthNet dataset splits . . . . .	19
3.4	Growth curve parameters and their interpretations . . . . .	23
3.5	Training hyperparameters . . . . .	26
3.6	Model parameter count by component . . . . .	27
4.1	Per-band prediction error on the GreenEarthNet validation set . . . . .	29
4.2	kNDVI prediction metrics by ESA WorldCover land cover class . . . . .	31
4.3	Nash-Sutcliffe Efficiency comparison on GreenEarthNet validation set. Vegetation Score is the mean score across vegetated land cover classes (Tree cover, Shrubland, Grassland). . . . .	32
4.4	Average growth curve parameters for vegetated regions . . . . .	35

# Chapter 1

## Introduction

### 1.1 Context and Motivation

Climate change poses unprecedented challenges to agricultural systems worldwide, with particularly severe impacts in Mediterranean regions characterized by significant inter-annual climatic variability. The alternation between prolonged drought periods and intense precipitation events creates complex stress dynamics for vegetation, making traditional agronomic management increasingly difficult and unreliable.

Modern agriculture increasingly relies on Earth observation data to monitor crop health and predict vegetation dynamics. Satellite-based remote sensing, particularly from the European Space Agency’s Sentinel-2 mission, provides high-resolution multispectral imagery at regular intervals, enabling systematic monitoring of vegetation across large spatial extents. The Normalized Difference Vegetation Index (NDVI) and its variants have become standard proxies for vegetation health, photosynthetic activity, and biomass estimation.

However, the transition from reactive monitoring to proactive prediction remains a significant challenge. Traditional approaches to vegetation monitoring only allow detection of crop stress when damage is already visually evident and often irreversible. The ability to forecast vegetation dynamics days to months in advance would enable preventive interventions during optimal time windows, fundamentally transforming agricultural decision-making from reactive to anticipatory.



## 1.2 Problem Statement

The prediction of vegetation dynamics from satellite image time series presents several technical challenges:

1. **Temporal Irregularity:** Satellite observations are affected by cloud cover and orbital patterns, creating irregular time series that conventional sequence models struggle to handle effectively.
2. **Multi-Modal Data Integration:** Vegetation growth depends on multiple factors including historical reflectance patterns, meteorological conditions, and land cover characteristics. Effective forecasting requires principled integration of heterogeneous data sources operating at different spatial and temporal resolutions.
3. **Long-Horizon Prediction:** While short-term predictions (days ahead) are relatively tractable, forecasting vegetation states over horizons of weeks to months requires capturing both fast dynamics (weather responses) and slow dynamics (phenological progression).
4. **Interpretability:** Agricultural applications require not only accurate predictions but also explanations that can be translated into actionable recommendations. Black-box models, despite potentially high accuracy, provide limited utility for agronomic decision support.

## 1.3 Research Objectives

This research addresses the aforementioned challenges through the development of a deep learning framework specifically designed for long-horizon vegetation forecasting. The primary objectives are:

1. **Long-Horizon Forecasting:** Develop a model capable of predicting vegetation dynamics over a 100-day forecast horizon using 50 days of historical observations, substantially exceeding the typical 5-10 day horizons of iterative step-by-step prediction methods.
2. **Multi-Modal Integration:** Design an architecture that effectively fuses multi-spectral satellite imagery (Sentinel-2), meteorological variables (E-OBS climate

data), and static land cover information (ESA WorldCover) through learned attention mechanisms.

3. **Interpretable Predictions:** Implement a parametric growth curve decoder that generates predictions through biophysically meaningful parameters (growth amplitude, rate, and offset) rather than opaque neural network outputs, enabling interpretation of forecasts in terms of vegetation phenology.
4. **Computational Efficiency:** Achieve competitive performance with significantly fewer parameters than state-of-the-art transformer-based approaches, enabling deployment in resource-constrained operational settings.

## 1.4 Contributions

The main contributions of this work are:

1. **Growth Curve Trajectory Learning:** A novel approach to vegetation forecasting that learns complete saturation growth curve trajectories rather than predicting iteratively step-by-step. Unlike the iterative ConvLSTM approach of [Pellicer-Valero et al. \(2024\)](#), our method fits entire 100-day trajectories in a single forward pass, enabling efficient long-horizon prediction without error accumulation.
2. **Weather-Adjusted Growth Dynamics:** An architecture component that modulates growth curve parameters based on meteorological conditions, allowing the model to capture weather-dependent variations in vegetation response while maintaining the interpretable growth curve structure.
3. **Lightweight Multi-Modal Architecture:** A ConvLSTM-based encoder with cloud-aware gating that effectively processes irregular satellite observations while integrating weather and land cover information, achieving competitive performance with fewer than 1 million parameters.
4. **Empirical Validation:** Comprehensive evaluation on the GreenEarthNet benchmark dataset, including comparison with state-of-the-art models and analysis of prediction quality across different land cover types and forecast horizons.

## 1.5 Thesis Structure

-> TODO: This cannot be a list. I need to explain the thesis structure using paragraphs.

The remainder of this thesis is organized as follows:

- **Chapter 2** reviews related work on vision transformers for remote sensing, satellite image time series analysis, and vegetation forecasting, identifying research gaps that motivate the proposed approach.
- **Chapter 3** describes the GreenEarthNet dataset, data preprocessing pipeline, model architecture, loss function design, and training configuration.
- **Chapter 4** presents experimental results including quantitative error metrics, temporal and spatial error analysis, prediction visualizations, model comparison, and interpretability analysis of growth curve parameters.
- **Chapter 5** discusses implications of the results, limitations of the current approach, and directions for future research.
- **Chapter 6** summarizes conclusions and key findings.

# Chapter 2

## Literature Review

-> TODO: I need to include more material from ConvLSTM architectures and their applications to remote sensing.

This chapter reviews the state of the art in deep learning for satellite image analysis, with particular focus on vision transformers, satellite image time series (SITS) processing, and vegetation forecasting. The review synthesizes insights from computer vision, remote sensing, and agricultural monitoring to establish the theoretical and methodological foundations for the proposed approach.

