



1 **Desert Model Intercomparison Project benchmark framework version 1.0 for assessing**
2 **land-surface dynamics and surface memory in monthly dust aerosol optical depth over**
3 **North Africa**

4 **Marzieh Mokarram¹, Mohammad Jafar Mokarram^{2*}, Huayu Lu^{3*}**

5 *¹Department of Geography, Faculty of Economics, Management and Social Sciences, Shiraz*
6 *University, Shiraz, Iran.*

7 *²School of Electronic Engineering and Intelligent Manufacturing, Anhui Xinhua University,*
8 *Hefei, China*

9 *³School of Geography and Ocean Science, Nanjing University, Nanjing 210023, China.*

10

11 **Corresponding author: Mohammad Jafar Mokarram and Huayu Lu (Email:*
12 *mjmokarram@axhu.edu.cn, huayulu@nju.edu.cn)*

13 **Abstract**

14 Deserts are the main global source of atmospheric mineral dust, yet large uncertainties remain
15 in the simulation of dust variability across space and time. Part of this uncertainty may reflect
16 the limited representation of dynamic land-surface states and antecedent surface conditions in
17 current dust-model formulations. Here, we develop an interpretable machine-learning
18 assessment framework, here termed the Desert Model Intercomparison Project (DesertMIP)
19 climate–surface–memory benchmark framework version 1.0, to quantify the added value of
20 land-surface dynamics and surface memory for monthly dust aerosol optical depth (AOD) over
21 North Africa during 2003–2020. Three hierarchical configurations were tested: Climate-only,
22 Climate+Surface, and Climate+Surface+Memory. Among the main explanatory models,
23 Random Forest gave the best overall performance. Test-set skill increased from $R^2 = 0.686$ in
24 the Climate-only case to 0.713 after adding surface variables and to 0.736 in the full
25 Climate+Surface+Memory case, while RMSE declined from 0.058 to 0.056 and then to 0.053.
26 The best model also gave MAE = 0.041 and Bias = 0.004. Residuals, defined here as observed
27 minus predicted, were centered close to zero, with a mean of -0.004 and a standard deviation
28 of 0.053, although residual spread increased at higher AOD values. Lag analysis showed a
29 persistent memory-sensitive signal, with the strongest negative associations near a three-month
30 lag for soil moisture ($r \approx -0.75$) and vegetation density (LAI; $r \approx -0.82$), whereas the highest
31 mean cross-validation skill among the tested memory windows occurred at six months. Despite
32 the gains in overall skill, severe dust outbreaks remained difficult to capture. For events above
33 the 95th percentile, the full model gave POD = 0.423, CSI = 0.269, FAR = 0.574, extreme
34 RMSE = 0.170, and extreme Bias = -0.060 . These results support a reproducible pilot
35 benchmark structure toward the DesertMIP.



36 **Keywords:** Dust AOD; North Africa; Land-surface dynamics; Surface memory; Dust-model
37 benchmarking; Interpretable machine learning; Extreme dust events; DesertMIP.

38 **1 Introduction**

39 Mineral dust is a major component of the Earth's climate system. It affects the radiation budget,
40 cloud processes, nutrient cycles, air quality, and many coupled atmospheric and environmental
41 processes (Kok et al., 2023; Mahowald et al., 2014). Deserts account for about 70–75 % of the
42 global atmospheric mineral aerosol source. Changes in the state and dynamics of desert regions
43 therefore have important consequences at regional and global scales (Kok et al., 2021). Dust
44 also influences climate through processes beyond direct radiative forcing. It modifies cloud
45 behavior, precipitation, and land–atmosphere interactions. These effects generate additional
46 feedbacks in the climate system (Wang et al., 2015).

47 Earth system models (ESMs) have been the main tools for representing and projecting dust
48 behavior over recent decades. Their development has advanced within major community
49 frameworks such as the Coupled Model Intercomparison Project, and their applications have
50 been central to climate-change assessment (Eyring et al., 2016). However, intercomparison
51 studies still report large uncertainties in dust emission, transport, and deposition (Aryal and
52 Evans, 2021; Huneus et al., 2011). These uncertainties are especially evident in arid regions
53 such as North Africa and the Middle East, where estimates of dust burden and aerosol optical
54 depth (AOD) differ substantially across models (Pu and Ginoux, 2018). Current model
55 evaluation often relies on bulk metrics such as coefficient of determination (R^2) and root mean
56 square error (RMSE). These metrics summarize overall predictive skill, but they do not by
57 themselves explain why errors occur or where they are concentrated. Recent evidence shows
58 that models may reproduce regional means while still exhibiting strong structural bias under
59 extreme conditions or within specific regions (Adebisi et al., 2023). A more diagnostic
60 evaluation framework is therefore needed. Such a framework should examine model behavior
61 across surface regimes, error structure, and possible memory effects, rather than relying on
62 domain-wide skill scores alone (Reichstein et al., 2019).

63 A related gap concerns the relative roles of atmospheric forcing and land-surface state in
64 shaping dust variability and model error. Wind speed is widely recognized as a primary driver
65 of aeolian erosion (Rohrman et al., 2013). Yet the effects of soil moisture, vegetation cover,
66 and surface roughness remain less systematically quantified, especially within nonlinear
67 diagnostic frameworks (Ginoux et al., 2012; Kim et al., 2024). These surface factors influence
68 dust emission through distinct physical pathways. Roughness reduces near-surface wind stress.
69 Soil moisture increases particle cohesion. Vegetation stabilizes exposed soil and limits sediment
70 entrainment (Zhang et al., 2022). At the same time, many climate and dust-model formulations
71 still depend on simplified or static descriptions of the land surface. Remote-sensing evidence



72 indicates that desert boundaries and land cover are highly dynamic. Seasonal and interannual
73 changes in vegetation, soil moisture, and land use can alter the potential for dust production
74 (Green et al., 2020; Van De Kerchove et al., 2021). Neglect of these dynamics can introduce
75 both spatial and temporal error. It can also reduce the ability of models to reproduce observed
76 patterns of dust activity.

77 An additional and less explored issue is surface memory. The present land surface does not
78 depend only on instantaneous conditions. It also reflects accumulated environmental history
79 over previous weeks and months. This history may include prolonged drought, gradual
80 vegetation decline, or persistent soil-moisture anomalies (Wen et al., 2025). Such processes can
81 produce delayed effects on sediment availability and erosion sensitivity. These considerations
82 motivate the hypothesis that antecedent surface states may provide additional explanatory
83 power beyond instantaneous surface conditions. However, surface-memory effects have rarely
84 been examined in dust-model diagnostics in a systematic and spatially explicit way.
85 Relationships between model bias, land-cover characteristics, and land-surface dynamics also
86 remain poorly resolved in structural terms (Ridley et al., 2014). A stronger focus on regime-
87 dependent error and lag-sensitive diagnostics may therefore help identify sensitive regions,
88 reveal missing controls, and guide future improvements in dust-model parameterization.

89 Interpretable machine learning provides a useful way to examine this problem because it can
90 represent nonlinear relationships while still allowing the contribution of different predictor
91 groups to be evaluated. In this study, machine learning is used mainly as a diagnostic tool rather
92 than only as a prediction method. The aim is to test how much additional information is gained
93 when dynamic surface conditions and antecedent surface states are added to climate-only
94 predictors.

95 The analysis is organized around three hierarchical configurations: Climate-only, Climate plus
96 Surface, and Climate plus Surface plus Memory. This structure allows the study to examine
97 whether land-surface variables improve dust AOD reconstruction, whether antecedent surface
98 conditions provide further information, and whether model errors differ across desert cores,
99 transition zones, and more vegetated environments. Extreme dust events are also evaluated
100 separately, because good mean-state performance does not necessarily imply good performance
101 in the upper tail of the distribution.

102 The main contribution of this study is therefore not simply the development of another
103 prediction model. Rather, it is the design of a benchmark-oriented diagnostic structure that links
104 climate forcing, land-surface dynamics, surface memory, regime-dependent error, and extreme-
105 event behavior. This structure provides a reproducible pilot framework for future dust-model
106 assessment over dynamic desert environments.

107 **2. Materials and Methods**



108 2.1. Overall Study Framework and Methodological Rationale

109 This study develops a benchmark-oriented framework to evaluate how much information is
110 added by dynamic land-surface conditions and antecedent surface states when reconstructing
111 monthly dust AOD over North Africa. The framework is referred to here as the Desert Model
112 Intercomparison Project (DesertMIP) benchmark framework version 1.0. In this structure, dust
113 AOD is treated as the response of a coupled climate–surface system rather than as the outcome
114 of atmospheric forcing alone.

115 The relationship was formulated as a nonlinear mapping between climatic drivers,
116 instantaneous surface conditions, memory terms, and dust AOD:

$$117 \quad Y_t = f(C_t, S_t, M_t) + \varepsilon_t \quad (1)$$

118 where Y_t denotes dust AOD at time t , C_t denotes climatic and atmospheric predictors, S_t denotes
119 instantaneous land-surface variables, M_t denotes antecedent or memory variables, and ε_t is the
120 residual term.

121 The unknown function was then approximated using a data-driven nonlinear model (Molnar,
122 2020; Reichstein et al., 2019):

$$123 \quad \hat{Y}_t = F_\theta(C_t, S_t, M_t) \quad (2)$$

124 where F_θ denotes a nonlinear model parameterized by θ . This modelling choice was used
125 because dust response may involve thresholds, nonlinear coupling, and interactions between
126 atmospheric forcing and surface resistance.

127 To isolate the contribution of each information block, three hierarchical configurations were
128 tested:

$$129 \quad \hat{Y}_t^{(1)} = F_\theta(C_t) \quad (3)$$

$$130 \quad \hat{Y}_t^{(2)} = F_\theta(C_t, S_t) \quad (4)$$

$$131 \quad \hat{Y}_t^{(3)} = F_\theta(C_t, S_t, M_t) \quad (5)$$

132 The first configuration uses only climatic predictors and serves as the baseline. The second adds
133 dynamic land-surface variables, and the third adds antecedent surface-memory terms.
134 Comparing these three configurations provides a direct way to assess whether surface state and
135 surface memory improve dust AOD reconstruction beyond climate-only forcing.

