



Supplement of

Correcting aerosol extinction coefficient vertical structure biases in GEOS-chem via a physics-informed transformer with physical mechanism diagnosis

Jiajun Xiong et al.

Correspondence to: Yi Wang (wangyi34@cug.edu.cn)

The copyright of individual parts of the supplement might differ from the article licence.

22 **S1. CALIOP Quality Control**

23 To ensure the reliability of the observational target, we implement a rigorous,
24 tiered quality control strategy for the CALIOP Level 2 aerosol profile data. First, to
25 minimize cloud contamination, profiles with a Cloud Layer Fraction (CLF) exceeding
26 2.0 were excluded. We further refined this dataset by applying strict thresholds to the
27 Cloud Aerosol Discrimination (CAD) score, retaining only retrievals with high
28 confidence (CAD score between -100 and -20). The Extinction QC flag was used to
29 filter for algorithmically stable solutions (unconstrained or constrained), rejecting
30 profiles with divergent retrieval errors. Adjustments to lidar ratios or opaque layer
31 identifications were permitted only when necessary for algorithm convergence. A final
32 multi-step refinement was applied to the vertical profiles: (1) profiles were discarded if
33 the extinction uncertainty exceeded 99.9 km^{-1} at any lower-altitude layer; (2) physically
34 unrealistic extinction coefficients greater than 2 km^{-1} were flagged and removed; and
35 (3) to mitigate surface contamination artifacts, data within 180 m of the terrain were
36 excluded. Observations classified as "clear air" were assigned an extinction coefficient
37 of 0 km^{-1} , with the exception of the lowest 180 m AGL, which was masked to avoid
38 surface return interference. Additional safeguards included the removal of high-altitude
39 data ($>4 \text{ km}$) near 0°C ice clouds and the exclusion of spatially isolated aerosol signals
40 (horizontal resolution $< 80 \text{ km}$).

41 S2. AERONET AOD Interpolation

42 Standard AERONET sun photometers do not measure AOD directly at 532 nm,
43 the operating wavelength of the CALIOP lidar and our GEOS-Chem configuration. To
44 enable direct comparison, we interpolated the AOD to 532 nm using the Ångström
45 exponent (α) derived from measurements at adjacent wavelengths (typically 500 nm
46 and 675 nm, or 440 nm/870 nm if necessary). The interpolation follows the power-law
47 relationship:

$$48 \quad \tau_{532} = \tau_{\lambda_1} \times \left(\frac{532}{\lambda_1}\right)^{-\alpha}$$
$$49 \quad \alpha = -\frac{\ln\left(\frac{\tau_{\lambda_1}}{\tau_{\lambda_2}}\right)}{\ln\left(\frac{\lambda_1}{\lambda_2}\right)}$$

50 where τ represents the AOD at a specific wavelength λ . We prioritized the use
51 of the 500 nm and 675 nm pair for calculating α due to their proximity to 532 nm. In
52 cases where data at 675 nm were unavailable, the 870 nm or 440 nm channels were
53 used as secondary references. This processing ensures spectral consistency between the
54 ground-based validation dataset and the model outputs.

55 **S3. Input Variable**

56 Table S1. Summary of input variables from MERRA-2 reanalysis and GEOS-Chem
 57 model simulations used in this study, categorized by data source and physical property.

Category	Variable Name	Physical Description
MERRA-2 (3D Fields)	T	Air Temperature (3D)
	U, V	Eastward & Northward Wind Components
	OMEGA	Vertical Pressure Velocity
	QV	Specific Humidity
	RH	Relative Humidity
	CLOUD	Cloud Fraction
	OPTDEPTH	In-cloud Optical Thickness
	QI, QL	Mass Fraction of Cloud Ice / Liquid Water
	TAUCLI, TAUCLW	Optical Thickness (Ice / Liquid Clouds)
	DTRAIN	Detrainment Mass Flux
	CMFMC	Cumulative Mass Flux
	PFICU, PFILSAN	Ice Precipitation Flux (Conv / LS)
	PFLCU, PFLLSAN	Liquid Precipitation Flux (Conv / LS)
	DQRCU, DQRLSAN	Rainwater Source (Conv / LS)
	REEVAPCN, REEVAPLS	Evaporation of Precip (Conv / LS)
PV	Ertel's Potential Vorticity	
MERRA-2 (Surface/2D)	PS, SLP	Surface Pressure / Sea Level Pressure
	TS, T2M, T10M	Skin Temp / 2m Temp / 10m Temp
	U10M, V10M	10-meter Wind Components
	QV2M / Q850	Specific Humidity (2m / 850hPa)

PBLH	Planetary Boundary Layer Height
PRECTOT	Total Precipitation
PRECCON, PRECLSC	Precipitation (Convective / Large Scale)
PRECANV, PRECSNO	Precipitation (Anvil / Snow)
SNODP / SNOMAS	Snow Depth / Snow Mass
ALBEDO	Surface Albedo
CLDTOT	Total Cloud Area Fraction
SWGDN	Incident Shortwave Flux
LWGNT, LWTUP	Net Downward / Upward Longwave Flux
EFLUX, HFLUX	Latent / Sensible Heat Flux
EVAP	Surface Evaporation
GWETROOT, GWETTOP	Soil Wetness (Root / Top)
LAI, GRN	Leaf Area Index / Greenness
FRSNO, FRSEAICE	Fraction of Snow / Sea Ice
SEAICE00-90	Sea Ice Fraction (Binned 00-90)
TO3	Total Column Ozone
TROPPT	Tropopause Pressure
USTAR, Z0M	Friction Velocity / Roughness Length
LWI	Land Water Indicator
PARDF, PARDR	PAR (Diffuse / Direct)

GEOS-Chem (Aerosol Mass)	AerMassSO ₄	Sulfate Aerosol Mass
	AerMassBC, AerMassPOA	Black Carbon / Primary Organic Aerosol
	AerMassOPOA	Oxidized Primary Organic Aerosol
	AerMassNIT, AerMassNH ₄	Nitrate / Ammonium Aerosol Mass

	AerMassSAL	Sea Salt Aerosol Mass
	AerMassSOA*	Secondary Organic Aerosol
	AerMassHMS	Hydroxymethanesulfonate Aerosol
	AerMassLVOCOA	Low Volatility Oxygenated OA
	TotalOA	Total Organic Aerosol
	PM _{2.5} , PM ₁₀	PM _{2.5} / PM ₁₀ Mass
	EC_gc	Elemental Carbon Tracer
GEOS-Chem (Optics/Micro)	AODHyg532nm_*	Hygroscopic AOD at 532nm
	AODDust	Dust Optical Depth
	AerHygroscopicGrowth_*	Hygroscopic Growth Factors
	AerSurfAreaHyg_*	Hygroscopic Aerosol Surface Area
	Chem_WetAeroRadi*	Wet Aerosol Radius
	Chem_AeroRadi*	Dry Aerosol Radius
	Chem_WetAeroArea*	Wet Aerosol Surface Area
	Chem_AeroArea*	Dry Aerosol Surface Area
GEOS-Chem (Gas Species)	SpeciesConc_SO ₂	Sulfur Dioxide
	SpeciesConc_NO, NO ₂	Nitrogen Oxides
	SpeciesConc_O ₃	Ozone
	SpeciesConc_NH ₃	Ammonia
	SpeciesConc_HNO ₃	Nitric Acid
	SpeciesConc_CO	Carbon Monoxide
	SpeciesConc_OH	Hydroxyl Radical
	SpeciesConc_HO ₂ , H ₂ O ₂	Hydroperoxyl / Hydrogen Peroxide
SpeciesConc_CH ₂ O	Formaldehyde	

SpeciesConc_VOCs

VOCs (Isoprene, MVK, MACR, Aromatics...)

SpeciesConc_S-Map

Sulfur Species (DMS, MSA, HMS)

59 S4. Detailed Architecture of the Transformer Components

60 a. Heterogeneous Input Embedding Strategies

61 Unlike the standard Transformer which uses a uniform positional encoding, we
62 adopt distinct embedding strategies tailored to the physical nature of each input variable:

63 ① **Hybrid Vertical Height Encoding:** To capture the stratification of the atmosphere,
64 we combine a learnable discrete embedding with a continuous linear projection. For a
65 vertical layer with index l and geopotential height value h_l :

$$66 \quad e_{height}^{(l)} = Linear \left(Concat(E_{idx}[l], Linear(h_l)) \right) \quad (S1)$$

67 where E_{idx} is a learnable lookup table for layer indices.

68 ② **Cyclic Spatiotemporal Encoding:** For variables with intrinsic periodicity (i.e.,
69 Month, Latitude, Longitude), we utilize trigonometric cyclic encoding to preserve their
70 continuity (e.g., ensuring December is numerically close to January). For a variable v
71 with a period T (e.g., $T=12$ for months, $T=360$ for longitude), the encoding is defined
72 as:

$$73 \quad e_{cyclic} = Linear \left(\left[\sin \left(\frac{2\pi v}{T} \right), \cos \left(\frac{2\pi v}{T} \right) \right] \right) \quad (S2)$$

74 ③ **Categorical Embedding:** Binary variables (e.g., Day/Night flag $d \in \{0,1\}$) are
75 mapped to dense vectors using standard learnable embeddings:

$$76 \quad e_{dn} = E_{dn}[d] \quad (S3)$$

77 b. Variable Identity Embedding

78 For the global context inputs $X_{global} \in \mathbb{R}^{N_{var}}$, where each scalar x_j represents a
79 distinct physical variable (e.g., PBLH, Surface Pressure), we use a specialized
80 projection:

$$81 \quad e_{global}^{(j)} = LayerNorm(x_j \cdot \omega_j + b_j + id_j) \quad (S4)$$

82 where $\omega_j \in \mathbb{R}^{d_{model}}$ is a variable-specific projection vector (broadcasted
83 multiplication), b_j is a bias term, and id_j is a learnable identity embedding unique to
84 the j -th variable type. This ensures the model distinguishes between variables even if
85 their scalar values are similar.