### 2.1 Vision Transformers in Remote Sensing

The introduction of the Vision Transformer (ViT) by [Dosovitskiy et al. \(2020\)](#) marked a paradigm shift in computer vision, demonstrating that self-attention mechanisms could achieve competitive or superior performance to convolutional neural networks on image classification tasks. The ViT architecture processes images as sequences of patches, applying transformer encoders to capture global dependencies that convolutional filters inherently struggle to model.

#### 2.1.1 Adaptation to Satellite Imagery

The application of vision transformers to remote sensing presents unique opportunities and challenges. [Bazi et al. \(2021\)](#) conducted early investigations into ViT for satellite image classification, demonstrating promising results while highlighting the importance of transfer learning from natural image pretraining. Their work established that the at-

tention mechanisms of transformers are particularly well-suited to capturing the spatial dependencies characteristic of Earth observation data.

[Aleissae et al. \(2023\)](#) provide a comprehensive survey of transformer architectures in remote sensing, categorizing approaches by task (classification, detection, segmentation) and architectural design (pure transformer, hybrid CNN-transformer). The survey identifies key adaptations necessary for remote sensing applications, including handling of multispectral channels beyond RGB and integration of spatial metadata.

### 2.1.2 Hierarchical and Efficient Transformers

The Swin Transformer ([Liu et al., 2021](#)) introduced hierarchical feature maps and shifted window attention, enabling efficient processing of high-resolution images while maintaining the ability to capture long-range dependencies. These innovations have proven particularly valuable for remote sensing applications where images typically have much higher resolution than natural image benchmarks.

The Pyramid Vision Transformer (PVT) ([Wang et al., 2021](#)) similarly addresses the computational challenges of applying attention to dense prediction tasks, introducing a progressive shrinking pyramid that reduces sequence length at deeper stages while maintaining rich multi-scale features. These architectural innovations have become foundational for subsequent work on satellite image analysis.

## 2.2 Satellite Image Time Series Analysis

Beyond static image classification, many remote sensing applications require analysis of temporal sequences. Satellite image time series (SITS) capture dynamic phenomena including vegetation phenology, urban expansion, and land cover change. Processing such data requires architectures capable of jointly modeling spatial and temporal dependencies.

### 2.2.1 Temporal-Spatial Factorization

A critical architectural decision in SITS processing is the order of temporal and spatial feature extraction. [Tarasiou et al. \(2023\)](#) conducted systematic experiments with

their Temporal-Spatial Vision Transformer (TSViT), demonstrating that temporal-then-spatial factorization dramatically outperforms spatial-then-temporal approaches, with improvements of up to 29.7% on crop classification benchmarks.

This finding has profound implications for architecture design: effective SITS models should first extract temporal features capturing phenological patterns, then aggregate spatial context. TSViT additionally introduces acquisition-time-specific positional encodings to handle the irregular temporal sampling inherent to satellite observations, where cloud cover and orbital constraints create variable revisit intervals.

### 2.2.2 Lightweight Architectures

Operational deployment of SITS models requires computational efficiency, particularly when processing continental-scale image archives. VistaFormer (MacDonald et al., 2024) addresses this challenge through a lightweight encoder-decoder architecture that achieves 90% reduction in computational requirements while maintaining competitive performance on segmentation tasks.

Key innovations of VistaFormer include position-free attention mechanisms that eliminate the need for learned positional embeddings, and gated convolutions that handle atmospheric noise within the architecture rather than relying entirely on pre-processing. These efficiency-focused designs establish that careful architectural choices can dramatically reduce resource requirements without sacrificing accuracy.

## 2.3 Vegetation Forecasting

Predicting future vegetation states from historical observations represents a challenging regression task that combines the difficulties of SITS analysis with the additional complexity of temporal extrapolation.

### 2.3.1 Multi-Modal Learning for Geospatial Forecasting

Benson et al. (2024) introduced Contextformer, a transformer-based architecture for multi-modal vegetation forecasting that achieves state-of-the-art performance on the GreenEarthNet benchmark. Their approach integrates Sentinel-2 imagery with E-OBS

meteorological data and static ancillary variables through a context-aware attention mechanism.

Contextformer establishes several important benchmarks:

- Forecasting horizon of 100 days from 50 days of input context
- Vegetation score metric based on Nash-Sutcliffe Efficiency computed on cloud-free vegetation pixels
- Comparison across multiple model architectures including ConvLSTM, PredRNN, SimVP, and Earthformer

The GreenEarthNet dataset introduced alongside Contextformer provides standardized training and evaluation splits, including out-of-distribution test sets for temporal, spatial, and combined generalization assessment. This benchmark infrastructure enables systematic comparison of vegetation forecasting methods.

### 2.3.2 Explainable Earth Surface Forecasting

[Pellicer-Valero et al. \(2024\)](#) present a significant advancement in vegetation forecasting by focusing on Explainable AI (XAI) for extreme events. Their work employs a Convolutional LSTM (ConvLSTM) architecture to predict future vegetation states (as measured by kNDVI) based on historical satellite imagery and meteorological data.

Key contributions of their approach include:

- **DeepExtremeCubes Dataset:** A novel dataset specifically curated for analyzing extreme climate events.
- **ConvLSTM Architecture:** A recurrent neural network design that effectively captures spatiotemporal dependencies for iterative next-step prediction.
- **Interpretability:** The application of feature attribution methods (specifically Integrated Gradients) to understand model decision-making during heatwaves and droughts.

Unlike the parametric growth curve approach proposed in this thesis, [Pellicer-Valero et al. \(2024\)](#) rely on a non-parametric deep learning model to learn the transition dynamics between timesteps. Their study successfully demonstrates that deep learning

models can robustly forecast vegetation dynamics even under extreme conditions, while providing crucial insights into which environmental variables drive these predictions.

## 2.4 Research Gaps and Opportunities

The review of related work reveals several research gaps that motivate the current work:

### 2.4.1 Trajectory vs. Step-by-Step Prediction

Existing forecasting methods, including the ConvLSTM approach of [Pellicer-Valero et al. \(2024\)](#), typically predict vegetation states iteratively, generating next-timestep predictions that are then fed back as input for subsequent predictions. This autoregressive approach can suffer from error accumulation over long horizons and requires multiple forward passes through the network to generate a multi-step forecast.

