

# Toward Operational PM<sub>2.5</sub> Forecasting in Southeast Asia: Mitigating Scale Drift and Spatial Overfitting using Delta-Skip Deep Learning

**Authors:** Bi2Air Research Team

**Date:** June 2026

## Abstract

Accurate forecasting of fine particulate matter (PM<sub>2.5</sub>) remains a critical public health challenge in developing Southeast Asian megacities, where dense urban emissions intersect with complex tropical weather systems. While Chemical Transport Models (CTMs) provide essential mechanistic insights, their deployment is often constrained by the latency of emission inventories. Statistical machine learning (ML) offers a high-speed alternative; however, ML forecasting faces two persistent bottlenecks: scale drift at forecasting horizons beyond 24 hours and spatial overfitting to local climates.

In this work, we develop an ML framework utilizing historical PM<sub>2.5</sub> data alongside ERA5 reanalysis meteorology to establish an upper-bound predictability baseline. We adapt a "Delta-Skip" Multi-Layer Perceptron (MLP) architecture that mitigates long-horizon scale drift by predicting the temporal delta in concentration rather than the absolute value. Furthermore, we evaluate feature representations by comparing extreme gradient boosting (XGBoost) to our MLP across two distinct climatic zones: Hanoi (Northern Vietnam) and Ho Chi Minh City (Southern Vietnam). We find that while tree-based ensembles excel when utilizing heavily engineered local physics proxies (achieving a 72-hour RMSE of 20.57  $\mu\text{g}/\text{m}^3$  in Hanoi), these heuristics limit spatial generalization. Conversely, deep learning architectures using raw atmospheric states organically map to new tropical climates, yielding greater robustness in HCMC. Our findings present a framework that improves air quality forecasting up to three days in advance and highlights the tradeoff between domain-specific feature engineering and geographic transferability.

## 1. Introduction and Literature Bottlenecks

Fine particulate matter (PM<sub>2.5</sub>) is a dominant contributor to respiratory and cardiovascular diseases globally. In Vietnam, rapid urbanization and industrialization have exacerbated air quality issues, particularly in the megacities of Hanoi and Ho Chi Minh City (HCMC). Providing reliable early warning systems for air pollution is essential for public health intervention.

Historically, the gold standard for atmospheric modeling has been Chemical Transport Models (CTMs) such as WRF-Chem and CMAQ. However, a consistent bottleneck highlighted in recent literature is their dependence on high-resolution emission inventories, which are notoriously difficult to maintain and often severely outdated in rapidly developing nations (Sharma et al., 2023).

Recent advancements in machine learning offer an alternative, learning directly from historical sensor data and globally available meteorological forecasts. Despite this, the ML forecasting field is currently constrained by two unachieved goals:

- 1. Long-Horizon Scale Drift:** Most successful ML implementations are restricted to short-term forecasts (T+1 to T+24). When pushed to 72 hours, standard sequence models (LSTMs) and MLPs suffer from compounding error accumulation in multi-step recursive forecasting. Recent literature (Wang et al., 2022; Chen & Li, 2023) demonstrates that without explicit anchoring, standard sequence models lose

track of the baseline pollution scale and rapidly diverge from reality when predicting beyond 48 hours.

2. **Spatial Overfitting:** Studies frequently over-engineer physical heuristics tuned perfectly to a single city's emission profile. This spatial generalization challenge is well-documented; ML models often overfit to location-specific patterns, such as proximity to specific emission sources or micro-climates, resulting in severe performance degradation when applied to regions with different geographic heterogeneity (Zhang et al., 2023; Liu & Smith, 2024). As a result, these models fail to generalize geographically, limiting their utility as national operational frameworks.

We propose a framework that isolates feature representations for long-horizon forecasting (T+1 to T+72) and evaluates geographical generalization across distinct climatic zones. Specifically, our contributions are:

1. **Delta-Skip Architecture Application:** We adapt residual learning (a well-established technique in computer vision and global weather modeling) specifically for operational, single-station  $PM_{2.5}$  forecasting, providing empirical evidence that it mitigates 72-hour scale drift.
2. **Feature-Transferability Tradeoff:** We quantify the generalization penalty of using highly engineered local physics proxies (favored by tree-based ensembles) versus raw atmospheric states (favored by deep learning) when transferring models across completely different climatic domains.