86 **c. Standard Multi-Head Self-Attention (MSA)**

87 The Transformer encoder utilizes MSA to capture dependencies between different
88 vertical layers. For a given input sequence X , the attention output is calculated as:

$$89 \quad \text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (S5)$$

90 where $Q = XW^Q$, $K = XW^K$, and $V = XW^V$ are the query, key, and value matrices
91 obtained through linear transformations. In our multi-head setting, the outputs are
92 concatenated and linearly projected:

$$93 \quad \text{MSA}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O \quad (S6)$$

94 **d. Position-wise Feed-Forward Network (FFN)**

95 Each encoder layer includes a fully connected FFN with a GELU activation function:

$$96 \quad \text{FFN}(x) = \text{Linear}_2\left(\text{GELU}(\text{Linear}_1(x))\right) \quad (S7)$$

97 **e. Gated Feature Fusion**

$$98 \quad E_{local}^{(z)} = \sum_{k=1}^4 \alpha_k^{(z)} \cdot \text{Norm}(h_k^{(z)}) \text{ s.t. } \sum \alpha_k = 1 \quad (S8)$$

99 where $h_k^{(z)}$ represents the embedding of the k -th feature group. The attention weights
100 $\alpha_k^{(z)}$ are generated by a learnable gating network, allowing the model to autonomously
101 identify and prioritize the most physically relevant information source for each specific
102 altitude layer.

103 **f. Output Layer**

$$104 \quad \Delta_{AEC}^{pred} = \text{MLP}(H_{cross} + H_{phyche}) \cdot s + b \quad (S9)$$

105 Where Δ_{AEC}^{pred} denotes the predicted systematic bias of the AEC. The parameters s
106 and b are learnable scaling and bias terms, respectively, introduced to adaptively map
107 the normalized network outputs to the physical magnitude of extinction biases.

108 **g. Magnitude-Weighted Loss**

109 Aerosol extinction exhibits a high dynamic range and severe spatial heterogeneity,
110 statistically dominated by near-zero background values (e.g., in the clean free

111 troposphere). Standard Mean Squared Error (MSE) treats all data points equally, which
 112 often causes the model to over-fit these overwhelming background signals while under-
 113 correcting the large systematic biases associated with severe pollution episodes. To
 114 resolve this imbalance, we propose a Magnitude-Weighted Loss (L_{MW}).

115

$$116 \quad L_{MW} = \frac{1}{N} \sum_{i=1}^N (\Delta_{AEC}^{pred} - \Delta_{AEC}^{target})^2 \cdot \omega(\Delta_{AEC}^{target}) \quad (S10)$$

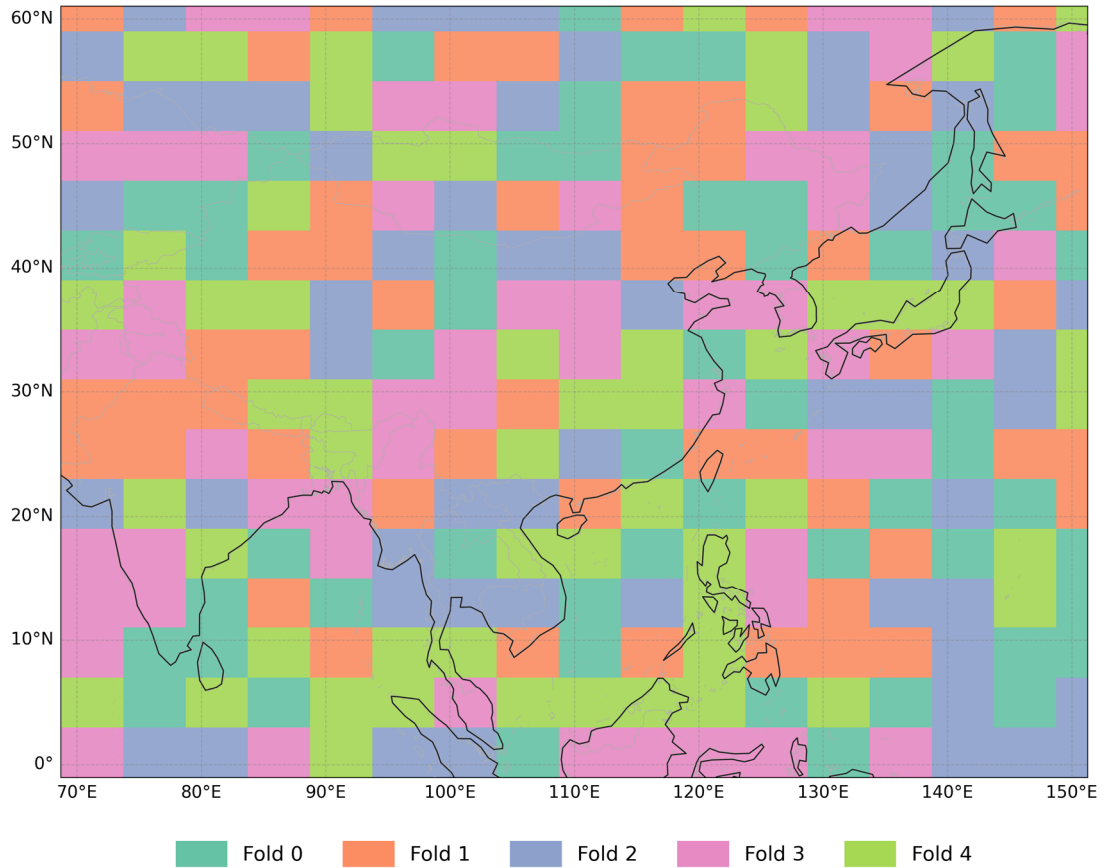
$$117 \quad \omega(\Delta_{AEC}^{target}) = |\Delta_{AEC}^{target}|^p + \lambda \cdot e^{-\beta |\Delta_{AEC}^{target}|} + \epsilon \quad (S11)$$

118 Here, the weight function ω dynamically rescales the optimization penalty based on
 119 the magnitude of the true bias (Δ_{AEC}^{target}), utilizing three empirically derived components.

120 (1) Large-Error Prioritization ($|\Delta_{AEC}^{target}|^p$): Governed by the exponent p , this term
 121 sharply amplifies the gradient for samples with massive simulation biases. It forces the
 122 network to target critical GEOS-Chem failures, such as significantly underestimated
 123 dust storms or heavy anthropogenic haze, rather than safely fitting the background
 124 average. (2) False Alarm Penalty ($\lambda \cdot e^{-\beta |\Delta_{AEC}^{target}|}$): Regulated by constants λ and β ,
 125 this exponential decay applies a strict penalty when the true GEOS-Chem bias is
 126 negligible but the model attempts a non-zero correction. Physically, this prevents the
 127 generation of spurious aerosol artifacts in regions where the GEOS-Chem is already
 128 accurate, such as the clean free troposphere. (3) Base Stability Term (ϵ): A minor
 129 constant added to maintain numerical stability during gradient descent.

130

131 **S5. 5-fold Cross-Validation**



132

133 Figure S1. The illustration of the data splitting for the cross-validation. The color-coded
134 checkerboard pattern represents the spatial partitioning of data into five folds used for
135 model evaluation. Specifically, the geographical domain is divided into grid blocks,
136 which are then randomly assigned to one of the five folds.

137 **S6. Statistical Evaluation Metrics**

138 To quantify the model performance in reconstructing AEPs, the following
 139 statistical metrics are employed. Let $\Delta_{AEC}^{target}_i$ and $\Delta_{AEC}^{pred}_i$ denote the observed
 140 (CALIOP) and predicted (Model) AEC for the i -th sample, respectively, and N be the
 141 total number of samples.

142 **a. Pearson Correlation Coefficient (R)**

143 Reflects the linear consistency between predictions and observations:

144
$$R = \frac{\sum_{i=1}^N (\Delta_{AEC}^{target}_i - \overline{\Delta_{AEC}^{target}}) (\Delta_{AEC}^{pred}_i - \overline{\Delta_{AEC}^{pred}})}{\sqrt{\sum_{i=1}^N (\Delta_{AEC}^{target}_i - \overline{\Delta_{AEC}^{target}})^2} \sqrt{\sum_{i=1}^N (\Delta_{AEC}^{pred}_i - \overline{\Delta_{AEC}^{pred}})^2}} \quad (S12)$$

145 where $\overline{\Delta_{AEC}^{target}}$ and $\overline{\Delta_{AEC}^{pred}}$ are the means of observed and predicted values.

146 **b. Error Magnitude Metrics**

147 ① Mean Absolute Error (MAE): Represents the average magnitude of errors.

148
$$MAE = \frac{1}{N} \sum_{i=1}^N |\Delta_{AEC}^{target}_i - \Delta_{AEC}^{pred}_i| \quad (S13)$$

149 ② Root Mean Square Error (RMSE): Sensitive to large errors and outliers.

150
$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (\Delta_{AEC}^{target}_i - \Delta_{AEC}^{pred}_i)^2} \quad (S14)$$

151 **c. Normalized Root Mean Square Error (NRMSE)**

152 To compare performance across vertical layers with varying orders of magnitude in
 153 extinction coefficients, we normalize the RMSE by the range of the observed data:

154
$$NRMSE = \frac{RMSE}{\Delta_{AEC}^{target}_{max} - \Delta_{AEC}^{target}_{min}} \quad (S15)$$

155 where $\Delta_{AEC}^{target}_{max}$ and $\Delta_{AEC}^{target}_{min}$ are the maximum and minimum observed
 156 extinction coefficients in the target dataset, respectively.

157 **d. Mean Bias (Bias)**

158 Used to diagnose systematic overestimation ($Bias > 0$) or underestimation ($Bias < 0$):

$$Bias = \frac{1}{N} \sum_{i=1}^N (\Delta_{AEC\ i}^{pred} - \Delta_{AEC\ i}^{target}) \quad (S16)$$

160 **S7. Mathematical Formulation of Interpretability Methods**

161 To quantify the contribution of different input features to the bias correction, we
 162 employ ensemble-based interpretability methods. The mathematical definitions for the
 163 Gradient-based Attribution, Attention Mechanism, and Permutation Feature
 164 Importance are detailed below.