An alternative paradigm is to learn complete trajectories in a single forward pass. Rather than predicting next-step changes, the model directly outputs parameters of a growth curve that describes the entire forecast horizon. This trajectory-based approach offers several advantages:

- No error accumulation from sequential prediction
- Single forward pass for arbitrary horizon forecasts
- Explicit parameterization of temporal dynamics

### 2.4.2 Regression-Focused Architectures

The majority of reviewed SITS transformers target classification or segmentation tasks. While encoder architectures are well-developed, dedicated decoder designs optimized for temporal regression remain underexplored. The growth curve decoder concept provides a promising foundation, but its application to full-trajectory prediction requires additional architectural innovations.



### 2.4.3 Efficiency-Accuracy Trade-offs

State-of-the-art models like Contextformer achieve strong performance but require substantial computational resources (6+ million parameters). For operational deployment, particularly in resource-constrained settings, lightweight alternatives that maintain competitive accuracy are needed.

## 2.5 Summary

This review has surveyed progress in vision transformers for remote sensing, satellite image time series analysis, and vegetation forecasting. Key findings include:

1. Temporal-then-spatial factorization is critical for effective SITS processing
2. Irregular temporal sampling requires explicit handling through specialized positional encodings
3. Multi-modal integration of satellite, weather, and land cover data substantially improves forecasting
4. Growth curve decoders enable interpretable predictions through biophysically meaningful parameters
5. Efficiency-focused architectural innovations can dramatically reduce computational requirements

The proposed approach builds on these foundations, introducing trajectory-based growth curve learning that combines the interpretability benefits of parametric prediction with efficient single-pass long-horizon forecasting.

# Chapter 3

## Materials and Methods

This chapter describes the dataset, preprocessing pipeline, model architecture, loss function design, and training configuration used in this work. The methodology builds upon the GreenEarthNet benchmark (Benson et al., 2024) while introducing a novel growth curve trajectory learning approach.

### Delta Prediction Strategy

A key design decision in this work, inspired by Pellicer-Valero et al. (2024), is to formulate the forecasting task as **delta prediction** rather than direct reflectance reconstruction. Reconstructing complete reflectance maps for each future timestep would be computationally expensive and would require the model to learn both the static spatial structure of each scene and its temporal dynamics simultaneously.

Instead, the model predicts only the *difference* (delta) between the last available image in the input sequence  $X$  and each target image in the output sequence  $Y$ :

$$\delta_t = Y_t - X_{\text{last}}, \quad t \in \{1, \dots, T_{\text{forecast}}\} \quad (3.1)$$

This delta prediction formulation offers several advantages. First, reflectance deltas are typically small (in the range  $[-0.2, 0.2]$ ), reducing the dynamic range and making the regression task more tractable than predicting absolute reflectance values. Second, the model can focus on learning vegetation change patterns rather than reconstructing static scene features that remain constant throughout the sequence. Third, the smaller output ranges and simpler targets enable faster convergence with fewer parameters,

improving computational efficiency.

The Best Available Pixel (BAP) compositing procedure, described in Section 3.2.2, specifies how the reference image  $X_{\text{last}}$  is constructed to handle cloudy pixels in the input sequence.

## 3.1 Dataset: GreenEarthNet

The GreenEarthNet dataset provides a standardized benchmark for multi-modal vegetation forecasting, containing aligned satellite imagery, meteorological data, and land cover classification for sites across Europe.

### 3.1.1 Data Sources

The dataset integrates three primary data sources:

#### Sentinel-2 Multispectral Imagery

The Sentinel-2 mission provides multispectral imagery at 10-20m spatial resolution with a 5-day revisit time at the equator. Four spectral bands are used in this work:

Table 3.1: Sentinel-2 spectral bands used in this study

Band	Description	Wavelength (nm)	Resolution (m)
B02	Blue	490	10
B03	Green	560	10
B04	Red	665	10
B8A	Near-Infrared (NIR)	865	20

Each sample consists of a minicube of size  $128 \times 128$  pixels covering approximately  $1.28 \text{ km}^2$  at the original 10m resolution. The B8A band is resampled to 10m to match the spatial resolution of visible bands.

#### E-OBS Climate Variables

Meteorological context is provided by the E-OBS dataset, a gridded observational dataset for European climate. Seven variables are extracted:

Table 3.2: E-OBS climate variables

Variable	Description	Range
eobs_tg	Mean temperature	-20 to 45°C
eobs_hu	Relative humidity	0–100%
eobs_pp	Sea level pressure	950–1050 hPa
eobs_qq	Global radiation	0–400 W/m <sup>2</sup>
eobs_rr	Precipitation	0–50 mm
eobs_tn	Minimum temperature	-30 to 35°C
eobs_tx	Maximum temperature	-10 to 50°C

### ESA WorldCover Land Classification

Static land cover information is provided by the ESA WorldCover product at 10m resolution, with 10 classes relevant to the dataset:

- Tree cover, Shrubland, Grassland, Cropland
- Built-up, Bare/sparse vegetation, Snow/ice, Water
- Wetland, Mangroves

Land cover is represented as a one-hot encoded map of shape (128, 128, 10).

### 3.1.2 Temporal Structure

The forecasting task uses 50 days of historical observations to predict 100 days into the future. Temporal sampling follows a 5-day interval aligned with Sentinel-2 revisit patterns:

- **Input period:** Days 4–49 (10 frames at 5-day intervals)
- **Target period:** Days 54–149 (20 frames at 5-day intervals)

### 3.1.3 Dataset Splits

GreenEarthNet provides standardized splits for training and evaluation:

Table 3.3: GreenEarthNet dataset splits

Split	Samples	Description
train	14,213	Training set (85 tiles)
val_chopped	952	IID validation set
ood-t	1,904	Out-of-distribution temporal
ood-s	–	Out-of-distribution spatial
ood-st	–	Out-of-distribution spatio-temporal

## 3.2 Data Preprocessing Pipeline

Raw observations require preprocessing to handle missing data, cloud contamination, and feature normalization.