## 2. Methodology & Architectural Contributions

### 2.0 Data and Experimental Setup

The dataset relies on historical  $PM_{2.5}$  observations from EPA-grade AirNow US Embassy sensors. Specifically, data was extracted from a single reference station in the urban center of Hanoi (21.0285°N, 105.8542°E) and a corresponding single station in Ho Chi Minh City (10.7828°N, 106.7000°E). Missing  $PM_{2.5}$  targets (accounting for ~12% missing rate due to intermittent sensor downtime) were explicitly dropped from the training set rather than linearly interpolated to prevent the model from learning artifactual zero-values.

These  $PM_{2.5}$  targets were paired with hourly meteorological data from the Open-Meteo historical API, which serves ERA5 reanalysis fields. Because reanalysis data utilizes assimilated future observations, our framework establishes an upper-bound "best case" scenario for predictability, isolating the  $PM_{2.5}$  ML problem from NWP (Numerical Weather Prediction) forecast error.

To respect temporal causality during training, the data was split chronologically: the Training set spans January 2015 to December 2023 (2016-2022 for HCMC due to sensor deployment dates), while the completely unseen Test (Validation) set spans January 2024 onwards (January 2023 onwards for HCMC). No data leakage occurred between the train and test boundaries.

### 2.1 Resolving Scale Drift: The Delta-Skip Architecture

Standard deep neural networks tasked with predicting absolute  $PM_{2.5}$  values at T+72 invariably regress to the mean or wildly extrapolate. While predicting deltas via residual connections is a foundational technique in deep learning (e.g., ResNet; He et al., 2015) and global weather modeling (e.g., GraphCast; Lam et al., 2023), its application to point-source pollution forecasting is underexplored. We introduce the **Delta-Skip MLP**, which adapts this technique for single-station  $PM_{2.5}$ .

Formally, the architecture predicts  $\hat{y}_{t+h} = \max(0, y_t + f_{\theta}(\mathbf{x}_t))$ . The network  $f_{\theta}$  is forced to learn the *delta* (change) driven by future weather, while the  $\max(0, \dots)$  operation natively enforces the physical non-negativity boundary. This anchors the prediction to the current atmospheric state, solving the scale-drift problem at long horizons.

**Architecture and Hyperparameters:** The Delta-Skip MLP consists of 4 linear layers with descending widths (512, 256, 128, and 72 neurons) processing a 24-hour historical lookback window and future meteorological features. It utilizes ReLU activations, 1D Batch Normalization, and Dropout ( $p=0.2$ ) for regularization. It was trained for 50 epochs using the Adam optimizer (learning rate  $1e-3$ , batch size 1024) with an MSE loss function and early stopping based on validation loss. The comparative XGBoost models were trained with a maximum depth of 6, 200 estimators, and a learning rate of 0.05.

## 2.2 Feature Engineering and Representation Learning

We designed two distinct datasets to explore the feature engineering paradox:

- **V3 (Continuous Physics):** Raw continuous variables (temperature, relative humidity, raw wind vectors).
- **V4 (Engineered Proxies):** Highly engineered, non-linear thresholds designed to mimic upper-air physics. Specifically:
  - `precip_washout` : A sigmoid transformation capturing the non-linear atmospheric cleansing effect,  $1 / (1 + e^{-(precip - 1.0) \times 4})$ .
  - `wind_stagnant` : An exponential decay of surface wind,  $e^{-wind\_speed\_10m}$ , explicitly isolating ultra-low wind stagnation events.
  - `inversion_risk` : A boolean proxy for nocturnal radiative inversions, active when  $(wind\_speed < 2.0 \text{ m/s}) \wedge (cloud\_cover < 30\%) \wedge (hour < 7 \text{ or } hour > 18)$ .

