# DengueNet: Dengue Prediction using Spatiotemporal Satellite Imagery for Resource-Limited Countries

**Kuan-Ting Kuo**[1] , **Dana Moukheiber**[2] , **Sebastian Cajas Ordonez**[3,4] , **David Restrepo**[2,5] , **Atika Rahman Paddo**[6] , **Tsung-Yu Chen**[1] , **Lama Moukheiber**[2] , **Mira Moukheiber**[2] , **Sulaiman Moukheiber**[7] , **Saptarshi Purkayastha**[6] , **Po-Chih Kuo**[1] and **Leo Anthony Celi**[2,3,8]

[1] National Tsing Hua Unversity, Taiwan
[2] Massachusetts Institute of Technology, USA
[3] Harvard University, USA
[4] University College Dublin, Ireland
[5] University of Cauca, Colombia
[6] Indiana University – Purdue University Indianapolis, USA
[7] Worcester Polytechnic Institute, USA
[8] Beth Israel Deaconess Medical Center, USA
{mimikuo365, lear1007}@gmail.com, {danamouk, davidres, lamam, miram, lceli}@mit.edu, apaddo@iu.edu, ulsordonez@unicauca.edu.co, swmoukheiber@wpi.edu, saptpurk@iupui.edu, kuopc@cs.nthu.edu.tw

## Abstract

Dengue fever presents a substantial challenge in developing countries where sanitation infrastructure is inadequate. The absence of comprehensive healthcare systems exacerbates the severity of dengue infections, potentially leading to life-threatening circumstances. Rapid response to dengue outbreaks is also challenging due to limited information exchange and integration. While timely dengue outbreak forecasts have the potential to prevent such outbreaks, the majority of dengue prediction studies have predominantly relied on data that impose significant burdens on individual countries for collection. In this study, our aim is to improve health equity in resource-constrained countries by exploring the effectiveness of high-resolution satellite imagery as a non-traditional and readily accessible data source. By leveraging the wealth of publicly available and easily obtainable satellite imagery, we present a scalable satellite extraction framework based on Sentinel Hub, a cloud-based computing platform. Furthermore, we introduce DengueNet[1], an innovative architecture that combines Vision Transformer, Radiomics, and Long Short-term Memory to extract and integrate spatiotemporal features from satellite images. This enables dengue predictions on an epidemiological-week basis. To evaluate the effectiveness of our proposed method, we conducted experiments on five municipalities in Colombia. We utilized a dataset comprising 780 high-resolution Sentinel-2 satellite images for training and evaluation. The performance of DengueNet was assessed using the mean absolute error (MAE) metric. Across the five municipalities, DengueNet achieved an average MAE of 43.92±42.19. Notably, the highest MAE was recorded in Cali at 113.65±0.08, whereas the lowest MAE was observed in Ibagué, amounting to 5.67±0.18. Our findings strongly support the efficacy of satellite imagery as a valuable resource for dengue prediction, particularly in informing public health policies within low- and middle-income countries. In these countries, where manually collected data of high quality is scarce and dengue virus prevalence is severe, satellite imagery can play a crucial role in improving dengue prevention and control strategies.

## 1 Introduction

Dengue, one of the most ubiquitous mosquito-borne viral infections, is the leading cause of hospitalization and death in many parts of the world, especially in tropical and subtropical countries [Cattarino *et al.*, 2020]. It is estimated that 129 countries [WHO, 2022] and 4 billion people [CDC, 2022] are at risk of dengue infection. In low- and middle-income countries (LMICs) where dengue fever is endemic, the prevalence of dengue outbreaks is exacerbated by multifarious factors such as barriers in the continuum of care, inequities in resource allocation, education levels, literacy, and income[Chaparro *et al.*, 2016]. Because there are no specific treatments available for the virus, dengue prevention is critical to reducing its infectious and fatality rate, particularly in hyperendemic regions in LMICs where dengue poses a significant public health predicament [Gutierrez-Barbosa *et al.*, 2020]. Therefore, the strategic utilization of viable early
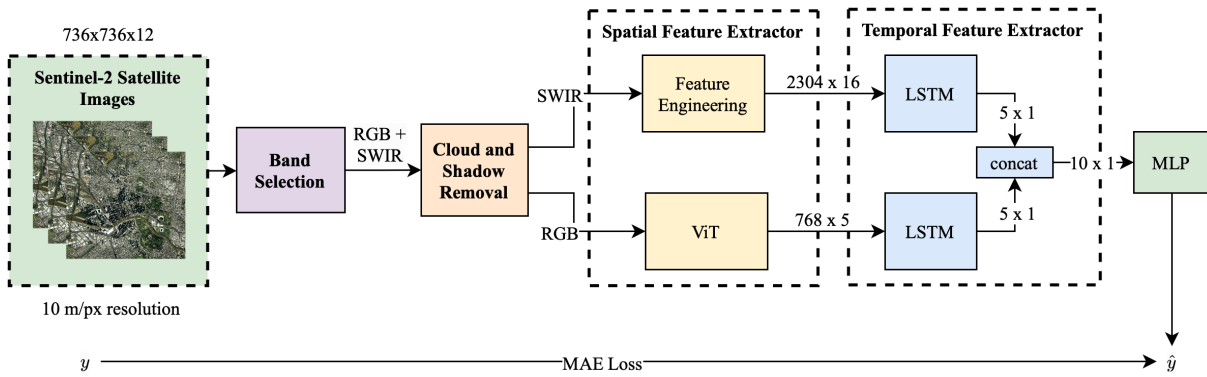
---

[1] https://github.com/mimikuo365/DengueNet-IJCAI

Figure 1: DengueNet model architecture takes in weekly satellite imagery and dengue cases $y$ as input for predicting $\hat{y}$ (m/px: meters per pixel; RGB: red, green and blue bands; SWIR: short wave infrared spectrum band; ViT: Vision Transformer; LSTM: Long Short-Term memory; MLP: Multilayer Perceptron). The LSTM module consists of three stacked standard LSTM layers.

detection approaches for dengue outbreaks in LMICs is not only imperative for promoting comprehensive well-being but also plays a crucial role in the pursuit of reducing health inequities. By employing these effective approaches, we can actively contribute to the realization of equitable healthcare access and outcomes, thereby fostering a more inclusive and just society.