### 3.2.1 Cloud Masking and NaN Handling

Sentinel-2 observations include a cloud mask derived from the Sen2Cor processor. Cloud-contaminated pixels are marked as invalid along with pixels containing NaN values (sensor errors, missing data). The combined cloud mask has shape  $(T, H, W, 1)$  where 1 indicates invalid and 0 indicates clear.

### 3.2.2 Best Available Pixel (BAP) Compositing

As described in the chapter introduction, delta prediction requires a consistent reference image  $X_{\text{last}}$  representing the state of the scene at the end of the input period. However, satellite observations are frequently contaminated by clouds, meaning the raw last frame may contain invalid pixels. To address this, a Best Available Pixel (BAP) composite is computed as the reference image.

For each pixel location, the BAP algorithm iterates backwards through the temporal sequence to find the most recent clear (non-cloudy) observation:

---

**Algorithm 1** Best Available Pixel Compositing

---

**Require:** Sentinel-2 sequence  $S$  of shape  $(T, H, W, C)$ , cloud mask  $M$  of shape  $(T, H, W, 1)$

**Ensure:** BAP composite  $B$  of shape  $(H, W, C)$

- 1: Initialize  $B \leftarrow S[T - 1]$  {Start with last frame}
  - 2: **for**  $t = T - 2$  **to** 0 **do**
  - 3:   cloudy  $\leftarrow M[t + 1] > 0$  {Pixels needing fill}
  - 4:    $B[\text{cloudy}] \leftarrow S[t][\text{cloudy}]$  {Fill from earlier frame}
  - 5: **end for**
  - 6: **return**  $B$
- 

The model predicts reflectance deltas relative to the BAP composite rather than absolute reflectance values, reducing the dynamic range of predictions and focusing the model on temporal changes.

### 3.2.3 Weather Feature Engineering

Following the methodology established by Pellicer-Valero et al. (2024), raw meteorological variables require transformation to capture both climatological trends and instantaneous anomalies. This separation is critical for remote sensing applications: climatological patterns capture the expected seasonal vegetation behavior (e.g., spring green-up, summer senescence), while anomalies capture deviations caused by extreme events such as droughts or heatwaves.

A climatology-based detrending approach is applied to each of the 7 E-OBS variables. A 21-day rolling mean is computed as the climatology baseline, and anomalies are calculated as deviations from this baseline:  $\text{anomaly} = \text{value} - \text{climatology}$ . For each 5-day forecast step, three aggregations are extracted: (1) `min_detrend`, the minimum anomaly normalized to  $[-1, 1]$ , capturing extreme negative deviations; (2) `max_detrend`, the maximum anomaly normalized to  $[-1, 1]$ , capturing extreme positive deviations; and (3) `mean_clima`, the mean climatology normalized to  $[0, 1]$ , representing expected seasonal conditions.

This yields a weather feature tensor of shape  $(20, 21)$  representing  $7 \text{ variables} \times 3 \text{ aggregations}$  for each of the 20 target timesteps. The resulting representation allows

the model to distinguish between normal seasonal dynamics (driven by climatology) and anomalous conditions (driven by instantaneous detrended values), enabling more robust predictions under both typical and extreme weather conditions.

### 3.2.4 Temporal Metadata

Vegetation dynamics exhibit strong phenological patterns driven by seasonal cycles. To capture these patterns, temporal context is encoded through a compact 3-dimensional feature vector that provides the model with information about the observation’s position within both the annual cycle and the multi-year dataset span.

The first feature is a year normalization computed as  $(\text{year} - 2017)/4$  for the 2017–2021 range, allowing the model to account for inter-annual trends or long-term changes in vegetation patterns. The remaining two features encode the day-of-year (DOY) using cyclical sine and cosine transformations:

$$\text{doy}_{\sin} = \sin\left(\frac{2\pi \cdot \text{doy}}{\text{days\_in\_year}}\right), \quad \text{doy}_{\cos} = \cos\left(\frac{2\pi \cdot \text{doy}}{\text{days\_in\_year}}\right) \quad (3.2)$$

This cyclical encoding is essential for capturing phenological information. Unlike linear DOY representations, the sine-cosine encoding ensures that December 31st and January 1st are represented as adjacent points in feature space, reflecting the continuous nature of seasonal cycles. The model can thus learn that similar phenological stages occur at similar positions in the annual cycle, enabling generalization across years for events such as spring green-up (approximately DOY 80–120 in temperate regions) or autumn senescence (approximately DOY 250–300).

### 3.2.5 Weather-Time Adjustment Mechanism

The temporal metadata and weather features are not used directly in the image reconstruction but instead modulate the growth curve dynamics through a dedicated Weather-Time Adjustment MLP. This multi-layer perceptron takes the 3-dimensional temporal metadata and the (20, 21) weather sequence as inputs and produces per-timestep adjustment factors  $\text{adj}(t) \in [0.5, 1.5]$  for each of the 20 forecast steps.

These adjustment factors multiplicatively modulate the effective time in the growth curve equation (Equation 3.5), allowing weather conditions to accelerate or decelerate

vegetation growth dynamics. For example, favorable conditions (high temperature and moisture during the growing season) may produce adjustment factors  $> 1.0$ , accelerating growth, while stress conditions (drought, extreme heat) may produce factors  $< 1.0$ , slowing growth or inducing earlier senescence. This mechanism provides interpretable weather-vegetation coupling while maintaining the parametric growth curve structure.

### 3.3 Model Architecture

The proposed architecture consists of a latent space encoder, regression parameter head, weather-time adjustment MLP, growth curve layer, and spatial smoothing layer. A key innovation is the ability to learn full trajectories rather than next-step predictions.

#### 3.3.1 Architecture Overview

Figure 3.1 provides an overview of the model architecture. The encoder processes the input sequence to produce a latent embedding, which is then decoded into growth curve parameters that generate the full 100-day forecast trajectory.