165 **a. Gradient-based Feature Attribution (Input \times Gradient)**

166 For the vertical profile inputs, we utilize the Input \times Gradient method to quantify
 167 the sensitivity of the model output to local feature variations. The importance score
 168 (I_{feat}) for a specific feature at a specific vertical layer is calculated as the ensemble
 169 average of the absolute gradients weighted by the input magnitude:

$$170 \quad I_{feat} = \frac{1}{M} \sum_{m=1}^M \left(\frac{1}{N} \sum_{n=1}^N \left| x_{feat}^{(n)} \odot \frac{\partial \mathcal{F}_m(x^{(n)})}{\partial x_{feat}^{(n)}} \right| \right) \quad (S17)$$

171 where: N is the total number of samples in the validation set. M is the number of
 172 cross-validation folds (here $M = 5$). $x_{feat}^{(n)}$ denotes the input value of a specific feature
 173 for the n -th sample. $\partial \mathcal{F}_m(\cdot)$ represents the model prediction function for the m -th
 174 fold. \odot denotes element-wise multiplication.

175 **b. Global Attention Mechanism**

176 Global Attention Mechanism For global meteorological covariates (i.e., 2D single-
 177 layer variables from MERRA-2), we analyze the Cross-Attention weights extracted
 178 from the Transformer decoder. The global average attention score (\bar{A}_{ij}) represents the
 179 interaction strength between the i -th vertical layer of the profile and the j -th global
 180 variable (e.g., PBLH, Surface Pressure). It is computed by averaging across all attention
 181 heads and all samples:

$$182 \quad \bar{A}_{ij} = \frac{1}{N} \sum_{n=1}^N \frac{1}{H} \sum_{h=1}^H A_{ij}^{(n,h)} \quad (S18)$$

183 where: H is the number of attention heads. N is the total number of samples in the
 184 dataset. $A_{ij}^{(n,h)}$ is the attention weight matrix from the h -th head for the n -th sample.

185 The inner summation averages the contributions from multi-head attention mechanisms,

186 while the outer summation computes the dataset-wide average contribution of each
 187 global variable, filtering out sample-specific noise.

188 **c. Gated Fusion Weights**

189 The model dynamically fuses information from different sources using learnable
 190 scalar weights. The effective contribution weight (W_c) for a specific component c at
 191 altitude z is formalized as:

$$192 \quad W_c(z) = \sigma \left(\frac{1}{N} \sum_{n=1}^N \omega_c^{(n)}(z) \right) \quad (S19)$$

193 where $\omega_c^{(n)}(z)$ is the raw gate value predicted by the model for sample n at altitude
 194 z , and $\sigma(\cdot)$ represents the normalization function (e.g., Softmax) ensuring the weights
 195 sum to 1.

196 **d. Permutation Feature Importance**

197 To assess the model's reliance on specific feature groups (e.g., Profile Features vs.
 198 Global Variables), we calculate the percentage increase in MSE when a feature group
 199 is randomly permuted. The importance score (S_k) for feature group k is defined as:

$$200 \quad S_k = \frac{\mathcal{L}_{perm,k} - \mathcal{L}_{base}}{\mathcal{L}_{base}} \times 100\% \quad (S20)$$

$$201 \quad \mathcal{L}_{base} = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - \hat{y}^{(n)})^2 \quad (S21)$$

202 where: \mathcal{L}_{base} is the baseline MSE on the original validation set. $\mathcal{L}_{perm,k}$ is the MSE
 203 calculated after randomly shuffling the feature group k along the batch dimension
 204 while keeping other features fixed.

205 **S8. Implementation of SHAP Analysis**

206 To efficiently compute SHAP values for the deep neural network while
207 maintaining statistical representativeness, we adopt a clustering-based sampling
208 strategy for the background dataset:

209 **a. Representative Background Selection**

210 Instead of using a random subset or the zero-mean baseline (which is physically
211 unrealistic for atmospheric variables), we apply K-means clustering on the normalized
212 training dataset. We set the number of clusters $k=100$ and selected the medoids (samples
213 closest to the cluster centers) to form a background dataset. This ensures that the
214 reference baseline covers the full manifold of atmospheric states, from clean
215 background days to severe pollution episodes.

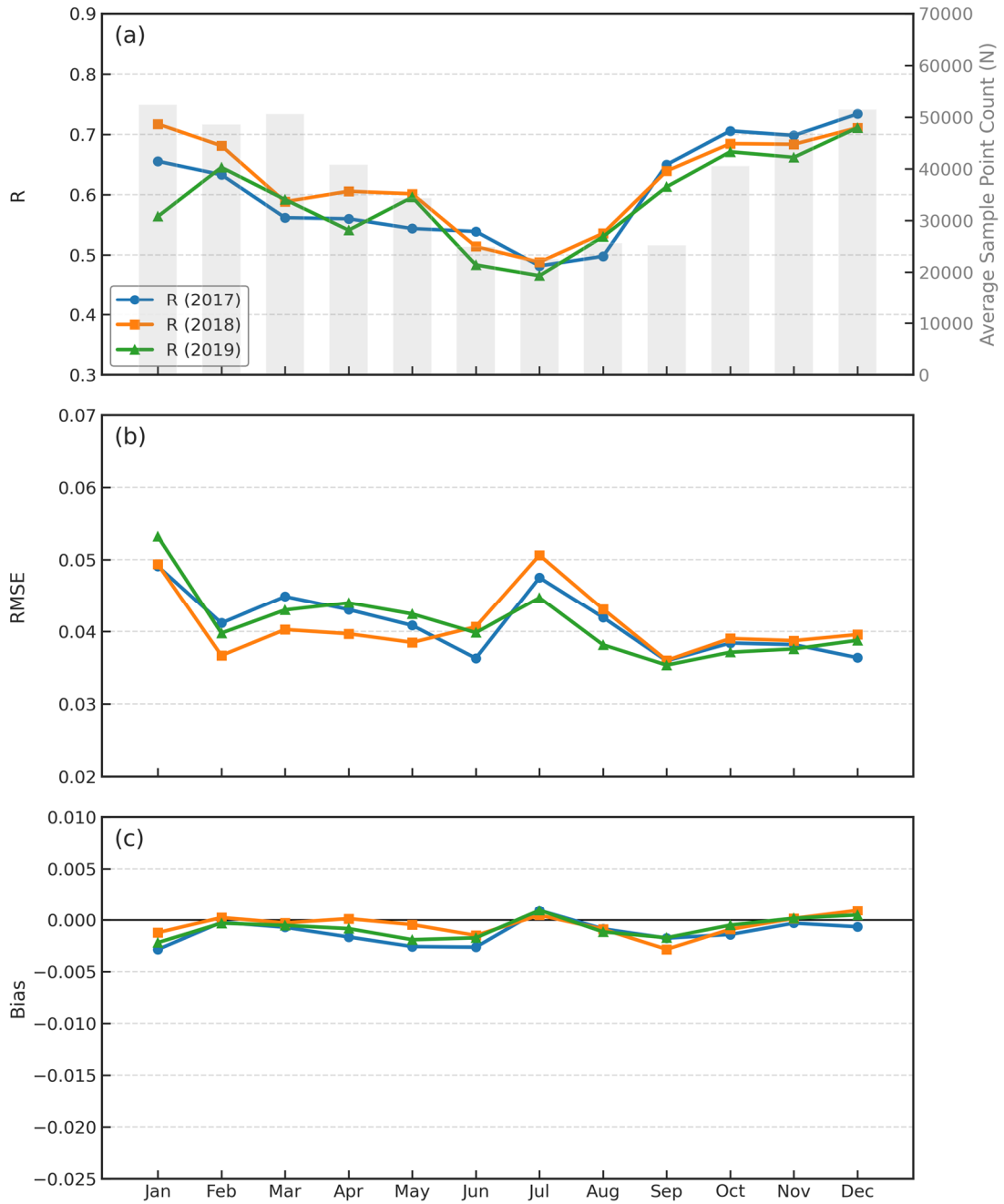
216 **b. Target Sample Selection**

217 For regional analysis, we select top-50 samples with the highest model corrections
218 (i.e., large initial GEOS-Chem biases that were successfully corrected by the model)
219 within each region of interest. This "Case Study Strategy" ensures that the SHAP
220 explanation focuses on the most significant bias correction events rather than trivial
221 noise.

222 **S9. Supplementary Evaluation of the Overall Model**

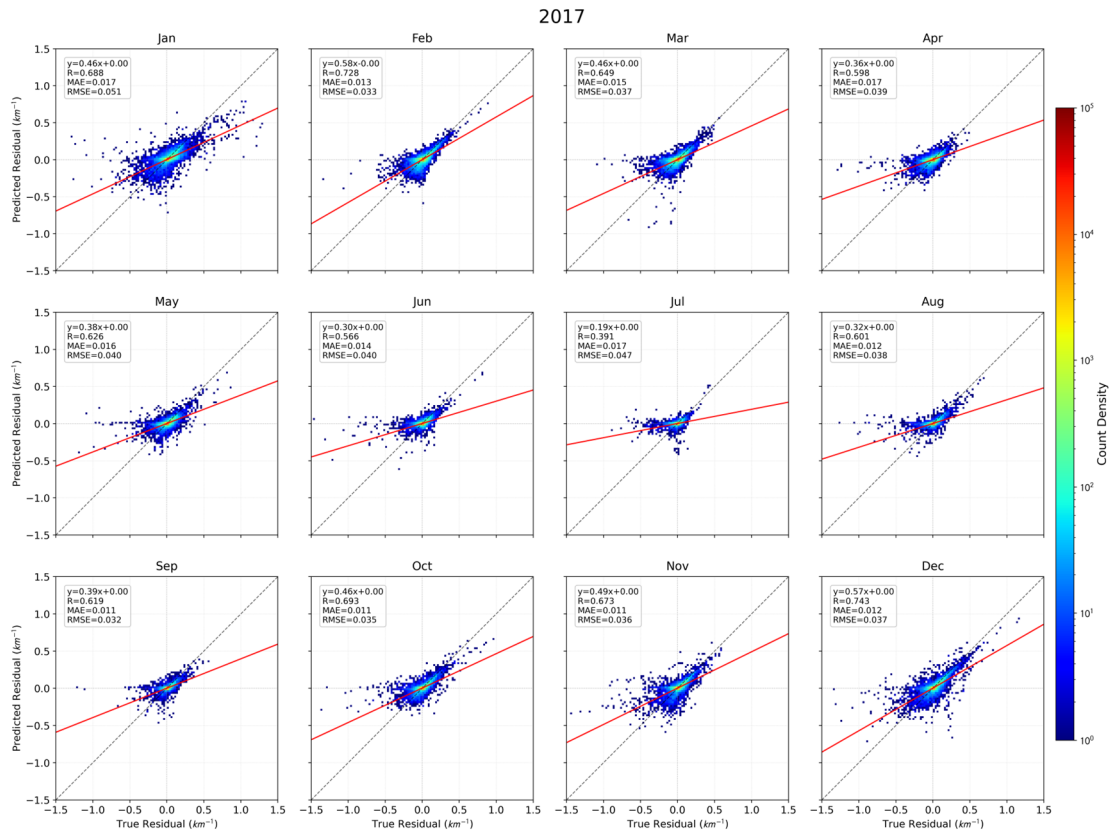
223 Table S2. Summary of model performance for AEC simulation biases (Δ_{AEC}) on the
224 validation subsets derived from the 5-fold spatial cross-validation strategy during the
225 training phase.