Prior research has demonstrated the potential for dengue forecasting utilizing pre-collected structural information like temperature and precipitation [Martheswaran *et al.*, 2022; Jain *et al.*, 2019]. However, conventional data collection techniques are both costly and difficult to scale. Therefore, seeking alternative resources, such as publicly available satellite imagery, is significant for LMICs where structured data is scarce and critical indicators remain lacking. Remote sensing satellite imagery can be a more cost-effective and efficient approach than alternative field survey methods and has shown potential correlation with weather variables [Ren *et al.*, 2021], which are one of the key factors behind dengue outbreaks. It also enables a higher revisit frequency and diverse resolutions of imagery over time than surveys where repeated measurements at a local level are limited [Lee *et al.*, 2017]. Furthermore, the development of surveillance systems that rely exclusively on satellite imagery to notify public health authorities of early dengue detection can cost-effectively enhance the response time to national crises in hyperendemic regions in LMICs.

This study employs recent advances in machine learning (ML) and proposes an ML-based approach for forecasting the incidence of dengue cases in five municipalities of Colombia using satellite imagery. This selection was made due to Colombia's persistent incidence of high levels of reported dengue outbreaks from 1978 until 2022 [National Institute of Health of Colombia, 2010]. As one of the top five countries in the Americas with the highest number of reported dengue cases, Colombia's dengue mortality rate is 4.84 times higher than that of other American countries [PAHO, 2022]. Below are the three principal contributions to this paper.

- We introduce a scalable data collection and processing framework to extract time-series data from the Sentinel-2 satellite.

- We propose a novel preprocessing pipeline that can effectively eliminate noises and extract spatiotemporal features from the collected satellite imagery.

- Our model, DengueNet, shows positive results, indicating dengue forecasting with time-series satellite imagery alone is a feasible approach for LMICs with limited resources.

## 2 Related Works
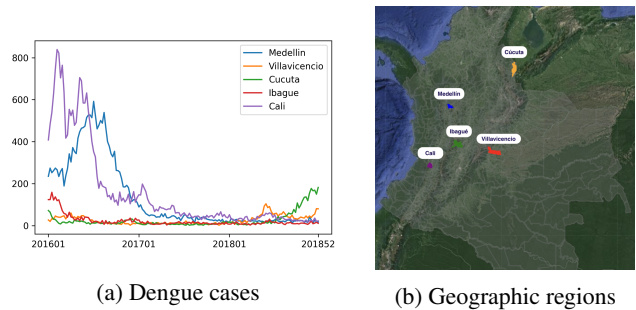


(a) Dengue cases      (b) Geographic regions

Figure 2: Municipality-level dengue case numbers and geographic locations. (a) Dengue cases from 2016 to 2018 were obtained from the SIVIGILA database for the top five affected municipalities in Colombia. (b) Geographic locations from satellite imagery for each municipality.

The epidemiology of dengue is influenced by multiple factors, including seasonal fluctuations in temperature and rainfall, socio-economic determinants such as education and household income [Morgan *et al.*, 2021; Watts *et al.*, 2020], and intra-strain genetic variability [Fontaine *et al.*, 2018]. To comprehend the determinants of dengue infection, studies have been conducted to evaluate the economic, societal, and other facets of dengue outbreaks worldwide. In terms of structured data, notable work by researchers has paired a boosted regression tree framework with longitudinal information and population surfaces to develop a risk map to understand the global distribution of dengue and improve disease
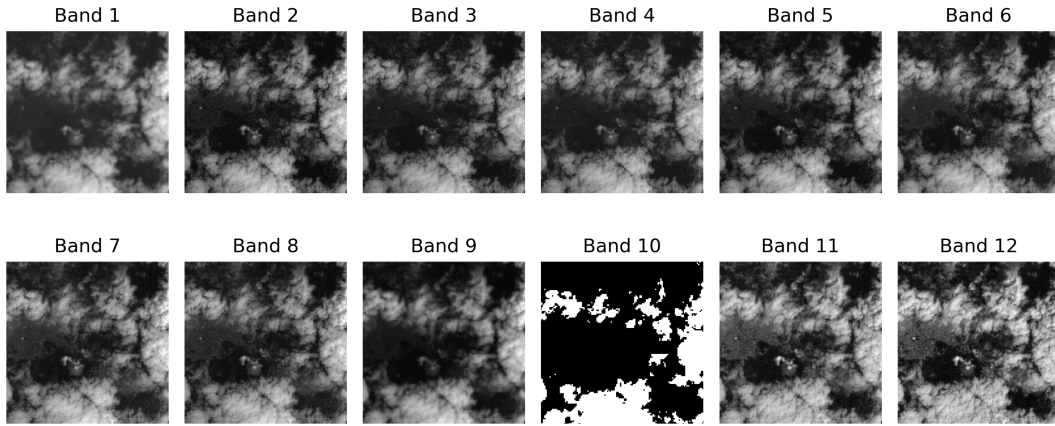
Figure 3: Gray-scale satellite band images captured by Sentinel-2 using different wavelengths.

management programs globally [Bhatt *et al.*, 2013]. Similar work has been established, which investigates the temporal and spatial distribution of dengue fever in India using Kulldorff's space-time permutation method [Mala and Jat, 2019]. Other work [Muñoz *et al.*, 2021] has also looked at the association of the local climate with dengue in Colombia using linear analysis tools and lagged crossed-correlations such as Pearson's test. Features highly associated with dengue, such as environmental, entomological, epidemiological, and human-related data, have been explored for dengue prediction [Roster and Rodrigues, 2021; Karim *et al.*, 2012; Guo *et al.*, 2017; Salim *et al.*, 2021]. Other studies have used human-related data like mobility [Datoc *et al.*, 2016], social media data [Livelo and Cheng, 2018], and distance to public transit [Shragai *et al.*, 2022] to build dengue early warning systems. In terms of unstructured data, studies compared street view and aerial images with different convolutional neural network architectures to estimate dengue rates [Andersson *et al.*, 2019].