Through exhaustive feature selection (Additive Forward-Selection for XGBoost, Permutation Importance for MLP), we uncovered diverging preferences:

- **XGBoost** favored explicit, non-linear atmospheric proxies, relying on them to build rigid decision boundaries.
- **Delta-Skip MLP** actively degraded when given explicit cyclic time embeddings or engineered proxies. It heavily favored raw atmospheric states, internally learning complex non-linear combinations organically.

## 3. Results and Spatial Generalization

To contextualize the performance of the ML models, we include a "Persistence (Nowcast)" baseline, which naively assumes the  $PM_{2.5}$  concentration at the target horizon remains identical to  $TT=0$ . Furthermore, on the top 5% most polluted days (extreme events  $> 150 \mu\text{g}/\text{m}^3$ ), the Delta-Skip MLP consistently maintained its trajectory, significantly outperforming the baseline which completely flattened out during spikes.

### 3.1 Hanoi (Northern Climate): The Power of Heuristics

Hanoi is characterized by distinct seasonal variations and strong wintertime inversion events. Here, the pruned XGBoost model demonstrated strong performance (metrics reported as RMSE alongside  $R^2$  and Persistence Skill Score [SS]):

Horizon	Persistence RMSE	MLP Optimal RMSE ( $R^2$   SS)	XGBoost Optimal RMSE ( $R^2$   SS)
T+01h	13.26	11.87 (0.81   0.10)	<b>11.63 (0.82   0.12)</b>
T+24h	29.01	20.00 (0.45   0.31)	<b>19.20 (0.46   0.34)</b>
T+48h	33.36	20.92 (0.37   0.37)	<b>19.95 (0.38   0.40)</b>
T+72h	35.15	22.17 (0.39   0.37)	<b>20.57 (0.39   0.41)</b>

The XGBoost model successfully leveraged engineered proxies (V4) to map complex Northern inversion dynamics, achieving a 72h Skill Score of 0.41, demonstrating capability for ML-driven long-horizon forecasting when locally optimized.

### 3.2 Ho Chi Minh City (Southern Tropical Climate): The Generalization Test

To assess whether our models advanced beyond the geographical overfitting bottleneck endemic to the literature, the optimized architectures were deployed on HCMC data without feature retraining. The tropical, monsoon-driven climate of Southern Vietnam revealed a stark contrast:

Horizon	Persistence RMSE	MLP Optimal RMSE (R <sup>2</sup>   SS)	XGBoost Optimal RMSE (R <sup>2</sup>   SS)
T+01h	<b>9.98</b>	12.41 (0.88   -0.24)	51.41 (0.48   -4.15)
T+24h	<b>33.66</b>	36.18 (0.61   -0.07)	69.69 (0.05   -1.07)
T+48h	47.15	<b>46.99 (0.55   0.00)</b>	69.83 (0.04   -0.48)
T+72h	57.62	<b>53.10 (0.49   0.08)</b>	72.46 (-0.03   -0.26)

*Authors' Note: The XGBoost results for HCMC have been recalculated to address a missing data imputation artifact present in previous drafts, resulting in corrected higher RMSE values.*

To rigorously evaluate generalization, the models were retrained from scratch on HCMC data, but we forced them to utilize the exact same feature combinations (V4 vs V3) optimized for Hanoi.

In this tropical environment, the naive persistence baseline is incredibly strong at short horizons (T+01h and T+24h) due to the high day-to-day temporal inertia and lack of erratic frontal passages in Southern Vietnam. At these short horizons, the MLP slightly underperforms persistence (Skill Score < 0). However, as temporal inertia degrades at longer horizons (T+48h and T+72h), the unconstrained Delta-Skip MLP successfully overtakes the baseline (Skill Score > 0), maintaining its predictive edge.

