

# GeoFormer: A Vision and Sequence Transformer-based Approach for Greenhouse Gas Monitoring

Madhav Khirwar  $\diamond$  Ankur Narang  $\spadesuit$   
 $\diamond$ madhavkhirwar49@gmail.com  $\spadesuit$ ankur.narang@fermionai.com

## Abstract

Air pollution represents a pivotal environmental challenge globally, playing a major role in climate change via greenhouse gas emissions and negatively affecting the health of billions. However predicting the spatial and temporal patterns of pollutants remains challenging. The scarcity of ground-based monitoring facilities and the dependency of air pollution modeling on comprehensive datasets, often inaccessible for numerous areas, complicate this issue. In this work, we introduce GeoFormer, a compact model that combines a vision transformer module with a highly efficient time-series transformer module to predict surface-level nitrogen dioxide (NO<sub>2</sub>) concentrations from Sentinel-5P satellite imagery. We train the proposed model to predict surface-level NO<sub>2</sub> measurements using a dataset we constructed with Sentinel-5P images of ground-level monitoring stations, and their corresponding NO<sub>2</sub> concentration readings. The proposed model attains high accuracy (MAE 5.65), demonstrating the efficacy of combining vision and time-series transformer architectures to harness satellite-derived data for enhanced GHG emission insights, proving instrumental in advancing climate change monitoring and emission regulation efforts globally.

Keywords: geo-spatial imaging, vision transformers, *ProbSparse* attention.

## Introduction

- The emission of greenhouse gases (GHGs), primarily from industrial and transportation activities, is a major contributor to the climate change crisis.
- NO<sub>2</sub> is associated with air contaminants like PM2.5 and is released alongside CO<sub>2</sub>, making it an effective indicator for CO<sub>2</sub> emissions.
- High-resolution satellite imagery, such as from the Sentinel-5P satellite, provides unprecedented monitoring opportunities for atmospheric pollutants.
- The TROPOMI device on Sentinel-5P enables detailed observation of NO<sub>2</sub> emissions on a global scale.

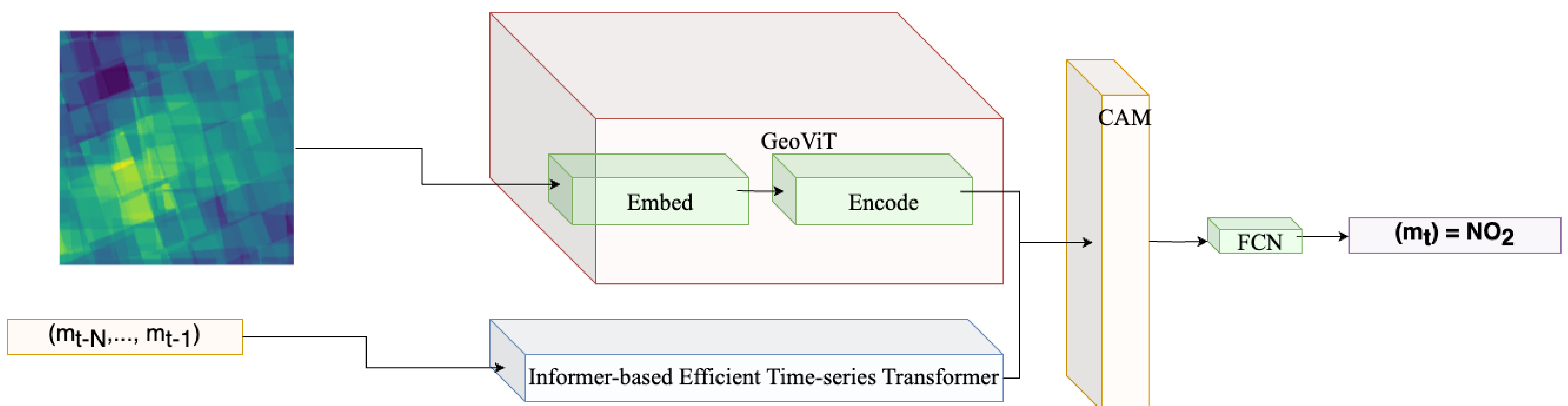
## Contributions

1. Introduced a comprehensive dataset pairing high-resolution Sentinel-5P imagery with surface-level NO<sub>2</sub> concentration measurements.
2. Proposed a novel, compact transformer-based model for efficient spatio-temporal analysis of NO<sub>2</sub> emissions.

## Methodology

### Vision Transformer Module

- Adapts images as sequences of patches for transformer processing.
- Applies linear transformation to project image patches into embeddings.



**Figure 1:** GeoFormer model architecture. Here,  $m_t$  represents an NO<sub>2</sub> prediction at timestamp  $t$ , and CAM represents the cross-attention module.

## References

- [1] L. Scheibenreif, M. Mommert, and D. Borth. Estimation of air pollution with remote sensing data: Revealing greenhouse gas emissions from space. *arXiv preprint arXiv:2108.13902*, 2021.

- Employs Multi-Head Self-Attention (MHSA) to capture spatial relationships:

$$\text{MHSA}(E) = \text{softmax} \left( \frac{EQE^TK}{\sqrt{d_k}} \right) EV \quad (1)$$

where  $E$  denotes patch embeddings, with  $Q$ ,  $K$ , and  $V$  as the query, key, and value projections, respectively.

### Efficient Time-series Transformer

- Processes sequences with a sparsity-enhanced self-attention mechanism for efficiency.
- Features a reduced time-complexity of  $O(N \log(N))$ , improving computational tractability:

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V \quad (2)$$

focusing on dominant queries for attention calculation, where  $Q$ ,  $K$ , and  $V$  are query, key, and value matrices from input data.

### Performance Summary

Model	MAE	MSE	Size (MB)
GeoViT*	5.84	58.9	850
CNN Backbone*	6.68	78.4	1964
GeoViT	6.69	72.70	65
CNN Backbone	6.49	67.25	<b>32</b>
<b>GeoFormer</b>	<b>5.65</b>	<b>56.95</b>	70

**Table 1:** Model comparison on GHG monitoring. Best Metrics are in **bold**. Models trained on the dataset introduced by [1] are marked with an asterisk (\*).

### Directions for Future Research:

- **Expansion:** Integrate GeoFormer with optical flow for enhanced sequence modeling from sequential satellite images.
- **Scalability:** Extend the model to monitor additional pollutants and environmental indicators.

GeoFormer sets a new standard for real-time and efficient GHG monitoring, opening pathways for significant advancements in climate change mitigation efforts.