Satellite imagery is often adopted with other statistical data to perform spatiotemporal tasks, such as weather forecasting, precipitation nowcasting [Moskolaï *et al.*, 2021; Son and Thong, 2017; de Witt *et al.*, 2020] and vector-borne disease case predictions [Rogers *et al.*, 2002; Li *et al.*, 2022a; Abdur Rehman *et al.*, 2019]. While LMICs lack access to reliable information systems for data collection and analysis [Ndabarora *et al.*, 2014; Kruk *et al.*, 2018; Fenech *et al.*, 2018], free sources of satellite imagery from cloud-based computing platforms, such as Google Earth Engine and Sentinel Hub, provide an alternative data asset for LMICs for early detection of dengue. In our work, we build a reproducible Sentinel-2 satellite data extraction framework leveraging Sentinel Hub and provide municipality-level predictions of dengue cases in Colombia per epi week. By solely adopting satellite imagery for dengue outbreak prediction, our model can focus on learning potential environmental information through difference in vegetation over time using time-series images to predict dengue cases [Moskolaï *et al.*, 2021].

## 3 Dataset

In this study, we collect satellite imagery and dengue incidences from 2016 to 2018 in five Colombian municipalities including Medellín, Ibagué, Cali, Villavicencio, and Cúcuta (Figure 2). These municipalities are chosen as they have reported relatively high dengue cases in Colombia. Sentinel Hub [Ltd, 2022] is used to collect and process Sentinel-2 satellite data. The regions of interest are pre-determined using the different municipalities' latitude and longitude square coordinates. Each area is sampled per epi week from Sentinel-2's launch date to the time frame before COVID-19, to create a time-series satellite imagery dataset. We focus on data before COVID-19, as studies show that COVID-19 has impacted dengue transmission [Lim *et al.*, 2020]. Our data is stored in a TIFF format and contains 12 bands from Sentinel-2 as shown in Figure 3. To account for differences in band resolution, we use nearest-neighbor interpolation to increase the resolution of all bands to a uniform 10 meters per pixel. Cloud inteferences are avoided using the LeastCC algorithm, which is configured using Sentinel Hub API to request the images with the least amount of clouds per epi week. We obtain weekly dengue incidences from the Colombian Public Health System (SIVIGILA). Satellite imagery is matched with dengue cases on an epi-week basis.

## 4 Methodology

### 4.1 Overview

To fully examine whether satellite imagery could be used to predict dengue cases, we introduce multiple modules in DengueNet (see Figure 1). The model components are designed to capture both the temporal and spatial information from satellite images for dengue outbreak forcasting. First, we conduct band correlation analysis to determine which satellite bands to select and use in our study. We then apply cloud and cloud shadow (CCS) removal on the selected bands to reduce noises in the satellite images. The preprocessed bands are then fed into two spatial feature extraction modules, the Feature-Engineering and the Vision-Transformer (ViT) feature extractors, respectively. The features extracted from
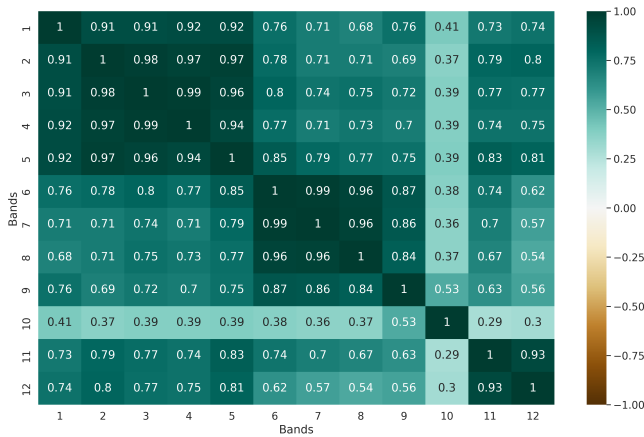
Figure 4: Average Pearson's correlation of the 12 bands for the Sentinel-2 satellite images across five Colombian municipalities in the training set from 2016 to 2018. The majority of correlations are statistically significant (p <0.001).
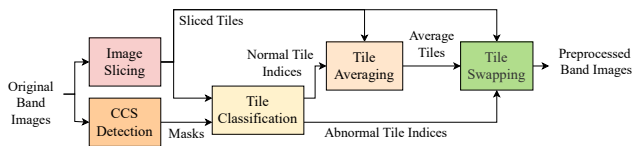


Figure 5: Stages involved in the cloud and cloud shadow removal module. The average tiles are generated using the normal tiles in the samples (CCS: cloud and cloud shadow).

the two modules are then fed into two multi-layer Long Short-term Memory (LSTM) networks that can extract temporal features, and eventually concatenated to a fully connected neural network for dengue case prediction.

## 4.2 Band Selection

Satellite imagery often contains multiple bands with different resolutions, central wavelengths, and channels. An example is shown in Figure 3. We aim to reduce the dimensionality of the input satellite images while preserving band variance. Thus, the band selection module contains two steps. We first compute the inter-band correlation matrix from the samples in the training set using Pearson's correlation coefficient (Figure 4). We then categorize the bands into different clusters and select the ones in different clusters.