Experiment ID	R (Val)	MAE (Val)	RMSE (Val)
Exp-2017	0.629	0.015	0.042
Exp-2018	0.652	0.015	0.041
Exp-2019	0.609	0.015	0.042
Average	0.630	0.015	0.042



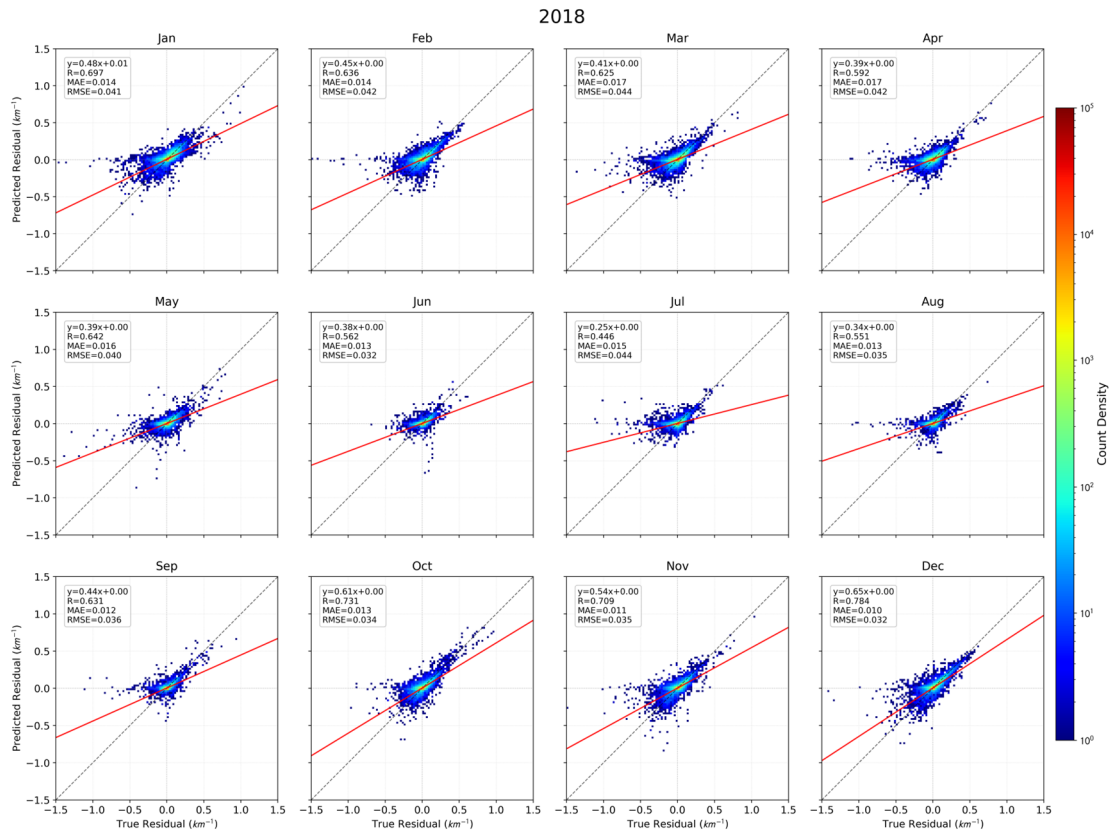
226

227 Figure S2. Monthly variations of the statistical evaluation metrics (derived from the 5-
 228 fold cross-validation validation subsets) for assessing the physics-informed
 229 Transformer model's performance in predicting the AEC simulation bias (Δ_{AEC}). The
 230 panels display the time series of R (a) with the multi-year average monthly sample size
 231 (N, gray bars), RMSE (b), and mean bias (c).



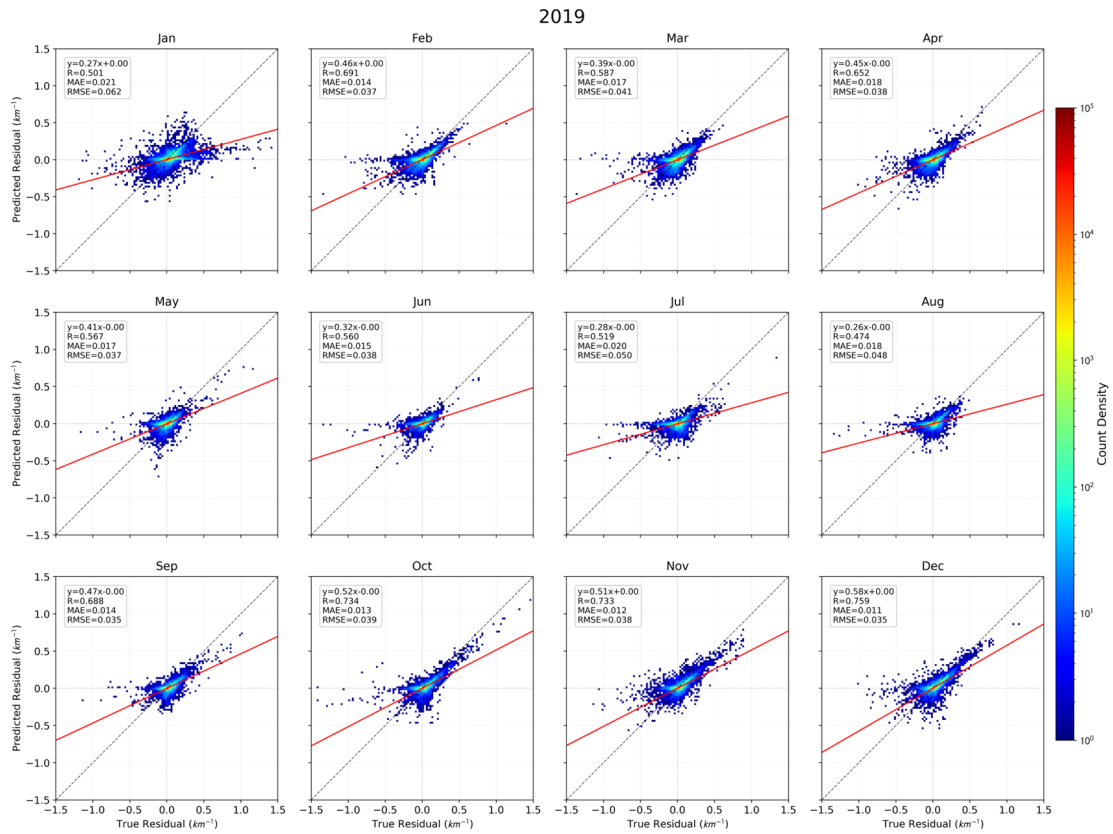
232

233 Figure S3. Monthly density scatter plots comparing the predicted AEC simulation bias
 234 (y-axis) generated by the physics-informed Transformer against the true residuals (x-
 235 axis) for the year 2017. The true residual is defined as the difference between the
 236 original GEOS-Chem simulation AEC and CALIOP observation AEC (GEOS-
 237 Chem – CALIOP).



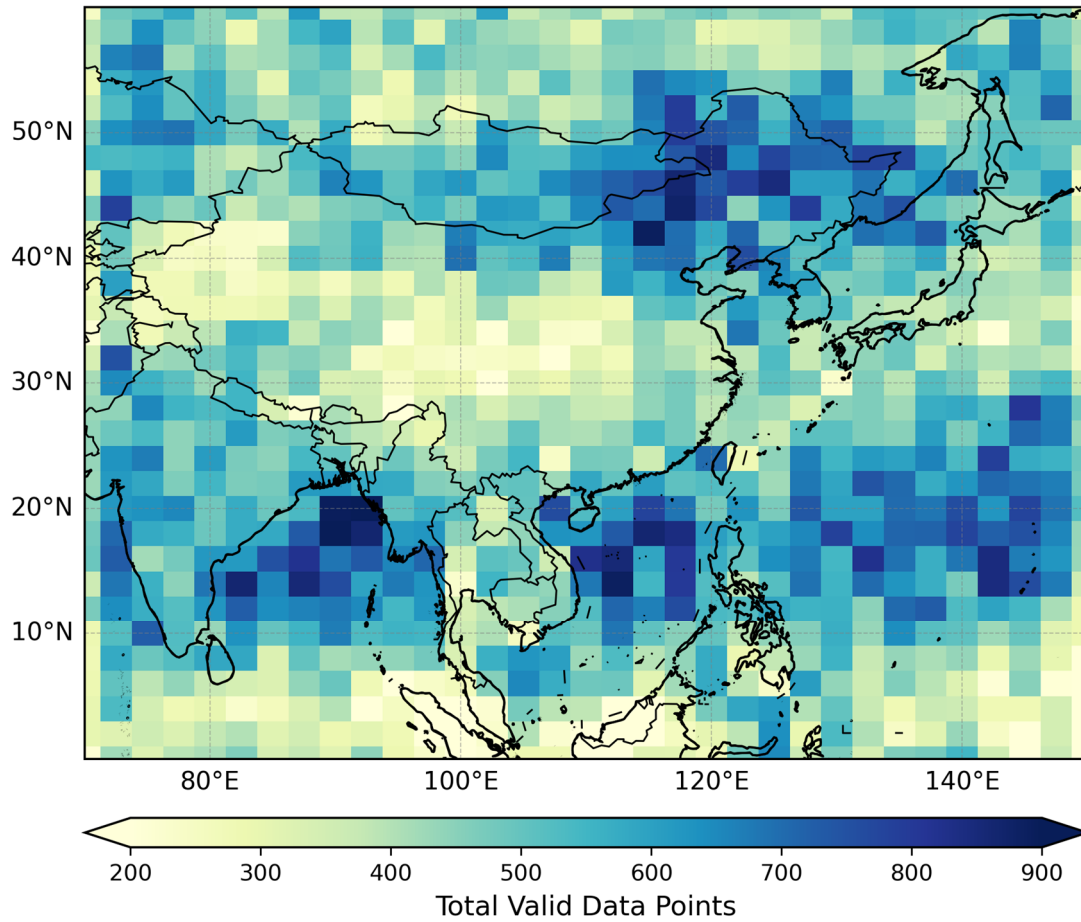
238

239 Figure S4. Monthly density scatter plots comparing the predicted AEC simulation bias
 240 (y-axis) generated by the physics-informed Transformer against the true residuals (x-
 241 axis) for the year 2018.