136 2.2. Spatial Domain, Temporal Coverage, and Spatio-Temporal Design

137 The study domain covers North Africa within 20°W–30° E and 10°N–25° N. This region
138 contains some of the world’s most important dust-source areas and shows strong gradients in



139 aridity, vegetation, and soil conditions (Kok et al., 2021; Pu and Ginoux, 2018). North Africa
140 is used here as a pilot benchmark domain for DesertMIP because it combines major dust sources
141 with strong ecological transition zones and pronounced land-surface contrasts. These
142 characteristics make it suitable for testing climate-only, surface-aware, and memory-aware
143 assessment designs.

144 The analysis period extends from 2003 to 2020. Spatial sampling used a regular grid at
145 $0.1^\circ \times 0.1^\circ$ resolution. The grid can be expressed as:

$$146 \quad G = \{(x_i, y_j) \mid x_i \in \lambda, y_j \in \phi\} \quad (6)$$

147 where x_i and y_j denote longitude and latitude coordinates, and λ and ϕ denote the full sets
148 of longitudes and latitudes in the study domain.

149 Daily data were aggregated to monthly values to reduce short-term noise and to ensure temporal
150 consistency across variables. The monthly value was defined as:

$$151 \quad X_m = \frac{1}{n} \sum_{t=1}^n X_t \quad (7)$$

152 where X_m is the monthly value, X_t is the daily value, and n is the number of days in the month.
153 Monthly sums were used for cumulative variables such as precipitation. Monthly means were
154 used for continuous variables such as temperature and humidity.

155 To characterize spatial heterogeneity of the land surface, a composite dryness index was defined
156 from soil moisture and vegetation state:

$$157 \quad D = 1 - (w_1 \cdot SM + w_2 \cdot NDVI) \quad (8)$$

158 where D denotes the dryness index, SM denotes soil moisture, $NDVI$ denotes the normalized
159 difference vegetation index, and w_1 and w_2 are weighting coefficients. Spatial transition
160 intensity was then estimated from the gradient magnitude of the dryness index:

$$161 \quad T = |\nabla D| \quad (9)$$

162 where T denotes transition intensity. Larger values of T indicate stronger ecological contrast
163 over space. Based on this indicator, the study area was stratified into three benchmark surface
164 regimes: desert core, transition zone, and vegetated region. This regime-based design forms an
165 important part of the benchmark logic because model skill is expected to differ across these
166 contrasting surface environments.

167 **2.3. Data Sources and Study Variables**

168 **2.3.1. Target Variable**

169 The target variable is dust AOD at 550 nm. Data was obtained from the Copernicus Atmosphere
170 Monitoring Service (CAMS). Dust AOD represents vertically integrated atmospheric dust load
171 and is widely used as an indicator of dust activity and dust burden (Kok et al., 2023).



172 A logarithmic transformation was applied to reduce skewness and stabilize model fitting:

$$173 \quad Y = \log(1 + AOD) \quad (10)$$

174 where AOD denotes dust AOD and Y denotes the transformed response variable.

175 Additional climatic and land-surface predictors were obtained from the ERA5 reanalysis
176 produced by the European Centre for Medium-Range Weather Forecasts (ECMWF) (Soci et al.,
177 2024). Data was accessed through the Copernicus Climate Data Store (CDS). The selected
178 observational and reanalysis products were chosen not only for the present analysis, but also
179 because they provide a publicly accessible and reproducible basis for future benchmark-
180 oriented experiments within a prototype DesertMIP structure. They also provide consistent
181 spatial and temporal coverage and are compatible with future multi-model comparison
182 workflows.

183 Quality-control procedures were applied before model development. Unrealistic values were
184 removed. Incomplete records were excluded or flagged during preprocessing. All variables
185 were then harmonized to a common monthly grid.

186 2.3.2. Climatic and Atmospheric Variables

187 Climatic predictors were selected to represent the main atmospheric controls on dust emission
188 and transport. Wind speed at 10 m is a primary driver of aeolian erosion (Ginoux et al., 2012).
189 Precipitation influences soil moisture and particle cohesion. Air temperature, relative humidity,
190 boundary layer height, surface pressure, and dew point temperature provide additional
191 information on atmospheric thermodynamic state and near-surface moisture conditions.

192 Wind speed was calculated from the zonal and meridional wind components:

$$193 \quad U = \sqrt{u_{10}^2 + v_{10}^2} \quad (11)$$

194 where U is wind speed in m s^{-1} , u_{10} is the zonal wind component, and v_{10} is the meridional
195 wind component.

196 Wind direction was calculated as:

$$197 \quad \theta = \text{atan2}(v_{10}, u_{10}) \quad (12)$$

198 where θ denotes wind direction and atan2 preserves the correct directional quadrant.

199 Precipitation was expressed in mm or mm month^{-1} , depending on the aggregation step.
200 Temperature was expressed in Kelvin or in degrees Celsius, depending on the source product
201 and preprocessing stage. These atmospheric variables define the climate-forcing block in the
202 hierarchical assessment framework.

203 2.3.3. Land-Surface Variables



204 Land-surface predictors were selected to represent the resistance of the surface to wind erosion
205 and its temporal variability. Soil moisture (SM) is a first-order control on sediment cohesion
206 and generally suppresses dust emission as it increases (Kim et al., 2024). The leaf area index
207 (LAI) represents vegetation density and affects roughness, sheltering, and sediment
208 stabilization.

209 Total vegetation state was represented as:

$$210 \quad LAI_{total} = LAI_{hv} + LAI_{lv} \quad (13)$$

211 where LAI_{hv} and LAI_{lv} denote high- and low-vegetation components, respectively.

212 Within the benchmark logic of this study, land-surface variables are not treated as secondary
213 descriptors. They form a distinct explanatory block that allows direct comparison between
214 climate-only and climate-plus-surface configurations. This distinction is important for
215 DesertMIP because it makes it possible to test whether model gains arise from atmospheric
216 forcing alone or from improved representation of dynamic surface state.

217 **2.3.4. Memory Variables and Lagged Effects**

218 Environmental memory effects were represented using antecedent windows of land-surface
219 variables. Two related lag treatments were used in this study. First, moving-average memory
220 predictors were constructed for use in the hierarchical explanatory models. Second, a broader
221 lag-correlation analysis was carried out to describe the temporal structure of land-surface
222 associations with dust AOD.

223 For the explanatory modeling framework, soil-moisture and vegetation memory were defined
224 as moving averages over antecedent months. The baseline memory formulation used in the
225 main Climate+Surface+Memory configuration was based on a 3-month moving average. Soil-
226 moisture memory was defined as

$$227 \quad SM_{memory,L}(t) = \frac{1}{L} \sum_{i=0}^{L-1} SM(t-i) \quad (14)$$

228 and vegetation memory was defined as

$$229 \quad LAI_{memory,L}(t) = \frac{1}{L} \sum_{i=0}^{L-1} LAI(t-i) \quad (15)$$

230 where L denotes the length of the memory window in months. In the primary explanatory
231 setting, $L = 3$. Additional lag windows of $L = 1$, $L = 3$, and $L = 6$ months were evaluated
232 in a sensitivity analysis to test how reconstruction skill changed with the duration of antecedent
233 surface memory. These lag windows were selected to represent short, intermediate, and
234 seasonally persistent memory scales.

235 In addition to the predictor-based memory formulations, lag-sensitive correlation diagnostics
236 were computed between dust AOD and key land-surface variables over lags from 0 to 6 months.



237 These diagnostics were used to characterize temporal association structure rather than to define
238 the main predictor set. This distinction is important because the lag-correlation analysis serves
239 as a benchmark descriptor of memory-sensitive behaviour, whereas the moving-average
240 memory predictors are part of the explanatory modeling framework.

241 These memory terms represent delayed land-surface effects that may reflect accumulated
242 drying, vegetation decline, or other forms of surface persistence. Environmental memory
243 effects have been reported to influence land–atmosphere interactions on seasonal timescales (Q.
244 Wen et al., 2025). In the present study, they are treated as observation-constrained descriptors
245 of antecedent surface state rather than as direct proof of causal memory mechanisms.

246 **2.3.5. Dryness Index and Nonlinear Interaction Features**

247 To represent the balance between atmospheric forcing and surface resistance, a process-oriented
248 dryness index was defined as:

$$249 \quad DI = \frac{U}{SM + \epsilon} \quad (16)$$

250 where DI denotes the dryness index, U denotes wind speed, SM denotes soil moisture, and
251 ϵ is a small constant that prevents division by zero.

252 Additional nonlinear interaction terms were also constructed to represent coupled controls
253 between wind forcing and land-surface state. These terms included $SM \times U$ and $LAI \times U$.
254 They were introduced to help the model capture situations in which the effect of wind depends
255 on surface moisture or vegetation conditions. This feature construction is consistent with the
256 physical expectation that dust emission is controlled by joint thresholds rather than by single
257 predictors in isolation.

258 From a benchmark perspective, these interaction variables provide a useful bridge between
259 purely statistical prediction and process-aware interpretation. They also support the broader
260 goal of DesertMIP, which is to assess dust-model behavior under a climate–surface–memory
261 framework rather than under climate forcing alone.

262 **2.4. Variables, Symbols, Units, and Model Roles**

263 Table 1 summarizes the main variable groups, symbols, units, data sources, temporal treatment,
264 and model roles used in this study. The table is structured to distinguish climate drivers,
265 instantaneous surface-state variables, surface-memory terms, process-aware interaction
266 features, regime descriptors, and extreme-event diagnostics. This grouping is also used as the
267 benchmark logic of the prototype DesertMIP framework, because it allows direct comparison
268 between climate-only, surface-aware, and memory-aware assessment tiers.