Figure 4 highlights three clusters in our data, each indicating the high correlation between the bands (bands 1-5, 6-9, and 11-12). We aim to select bands from different clusters for the two feature extraction modules to preserve band variance. Since bands 11 and 12 correspond to the Short Wave Infrared (SWIR) spectrum, which is mainly used for measuring soil and vegetation moisture content as it provides good contrast between different vegetation types, we intend to select bands from this cluster for the Feature-Engineering pipeline. Given that both bands show a high correlation, we select band 12 for its relatively lower correlation coefficient against the other satellite bands (bands 1-10) to avoid multicollinearity. For the



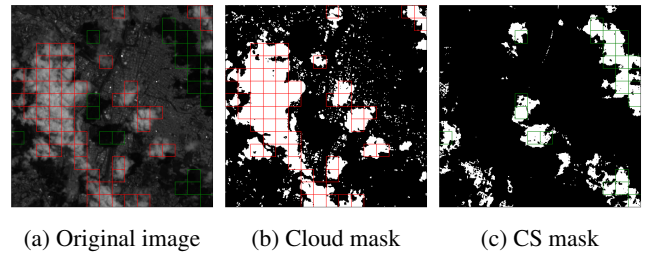(a) Original image  (b) Cloud mask  (c) CS mask

Figure 6: Cloud and cloud shadow masks generated in the CCS detection stage in Figure 5. (a) Original image where abnormal tiles will be swapped with the average of normal tiles. (b) Cloud mask with detected abnormal cloudy pixels in white and normal pixels in black. Abnormal tiles detected by the cloud mask are highlighted in red. (c) Cloud shadow (CS) mask with detected abnormal shadowy pixels in white and normal pixels in black. Abnormal tiles detected by the shadow mask are highlighted in green.

ViT feature extraction module, to preserve band diversity and match channels with the pre-training image set, we use bands 2, 3, and 4, which correspond to the Red, Green, and Blue channels.

## 4.3 Cloud and Cloud Shadow Removal

The cloud and cloud shadow removal (CSR) module is used to remove the cloud and cloud shadow from the selected satellite bands by performing CCS detection, image slicing, tile classification, tile averaging, and tile swapping (see Figure 5).

As satellite imagery often contains many cloud and cloud shadow noises, CCS detection [Li *et al.*, 2022b] is an essential stage for reducing noises. To identify noisy pixels caused by cloud or cloud shadow coverage, two thresholds are utilized to determine whether a pixel is considered noisy due to the often extreme pixel values in the affected areas. To establish thresholds for detecting cloud and cloud shadow, we evaluate the effectiveness of using pixel value percentiles from the training set and compare their performance. Through testing percentiles ranging from the 5th to 95th percentile at 5 percentile intervals, we choose two percentiles as the detection thresholds for cloud and cloud shadow, respectively. These thresholds are then used to generate the corresponding masks for cloud and cloud shadow (see Figure 6).

After obtaining the two masks, we slice each satellite band image into 16×16 tiles. With the sliced tiles and the cloud and cloud shadow masks, tiles are classified into abnormal and normal tiles, where an abnormal tile indicates more than 50 percent of pixels in the tile are marked as noise in either mask. For each tile in a different position in the images, we calculate the average tile of that position using the normal tiles By replacing the abnormal tiles in each sample with the corresponding average tiles, we generate noise-eliminated images. These average tiles are obtained by computing the average of normal tiles for a specific position in the images.

## 4.4 Spatial Feature Extractors

We adopt two feature extractors to extract different types of spatial features from the satellite images. In the Feature-Engineering feature extractor, we extract statistical pixel-based features from the SWIR band to obtain the texture in-

formation. Nine features from both first-order and higher-order features, such as Skewness and Joint Average, are collected using the PyRadiomics library [Van Griethuysen *et al.*, 2017]. The details can be found in the GitHub repository. For the ViT module, we adopt transfer learning to overcome the limited number of real-world satellite imagery in our dataset. We utilize a ViT [Wu *et al.*, 2020] pre-trained on ImageNet [Deng *et al.*, 2009] to collect deep learning-based features from the RGB bands. The RGB bands are downscaled from $736 \times 736$ to $224 \times 224$ to fit the model.

### 4.5 Model

The spatial feature extractors are both concatenated to a multi-layer LSTM module for extracting the temporal characteristics. To mitigate overfitting, a dropout layer is added after each LSTM layer in the module. The last LSTM layers are then concatenated to a multilayer perceptron (MLP) with one dense layer and one neuron as the final layer. We chose Leaky ReLu [Maas *et al.*, 2013] as the activation function to add non-linearity to the regression task. All models are trained for 100 epochs with an adaptive learning rate starting from 0.0001.

In this work, we train and evaluate the proposed structure on each municipality individually. This is because, with limited amount of training data, the model may prioritize learning the geographic meaning of different tile positions, within the same municipality. Since historical dengue cases are commonly used for dengue prediction, we evaluate the effectiveness of satellite imagery with dengue cases. To do so, we use the same multi-layer LSTM structure to create a LSTM model which takes cases as the model inputs. We also explore model performance with both satellite images and cases as inputs by concatenating the two LSTM modules from DengueNet with the LSTM module from the case model, resulting in a $10 \times 1$ dimension input to the MLP.

### 4.6 Evaluation and Performance Metrics

For each municipality, we use the first 80 percent of the data for training, the next 10 percent of the data for validation, and the last 10 percent for testing. We evaluate the proposed model structure using Mean Absolute Error (MAE), Symmetric Mean Absolute Percentage Error (sMAPE), and Root-Mean-Square Error (RMSE) metrics. sMAPE computes the percentage error between the actual value and the predicted value. We choose to use sMAPE over MAPE because the dengue cases in our dataset have relatively low actual values. RMSE penalizes the cases where the difference between actual and the predicted value is the greatest.

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |\hat{y}_i - y_i|, \tag{1}$$

$$sMAPE = \frac{100\%}{n} \sum_{i=1}^{n} \frac{2 \times |\hat{y}_i - y_i|}{(|\hat{y}_i| + |y_i|)} \tag{2}$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2}, \tag{3}$$

Refering to Equations 1,2,3, $n$ is the total number of samples to evaluate in the test set, and $i$ represents the sample number. $\hat{y}_i$ represents the predicted value from the model, and $y_i$ represents the actual value from the test set for each sample starting from $(i = 1)$ to $(i = n)$.

## 5 Results

Table 1 presents the performance evaluation of DengueNet in forecasting dengue cases using a time-series of satellite imagery with a window size of five weeks. Among the five municipalities assessed, Ibagué exhibits the most favorable performance across all metrics, while Cúcuta reports the least favorable performance. These results are anticipated. In Ibagué, apart from an initial peak, the dengue trend is comparatively more stable than in other municipalities. While the number of dengue cases in Cali appears stable, the high baseline number of cases results in an increase in the MAE. In the case of Cúcuta, given that the training set has relatively low occurrences of dengue, it is reasonable that the model fails to accurately reflect the actual trend of dengue cases for Cúcuta during the testing period. A notable observation is that while the three metrics have different values within one municipality, they report similar results acros municipalities, indicating that DengueNet exhibits relatively stable performance across different metrics.