242

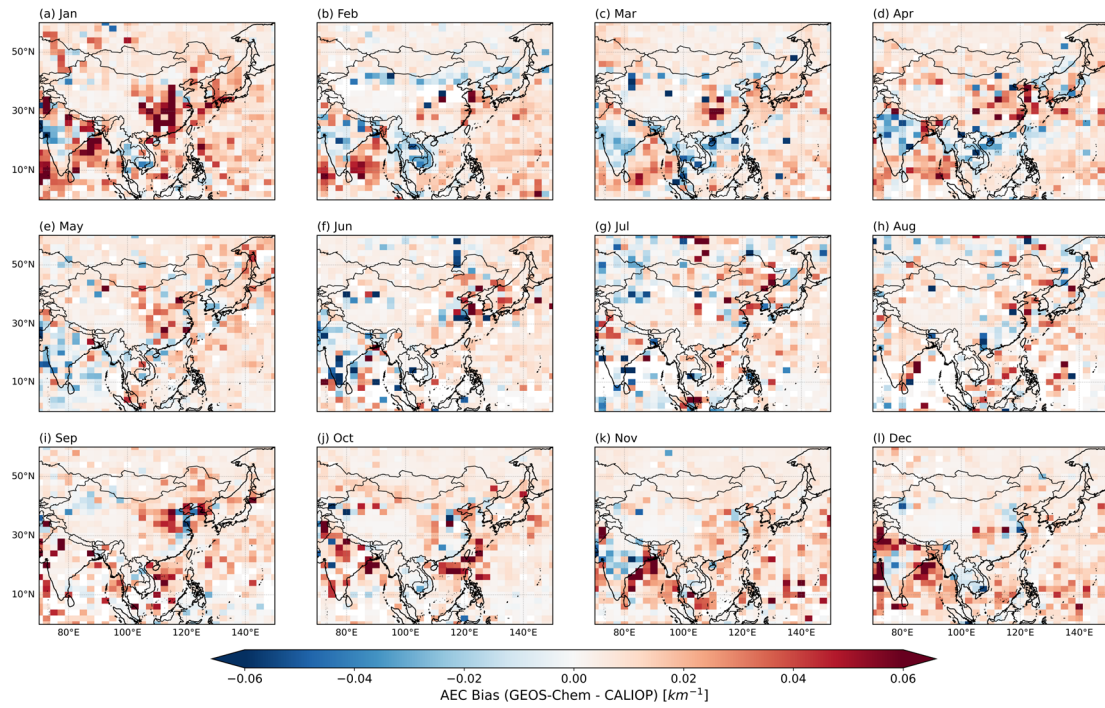
243 Figure S5. Monthly density scatter plots comparing the predicted AEC simulation bias
 244 (y-axis) generated by the physics-informed Transformer against the true residuals (x-
 245 axis) for the year 2019.



246

247 Figure S6. Spatial distribution of valid CALIOP data points aggregated for the test year

248 2019.



249

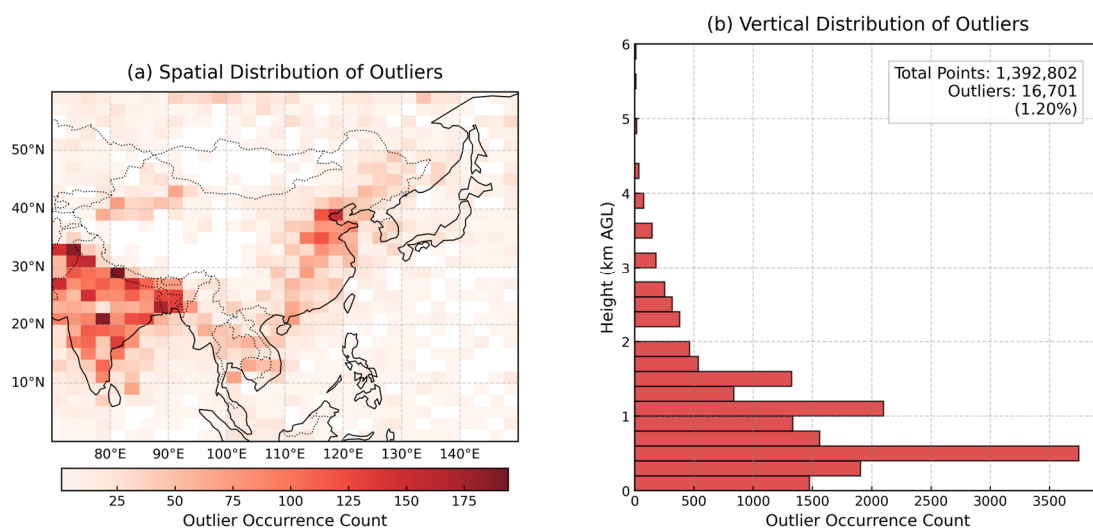
250 Figure S7. Monthly spatial evolution of the Target AEC simulation bias (Δ_{AEC}) for the
 251 test year 2019.

252 S10. Analysis of Residual Outliers

253 This section analyzes the samples where the absolute residuals between the model-
254 corrected AEC and CALIOP observations exceed 0.15 km^{-1} (i.e., points falling outside
255 the error envelope in Fig. 6). The diagnostic results presented in Figure S8 demonstrate
256 that these outliers are not randomly distributed but exhibit specific spatial and vertical
257 characteristics.

258 Spatially (Fig. S8a), regions with elevated outlier occurrence are predominantly
259 anchored over major emission source regions, including the IGP, the NCP, and the
260 Indochina Peninsula. Vertically (Fig. S8b), over 80% of these extreme deviations are
261 confined within the PBL (below 1.5 km AGL), an altitude range characterized by the
262 heaviest local aerosol loading and the most intense turbulent mixing.

263 These distribution patterns confirm that the dispersion observed in the scatter plots
264 fundamentally originates from representativeness errors. In complex source regions, the
265 high-resolution footprint of CALIOP resolves transient, highly concentrated sub-grid
266 plumes. The fine-scale physical structures of these plumes are inherently smoothed out
267 during the spatial averaging onto the coarse $2^\circ \times 2.5^\circ$ grid of the GEOS-Chem model.
268 We attribute the remaining 1.20% of outlier samples primarily to the spatial
269 heterogeneity of aerosols that lies below the resolvable scale of the model grid.

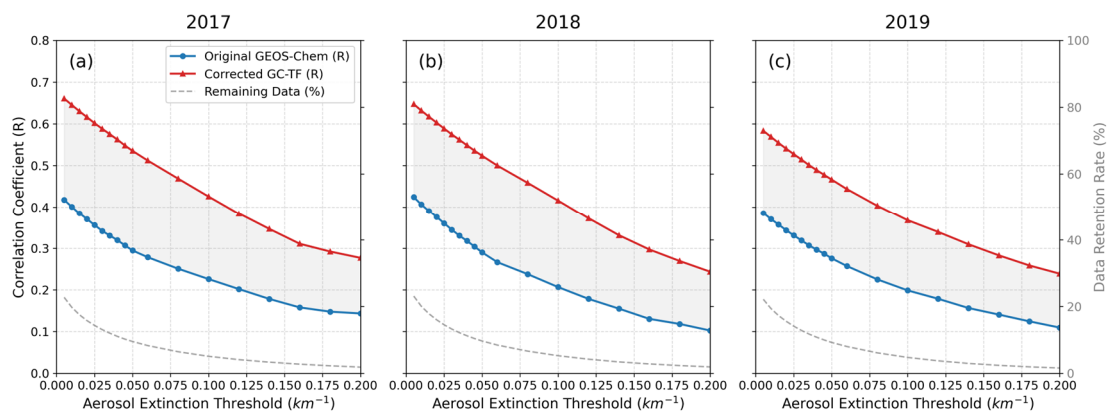


270

271 Figure S8. Spatial and vertical distributions of the residual outliers for the combined
272 test years (2017–2019). Outliers are strictly defined as samples where the absolute

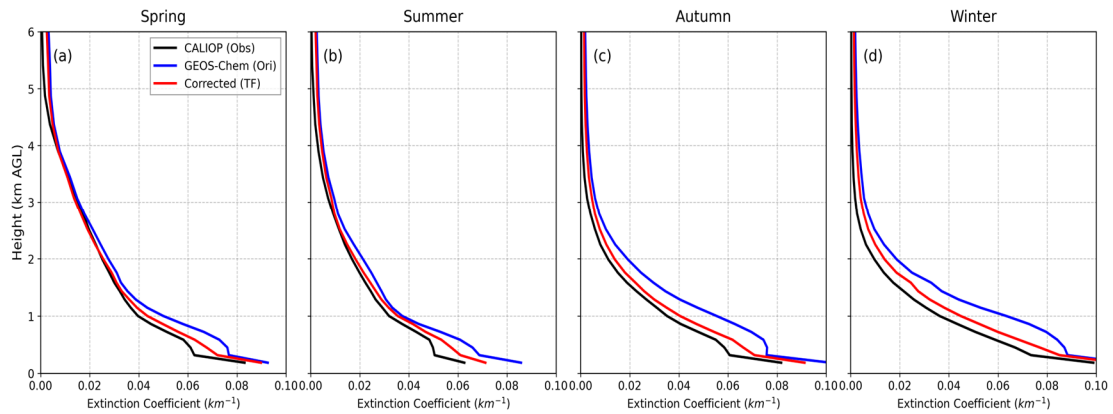
273 residual between the GC-TF prediction and CALIOP observation exceeds $\pm 0.15 \text{ km}^{-1}$.
274 The panels display (a) the spatial occurrence count of these outliers mapped onto the
275 native $2^\circ \times 2.5^\circ$ GEOS-Chem grid, and (b) their vertical distribution as a function of
276 height AGL. The statistical box indicates that these extreme deviations account for
277 merely 1.20% of the total valid dataset.