269 **2.5. Data Preprocessing and Analytical Dataset Construction**



270 **2.5.1. Spatio-Temporal Harmonization**

271 All variables were harmonized onto a common $0.1^\circ \times 0.1^\circ$ grid and a monthly temporal
 272 resolution. This step was required to ensure consistency across products and to allow direct
 273 comparison among the three hierarchical predictor configurations. The harmonized value was
 274 defined as:

275
$$X_{i,j,m}^* = A\left(X_t^{(r)}\right) \tag{17}$$

276 where $X_t^{(r)}$ denotes the original variable from source r at time t , $X_{i,j,m}^*$ denotes the
 277 harmonized value at grid cell (i, j) and month m , and A denotes the harmonization operator.
 278 This operator includes spatial remapping and temporal aggregation. The use of a common grid
 279 is also important from a benchmark perspective because future DesertMIP-style experiments
 280 require transferable inputs and directly comparable outputs across datasets and models.

281 **2.5.2. Quality Control and Missing Data**

282 Quality-control procedures were applied before model development. Values outside physically
 283 realistic ranges were removed using variable-specific limits. Missing values were imputed
 284 using the median of the training data for each predictor. This choice reduced sensitivity to
 285 outliers and avoided leakage of information from the test period into the training workflow.
 286 Feature scaling was applied only where required by the modeling algorithm. Tree-based models
 287 were trained on the original feature scale, whereas scale-sensitive models were standardized
 288 before fitting.

289 **Table 1.** Proposed DesertMIP benchmark variable groups

Group	Variable	Role in dust system	Data source	Temporal treatment	Diagnostic use in DesertMIP
Climate drivers	Dust AOD	Target variable for dust variability	CAMS	Daily to monthly; log-transformed for modeling	Benchmark target for all tiers
Climate drivers	$u_{10}, v_{10}, U, \theta$	Wind forcing for erosion and transport	ERA5	Daily to monthly means	Core predictors in DMP-1
Climate drivers	$P_{or}, TP, T, RH, BLH, SP, DPT$	Moisture, thermodynamic state, and boundary-layer controls	ERA5	Monthly means or sums	Core predictors in DMP-1
Instantaneous surface state	$SM, LAI_{total}, NDVI$	Cohesion, roughness, vegetation shielding, and surface resistance	ERA5 and harmonized land-surface	Monthly means	Added explanatory block in DMP-2



Group	Variable	Role in dust system	Data source	Temporal treatment	Diagnostic use in DesertMIP
Surface memory terms	SM_{memory} , LAI_{memory}	Antecedent surface persistence and delayed effects	Derived from products monthly surface variables	Moving averages and lag windows	Added explanatory block in DMP-3
Process-aware derived features	DI , $SM \times U$, $LAI \times U$	Coupled forcing–resistance effects and nonlinear thresholds	Derived variables	Monthly	Process-aware interpretation and sensitivity analysis
Regime descriptors	D , T , regime class	Dryness structure and ecological transition intensity	Derived from harmonized surface variables	Spatial diagnostics on monthly grid	Regime-based benchmarking
Extreme-event diagnostics	Q95 threshold, POD, CSI, FAR, extreme RMSE, extreme Bias	Tail behavior and severe-event performance	Derived from observed and predicted AOD	Evaluated on test period	Mandatory tail diagnostics in prototype DesertMIP

290

291 **2.5.3. Feature Matrix Construction**

292 The analytical feature matrix combined climatic predictors, surface-state variables, memory
293 terms, and derived interaction features. It was defined as:

294
$$X = [C_t, S_t, M_t, DI, SM \times U, LAI \times U, \dots] \tag{18}$$

295 The response vector was defined as:

296
$$Y = [Y_1, Y_2, \dots, Y_n] \tag{19}$$

297 where C_t denotes climatic predictors, S_t denotes instantaneous surface variables, and
298 M_t denotes memory variables. In the benchmark logic of this study, these feature blocks
299 correspond directly to the progressive information content of the three explanatory tiers. This
300 structure allows the added value of surface and memory information to be quantified in a
301 transparent and transferable way.

302 For the primary Climate+Surface+Memory configuration, the memory block used the 3-month
303 moving-average formulation, while additional 1-, 3-, and 6-month windows were examined in
304 dedicated sensitivity tests.

305 **2.6. Data Partitioning and Validation Strategy**

306 Temporal ordering was preserved throughout model development in order to reduce
307 information leakage across time. The first 80 % of the time series was used for model training
308 and the remaining 20 % was reserved as an independent temporally out-of-sample test set:



309 $D_{train} = \{t_1, t_2, \dots, t_k\}$ (20)

310 $D_{test} = \{t_{k+1}, \dots, t_n\}$ (21)

311 Hyperparameter tuning was performed within the training period using blocked time-series
312 cross-validation. This design was chosen to assess temporal generalization while preventing
313 future information from entering model fitting. The present validation should therefore be
314 interpreted as a pilot temporal benchmark rather than a fully independent regional benchmark.
315 Because the split was based primarily on time, reported test performance may still be influenced
316 by spatial autocorrelation among neighboring grid cells. A stricter evaluation of generalization
317 would require spatial blocking or leave-region-out experiments. In a future DesertMIP
318 implementation, these forms of validation should be treated as mandatory components of the
319 protocol.

320 2.7. Modeling Strategy and Scenario Design

321 The modeling strategy combined linear and nonlinear algorithms in order to compare simple
322 baselines with more flexible learners. Ridge regression was used as a linear reference model.
323 Random Forest, HistGradientBoosting, and XGBoost were used to represent nonlinear tree-
324 based approaches that can capture threshold behavior and interactions between atmospheric
325 forcing and land-surface state (Reichstein et al., 2019). Algorithm selection was guided by two
326 criteria. The first was the ability to represent nonlinear responses. The second was relative
327 robustness in the presence of correlated predictors.

328 Three primary explanatory scenarios were evaluated. Scenario 1 corresponds to the Climate-
329 only configuration defined in Eq. (3). Scenario 2 corresponds to the Climate+Surface
330 configuration defined in Eq. (4). Scenario 3 corresponds to the Climate+Surface+Memory
331 configuration defined in Eq. (5). For clarity and future protocol transferability, these three
332 explanatory scenarios are also referred to here as pilot DesertMIP tiers: DMP-1 (climate-only),
333 DMP-2 (climate plus surface), and DMP-3 (climate plus surface plus memory). This tiered
334 structure converts the present analysis into a benchmark architecture rather than a single-model
335 exercise. Within DMP-3, the baseline memory predictors were defined using a 3-month moving
336 window, and alternative 1-, 3-, and 6-month formulations were evaluated separately to assess
337 the sensitivity of benchmark performance to memory-window length.

338 An ablation analysis was used to quantify the incremental contribution of each predictor block.
339 This analysis tested how model skill changed after removal of surface-state variables,
340 vegetation terms, soil-moisture terms, or memory features. From a DesertMIP perspective, this
341 step provides a diagnostic link between aggregate model performance and the information
342 content of specific variable groups.

343 In addition to the three explanatory tiers, an auxiliary autoregressive benchmark was also
344 evaluated by adding lagged dust AOD (*dust_lag1*) as a predictor. This benchmark is referred



345 to here as DMP-AR. It was not treated as part of the primary explanatory comparison because
346 it includes prior information from the target variable itself. It therefore provides a statistical
347 upper-bound reference rather than a directly comparable process-oriented scenario. Its purpose
348 is to show how much predictive skill can be gained when persistence in the target series is used
349 explicitly.

350 **2.8. Model Performance Metrics**

351 Model evaluation was organized around two diagnostic levels. The first level includes **core**
352 **diagnostics** that quantify overall predictive skill and benchmark performance across regimes
353 and event classes. The second level includes **extended diagnostics** that describe how and where
354 model behavior changes across the predictor hierarchy. This distinction is useful for the present
355 analysis and also provides a transferable metric structure for a prototype DesertMIP benchmark.

356 The core diagnostics include the coefficient of determination (R^2), root mean square error
357 ($RMSE$), mean absolute error (MAE), mean bias, regime-specific RMSE, and extreme-event
358 detection metrics. The coefficient of determination was defined as

$$359 \quad R^2 = 1 - \frac{\sum_{t=1}^n (Y_t - \hat{Y}_t)^2}{\sum_{t=1}^n (Y_t - \bar{Y})^2} \quad (22)$$

360 The root mean square error was defined as

$$361 \quad RMSE = \sqrt{\frac{1}{n} \sum_{t=1}^n (Y_t - \hat{Y}_t)^2} \quad (23)$$

362 The mean absolute error was defined as

$$363 \quad MAE = \frac{1}{n} \sum_{t=1}^n |Y_t - \hat{Y}_t| \quad (24)$$

364 To ensure a consistent sign convention across the manuscript, mean bias was defined as
365 predicted minus observed:

$$366 \quad Bias = \frac{1}{n} \sum_{t=1}^n (\hat{Y}_t - Y_t) \quad (25)$$

367 Under this definition, positive bias indicates overestimation and negative bias indicates
368 underestimation.

369 Extreme dust events were defined using the 95th percentile of the observed AOD distribution
370 in the evaluation dataset. Event-detection skill was assessed using probability of detection
371 (POD),

$$372 \quad POD = \frac{Hits}{Hits + Misses} \quad (26)$$

373 critical success index (CSI),



374
$$CSI = \frac{Hits}{Hits+Misses+FalseAlarms} \quad (27)$$

375 and false alarm ratio (*FAR*),

376
$$FAR = \frac{FalseAlarms}{Hits+FalseAlarms} \quad (28)$$

377 Core diagnostics were computed for each model configuration and for each benchmark surface
378 regime. The extended diagnostics include conditional-bias curves, lag-dependence diagnostics,
379 spatial error maps, and variable-group contribution analysis. Within the prototype DesertMIP
380 logic, the core diagnostics define minimum benchmark outputs, whereas the extended
381 diagnostics provide process-oriented interpretation of performance differences across tiers.

382 2.9. Explainability Analysis and Variable Contribution

383 Explainability analysis was used to quantify how predictor importance and response structure
384 changed across the three hierarchical configurations. Global predictor importance was
385 estimated from the fitted model using error-based importance scores, and complementary
386 analyses were used to compare the relative role of climatic, surface, and memory variables
387 (Molnar, 2020). These diagnostics were intended to support interpretation of model behavior
388 rather than to establish direct physical causality.