Figure 7 depicts the forecasted dengue cases for five municipalities utilizing a diverse set of input data, including features extracted from satellite imagery and historical dengue cases. Comparative analysis is conducted against actual dengue incidences, an LSTM model relying solely on historical cases, and a combined model incorporating both satellite images and cases as input. Upon examination of the figures, it is evident that DengueNet demonstrates the capability to accurately predict most trends, even in the case of Villavicencio (refer to Figure 7c), which exhibits greater fluctuations in dengue cases over time. This observation substantiates the effectiveness of DengueNet in forecasting outbreak patterns within the majority of municipalities, relying solely on satellite images as input. Furthermore, our model exhibits robust predictive capabilities not only for short-term trends, while performing slightly less worse compared to the LSTM model that solely relies on historical case data, but also demonstrates adaptability by easily incorporating historical case data when available, thus enhancing prediction accuracy.

## 6 Ablation Studies

For the ablation studies, we evaluate the usage of the two feature extraction modules as shown in Figure 1, and the CSR module as presented in Table 2. As we observe a high degree of similarity among the MAE, sMAPE, and RMSE metrics in Table 1, our analysis focuses on examining the differences between the MAE with and without the inclusion of these three modules. For the Feature-Engineering module, four municipalities result in improved MAE, with Medellín having the most significant MAE improvement when paired with the CSR module. On the other hand, the CSR module has less impact on the ViT module, with only one municipality showing improved MAE. However, after combining both

| Metrics | Villavicencio | Medellín | Cúcuta | Ibagué | Cali | Average |
|---|---|---|---|---|---|---|
| MAE | 25.54±0.06 | 50.96±0.34 | **113.65±0.08** | **5.67±0.18** | 23.77±0.95 | 43.92±42.19 |
| sMAPE | 72.90±0.27 | 92.02±0.33 | **162.91±0.25** | **40.06±0.83** | 56.16±1.15 | 84.81±47.74 |
| RMSE | 30.62±0.03 | 67.86±0.40 | **120.57±0.07** | **7.45±0.22** | 31.80±1.46 | 51.66±44.17 |

Table 1: DengueNet evaluation across five municipalities. All experiments are repeated three times, with the average value reported with the standard deviation. The scores for the municipalities with the best and worst scores are indicated.

| ViT | FEng | CSR | Villavicencio | Medellín | Cúcuta | Ibagué | Cali |
|---|---|---|---|---|---|---|---|
| ✓ | | ✓ | **24.67±0.26** | 45.48±5.56 | 113.10±0.08 | 13.46±0.08 | 58.10±1.27 |
| ✓ | | | 26.25±0.00 | **44.77±0.79** | **109.31±0.00** | **6.21±0.13** | **33.42±0.42** |
| | ✓ | ✓ | **24.00±0.05** | **80.46±0.03** | **113.46±0.08** | **3.52±0.06** | 96.71±0.08 |
| | ✓ | | 27.21±0.29 | 111.15±0.19 | 113.58±0.03 | 6.96±0.16 | **48.15±0.31** |
| ✓ | ✓ | ✓ | 25.54±0.06 | 50.96±0.34 | **113.65±0.08** | **5.67±0.18** | 23.77±0.95 |
| ✓ | ✓ | | **24.40±0.06** | **42.48±0.96** | 114.19±0.09 | 7.25±0.09 | 42.35±0.81 |

Table 2: MAE scores with or without the cloud shadow removal (CSR) module combined with different feature extractors across five municipalities. ViT indicates only features extracted from the ViT module are used. FEng indicates only features extracted from the feature-engineering module are used. All experiments are repeated three times. Average values are reported ± the standard deviation. The best scores are highlighted.



(a) Medellín
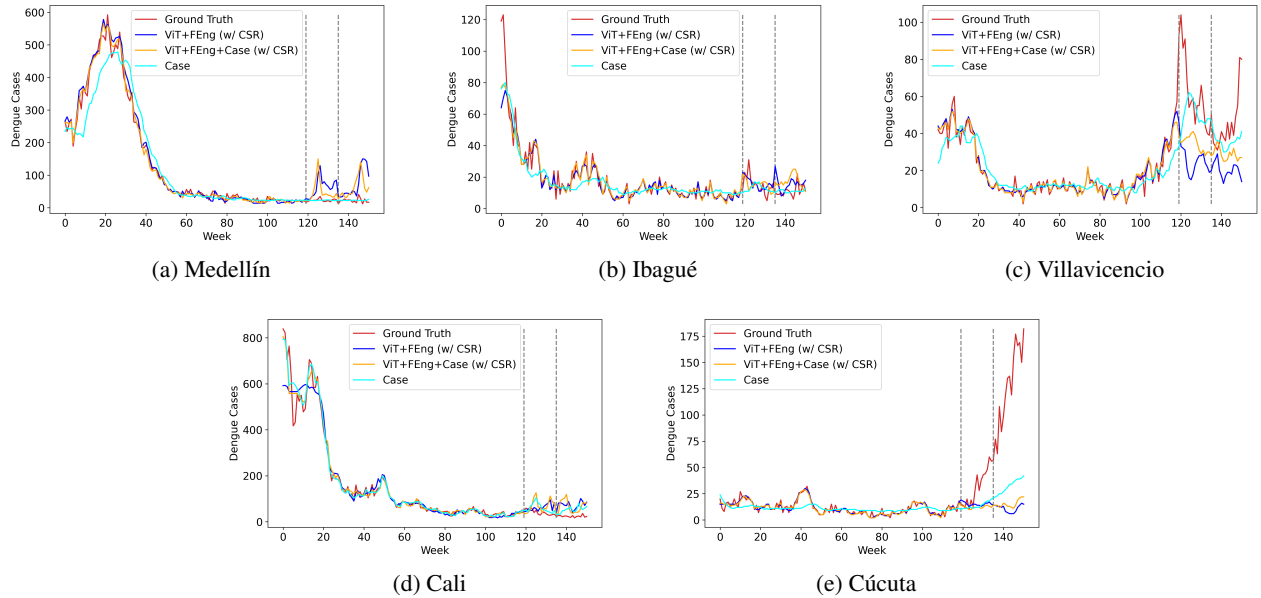
(b) Ibagué

(c) Villavicencio

(d) Cali

(e) Cúcuta

Figure 7: Dengue case prediction was performed for five municipalities per epidemiological week from 2016 to 2018. Three approaches were evaluated: using satellite imagery features (ViT+FEng), case data (Case), and a combination of both (ViT+FEng+Case). The Ground Truth label represents the actual number of dengue cases per week. The grey vertical dashed lines indicate the starting weeks of the validation and testing sets.