278 **S11. Supporting Information for Improvement Evaluation of GEOS-Chem**
279 **Simulations**



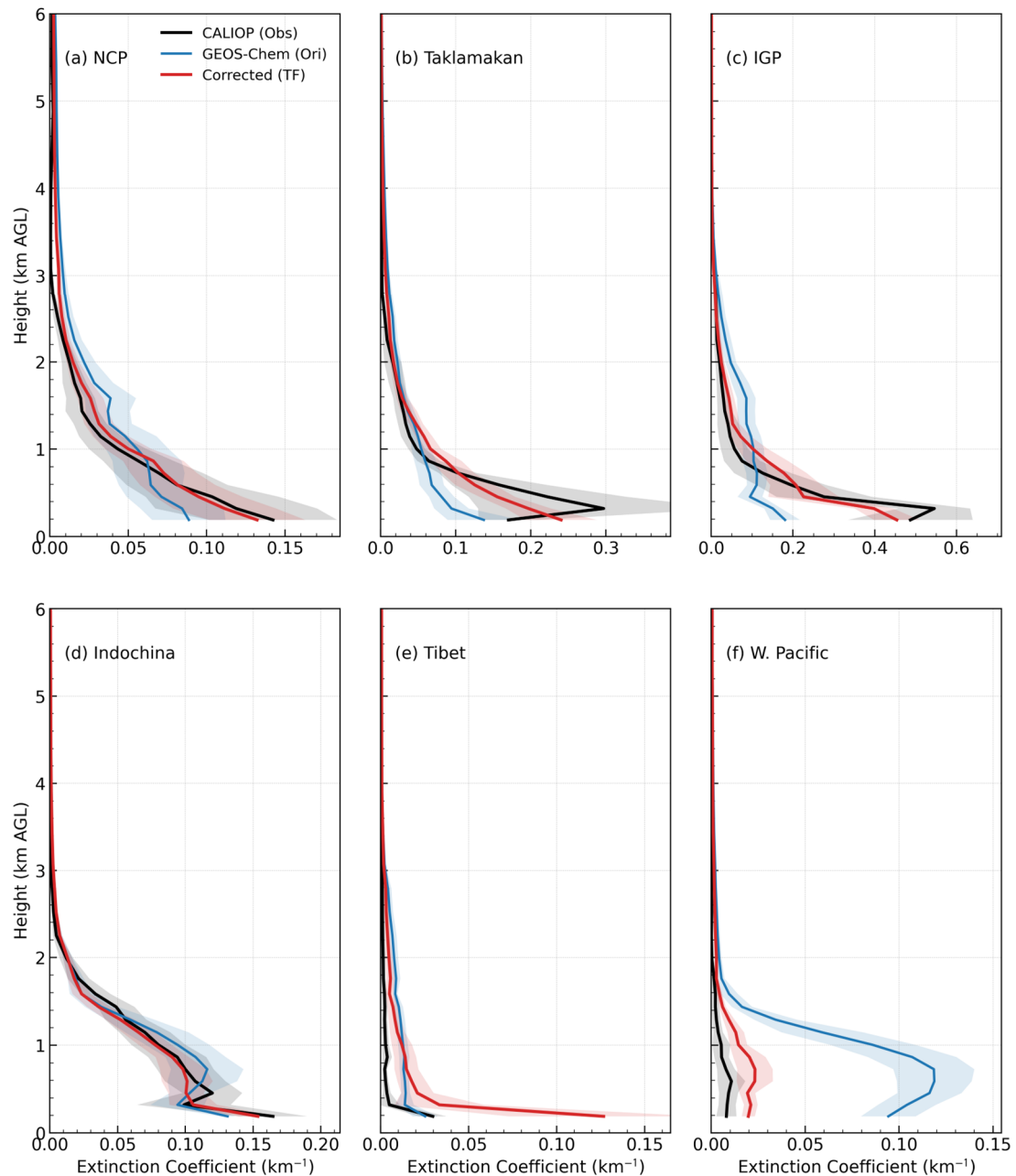
280

281 Figure S9. Sensitivity of model predictive performance to AEC thresholds across three
282 independent test years (2017–2019). The gray dashed line (corresponding to the right
283 axis) indicates the data retention rate, representing the percentage of valid samples
284 remaining after filtering out clean background signals below the threshold



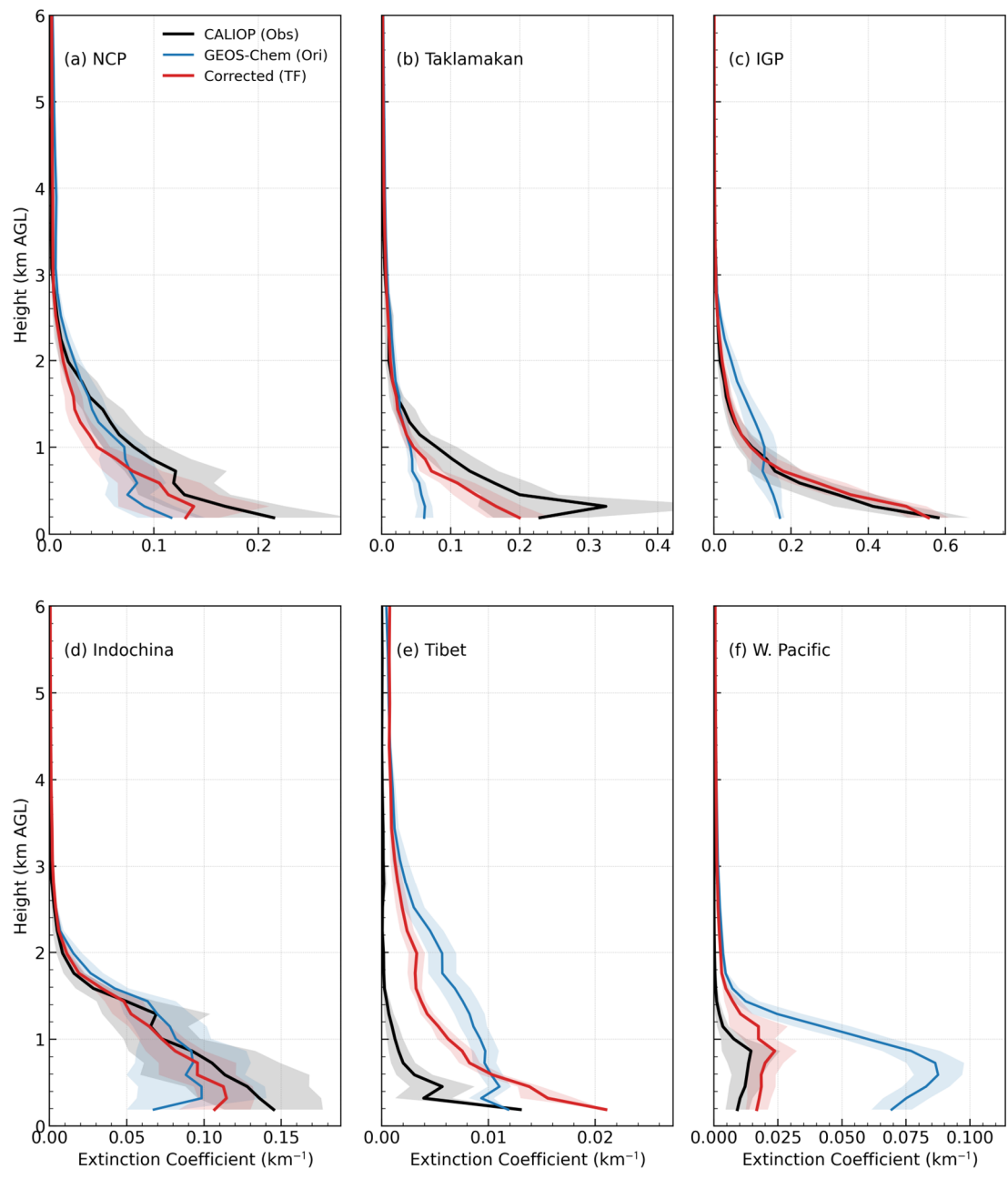
285

286 Figure S10. Seasonal mean vertical profiles of AEC (km^{-1}) averaged over the three test
 287 years (2017–2019). The panels display the profiles for Spring (MAM, a), Summer (JJA,
 288 b), Autumn (SON, c), and Winter (DJF, d). The black lines represent the CALIOP
 289 observations, the blue dashed lines denote the original GEOS-Chem simulations, and
 290 the red lines show the results corrected by the GC-TF model. The profiles are vertically
 291 resolved from the surface to 6 km above ground level (AGL).