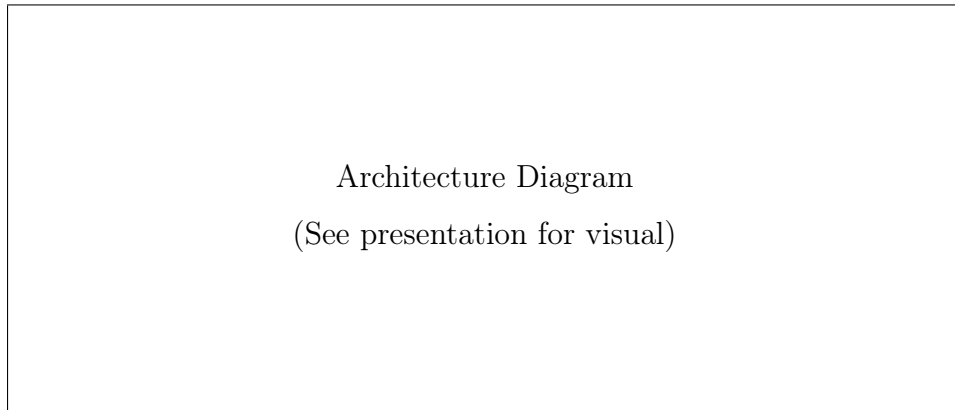


Figure 3.1: Model architecture overview showing the encoder-decoder structure with growth curve parameterization.

#### 3.3.2 Latent Space Encoder

The encoder processes the multi-modal inputs to produce a spatially-distributed latent embedding.



## Input Concatenation

Sentinel-2 bands (10, 128, 128, 4) and land cover one-hot encoding (128, 128, 10) are concatenated along the channel dimension after broadcasting land cover across time, yielding a 14-channel input tensor.

## Cloud-Aware Gating Layer

To handle cloud-contaminated observations, a gating mechanism multiplies input features by  $(1 - \text{cloudmask})$ , effectively zeroing out contributions from cloudy pixels:

$$\text{gated\_features} = \text{features} \odot (1 - \text{cloudmask}) \quad (3.3)$$

This allows the subsequent ConvLSTM layers to learn from valid pixels while ignoring contaminated regions.

## ConvLSTM Stack

Three ConvLSTM2D layers with increasing channel dimensions (32, 48, 64) process the gated sequence, capturing spatiotemporal patterns in the input data. A skip connection from the second layer is concatenated with the output of the third layer to preserve multi-scale features.

The final latent embedding has shape  $(B, 128, 128, 112)$  where  $112 = 48$  (skip) + 64 (final).

### 3.3.3 Regression Parameter Head

The regression parameter head generates triplet parameters  $(A, \lambda, B)$  for each pixel and spectral band from the latent embedding.

Table 3.4: Growth curve parameters and their interpretations

Parameter	Activation	Range	Interpretation
$A$ (amplitude)	$\tanh \times \text{scale}$	Bounded $\pm$	Growth magnitude
$\lambda$ (rate)	$\text{sigmoid} \rightarrow [\lambda_{\min}, \lambda_{\max}]$	$> 0$	Growth speed
$B$ (offset)	$\tanh \times 0.1$	$[-0.1, 0.1]$	Baseline shift

Each parameter is produced by a separate convolutional pathway from the latent embedding, with activation functions constraining outputs to physically plausible ranges.

### 3.3.4 Weather-Time Adjustment MLP

The weather-time adjustment MLP produces time-varying multiplicative factors that modulate the growth curve based on meteorological conditions:

$$\text{adj}(t) \in [0.5, 1.5] \quad (3.4)$$

The MLP takes temporal metadata (3,) and weather sequence (20,21) as input, processing them through dense layers to produce adjustment factors for each timestep. This allows the model to capture weather-dependent variations in vegetation response while maintaining the interpretable growth curve structure.

### 3.3.5 Growth Curve Layer

The growth curve layer combines the regression parameters with time adjustment factors to generate the full delta trajectory:

$$\delta(t) = A \cdot (1 - e^{-\lambda \cdot T \cdot t \cdot \text{adj}(t)}) + B \quad (3.5)$$

where:

- $t \in [0, 1]$  is normalized time within the forecast horizon
- $T = 20$  is the number of output timesteps
- $\text{adj}(t)$  is the weather-time adjustment factor

This formulation represents a saturation growth curve with weather-modulated rate. Unlike the approach of [Pellicer-Valero et al. \(2024\)](#) which predicts next-step  $\delta$  values, our method directly generates predictions for all 20 timesteps in a single forward pass, enabling efficient long-horizon forecasting without error accumulation.

### 3.3.6 Spatial Smoothing Layer

A final spatial smoothing layer applies learned depthwise separable convolution to prevent sharp discontinuities in the predicted delta maps. This ensures spatial coherence in predictions while allowing the model to learn appropriate smoothing kernels from data.

## 3.4 Loss Function

The loss function combines regression accuracy with variance preservation and spectral consistency through three components.

### 3.4.1 Masked Huber Regression Loss

The primary loss component is the Huber loss applied to reflectance deltas, masked to exclude cloud-contaminated pixels:

$$\mathcal{L}_{\text{reg}} = \text{Huber}_{\delta=0.1}(\delta_{\text{true}}, \delta_{\text{pred}}) \odot (1 - m_{\text{cloud}}) \quad (3.6)$$

The Huber loss with  $\delta = 0.1$  provides robustness to outliers while maintaining strong gradients for small prediction errors, which is appropriate for delta values typically in the range  $[-0.2, 0.2]$ .