389 Partial-dependence analysis was used to summarize average response patterns for selected
390 predictors:

391
$$PD(X_j) = E_{X_{-j}}[F_{\theta}(X)] \quad (29)$$

392 where $PD(X_j)$ denotes the partial dependence of predictor X_j , and X_{-j} denotes the set of all
393 remaining predictors.

394 To examine whether the addition of surface and memory variables changed systematic error,
395 scenario-wise differences in mean bias were defined as

396
$$\Delta Bias_S = Bias(C) - Bias(C + S) \quad (30)$$

397
$$\Delta Bias_M = Bias(C + S) - Bias(C + S + M) \quad (31)$$

398 where C denotes the Climate-only configuration, $C + S$ denotes the Climate+Surface
399 configuration, and $C + S + M$ denotes the Climate+Surface+Memory configuration. Positive
400 values of $\Delta Bias$ indicate reduced absolute bias after the addition of the corresponding predictor
401 block.

402 These diagnostics were examined at the full-domain scale and across benchmark surface
403 regimes. This design allows the analysis to test whether soil moisture, vegetation state, and
404 antecedent memory terms contribute differently across desert cores, transition zones, and more
405 vegetated environments. Within a future DesertMIP framework, these explainability
406 diagnostics are intended as standardized interpretive outputs rather than as stand-alone causal



407 evidence.

408 **2.10. Computational Environment and Reproducibility**

409 All analyses were conducted in the Python scientific-computing environment. Core libraries
410 included xarray, netCDF4, NumPy, pandas, scikit-learn, and XGBoost where required.
411 Reproducibility was supported through fixed random seeds, explicit preprocessing rules, and a
412 documented workflow from data harmonization to model evaluation.

413 The computational pipeline included data acquisition, spatial and temporal harmonization,
414 quality control, construction of base predictors, generation of interaction and memory features,
415 assembly of the analytical feature matrix, temporal partitioning, model fitting, hyperparameter
416 tuning, benchmark evaluation, and explainability analysis. To support reproducibility beyond
417 the present manuscript, the archived workflow includes the exact repository version, the exact
418 Zenodo DOI, experiment configuration files for all benchmark tiers, region masks for the three
419 benchmark regimes, plotting scripts, and benchmark summary tables. Together, these archived
420 materials allow the main results of the pilot benchmark to be reproduced from the documented
421 inputs and configuration settings. This level of documentation is especially important for a
422 benchmark-oriented study because future DesertMIP-style applications require transferable
423 experiment definitions rather than only narrative method descriptions.

424 **2.11. Prototype DesertMIP Benchmark Logic and Protocol**

425 This study defines a prototype DesertMIP benchmark logic and protocol for transferable dust-
426 model assessment. The proposed framework should be interpreted as a pilot proof of concept
427 rather than as a finalized community intercomparison standard. Its main objective is to diagnose
428 the added value of dynamic land-surface information and antecedent surface memory in dust-
429 model assessment within a structured and reproducible benchmarking architecture.

430 The pilot benchmark domain is North Africa at monthly resolution over 2003–2020. This
431 domain was selected because it combines strong dust activity with pronounced gradients in
432 aridity, vegetation cover, and ecological transition intensity, making it a suitable test bed for
433 climate–surface–memory benchmarking. The benchmark is organized around three primary
434 explanatory tiers: DMP-1, corresponding to the Climate-only configuration; DMP-2, which
435 adds instantaneous land-surface variables; and DMP-3, which adds antecedent surface-memory
436 terms. In addition, DMP-AR is included as an optional autoregressive reference that provides
437 a statistical upper-bound benchmark but is not treated as part of the main explanatory hierarchy.

438 Mandatory benchmark outputs include overall skill metrics, regime-specific skill metrics,
439 extreme-event diagnostics, lag-sensitive diagnostics, and spatial bias maps. Benchmark
440 evaluation is performed across three required reporting strata: desert core, transition zone, and
441 vegetated region. These regimes are treated as essential benchmark classes because model



442 behavior is expected to differ substantially across them.

443 In this sense, the present implementation provides a reproducible pilot benchmark structure for
444 climate–surface–memory assessment of dust models. It does not yet constitute a community-
445 finalized DesertMIP protocol, but it offers a structured blueprint that can be extended in future
446 developments through additional domains, stricter spatial validation strategies, and direct
447 comparison with process-based dust models.

448 **3. Results**

449 **3.1. Overall Evaluation of Model Performance in Reconstructing Dust AOD**

450 **3.1.1. General performance across the three explanatory configurations**

451 The three explanatory configurations show a clear and ordered increase in reconstruction skill
452 as additional predictor groups are introduced. Across most algorithms, performance improves
453 from the Climate-only configuration to Climate+Surface and then to Climate+Surface+Memory.
454 This pattern indicates that dynamic land-surface information and antecedent memory terms add
455 useful explanatory information beyond climate-only forcing.

456 Within the Climate-only configuration, XGBoost achieved the highest test performance, with
457 $R^2 = 0.700$, $MAE = 0.043$, and $RMSE = 0.057$. Random Forest followed closely, with
458 $R^2 = 0.686$ and $RMSE = 0.058$. After surface variables were added, Random Forest became
459 the best-performing model, with $R^2 = 0.713$, $MAE = 0.042$, and $RMSE = 0.056$. In the
460 full Climate+Surface+Memory configuration, Random Forest again provided the highest
461 overall skill, with $R^2 = 0.736$, $MAE = 0.041$, and $RMSE = 0.053$. This was the best test
462 performance among all explanatory models examined in this study.

463 The same comparison also shows that tree-based models consistently outperformed the linear
464 Ridge model. This result suggests that the relationship between dust AOD and the predictor set
465 is not well represented by a purely linear structure. Instead, the reconstruction problem appears
466 to benefit from methods that can represent nonlinear responses and interactions.

467 Although XGBoost slightly outperformed Random Forest in the Climate-only configuration,
468 Random Forest was selected as the reference model for subsequent interpretation. Its
469 performance was strongest in the two surface-aware configurations that are most relevant to the
470 climate–surface–memory framework. It also provided the best overall result in the full
471 explanatory setting. For this reason, Random Forest was used as the main interpretive model in
472 the analyses that follow (Table 2).

473 **Table 2.** Performance metrics of machine-learning models for dust AOD reconstruction under
474 three predictor configurations.



Scenario	Model	Train (R ²)	Test (R ²)	Test (MAE)	Test (RMSE)	Bias
Climate-only	XGBoost	0.859	0.7	0.043	0.057	0.007
Climate-only	RandomForest	0.979	0.686	0.044	0.058	0
Climate-only	HistGB	0.948	0.629	0.048	0.063	0.021
Climate-only	Ridge	0.601	0.594	0.052	0.066	0.005
Climate+Surface	RandomForest	0.976	0.713	0.042	0.056	-0.002
Climate+Surface	XGBoost	0.859	0.706	0.043	0.056	-0.001
Climate+Surface	HistGB	0.945	0.675	0.045	0.059	-0.011
Climate+Surface	Ridge	0.624	0.614	0.051	0.064	0.004
Climate+Surface+Memory	RandomForest	0.983	0.736	0.041	0.053	0.004
Climate+Surface+Memory	XGBoost	0.857	0.721	0.042	0.055	0.003
Climate+Surface+Memory	Ridge	0.678	0.676	0.046	0.059	0.003
Climate+Surface+Memory	HistGB	0.951	0.672	0.045	0.059	-0.004

475

476 **3.1.2. Comparative performance assessment across configurations**

477 The comparison among the three explanatory configurations shows a gradual improvement in
 478 model skill as additional information is introduced. For the Random Forest model, test R²
 479 increased from 0.686 in the Climate-only configuration to 0.713 in Climate+Surface and to
 480 0.736 in Climate+Surface+Memory. Over the same sequence, test RMSE decreased from 0.058
 481 to 0.056 and then to 0.053. Other nonlinear models showed broadly similar behavior, although
 482 the size of the improvement differed among algorithms.

483 This stepwise improvement suggests that land-surface variables add information that is not fully
 484 captured by atmospheric predictors alone. The additional gain after including memory terms is
 485 smaller, but it still indicates that antecedent surface conditions contain useful information
 486 beyond the instantaneous surface state. In the proposed benchmark hierarchy, these results
 487 define the basic skill progression from DMP-1 to DMP-3.

488 **3.1.3. Relationship between observed and predicted values**

489 Figure 1b compares observed and predicted dust AOD for the best-performing explanatory
 490 model, Random Forest under the Climate+Surface+Memory configuration. Most points are
 491 concentrated near the 1:1 line, showing that the model captures much of the observed variability.
 492 The agreement is strongest for low and moderate AOD values, where most observations are
 493 located.

494 The spread becomes larger at higher AOD values. Several high observed AOD cases fall below
 495 the 1:1 line, indicating underestimation of intense dust conditions. This behavior shows that the
 496 model performs well for the dominant range of the data but becomes less reliable toward the
 497 upper tail of the distribution.



498 **3.1.4. Distribution of model residuals**

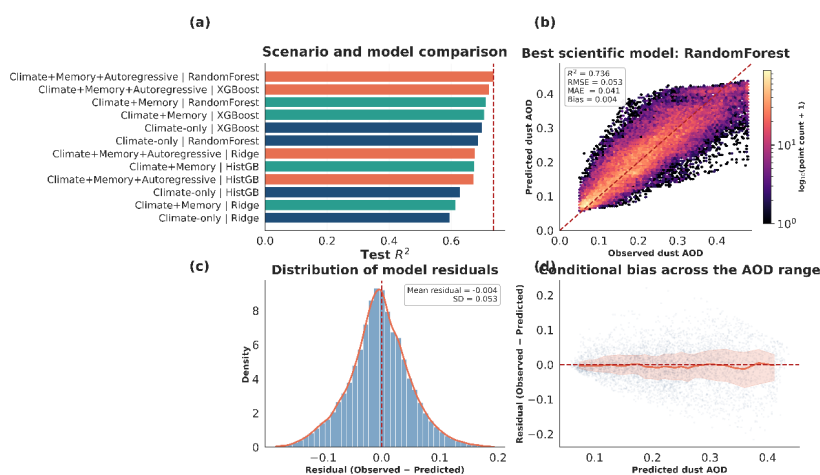
499 Figure 1c shows the residual distribution for the best-performing explanatory model. The
 500 residuals are centered close to zero and are approximately symmetric, with a mean of -0.004
 501 and a standard deviation of 0.053 . This indicates that the model has little average bias across
 502 the full evaluation set.