| Models | MAE | sMAPE | RMSE |
|---|---|---|---|
| ViT (w/ CSR) | 50.96 | 97.66 | 60.20 |
| FEng (w/ CSR) | 63.63 | 99.24 | 74.02 |
| ViT+FEng (w/ CSR) | **43.92** | **84.81** | **51.66** |

Table 3: Performance comparison of different feature extractors with the cloud and shadow removal module (w/ CSR). All experiments are repeated three times and average values are reported. The best scores are highlighted.

spatial feature extraction modules as inputs, the CSR module improves the performance across three municipalities, and the average MAE across five municipalities also decreases from 54.14 to 51.66.

The effectiveness of having both spatial feature extractors is also analyzed in Table 3. With a single feature extractor, the ViT feature extractor performs slightly better than the Feature-Engineering extractor. However, the lowest average MAE, sMAPE, and RMSE are observed when both feature extractors are used. This finding is reasonable as the two feature extractors retrieve different types of information from the satellite imagery. This model architecture design enables DengueNet to maintain high performance even if one of the feature extraction modules fails to extract crucial features, as the other feature extractor can compensate for it.

## 7 Discussion

This study introduces a robust and efficient approach for extracting satellite data and presents DengueNet, a novel architecture for predicting dengue outbreaks using satellite imagery. The experimentation phase involves the analysis of satellite images and dengue cases spanning from 2016 to 2018, focusing specifically on five municipalities in Colombia, a country significantly affected by the prevalence of dengue fever. The proposed model combines ViTs with concatenated multi-layer LSTMs to effectively extract both spatial and temporal information from a series of satellite imagery, resulting in comparable dengue case predictions.

To address the challenges posed by the dimensionality of satellite images, the study incorporates band selection based on band-to-band Pearson's correlation, enabling a comprehensive assessment of Sentinel-2 satellite images. The selected bands undergo feature extraction through the use of both the feature-engineering and ViT modules. The feature-engineering pipeline involves dividing satellite images into tiles and employing CCS detection to minimize the presence of environmental noise artifacts, allowing for the extraction of noise-free pixel features. On the other hand, the ViT module utilizes transfer learning from a pre-trained ViT model to extract features. These extracted features from both modules are subsequently integrated into a concatenated LSTM-based model for predicting dengue cases.

Incorporating freely accessible satellite imagery into our DengueNet model holds significant potential for making a substantial impact on public health legislation and fairness in health. Over the past two decades, dengue fever has emerged as a prevalent epidemic in tropical developing countries, necessitating the establishment of an effective early warning system for preventing and monitoring outbreaks. The feasibility of DengueNet for predicting dengue outbreaks has been successfully demonstrated in five municipalities, showcasing its potential for transferability to other geographical regions. Moreover, the computational requirements of the model are relatively low, and its deployment only requires minimal resources, making it an accessible alternative for resource-constrained developing countries.

The proposed approach is further reinforced by the inclusion of a dockerized version of the satellite extraction framework, leveraging Sentinel Hub, which ensures data reproducibility and scalability [Alberto *et al.*, 2023]. This empowers LMICs to leverage higher quality and more frequently updated satellite data, overcoming the limitations of field data collection characterized by irregular revisit rates and varying data quality. The utilization of such information can significantly contribute to informed policy decisions and strategies at the municipality level, enabling early containment of the dengue virus. Ultimately, the proposed method holds immense potential to enhance the prevention and control of dengue fever outbreaks in developing countries, thereby advancing public health outcomes and promoting health equity.

## 8 Conclusion

The dockerized satellite extraction framework and lightweight DengueNet model presented in this work present a viable alternative for LMICs, where data collection and preprocessing pose substantial challenges. The performance of DengueNet, which leverages publicly accessible satellite imagery, exhibits comparable performance to that of a straightforward LSTM model that relies exclusively on dengue cases for dengue prediction. This approach takes us closer to the democratization of data access and the implementation of machine learning models globally, thereby aiding in the formulation of informed public health policies and strategies for early warning systems. To ensure safe and responsible integration of satellite imagery and DengueNet, future work should understand and mitigate the sources of bias inherent in machine learning models[Celi *et al.*, 2022; Nazer *et al.*, 2023] to promote fairness and reduce disparities in public health across diverse populations.

## Acknowledgments

# References

[Abdur Rehman *et al.*, 2019] Nabeel Abdur Rehman, Umar Saif, and Rumi Chunara. Deep landscape features for improving vector-borne disease prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 44–51, 2019.

[Alberto *et al.*, 2023] Isabelle Rose I Alberto, Nicole Rose I Alberto, Arnab K Ghosh, Bhav Jain, Shruti Jayakumar, Nicole Martinez-Martin, Ned McCague, Dana Moukheiber, Lama Moukheiber, Mira Moukheiber, et al. The impact of commercial health datasets on medical research and health-care algorithms. *The Lancet Digital Health*, 5(5):e288–e294, 2023.

[Andersson *et al.*, 2019] Virginia Ortiz Andersson, Cristian Cechinel, and Ricardo Matsumura Araujo. Combining street-level and aerial images for dengue incidence rate estimation. In *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2019.