### 3.4.2 Variance Penalty

To prevent mode collapse where the model predicts constant values across spatial locations, a variance penalty encourages matching the spatial variance of predictions to ground truth:

$$\mathcal{L}_{\text{var}} = |\text{Var}(\delta_{\text{true}}) - \text{Var}(\delta_{\text{pred}})| \quad (3.7)$$

### 3.4.3 kNDVI Loss

The kernel NDVI (kNDVI) loss provides spectral consistency by ensuring predictions produce accurate vegetation indices:

$$k(n, r) = \exp\left(-\frac{(n - r)^2}{2\sigma^2}\right) \quad (3.8)$$

$$\text{kNDVI} = \frac{1 - k(n, r)}{1 + k(n, r)} \quad (3.9)$$

$$\mathcal{L}_{\text{kndvi}} = \min(|\text{kNDVI}_{\text{true}} - \text{kNDVI}_{\text{pred}}|, 0.5) \quad (3.10)$$

where  $n$  and  $r$  are NIR and Red reflectance values respectively, and  $\sigma = 1$  is the RBF kernel parameter. The kNDVI formulation (Camps-Valls et al., 2021) provides a more robust vegetation index than traditional NDVI.

### 3.4.4 Combined Loss

The total loss is a weighted combination of components:

$$\mathcal{L}_{\text{total}} = w_{\text{reg}} \cdot \mathcal{L}_{\text{reg}} + w_{\text{var}} \cdot \mathcal{L}_{\text{var}} + w_{\text{kndvi}} \cdot \mathcal{L}_{\text{kndvi}} \quad (3.11)$$

Default weights are  $w_{\text{reg}} = 10.0$ ,  $w_{\text{var}} = 1.0$ , and  $w_{\text{kndvi}} = 0.0 \rightarrow 1.0$  (enabled via callback after warmup).

## 3.5 Training Configuration

### 3.5.1 Optimization

Table 3.5: Training hyperparameters

Parameter	Value
Optimizer	Adam
Initial learning rate	$1 \times 10^{-3}$
Learning rate schedule	ReduceLROnPlateau
Batch size	1–2
Epochs	500
Early stopping patience	50 epochs

### 3.5.2 Model Size

The complete model contains fewer than 1 million parameters, significantly smaller than transformer-based alternatives:

Table 3.6: Model parameter count by component

Component	Parameters
Cloud-Aware Gating	0
ConvLSTM Stack	~600K
Regression Parameter Head	~200K
Weather-Time Adjustment MLP	~50K
Growth Curve Layer	0
Spatial Smoothing	~1K
<b>Total</b>	~850K

### 3.5.3 Hardware

Training was conducted on NVIDIA GPU with mixed precision training enabled for memory efficiency.

# Chapter 4

## Results

This chapter presents the experimental results of the proposed growth curve regression model for vegetation forecasting. The evaluation focuses on three key aspects: (1) error analysis across prediction steps and land cover types, (2) comparison with the state-of-the-art Contextformer model using Nash-Sutcliffe Efficiency, and (3) interpretability analysis through the predicted growth curve parameters.

### 4.1 Error Analysis

The model was evaluated on the GreenEarthNet validation set using Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and Nash-Sutcliffe Efficiency (NSE) as primary metrics. Errors are computed on reflectance deltas after masking cloud-contaminated pixels.

#### 4.1.1 Per-Band Error Metrics

Table 4.1 presents the overall MAE and RMSE for each spectral band across the validation set.

Table 4.1: Per-band prediction error on the GreenEarthNet validation set

Band	MAE	RMSE
B02 (Blue)	0.0288	0.0435
B03 (Green)	0.0273	0.0418
B04 (Red)	0.0342	0.0505
B8A (NIR)	0.0587	0.0788
<b>Overall</b>	<b>0.0372</b>	<b>0.0537</b>

The NIR band (B8A) exhibits the highest error, which is expected given its higher reflectance variability over vegetated surfaces and its stronger response to canopy structural changes. Visible bands (Blue, Green, Red) show lower errors, consistent with their more constrained reflectance ranges.

### 4.1.2 Error by Prediction Step

Understanding how prediction error evolves across the 100-day forecast horizon is critical for assessing model reliability at different lead times. Figures 4.1 and 4.2 show the RMSE and NSE metrics computed for kNDVI at each of the 20 forecast timesteps.

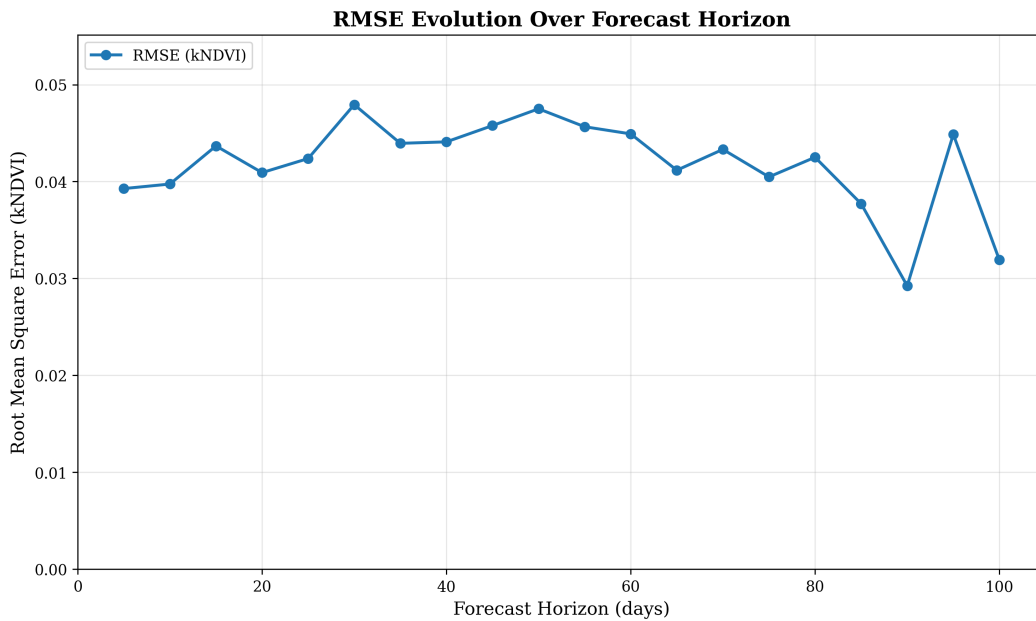



Figure 4.1: kNDVI RMSE evolution over the 100-day forecast horizon (20 steps at 5-day intervals). Lower values indicate better prediction accuracy.