503 However, the distribution also has elongated tails. These larger residuals occur under a smaller
 504 subset of conditions and are consistent with the increased spread observed at higher AOD values
 505 in the observed–predicted comparison. For this reason, average error statistics alone are not
 506 sufficient to describe model performance across the full dust range.

507 **3.1.5. Conditional bias across the AOD range**

508 The mean residual provides a useful summary, but it does not show how error changes across
 509 the AOD distribution. Figure 1d therefore presents conditional bias as a function of predicted
 510 AOD. In the low-to-moderate range, conditional bias remains close to zero, indicating relatively
 511 stable calibration for the most common conditions.

512 At higher predicted AOD values, conditional bias becomes increasingly negative. Because
 513 residuals are defined as observed minus predicted in this diagnostic, negative values indicate
 514 underestimation in the upper part of the distribution. The model therefore reconstructs typical
 515 dust conditions more accurately than severe dust episodes. This range-dependent behavior
 516 motivates the separate analysis of extreme-event performance in the following sections.



517

518 **Figure 1.** Overall evaluation of model performance in reconstructing dust AOD. (a) Test-set
 519 R^2 across algorithms and explanatory configurations. (b) Observed versus predicted AOD for
 520 the best-performing explanatory model, Random Forest under the Climate+Surface+Memory
 521 configuration. The dashed line marks the 1:1 reference. (c) Distribution of residuals, defined as



522 observed minus predicted. (d) Conditional bias as a function of predicted AOD, showing how
523 model error changes across the response range.

524 **3.2. Added Value of Land-Surface Information**

525 **3.2.1. Improvement in skill after adding surface variables**

526 Adding land-surface variables to the climate-only predictor set improved dust AOD
527 reconstruction in a consistent way. At the domain scale, the Random Forest model improved
528 from $R^2 = 0.686$ in the Climate-only configuration to $R^2 = 0.713$ in the Climate+Surface
529 configuration. Over the same transition, *RMSE* decreased from 0.058 to 0.056 and
530 *MAE* decreased from 0.044 to 0.042. These changes correspond to an absolute gain of 0.027 in
531 R^2 , an absolute reduction of 0.002 in *RMSE*, and an absolute reduction of 0.002 in *MAE*. Bias
532 changed only slightly, from 0.000 to -0.002 .

533 These results indicate that instantaneous land-surface information adds explanatory power
534 beyond atmospheric predictors alone. The gain is physically plausible because climate-only
535 predictors do not directly represent surface resistance, vegetation shielding, or moisture-related
536 changes in sediment availability. Surface variables provide complementary information on the
537 state of erosion thresholds and the spatial structure of erodibility. In the benchmark logic of this
538 study, this step defines the added value of the DMP-2 tier relative to DMP-1.

539 **3.2.2. Spatial pattern of error reduction**

540 The gain associated with surface variables is not spatially uniform. Figure 2 shows the
541 difference in *RMSE* between the Climate-only and Climate+Surface configurations. Positive
542 values indicate reduced error after the inclusion of surface predictors. The largest improvements
543 are concentrated in transition areas and in regions with stronger vegetation influence. Smaller
544 gains are found in more homogeneous desert interiors.

545 This spatial pattern suggests that surface information is most useful where soil moisture and
546 vegetation vary more strongly over time and space. In such regions, climate-only predictors are
547 less able to capture changes in effective dust-source conditions. By contrast, bare desert
548 interiors tend to show more stable surface states. In those areas, the additional information
549 carried by surface variables is smaller. From a benchmark perspective, these results show that
550 DMP-2 is especially informative in regions where ecological transitions are active and surface-
551 state dynamics are strong.

552 **3.2.3. Regime-specific changes in error and bias**

553 The regime-based results confirm that the impact of surface variables differs across benchmark
554 environments. Relative *RMSE* reduction was largest in the vegetated zone, where error
555 decreased by 2.89 %. The corresponding reductions were 1.67 % in the transition zone and



556 0.80 % in the core desert. Bias behavior was less uniform. The vegetated zone showed a 2.56 %
557 reduction in bias, whereas the transition zone and core desert showed small bias degradations
558 of -0.70 % and -1.57 %, respectively. These results suggest that surface variables mainly
559 reduce random error across all regimes, but their effect on systematic bias is more region-
560 dependent. In vegetated and semi-vegetated areas, surface predictors improve both overall fit
561 and consistency. In very dry regions, they still reduce RMSE, but they may also introduce small
562 residual bias. This behavior may reflect measurement limitations over bright desert surfaces or
563 unresolved microtopographic effects.

564

565

566

567

568

569 **Table 3.** Relative changes in RMSE and Bias after adding instantaneous surface variables to
570 the Climate-only configuration. Positive values indicate error reduction. Negative values
571 indicate performance degradation.

Surface regime	RMSE change (%)	Bias change (%)
Vegetated zone	2.89	2.56
Transition zone	1.67	-0.70
Core desert	0.80	-1.57

572

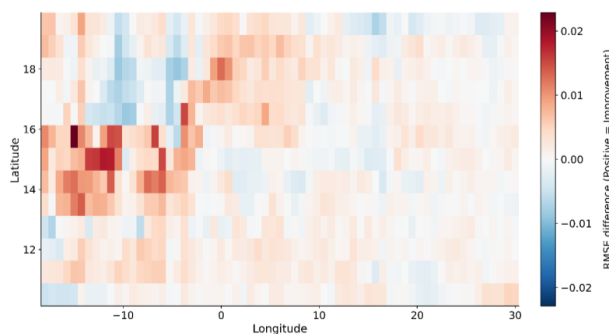
573 3.2.4. Role of soil moisture and vegetation indicators

574 The spatial and regime-dependent gains are consistent with the physical role of soil moisture
575 and vegetation in dust emission. Soil moisture increases particle cohesion and raises the
576 effective erosion threshold. Vegetation increases roughness and sheltering, and it can also
577 reduce sediment availability at the surface. These effects act directly on the susceptibility of the
578 land surface to wind erosion.

579 The larger gains observed in vegetated and transition regions indicate that these controls are
580 especially relevant where surface conditions are seasonally variable. They also help explain
581 why nonlinear tree-based models performed better once surface predictors were included. The
582 effect of wind depends on moisture and vegetation state. A model that can represent interactions



583 and threshold behavior is therefore better suited to this part of the problem.



584

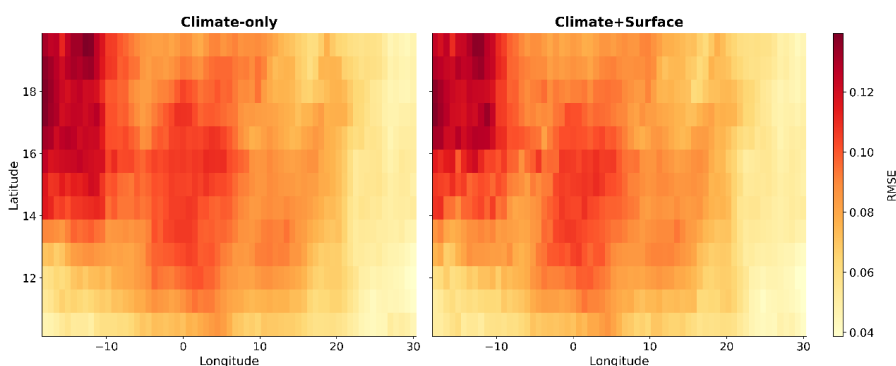
585 **Figure 2.** Difference in RMSE between the Climate-only and Climate+Surface configurations.
586 Positive values indicate reduced error after the inclusion of instantaneous surface variables. The
587 strongest gains occur in transition zones and in regions with stronger vegetation influence.

588

589 3.2.5. Benchmark interpretation across surface regimes

590 Taken together, the results show that the added value of surface variables is benchmark-region
591 dependent. Vegetated and semi-vegetated zones show the strongest gains. Transition zones
592 show intermediate but still meaningful improvement. Core desert regions show the smallest
593 gain, which is consistent with their more stable surface conditions. This hierarchy is visible
594 both in the regime statistics and in the spatial RMSE maps shown in Figure 3.

595 This pattern has a clear implication for the benchmark structure proposed in this study. Surface-
596 state information should not be treated as an optional refinement that only improves local detail.
597 It changes model skill in a systematic and regionally interpretable way. These findings justify
598 the inclusion of an explicit surface-state experiment tier in a future DesertMIP protocol,
599 particularly for transition and semi-vegetated regimes where climate-only predictors are least
600 adequate.



601



602 **Figure 3.** Spatial distribution of RMSE for the Climate-only configuration (a) and the
 603 Climate+Surface configuration (b). Error reduction is concentrated in benchmark regions with
 604 stronger land-surface variability, especially semi-arid and transition environments.

605 3.3. Role of Surface Memory in Improving Dust Reconstruction

606 3.3.1. Improvement in skill after adding memory components

607 Adding surface-memory terms produced a further improvement in dust AOD reconstruction
 608 after instantaneous surface variables had already been included. The increase in skill was
 609 modest, but it was consistent with the idea that dust activity is influenced not only by current
 610 land-surface conditions, but also by surface states inherited from previous months. Antecedent
 611 soil moisture and vegetation conditions may reflect accumulated drying, delayed vegetation
 612 response, or persistent exposure of erodible surfaces.

613 The improvement should not be overstated. The differences among the memory configurations
 614 are small, and they do not imply that a single lag window controls dust variability everywhere.
 615 Nevertheless, the direction of the change supports the inclusion of a separate memory-aware
 616 tier in the proposed benchmark hierarchy.