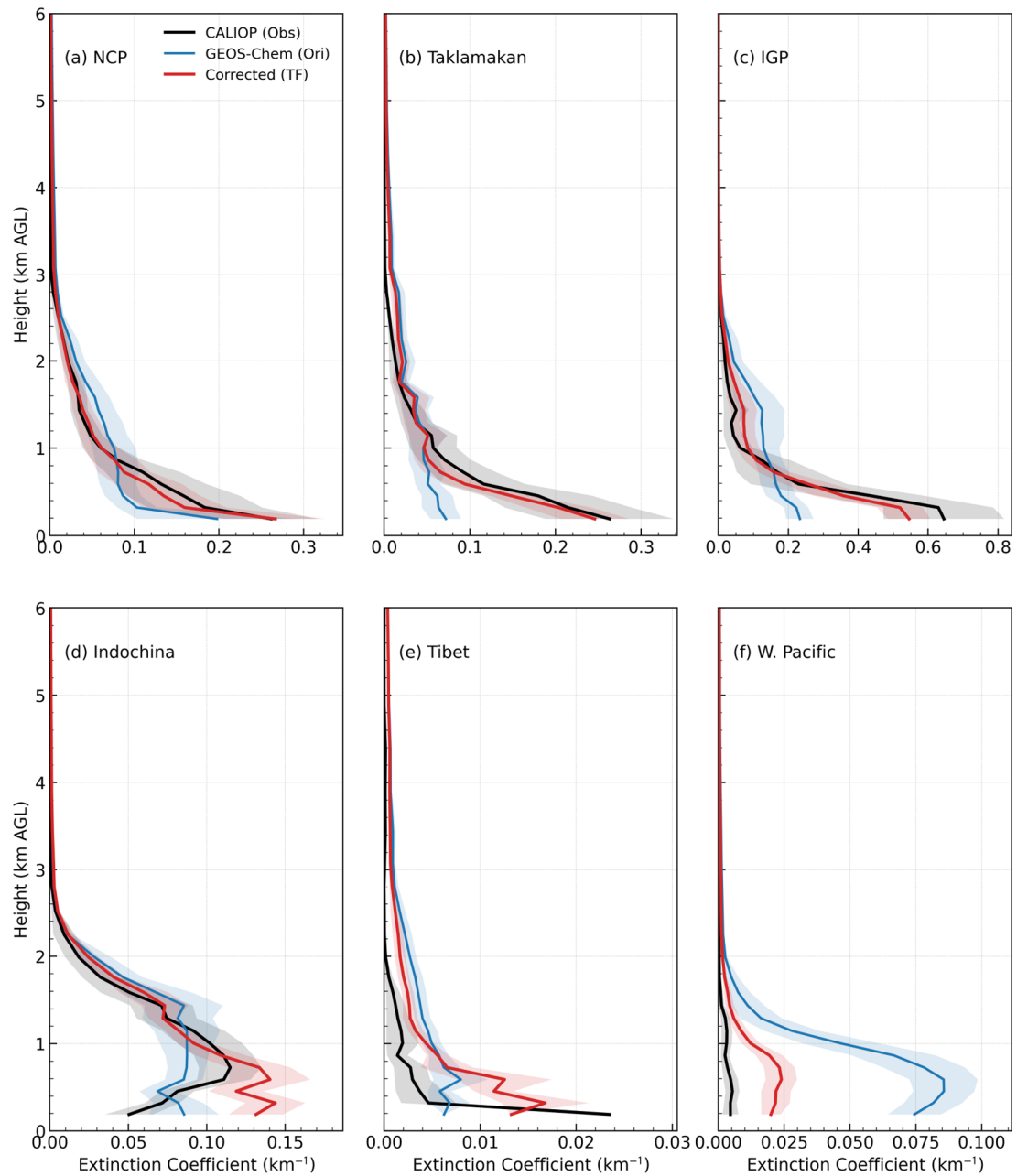
292

293 Figure S11. Regional annual mean vertical profiles of AEC over six representative sub-
 294 regions for the test year 2017. The panels display the profiles for NCP (a), Taklamakan
 295 Desert (b), IGP (c), Indochina (d), Tibetan Plateau (e), and Western Pacific (f). The
 296 black lines represent CALIOP observations, the blue dashed lines denote the original
 297 GEOS-Chem simulations, and the red lines show the results corrected by the GC-TF
 298 model. The shaded areas indicate the regional standard deviation ($\pm 0.5\sigma$), representing
 299 the spatial variability within each sub-region.



300

301 Figure S12. Regional annual mean vertical profiles of AEC over six representative sub-
 302 regions for the test year 2018.



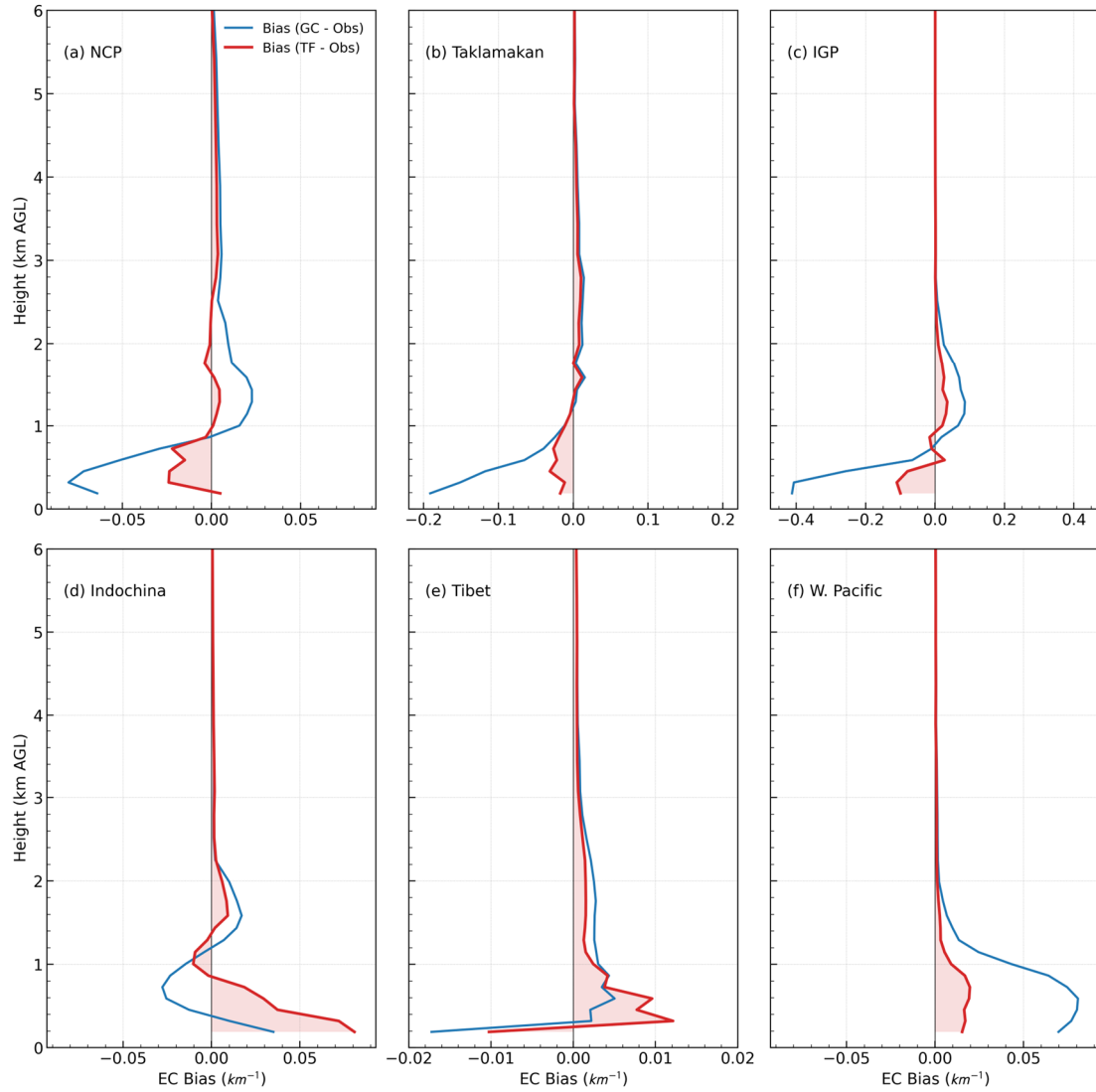
303

304

Figure S13. Regional annual mean vertical profiles of AEC over six representative sub-

305

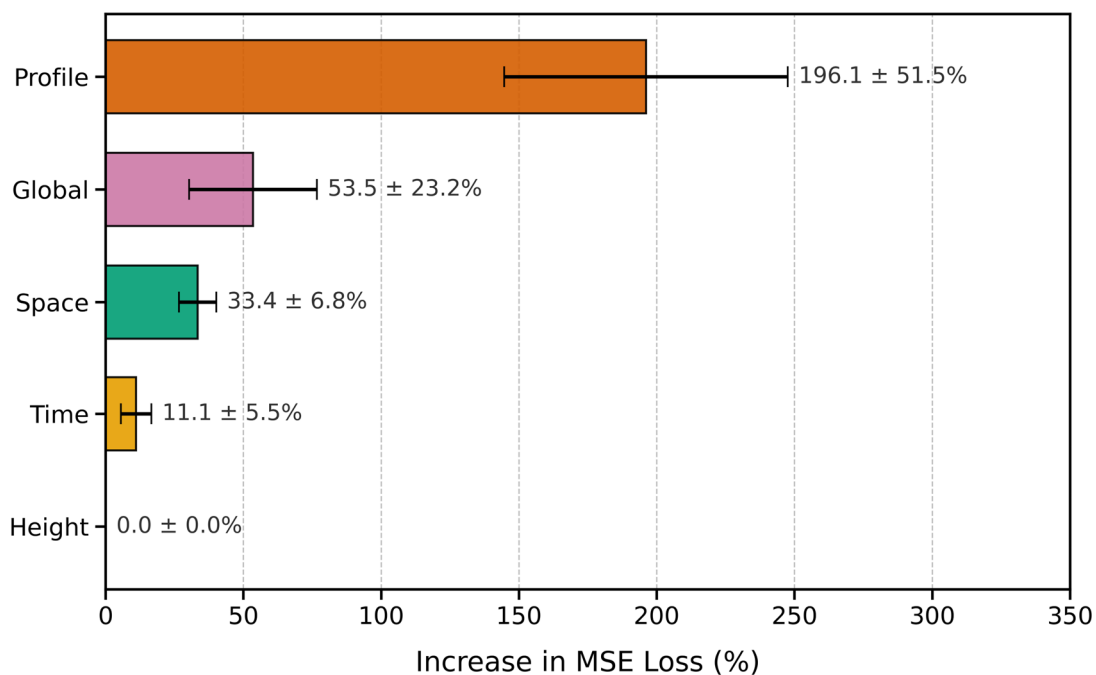
regions for the test year 2019.