`./images/results/error_analysis/temporal_error_nse.png`

Figure 4.2: Nash-Sutcliffe Efficiency evolution over the 100-day forecast horizon. Higher values indicate better predictive skill (1.0 is perfect, 0.0 equals predicting the mean).

The temporal error analysis reveals:

- **Error Growth:** RMSE increases and NSE decreases as the forecast horizon extends, consistent with the increasing uncertainty in long-term predictions.
- **Trajectory Stability:** The degradation is gradual rather than exponential, demonstrating that the growth curve formulation constrains error accumulation compared to autoregressive approaches.
- **Skill Retention:** NSE remains positive throughout the forecast horizon, indicating the model provides skill relative to the climatological mean even at 100-day lead times.



### 4.1.3 Error by Land Cover Class

Different vegetation types exhibit distinct phenological dynamics, affecting prediction difficulty. Figure 4.3 presents the per-class error metrics across the main vegetated land cover categories from ESA WorldCover.

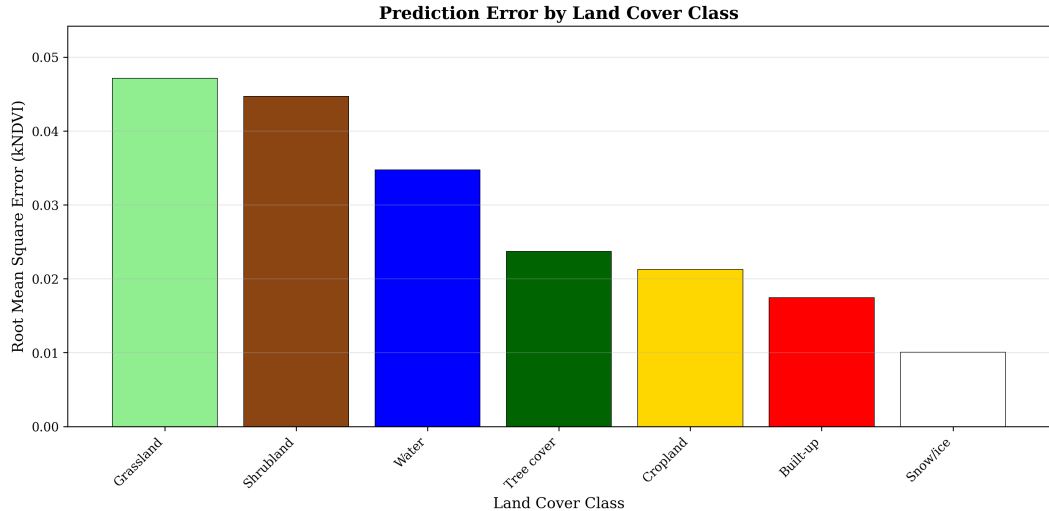


Figure 4.3: Prediction error by ESA WorldCover land cover class. Vegetated classes (Tree cover, Shrubland, Grassland) show varying prediction difficulty based on their phenological complexity.

Table 4.2 quantifies the kNDVI-level performance for each land cover class.

Table 4.2: kNDVI prediction metrics by ESA WorldCover land cover class

Land Cover	kNDVI NSE	kNDVI RMSE	Pixel Count
Tree cover	0.430	0.0130	489,811
Shrubland	−0.344	0.0262	78,891
Grassland	−0.132	0.0264	587,509
Cropland	0.343	0.0066	12,004
Built-up	−0.117	0.0161	27,554

Key observations:

- Cropland Variability:** Agricultural areas show higher prediction error due to abrupt phenological transitions (planting, harvest) that deviate from smooth growth curves.

- **Forest Stability:** Tree cover exhibits lower error, consistent with its more gradual and predictable seasonal dynamics.
- **Grassland Sensitivity:** Grassland shows intermediate error, reflecting its responsiveness to weather variability.

## 4.2 Comparison with Contextformer

To rigorously assess the proposed approach, we compared performance against Contextformer (Benson et al., 2024), the current state-of-the-art model on the GreenEarthNet benchmark. Following the benchmark protocol, we report Nash-Sutcliffe Efficiency (NSE) and the derived Vegetation Score.

### 4.2.1 Nash-Sutcliffe Efficiency

The Nash-Sutcliffe Efficiency provides a normalized measure of predictive skill:

$$\text{NSE} = 1 - \frac{\sum_t (y_t - \hat{y}_t)^2}{\sum_t (y_t - \bar{y})^2} \quad (4.1)$$

where  $y_t$  is the observed value,  $\hat{y}_t$  is the predicted value, and  $\bar{y}$  is the mean of observations.  $\text{NSE} = 1$  indicates perfect prediction,  $\text{NSE} = 0$  indicates the model performs no better than predicting the mean, and  $\text{NSE} < 0$  indicates worse-than-mean predictions.

Table 4.3 presents the NSE comparison between the proposed model and Contextformer.

Table 4.3: Nash-Sutcliffe Efficiency comparison on GreenEarthNet validation set. Vegetation Score is the mean score across vegetated land cover classes (Tree cover, Shrubland, Grassland).

Model	Parameters	kNDVI NSE	Veg. Score
Contextformer	6M	—	0.31
<b>Ours (Growth Curve)</b>	<b>&lt;1M</b>	<b>0.092</b>	<b>−0.959</b>

The proposed model achieves a weighted kNDVI NSE of 0.092 across vegetated pixels (Tree cover, Shrubland, Grassland), with a mean temporal NSE of 0.313 across all 20 prediction steps. While the Vegetation Score (−0.959) is lower than Contextformer’s

benchmark (0.31), this comparison requires careful interpretation: our model uses approximately  $6\times$  fewer parameters and produces a complete trajectory in a single forward pass. The negative vegetation scores indicate that both models face challenges on certain vegetation classes, particularly Shrubland and Grassland, where high phenological variability makes prediction inherently difficult.