617 3.3.2. Sensitivity to lag length

618 Sensitivity analysis was carried out for lag windows of 0, 1, 3, and 6 months. All lagged-
 619 memory cases performed slightly better than the zero-lag case in terms of mean R^2 . The zero-
 620 lag configuration gave mean $R^2 = 0.671$ and mean $RMSE = 0.0576$. The 1-month and 3-
 621 month lags both yielded mean $R^2 = 0.678$ and mean $RMSE = 0.057$. The 6-month lag gave
 622 the highest mean $R^2 = 0.679$ and the lowest standard deviation of $R^2 = 0.0699$. It also gave
 623 the lowest standard deviation of $RMSE = 0.0043$ (Figure 4).

624 These results indicate that lagged predictors improve both mean skill and stability relative to
 625 the zero-lag baseline. At the same time, the differences among the lagged cases remain small.
 626 Importantly, the six-month lag gives the best predictive cross-validation performance in the
 627 present setup, whereas Sect. 3.7 identifies the strongest lagged associations near three months.
 628 This pattern supports the interpretation that antecedent surface conditions contain useful
 629 information, but it does not justify a strong claim that one lag window is uniquely optimal under
 630 all settings (Table 4).

631 **Table 4.** Cross-validation performance for different surface-memory lag lengths.

Lag (months)	Mean R^2	SD(R^2)	Mean RMSE	SD(RMSE)
0	0.671	0.0766	0.0576	0.0047



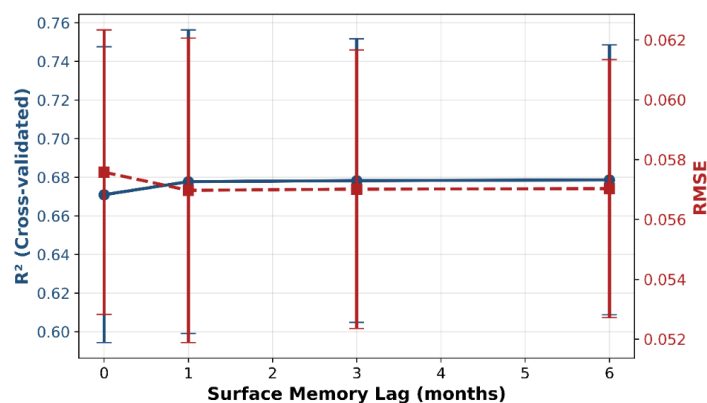
Lag (months)	Mean R ²	SD(R ²)	Mean RMSE	SD(RMSE)
1	0.678	0.0786	0.0570	0.0051
3	0.678	0.0734	0.0570	0.0046
6	0.679	0.0699	0.0570	0.0043

632

633 3.3.3. Variation of model performance across lag windows

634 The lag-sensitivity results suggest that memory terms may help stabilize model performance.
 635 When lagged surface variables were included, mean skill increased slightly and the spread of
 636 performance across cross-validation folds generally declined, especially for the 3-month and 6-
 637 month windows. This behavior is consistent with the role of memory variables as temporal
 638 summaries of recent surface conditions, rather than as direct replacements for current-month
 639 predictors.

640 At the same time, the effect remains moderate. The results do not show a sharp separation
 641 among the tested lag windows. Instead, they indicate that antecedent surface information is
 642 useful within a short-to-seasonal time range. Together with the lag-association analysis in Sect.
 643 3.7, these findings justify treating surface memory as a distinct diagnostic tier, while avoiding
 644 a strong claim that one lag length is universally optimal.



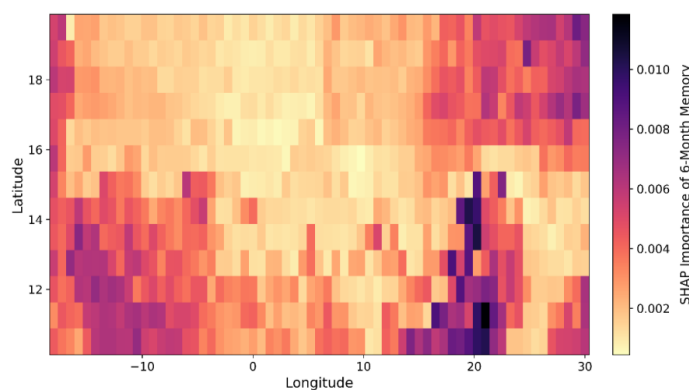
645

646 **Figure 4.** Cross-validation performance as a function of the lag applied to land-surface
 647 variables. Lagged-memory formulations show slightly higher mean R² and lower performance
 648 spread than the zero-lag baseline. The 6-month lag gives the best mean predictive performance
 649 in the present setup, although differences among lag windows remain modest. This should be
 650 distinguished from the lag-association analysis in Sect. 3.7, which peaks near three months.



651 **3.3.4. Spatial pattern of memory contribution**

652 The contribution of memory terms is not spatially uniform. Feature-importance analysis based
653 on SHAP indicates that memory effects are strongest in dry and semi-dry parts of the study
654 region. These areas are more likely to retain multi-month signatures of drying, vegetation
655 decline, and persistent surface exposure. By contrast, regions with denser vegetation and wetter
656 surface conditions show weaker memory contributions. This spatial contrast is consistent with
657 the expectation that antecedent surface state matters most where erosion sensitivity evolves
658 gradually through time rather than responding only to the current month.



659

660 **Figure 5.** Spatial importance of land-surface memory based on SHAP values for the 6-month
661 lag formulation. Higher values indicate regions where antecedent surface state contributes more
662 strongly to the reconstruction of dust AOD. Memory effects are strongest in dry and semi-dry
663 environments and weaker in more densely vegetated areas.

664 **3.4. Model Explainability and Variable Importance**

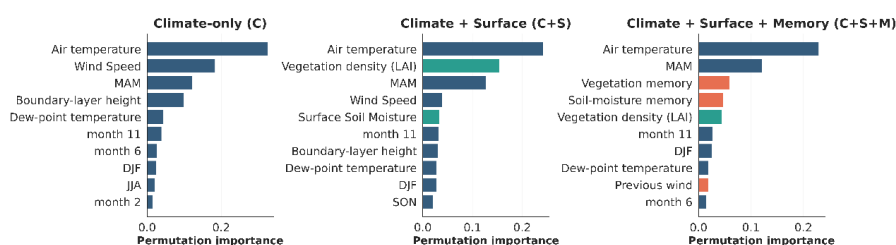
665 **3.4.1. Global importance structure across model configurations**

666 The explainability analysis shows that dust AOD variability is controlled by a combination of
667 atmospheric forcing, instantaneous land-surface state, and antecedent surface memory. Figure
668 6 summarizes the importance structure across the three explanatory configurations. In the
669 Climate-only configuration, atmospheric variables dominate by construction. Among these
670 predictors, air temperature shows the highest importance, followed by wind-related terms and
671 the spring seasonal indicator. This pattern suggests that the large-scale thermodynamic and
672 dynamical state of the atmosphere exerts a strong influence on the background conditions that
673 favor dust activity.

674 The inclusion of land-surface variables in the Climate+Surface configuration changes the
675 structure of predictor importance. In particular, vegetation-related predictors such as leaf area
676 index (LAI) become prominent. This shift indicates that part of the explained variability is



677 transferred from purely atmospheric drivers to variables that represent surface resistance and
 678 vegetation control. After memory terms are added, the importance structure expands further. In
 679 the full configuration, both instantaneous surface variables and memory-based predictors
 680 contribute to the explanation of dust variability. These results indicate that atmospheric forcing
 681 alone does not capture the full information content of the reconstruction problem.



682

683 **Figure 6.** Permutation-based variable importance across the three explanatory configurations.
 684 Atmospheric variables dominate in the Climate-only configuration, while land-surface and
 685 memory variables become more important in the Climate+Surface and
 686 Climate+Surface+Memory configurations.

687 3.4.2. Shifts in variable-group contribution

688 The transition from Climate-only to Climate+Surface and then to Climate+Surface+Memory is
 689 associated with a redistribution of importance across predictor groups. In the Climate-only case,
 690 all explanatory power is assigned to atmospheric variables. Once surface variables are
 691 introduced, part of that importance shifts toward predictors that describe vegetation and
 692 moisture state. In the full configuration, memory variables account for an additional share of
 693 the explanatory structure. This result suggests that antecedent land-surface conditions contain
 694 information that is not fully captured by climate forcing or instantaneous surface state alone.

695 The grouped contribution estimates from the importance analysis are reported here only as
 696 relative importance shares, not as a second set of model-performance metrics. In the
 697 Climate+Surface configuration, climate variables account for about 74.45 % of the total
 698 importance, whereas surface variables account for about 25.54 %. In the
 699 Climate+Surface+Memory configuration, the corresponding shares are about 71.14 % for
 700 climate variables, 9.15 % for instantaneous surface variables, and 19.69 % for memory
 701 variables. These percentages indicate how the explanatory structure of the model changes as
 702 additional information is introduced. They should therefore be interpreted as diagnostic
 703 indicators of model dependence on different predictor groups, rather than as causal attribution
 704 or as substitutes for the performance metrics reported in Table 2.