[Bhatt *et al.*, 2013] Samir Bhatt, Peter W Gething, Oliver J Brady, Jane P Messina, Andrew W Farlow, Catherine L Moyes, John M Drake, John S Brownstein, Anne G Hoen, Osman Sankoh, et al. The global distribution and burden of dengue. *Nature*, 496(7446):504–507, 2013.

[Cattarino *et al.*, 2020] Lorenzo Cattarino, Isabel Rodriguez-Barraquer, Natsuko Imai, Derek AT Cummings, and Neil M Ferguson. Mapping global variation in dengue transmission intensity. *Science translational medicine*, 12(528):eaax4144, 2020.

[CDC, 2022] CDC. Dengue. https://www.cdc.gov/dengue/index.html, 2022. Accessed: 2023-01-15.

[Celi *et al.*, 2022] Leo Anthony Celi, Jacqueline Cellini, Marie-Laure Charpignon, Edward Christopher Dee, Franck Dernoncourt, Rene Eber, William Greig Mitchell, Lama Moukheiber, Julian Schirmer, Julia Situ, et al. Sources of bias in artificial intelligence that perpetuate healthcare disparities—a global review. *PLOS Digital Health*, 1(3):e0000022, 2022.

[Chaparro *et al.*, 2016] P Chaparro, W León, and CA Castañeda. Comportamiento de la mortalidad por dengue en colombia entre 1985 y 2012. *Biomédica*, 36(Supl 2):125–34, 2016.

[Datoc *et al.*, 2016] Hillary Ingrid Datoc, Romeo Caparas, and Jaime Caro. Forecasting and data visualization of dengue spread in the philippine visayas island group. In *2016 7th International Conference on Information, Intelligence, Systems & Applications (IISA)*, pages 1–4. IEEE, 2016.

[de Witt *et al.*, 2020] Christian Schroeder de Witt, Catherine Tong, Valentina Zantedeschi, Daniele De Martini, Freddie Kalaitzis, Matthew Chantry, Duncan Watson-Parris, and Piotr Bilinski. Rainbench: towards global precipitation forecasting from satellite imagery. *arXiv preprint arXiv:2012.09670*, 2020.

[Deng *et al.*, 2009] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.

[Fenech *et al.*, 2018] Matthew Fenech, Nika Strukelj, and Olly Buston. Ethical, social, and political challenges of artificial intelligence in health. *London: Wellcome Trust Future Advocacy*, 12, 2018.

[Fontaine *et al.*, 2018] Albin Fontaine, Sebastian Lequime, Isabelle Moltini-Conclois, Davy Jiolle, Isabelle Leparc-Goffart, Robert Charles Reiner Jr, and Louis Lambrechts. Epidemiological significance of dengue virus genetic variation in mosquito infection dynamics. *PLoS pathogens*, 14(7):e1007187, 2018.

[Guo *et al.*, 2017] Pi Guo, Tao Liu, Qin Zhang, Li Wang, Jianpeng Xiao, Qingying Zhang, Ganfeng Luo, Zhihao Li, Jianfeng He, Yonghui Zhang, et al. Developing a dengue forecast model using machine learning: A case study in china. *PLoS neglected tropical diseases*, 11(10):e0005973, 2017.

[Gutierrez-Barbosa *et al.*, 2020] Hernando Gutierrez-Barbosa, Sandra Medina-Moreno, Juan C Zapata, and Joel V Chua. Dengue infections in colombia: epidemiological trends of a hyperendemic country. *Tropical Medicine and Infectious Disease*, 5(4):156, 2020.

[Jain *et al.*, 2019] Raghvendra Jain, Sra Sontisirikit, Sopon Iamsirithaworn, and Helmut Prendinger. Prediction of dengue outbreaks based on disease surveillance, meteorological and socio-economic data. *BMC infectious diseases*, 19(1):1–16, 2019.

[Karim *et al.*, 2012] Md Nazmul Karim, Saif Ullah Munshi, Nazneen Anwar, and Md Shah Alam. Climatic factors influencing dengue cases in dhaka city: a model for dengue prediction. *The Indian journal of medical research*, 136(1):32, 2012.

[Kruk *et al.*, 2018] Margaret E Kruk, Anna D Gage, Catherine Arsenault, Keely Jordan, Hannah H Leslie, Sanam Roder-DeWan, Olusoji Adeyi, Pierre Barker, Bernadette Daelmans, Svetlana V Doubova, et al. High-quality health systems in the sustainable development goals era: time for a revolution. *The Lancet global health*, 6(11):e1196–e1252, 2018.

[Lee *et al.*, 2017] Jung-Seok Lee, Mabel Carabali, Jacqueline K Lim, Victor M Herrera, Il-Yeon Park, Luis Villar, and Andrew Farlow. Early warning signal for dengue outbreaks and identification of high risk areas for dengue fever in colombia using climate and non-climate datasets. *BMC Infectious Diseases*, 17(1):1–11, 2017.

[Li *et al.*, 2022a] Zhichao Li, Helen Gurgel, Lei Xu, Linsheng Yang, and Jinwei Dong. Improving dengue forecasts by using geospatial big data analysis in google earth engine and the historical dengue information-aided long short term memory modeling. *Biology*, 11(2):169, 2022.

[Li *et al.*, 2022b] Zhiwei Li, Huanfeng Shen, Qihao Weng, Yuzhuo Zhang, Peng Dou, and Liangpei Zhang. Cloud and cloud shadow detection for optical satellite imagery: Features, algorithms, validation, and prospects. *ISPRS Journal of Photogrammetry and Remote Sensing*, 188:89–108, 2022.

[Lim *et al.*, 2020] Jue Tao Lim, Borame Sue Lee Dickens, Lawrence Zheng Xiong Chew, Esther Li Wen Choo, Joel Ruihan Koo, Joel Aik, Lee Ching Ng, and Alex R Cook. Impact of sars-cov-2 interventions on dengue transmission. *PLoS neglected tropical diseases*, 14(10):e0008719, 2020.

[Livelo and Cheng, 2018] Evan Dennison Livelo and Charibeth Cheng. Intelligent dengue infoveillance using gated recurrent neural learning and cross-label frequencies. In *2018 IEEE International Conference on Agents (ICA)*, pages 2–7. IEEE, 2018.