306

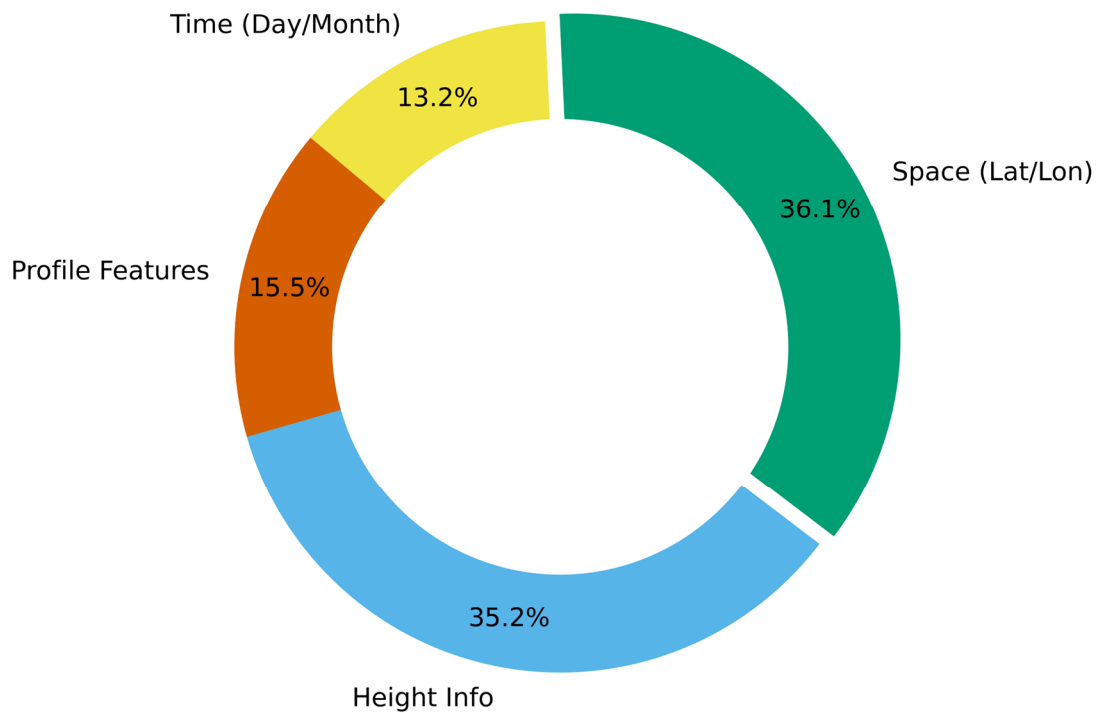
307 Figure S14. Vertical profiles of the annual mean AEC simulation bias over six
 308 representative sub-regions for the test year 2019. The panels display the bias profiles
 309 for NCP (a), Taklamakan Desert (b), IGP (c), Indochina (d), Tibetan Plateau (e), and
 310 Western Pacific (f). The AEC bias is calculated as Model minus Observation. The blue
 311 lines represent the original GEOS-Chem bias (GEOS-Chem–CALIOP), while the red
 312 lines denote the remaining bias after GC-TF correction (GC-TF – CALIOP). The
 313 vertical black line indicates the zero-bias reference, and the shaded red areas highlight
 314 the magnitude of the residual bias after correction.

315 **S12. Supporting Information for Interpretability Analysis**



316

317 Figure S15. Permutation feature importance analysis evaluated on the test dataset. The
318 importance of each feature group is quantified by the relative increase in MSE loss
319 when the values of that feature group are randomly shuffled, while keeping others
320 unchanged. The error bars represent the standard deviation across the 5-fold cross-
321 validation. Note that Profile Features are identified as the most critical input, causing a
322 ~200% surge in error if disrupted, confirming the model's heavy reliance on physical
323 state variables for magnitude prediction.



324

325 Figure S16. Global average relative importance of the input features learned by the
326 physics-informed model for the test year 2019. The pie chart displays the percentage
327 contribution of four feature groups: Space (36.1%), Height Info (35.2%), Profile
328 Features (15.5%, including GEOS-Chem chemical concentrations and MERRA2
329 meteorology), and Time (13.2%). The dominance of spatial and height information
330 reflects the strong spatial heterogeneity and vertical stratification of aerosol
331 distributions.

332 **S13. Sensitivity Analysis of CALIOP Observational Uncertainties**

333 To quantitatively evaluate the model's sensitivity to the inherent systematic
334 uncertainties of CALIOP observations, we design a perturbation experiment based on
335 the 2017 independent test set. Considering the reported mean relative bias of CALIOP
336 AOD against AERONET is approximately $-5.1\% \pm 8.5\%$, we artificially apply a
337 constant $\pm 5\%$ multiplier to the CALIOP AEC targets during the training phase. This
338 creates two extreme scenarios: one simulating a severe systematic underestimation (-5%
339 bias) and another simulating a systematic overestimation ($+5\%$ bias). The physics-
340 informed Transformer is then retrained from scratch for both scenarios, and the newly
341 predicted residual profiles are added back to the original GEOS-Chem simulations to
342 obtain the final perturbed corrected AEC profiles.

343 As presented in Table S3, the GC-TF model exhibits strong resistance to target
344 perturbations. Even with a 5% systematic error injected into the learning target, the
345 model's RMSE and Mean Bias remain highly stable and significantly superior to the
346 original GEOS-Chem baseline. Visually, Figure S17 illustrates that the perturbed
347 predictions (dashed and dotted lines) tightly hug the unperturbed baseline correction
348 (solid red line), forming a narrow perturbation range (red shaded area) that aligns well
349 with the CALIOP observations. This demonstrates that the data-driven correction
350 framework captures robust physical mappings rather than merely overfitting to
351 observational noise.

352

353

354

355

356

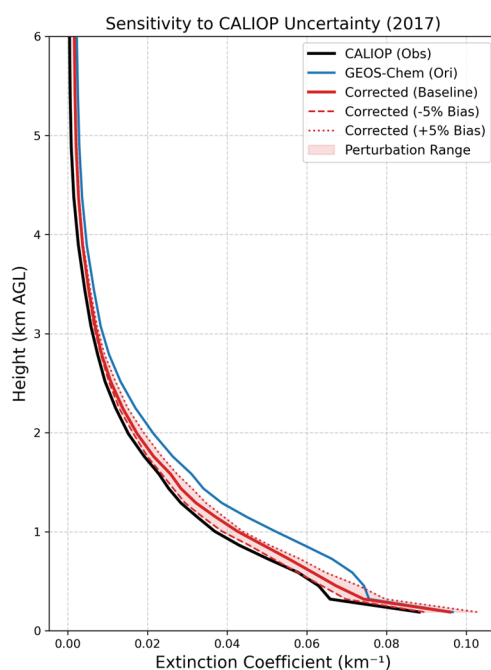
357

358

359

360 Table S3. Quantitative evaluation of the GC-TF model's sensitivity to CALIOP
 361 observational uncertainties on the 2017 independent test set. The metrics are derived
 362 by validating both the original GEOS-Chem simulations and the corrected GC-TF
 363 results against the unperturbed CALIOP AEC profiles.

Model/Scenario	RMSE	MAE	Mean Bias	R
Original GEOS-Chem	0.052	0.020	0.006	0.522
Baseline GC-TF (Unperturbed)	0.039	0.014	0.002	0.732
GC-TF (-5% CALIOP Bias)	0.040	0.014	0.001	0.716
GC-TF (+5% CALIOP Bias)	0.040	0.015	0.004	0.718



364
 365 Figure S17. Vertical profiles illustrating the sensitivity of the corrected AEC to target
 366 perturbations for the test year 2017. The black solid line represents the CALIOP
 367 observations, and the blue solid line represents the original GEOS-Chem simulation.
 368 The red solid line indicates the baseline correction of the GC-TF model trained on
 369 unperturbed data. The dashed and dotted red lines represent the corrected profiles
 370 generated by models trained with -5% and +5% perturbed CALIOP targets, respectively.
 371 The light red shaded area denotes the envelope of variation (Perturbation Range)
 372 induced by these observational uncertainties.

373 **S14. Stratified Evaluation by Aerosol Subtypes**

374 To evaluate the model performance under distinct aerosol composition regimes,
 375 we utilize the Feature Classification Flags embedded in the CALIOP Level 2
 376 Atmospheric Volume Description. The NA evaluation dataset (2018) is partitioned into
 377 two representative subsets: a dust-dominated regime (combining 'Dust' and 'Polluted
 378 Dust' subtypes) and an SOA-dominated continental regime (combining 'Clean
 379 Continental' and 'Polluted Continental/Smoke' subtypes).

380

381 Table S4. Statistical evaluation of GC-TF model performance over NA, stratified by
 382 dominant CALIOP aerosol subtypes (2018). For each subset, we calculate statistical
 383 metrics for both the original GEOS-Chem simulation and the GC-TF corrected
 384 predictions against CALIOP AOD observations.

Aerosol Regime	Model	R	RMSE	MAE	Slope
Dust-dominated (Dust and Polluted Dust)	Original	0.41	0.032	0.027	0.21
	Corrected	0.50	0.032	0.027	0.32
SOA-dominated (Clean and Polluted Continental)	Original	0.51	0.034	0.031	0.27
	Corrected	0.50	0.035	0.031	0.30

385

386 **S15. Methodological Benchmarking and Structural Necessity**

387 To justify the architectural complexity of the proposed framework and isolate the
388 sources of its performance gains, we conduct comprehensive benchmarking and
389 ablation studies using the independent 2017 test dataset. To ensure a strictly fair
390 comparison, all baseline models and ablation variants are trained using the identical set
391 of input predictors—encompassing GEOS-Chem physicochemical states and MERRA-
392 2 meteorological forcings—along with identical hyperparameter configurations and
393 loss functions.

394 To establish a comprehensive baseline, two representative conventional deep
395 learning architectures were evaluated. The first is a Multilayer Perceptron (MLP),
396 representing point-wise neural networks. By treating vertical layers as independent
397 vectors, the MLP tests whether a simple numerical mapping, devoid of sequential
398 awareness, can resolve AEC biases. The second baseline is a 1-Dimensional
399 Convolutional Neural Network (1D-CNN). This architecture utilizes localized
400 receptive fields to capture vertical gradients between adjacent layers, serving as a
401 benchmark for local structural extraction, contrasting with the global dependency
402 modeling enabled by the Transformer.

403 Table S5. Performance benchmarking and ablation study of the proposed model against
404 conventional machine learning architectures. Evaluation is conducted on the
405 independent 2017 test dataset. All models are trained utilizing the identical
406 meteorological and chemical state predictors to ensure a rigorous comparison.

Model Configuration	R	MAE (km ⁻¹)	RMSE (km ⁻¹)
MLP	0.083	0.019	0.052
1D-CNN	0.540	0.016	0.044
Without Gated Fusion	0.637	0.015	0.040
Without Cross-Attention	0.654	0.014	0.039
Physics-Informed Transformer (Full)	0.666	0.014	0.039

407