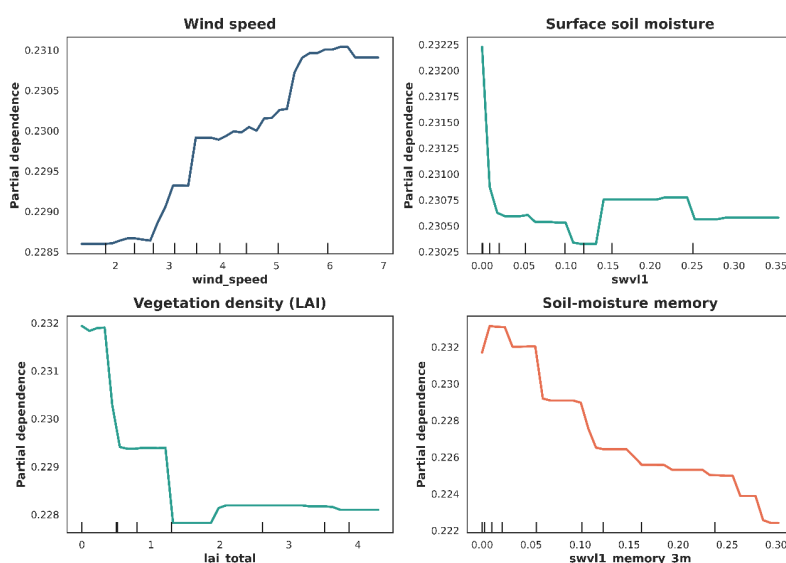
705 3.4.3. Response-shape diagnostics for key predictors



706 Partial-dependence plots provide additional insight into how the reference model responds to
 707 key predictors. Figure 7 shows that several predictors are associated with nonlinear and
 708 threshold-like behavior. Wind speed is associated with a rapid rise in predicted AOD at low to
 709 moderate values. The response then tends to flatten near about 5–6 m s⁻¹. This pattern is
 710 consistent with a transition from a strongly wind-limited regime to a more supply-limited
 711 regime.

712 Surface soil moisture shows a negative and nonlinear relationship with predicted AOD. The
 713 strongest decline occurs at very low moisture levels, roughly within the 0–0.05 range. At higher
 714 moisture values, the response becomes flatter. This pattern is physically plausible because small
 715 increases in moisture at the dry end can strongly increase particle cohesion, while further
 716 increases have less incremental effect on already stabilized surfaces.

717 Vegetation and memory predictors show similar response-shape behavior. Higher LAI values
 718 are associated with lower predicted AOD. The three-month soil-moisture memory term also
 719 shows a decreasing response. This result suggests that wetter antecedent conditions reduce the
 720 likelihood of dust activity during the current month. These plots should be interpreted as
 721 standardized response summaries of the fitted model rather than as direct causal evidence.



722

723 **Figure 7.** Partial-dependence plots for selected predictors, including wind speed, soil moisture,
 724 LAI, and soil-moisture memory. The plots show nonlinear response shapes and threshold-like
 725 behavior in the fitted model.

726 3.4.4. Interaction diagnostics in the erosion response

727 The interaction between wind speed and soil moisture provides a useful example of coupled



728 control in the fitted model. Figure S1 shows that high predicted AOD values occur mainly under
729 the joint condition of strong winds and very low soil moisture. As soil moisture increases,
730 predicted AOD declines even when wind remains relatively strong. This pattern is consistent
731 with the idea that atmospheric forcing alone is not sufficient to explain erosion response. The
732 land surface modulates whether strong wind can translate into high dust activity.

733 This interaction should be interpreted as a model-based diagnostic of coupled sensitivity rather
734 than as a direct process proof. Even so, it provides useful interpretive support for the climate–
735 surface–memory framework used in this study. Within a prototype DesertMIP framework,
736 variable-group attribution and response-shape diagnostics are proposed as complementary
737 interpretive products to accompany standard error metrics.

738 **3.5. Regime-based Analysis of Error and Bias**

739 Regime-based evaluation is treated here as a core DesertMIP diagnostic because dust-model
740 behaviour is expected to differ systematically between desert cores, ecological transition zones,
741 and more vegetated margins. In this section, mean bias is defined consistently as predicted
742 minus observed. Under this convention, positive bias indicates overestimation and negative bias
743 indicates underestimation.

744 **3.5.1. Error distribution across core desert, transition, and vegetated regimes**

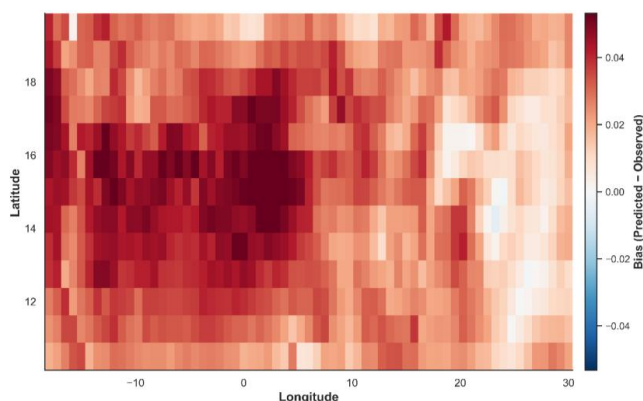
745 Evaluation of the final Climate+Surface+Memory configuration shows that prediction error
746 varies across benchmark surface regimes. The highest RMSE occurs in the transition zone
747 (0.078), followed by the core desert (0.076), whereas the vegetated zone shows the lowest error
748 (0.067). This pattern indicates that transition environments are the most difficult regime for dust
749 AOD reconstruction. The elevated uncertainty in transition zones is consistent with their
750 stronger environmental heterogeneity. These regions combine sharper gradients in vegetation,
751 moisture, and surface resistance than the more homogeneous desert core. By contrast, the
752 vegetated zone shows the lowest error level, which suggests that model behaviour is more stable
753 where land-surface conditions are comparatively less dust-prone and more structured. The core
754 desert occupies an intermediate position, with error levels lower than in the transition zone but
755 higher than in the vegetated regime.

756 **3.5.2. Spatial concentration of bias**

757 The regional bias map for the final configuration is shown in Figure 8. Because bias is defined
758 here as predicted minus observed, positive values indicate overestimation. Red areas therefore
759 mark regions where the model tends to predict dust AOD values that are too high. The map
760 shows that positive bias is spatially concentrated in transition environments and along desert
761 margins. Mean bias values are positive in all three regimes, but they differ only modestly: the
762 transition zone shows the largest mean bias (0.015), followed by the core desert (0.014) and the



763 vegetated zone (0.013). These results indicate that transition environments are not only the most
764 uncertain in terms of RMSE, but also the most sensitive to systematic error. At the same time,
765 the relatively small differences among the three regimes suggest that regime contrasts are more
766 pronounced for error magnitude than for mean bias alone. This behaviour is consistent with the
767 stronger ecological and climatic gradients that characterize transition regions.



768

769 **Figure 8.** Spatial distribution of model bias for the Climate+Surface+Memory configuration.
770 Bias is defined as predicted minus observed. Positive values (red) indicate overestimation, and
771 negative values (blue) indicate underestimation. Overestimation is concentrated mainly in
772 transition environments and along desert margins.

773 3.5.3. Relationship between bias and the Desert Proxy Index

774 Figure 9a examines the relationship between model bias and the Desert Proxy Index. The
775 overall association is weak. The Pearson correlation coefficient is $r = -0.02$, and the
776 Spearman coefficient is $\rho = -0.08$. These values indicate that there is no strong monotonic
777 relationship between mean bias and increasing desert-like conditions across the full domain.

778 Even so, the hexbin distribution and the median trend line show greater dispersion of positive
779 bias values within the Desert Proxy Index range of about 0.6 to 0.9. This pattern suggests that
780 highly arid environments can still be associated with less stable model behaviour, even though
781 the relationship is not strong in a domain-wide linear sense. The result therefore points to
782 localized sensitivity rather than to a simple large-scale bias gradient.

783 3.5.4. Relationship between RMSE and transition intensity

784 Figure 9b shows the relationship between RMSE and transition intensity. The association is
785 again weak, with a Pearson correlation coefficient of $r = -0.06$. Mean RMSE values across
786 transition-intensity classes are also very similar. The mean RMSE is 0.081 for low intensity,
787 0.078 for moderate intensity, and 0.079 for high intensity. These results indicate that transition
788 intensity alone does not explain error magnitude in a simple linear way. At the same time, the



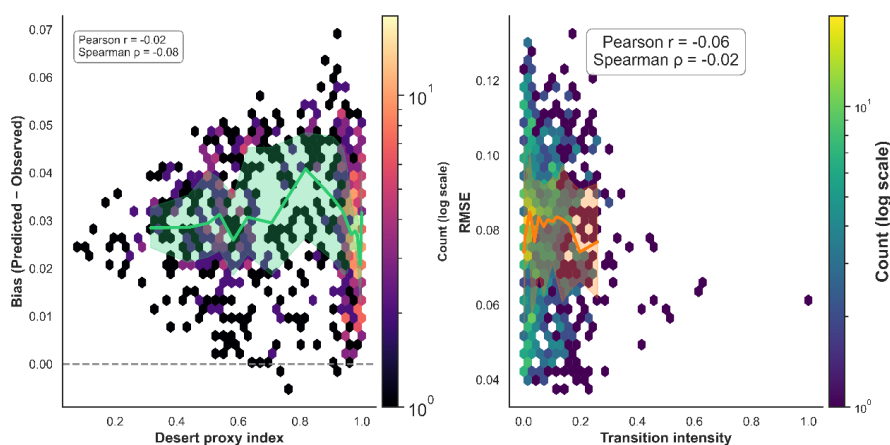
789 regime-based comparison in Sect. 3.5.1 still shows that transition zones have the highest overall
790 RMSE. The implication is that transition-related uncertainty is better expressed at the regime
791 level than through a simple pointwise correlation with transition intensity.

792

793 3.5.5. Changes in spatial error patterns across model configurations

794 The comparison across model configurations shows that adding surface and memory predictors
795 improves the spatial structure of model error, although the improvement is not uniform across
796 regimes. In the final Climate+Surface+Memory configuration, RMSE is 0.076 in the core desert,
797 0.078 in the transition zone, and 0.067 in the vegetated zone (Table 7). The transition zone
798 therefore remains the most difficult environment to reconstruct, while the vegetated zone has
799 the lowest residual error.

800 This ordering is consistent with the idea that desert margins and transition environments are
801 more difficult to model because they combine sharp gradients in vegetation, soil moisture, and
802 surface resistance. Figure S2 further shows that positive-bias clusters are stronger in the
803 Climate-only configuration and become weaker after surface and memory information is
804 included. Thus, the full configuration improves not only the average skill, but also the regional
805 coherence of the error field.



806

807 **Figure 9.** Relationships between model error metrics and surface indices. (a) Relationship
808 between model bias, defined as predicted minus observed, and the Desert Proxy Index. The
809 hexbin shading shows point density, and the thick line shows median bias across index intervals.
810 (b) Relationship between RMSE and transition intensity. Mean RMSE varies only weakly
811 across transition-intensity classes.