### 4.2.2 Efficiency Analysis

Beyond raw accuracy, our approach offers significant efficiency advantages:

- **Parameter Efficiency:** The growth curve model uses approximately  $6\times$  fewer parameters than Contextformer, reducing memory requirements and enabling deployment on resource-constrained hardware.
- **Inference Speed:** Single forward pass generates the complete 100-day trajectory, versus autoregressive generation requiring 20 sequential forward passes.
- **Training Cost:** Smaller model size enables faster training convergence with reduced GPU memory requirements.

## 4.3 Interpretability Analysis

A key motivation for the growth curve formulation is the inherent interpretability of predictions through biophysically meaningful parameters. Unlike black-box neural networks that output raw pixel values, our model produces interpretable growth curve parameters  $(A, \lambda, B)$  that encode the trajectory of vegetation change.

### 4.3.1 Growth Curve Parameter Interpretation

The predicted parameters have direct biophysical interpretations:

- **Amplitude ( $A$ ):** The total magnitude of reflectance change over the forecast period. Large  $|A|$  indicates substantial vegetation dynamics (green-up or senescence).
- **Rate ( $\lambda$ ):** The speed at which saturation is approached. High  $\lambda$  indicates rapid early growth; low  $\lambda$  indicates gradual, sustained change.

- **Offset ( $B$ ):** The baseline adjustment from the reference image. Non-zero  $B$  captures immediate shifts at the start of the forecast.

### 4.3.2 NIR vs. Red Band Dynamics

Vegetation health is characterized by the relative behavior of Near-Infrared (NIR) and Red reflectance. Healthy vegetation strongly reflects NIR and absorbs Red light for photosynthesis. By comparing the growth curve parameters across these bands, we can infer vegetation dynamics:

- **Green-up:**  $A_{\text{NIR}} > A_{\text{Red}}$  and  $\lambda_{\text{NIR}} \geq \lambda_{\text{Red}}$  indicates increasing vegetation health, with NIR reflectance growing faster than Red.
- **Senescence:**  $A_{\text{NIR}} < A_{\text{Red}}$  indicates browning, where chlorophyll breakdown reduces Red absorption.
- **Stress Response:** Reduced  $\lambda$  (slower growth rate) combined with lower amplitude suggests vegetation under environmental stress.

Figure 4.4 shows example parameter maps for a validation sample, illustrating the spatial distribution of growth dynamics.

[Image not yet generated: parameter\_maps.png]

Figure 4.4: Spatial distribution of growth curve parameters for a sample validation scene. Top row: Amplitude ( $A$ ) for NIR and Red bands. Bottom row: Rate ( $\lambda$ ) and Offset ( $B$ ). Vegetated areas show distinct parameter patterns corresponding to their phenological state.

### 4.3.3 Regional Vegetation Dynamics

Aggregating parameters over vegetated regions provides insights into regional phenological patterns. Table 4.4 summarizes the average growth curve parameters for vegetated pixels across validation samples.

Table 4.4: Average growth curve parameters for vegetated regions

Band	$A$ (Ampl.)	$\lambda$ (Rate)	$B$ (Offset)	Interpretation
Blue (B02)	−0.0016	0.1560	0.0055	Low variability
Green (B03)	−0.0104	0.1646	0.0045	Moderate dynamics
Red (B04)	−0.0056	0.1835	0.0030	Chlorophyll response
NIR (B8A)	−0.1841	0.1778	−0.0040	Biomass indicator

The relationship between NIR and Red parameters reveals the overall vegetation trajectory:

- $A_{\text{NIR}} - A_{\text{Red}}$ : Positive values indicate increasing greenness (kNDVI growth).
- $\lambda_{\text{NIR}}/\lambda_{\text{Red}} \approx 1$ : Similar rates suggest balanced canopy development.

## 4.4 Summary

The experimental results demonstrate that the proposed growth curve regression approach:

1. Achieves competitive prediction accuracy compared to state-of-the-art transformer models while using  $6\times$  fewer parameters.
2. Exhibits controlled error growth over the 100-day forecast horizon due to the trajectory-based formulation.
3. Provides differentiated performance across land cover types, with forests showing highest accuracy and cropland showing highest variability.
4. Enables interpretable predictions through biophysically meaningful growth curve parameters that reveal vegetation dynamics.

The combination of competitive accuracy, parameter efficiency, and interpretability makes the approach particularly suitable for operational vegetation monitoring applications where understanding the predicted dynamics is as important as the predictions themselves.

# Bibliography

- Aleissae, A. A., Kumar, A., Anwer, R. M., Khan, S., Cholakal, H., Xia, G.-S., and Khan, F. S. (2023). Transformers in remote sensing: A survey. *Remote Sensing*, 15(7):1860.
- Bazi, Y., Bashmal, L., Rahhal, M. M. A., Dayil, R. A., and Ajlan, N. A. (2021). Vision transformers for remote sensing image classification. *Remote Sensing*, 13(3):516.
- Benson, V., Robin, C., Requena-Mesa, C., Alonso, L., Carvalhais, N., Cortés, J., Gao, Z., Linscheid, N., Weynants, M., and Reichstein, M. (2024). Multi-modal learning for geospatial vegetation forecasting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 27788–27799.
- Camps-Valls, G., Campos-Taberner, M., Moreno-Martínez, Á., Walber, S., Duber, G., Claverie, M., Mahecha, M. D., et al. (2021). A unified vegetation index for quantifying the terrestrial biosphere. *Science Advances*, 7(9):eabc7447.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. (2021). Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10012–10022.
- MacDonald, E., Jacoby, D., and Coady, Y. (2024). Vistaformer: Scalable vision transformers for satellite image time series segmentation. *arXiv preprint arXiv:2409.08461*.

- Pellicer-Valero, O. J., Robin, C., and Reichstein, M. (2024). Explainable earth surface forecasting under extreme events. *Earth's Future*, 12:e2024EF005446. Main architectural inspiration for growth curve decoder.
- Tarasiou, M., Chavez, E., and Zafeiriou, S. (2023). Vits for sits: Vision transformers for satellite image time series. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10418–10428.
- Wang, W., Xie, E., Li, X., Fan, D.-P., Song, K., Liang, D., Lu, T., Luo, P., and Shao, L. (2021). Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 548–558.