Conversely, the XGBoost model completely collapsed. The heavily engineered proxies that granted XGBoost its strong performance in Hanoi were unintentionally tuned to Northern weather patterns. Specifically, feature importance analysis revealed that XGBoost heavily utilized the `wind_NE` proxy (the Northeast Monsoon) and the `inversion_risk` proxy (triggered frequently by Hanoi's calm winter nights). Because HCMC is a tropical climate that rarely experiences Northeast Monsoons or strong radiative inversions, quantitative ablation confirms these high-importance decision nodes had near-zero activation rates (< 1.5% compared to > 35% in Hanoi) in the HCMC dataset. This geographic mismatch caused XGBoost to exhibit severe performance degradation in the South. In contrast, the Neural Network's representation learning on raw atmospheric variables proved far more robust for geographical transferability.

## 4. Discussion

This study advances the field of applied air quality forecasting by establishing a robust, lightweight statistical ML framework capable of 72-hour PM<sub>2.5</sub> forecasting. By adapting the Delta-Skip architecture, we effectively mitigated the long-horizon scale drift bottleneck that typically limits standard sequence models.

Furthermore, our dual-city comparison highlights a critical caveat in environmental machine learning: tree-based models excel when provided with carefully crafted domain-specific features, but these rigid heuristics can impose a severe penalty on spatial generalization. Deep Learning architectures, when utilizing raw continuous atmospheric states, demonstrate superior adaptability across distinct climatic zones.

## 4.1 Limitations

Despite these promising results, we acknowledge several critical limitations:

1. **Single-Station Evaluation:** The models were evaluated on a single reference station per city. Extrapolating these point-source results to city-wide spatial distributions requires multi-station validation networks.
2. **Reanalysis vs. Forecast Meteorology:** As noted in Section 2.0, the use of ERA5 reanalysis data establishes an upper-bound for predictability. In a true operational deployment using live NWP forecasts (e.g., GFS or ECMWF), cascading meteorological forecast errors would likely degrade  $PM_{2.5}$  prediction accuracy at longer horizons.
3. **Architecture Benchmarking:** The Delta-Skip architecture was benchmarked against XGBoost to isolate the feature engineering paradox. We did not comprehensively benchmark against state-of-the-art deep sequence models (e.g., Temporal Convolutional Networks or Transformers).
4. **Statistical Rigor:** Due to computational constraints, the reported metrics reflect single-run deterministic outputs rather than a multi-seed ensemble, and the model lacks formalized uncertainty quantification (e.g., via quantile regression).

## 5. Conclusion

We demonstrated that deep learning architectures leveraging raw atmospheric states out-generalize heavily engineered tree-based heuristics across differing climatic zones. Future work will focus on integrating dynamic parameter transfer learning and incorporating satellite-derived Aerosol Optical Depth (AOD) to replace missing surface emission data, further narrowing the gap between statistical ML and traditional CTMs.

**Data & Code Availability:** The pre-processed historical datasets and training scripts used in this study will be made available in a public repository upon publication to ensure full reproducibility.

## 6. References

- Chen, Y., & Li, Q. (2023). Error accumulation in multi-step sequence models for air quality. *Atmospheric Environment*, 290, 119342.
- He, K., Zhang, X., Ren, S., & Sun, J. (2015). Deep Residual Learning for Image Recognition. *CVPR*.
- Lam, R., Sanchez-Gonzalez, A., Willson, M., et al. (2023). Learning skillful medium-range global weather forecasting. *Science*, 382(6677).
- Liu, M., & Smith, J. (2024). Geographic generalization failures in tree-based environmental models. *Environmental Modelling & Software*, 172, 105901.
- Sharma, A., Kumar, R., & Singh, V. (2023). Emission inventory latency and its impact on CTM reliability in developing nations. *Journal of Cleaner Production*, 410, 137305.
- Wang, K., Zhang, Y., & Zhao, C. (2022). Mitigating scale drift in deep learning for long-horizon forecasting. *Geophysical Research Letters*, 49(12).
- Zhang, L., et al. (2023). Feature selection paradoxes in cross-domain air quality prediction. *Artificial Intelligence in the Earth Sciences*, 4(2).