[Ltd, 2022] Sinergise Ltd. Sentinel-2 L2A about sentinet-2 l2a data. https://www.sentinel-hub.com/, 2022. Accessed: 2022-08-13.

[Maas *et al.*, 2013] Andrew L Maas, Awni Y Hannun, Andrew Y Ng, et al. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, page 3. Atlanta, Georgia, USA, 2013.

[Mala and Jat, 2019] Shuchi Mala and Mahesh Kumar Jat. Geographic information system based spatio-temporal dengue fever cluster analysis and mapping. *The Egyptian Journal of Remote Sensing and Space Science*, 22(3):297–304, 2019.

[Martheswaran *et al.*, 2022] Tarun Kumar Martheswaran, Hamida Hamdi, Amal Al-Barty, Abeer Abu Zaid, and Biswadeep Das. Prediction of dengue fever outbreaks using climate variability and markov chain monte carlo techniques in a stochastic susceptible-infected-removed model. *Scientific Reports*, 12(1):5459, 2022.

[Morgan *et al.*, 2021] Jasmine Morgan, Clare Strode, and J Enrique Salcedo-Sora. Climatic and socio-economic factors supporting the co-circulation of dengue, zika and chikungunya in three different ecosystems in colombia. *PLoS Neglected Tropical Diseases*, 15(3):e0009259, 2021.

[Moskolaï *et al.*, 2021] Waytehad Rose Moskolaï, Wahabou Abdou, and Albert Dipanda. Application of deep learning architectures for satellite image time series prediction: A review. *Remote Sensing*, 13(23):4822, 2021.

[Muñoz *et al.*, 2021] Estefanía Muñoz, Germán Poveda, M Patricia Arbeláez, and Iván D Vélez. Spatiotemporal dynamics of dengue in colombia in relation to the combined effects of local climate and enso. *Acta Tropica*, 224:106136, 2021.

[National Institute of Health of Colombia, 2010] National Institute of Health of Colombia. Comportamiento epidemiológico del dengue en colombia año 2010. http://www.ins.gov.co/buscador-eventos/Paginas/Info-Evento.aspx, 2010. Accessed: 2022-08-13.

[Nazer *et al.*, 2023] Lama H Nazer, Razan Zatarah, Shai Waldrip, Janny Xue Chen Ke, Mira Moukheiber, Ashish K Khanna, Rachel S Hicklen, Lama Moukheiber, Dana Moukheiber, Haobo Ma, et al. Bias in artificial intelligence algorithms and recommendations for mitigation. *PLOS Digital Health*, 2(6):e0000278, 2023.

[Ndabarora *et al.*, 2014] Eleazar Ndabarora, Jennifer A Chipps, and Leana Uys. Systematic review of health data quality management and best practices at community and district levels in lmic. *Information Development*, 30(2):103–120, 2014.

[PAHO, 2022] PAHO. Dengue. https://www.paho.org/en/topics/dengue, 2022. Accessed: 2022-08-13.

[Ren *et al.*, 2021] Xiaoli Ren, Xiaoyong Li, Kaijun Ren, Junqiang Song, Zichen Xu, Kefeng Deng, and Xiang Wang. Deep learning-based weather prediction: a survey. *Big Data Research*, 23:100178, 2021.

[Rogers *et al.*, 2002] David J Rogers, Sarah E Randolph, Robert W Snow, and Simon I Hay. Satellite imagery in the study and forecast of malaria. *Nature*, 415(6872):710–715, 2002.

[Roster and Rodrigues, 2021] Kirstin Roster and Francisco A Rodrigues. Neural networks for dengue prediction: a systematic review. *arXiv preprint arXiv:2106.12905*, 2021.

[Salim *et al.*, 2021] Nurul Azam Mohd Salim, Yap Bee Wah, Caitlynn Reeves, Madison Smith, Wan Fairos Wan Yaacob, Rose Nani Mudin, Rahmat Dapari, Nik Nur Fatin Fatihah Sapri, and Ubydul Haque. Prediction of dengue outbreak in selangor malaysia using machine learning techniques. *Scientific reports*, 11(1):1–9, 2021.

[Shragai *et al.*, 2022] Talya Shragai, Juliana Pérez-Pérez, Marcela del Pilar Quimbayo-Forero, Raúl Rojo, Laura C. Harrington, and Guillermo Rúa-Uribe. Distance to public transit predicts spatial distribution of dengue virus incidence in medellín, colombia. *Scientific Reports*, 12(1):8333, May 2022.

[Son and Thong, 2017] Le Hoang Son and Pham Huy Thong. Some novel hybrid forecast methods based on picture fuzzy clustering for weather nowcasting from satellite image sequences. *Applied Intelligence*, 46(1):1–15, 2017.

[Van Griethuysen *et al.*, 2017] Joost JM Van Griethuysen, Andriy Fedorov, Chintan Parmar, Ahmed Hosny, Nicole Aucoin, Vivek Narayan, Regina GH Beets-Tan, Jean-Christophe Fillion-Robin, Steve Pieper, and Hugo JWL Aerts. Computational radiomics system to decode the radiographic phenotype. *Cancer research*, 77(21):e104–e107, 2017.

[Watts *et al.*, 2020] Matthew J Watts, Panagiota Kotsila, P Graham Mortyn, Victor Sarto i Monteys, and Cesira Urzi Brancati. Influence of socio-economic, demographic and climate factors on the regional distribution of dengue in the united states and mexico. *International journal of health geographics*, 19(1):1–15, 2020.

[WHO, 2022] WHO. Dengue and severe dengue. https://www.who.int/news-room/fact-sheets/detail/dengue-and-severe-dengue, 2022. Accessed: 2022-08-13.

[Wu *et al.*, 2020] Bichen Wu, Chenfeng Xu, Xiaoliang Dai, Alvin Wan, Peizhao Zhang, Zhicheng Yan, Masayoshi Tomizuka, Joseph Gonzalez, Kurt Keutzer, and Peter Vajda. Visual transformers: Token-based image representation and processing for computer vision, 2020.