812



813 **3.6. Model Performance in the Reconstruction of Extreme Dust Events**

814 **3.6.1. Evaluation of extreme-event detection**

815 Model performance under extreme conditions was evaluated using the upper 5 % of the AOD
816 distribution. This subset corresponds to 2,231 samples in the test dataset and represents severe
817 dust events. Detection metrics for these events are summarized in Table 5. The results show a
818 clear reduction in skill relative to the overall reconstruction problem. In particular, the
819 probability of detection (POD) and the critical success index (CSI) are lower than would be
820 expected from the aggregate skill scores alone. This pattern indicates that severe dust outbreaks
821 are more difficult to reconstruct than moderate or typical conditions.

822 This result is physically and statistically plausible. Extreme events are rare. They also tend to
823 occupy the upper tail of the response distribution, where regression models often show stronger
824 shrinkage toward central values. As a result, acceptable mean-state performance does not
825 guarantee accurate representation of severe dust outbreaks.

826 **3.6.2. Comparison of model configurations for extreme events**

827 A comparison of the three explanatory configurations shows that the inclusion of surface and
828 memory predictors does not improve extreme-event detection. Instead, detection metrics
829 decline slightly from the Climate-only configuration to Climate+Surface and then to
830 Climate+Surface+Memory. POD decreases from 0.437 in the Climate-only case to 0.431 in
831 Climate+Surface and to 0.423 in Climate+Surface+Memory. CSI decreases from 0.276 to 0.272
832 and then to 0.269. FAR remains high in all cases, with values between 0.572 and 0.576.

833 These results indicate that the added value of surface and memory predictors is concentrated in
834 the reconstruction of the overall distribution and mean-state structure rather than in the
835 detection of the most severe events. From a DesertMIP perspective, this result is valuable
836 because it identifies tail performance as a separate diagnostic target rather than assuming that
837 improvements in overall skill automatically transfer to severe dust outbreaks.

838 **Table 5.** Detection metrics for extreme dust events across the three explanatory configurations.
839 Extreme events are defined as observations above the 95th percentile of the AOD distribution.

Scenario	POD	CSI	FAR
C	0.437	0.276	0.572
C+S	0.431	0.272	0.576
C+S+M	0.423	0.269	0.574

840



841 **3.6.3. Error amplification during extreme events**

842 The error structure changes markedly under extreme conditions. For the full
843 Climate+Surface+Memory configuration, the extreme-event RMSE is 0.170 and the extreme-
844 event MAE is 0.141, with an extreme bias of -0.060 . Under the bias convention used here, this
845 negative value indicates systematic underestimation of extreme dust intensity. Similar behavior
846 is also found in the Climate-only and Climate+Surface configurations, which yield extreme
847 bias values of -0.058 and -0.059 , respectively.

848 Taken together, these results show that the upper tail of the AOD distribution remains difficult
849 to reconstruct across all configurations. While surface and memory information improve the
850 physical realism of the framework and support better process representation, they do not
851 eliminate the tendency toward underestimation in the most severe events. In other words, the
852 main remaining challenge lies less in identifying where extreme events occur than in
853 reproducing their full magnitude.

854 **Table 6.** Quantitative error statistics for extreme dust events across the three explanatory
855 configurations

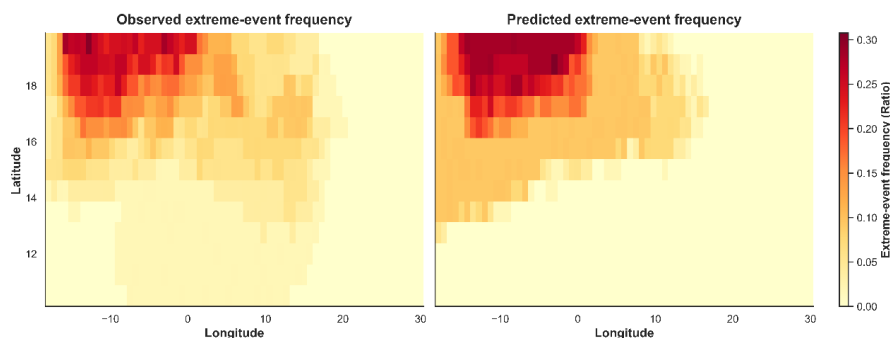
Scenario	n extreme	Extreme RMSE	Extreme MAE	Extreme Bias
C	2231	0.161	0.133	-0.058
C+S	2231	0.163	0.134	-0.059
C+S+M	2231	0.170	0.141	-0.060

856

857 **3.6.4. Implications for benchmark design and future model development**

858 Figure 10 shows that the full model reproduces the broad geographic pattern of extreme dust
859 occurrence, but it still underestimates event intensity. This distinction is important because
860 spatial placement and event magnitude are not the same measure of skill. A model may identify
861 the regions where severe dust is likely to occur while still failing to reproduce the strength of
862 the most intense events.

863 For this reason, severe-event evaluation should be treated separately from mean-state
864 reconstruction. The present results do not show that surface-memory terms improve extreme-
865 event detection. Future model development may therefore require methods that are more
866 sensitive to rare events, such as tail-aware loss functions, resampling strategies, or specialized
867 extreme-value diagnostics.



868

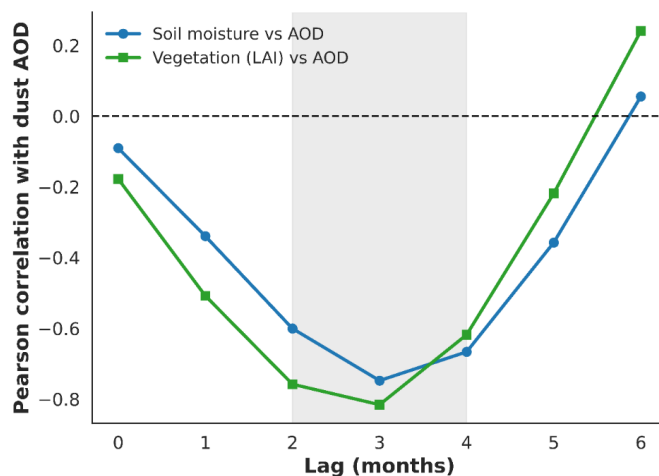
869 **Figure 10.** Spatial distribution of extreme dust AOD events for the observed data and for the
870 Climate+Surface+Memory reconstruction. The model reproduces the broad spatial structure of
871 extreme-event occurrence but underestimates event intensity, especially in the most severe
872 cases.

873 3.7. Surface Memory Analysis and Lagged Relationships

874 3.7.1. Lag-dependent associations between surface variables and AOD

875 To further examine the temporal structure of surface influence, Pearson correlations were
876 calculated between dust AOD and two key land-surface variables, soil moisture (SM) and leaf
877 area index (LAI), across lags from 0 to 6 months. The results are shown in Figure 11.
878 Correlation strength varies systematically with lag. Contemporaneous correlations are
879 relatively weak, with values of about -0.09 for soil moisture and -0.18 for LAI at lag 0.
880 Stronger negative correlations emerge at longer lags, and the largest magnitudes occur near a
881 lag of about 3 months. Peak values reach about -0.75 for soil moisture and -0.82 for LAI.
882 Correlation strength then weakens after the fourth month.

883 These results indicate that lagged surface associations are stronger than contemporaneous
884 associations in the present dataset. Importantly, these association peaks near three months
885 should be distinguished from the predictive-skill comparison in Sect. 3.3, where the six-month
886 lag yielded the highest mean cross-validation performance, albeit by a small margin. However,
887 these patterns should be interpreted with caution. They provide evidence for memory-sensitive
888 benchmarking, not direct proof of causal memory mechanisms. These lag diagnostics are
889 proposed as benchmark descriptors within the prototype DesertMIP framework, while
890 recognizing that part of the signal may reflect seasonality and temporal autocorrelation.



891

892 **Figure 11.** Pearson correlation between dust AOD and key land-surface variables, soil moisture
 893 (SM) and leaf area index (LAI), across temporal lags from 0 to 6 months. The strongest negative
 894 associations occur within the two- to four-month window, with a peak near three months. This
 895 lag of maximum association is distinct from the predictive-skill analysis in Sect. 3.3, in which
 896 the six-month lag produced the highest mean cross-validation score. These lag patterns are
 897 presented as benchmark descriptors of memory-sensitive behavior rather than as stand-alone
 898 proof of causal memory mechanisms.

899 3.7.2. Comparison of vegetation and soil-moisture behavior

900 The lag structure differs between LAI and soil moisture. LAI shows slightly stronger negative
 901 associations with dust AOD than soil moisture across the main lag window. Its maximum
 902 correlation reaches about -0.82 , whereas the maximum for soil moisture is about -0.75 . The
 903 decline in correlation strength at longer lags is also more gradual for LAI than for soil moisture.

904 This contrast is consistent with the idea that vegetation-related controls may persist longer than
 905 moisture-related controls in the present dataset. Vegetation affects roughness, sheltering, and
 906 surface stabilization. Soil moisture mainly affects particle cohesion. These differences may help
 907 explain why LAI retains a stronger lagged signal over a wider temporal window. Even so, the
 908 present results should be read as comparative lag associations rather than as direct evidence
 909 that one process is mechanistically dominant.

910 3.7.3. Temporal interpretation of lagged effects

911 The lag structure can be summarized conceptually by writing dust AOD as a function of both
 912 present and antecedent surface states:

$$913 \quad AOD_t = f(S_t, S_{t-1}, S_{t-2}, \dots) \quad (32)$$



914 where S denotes land-surface state. In the present dataset, the strongest associations occur
915 within the two- to four-month window prior to the target month, with the most pronounced
916 signal near three months. This association-based result is compatible with, but not identical to,
917 the predictive experiment reported in Section 3.3, in which the six-month lag provided the
918 highest mean cross-validation skill. The key point here is not that a unique causal lag has been
919 identified. Rather, the results indicate that benchmark diagnostics based on antecedent surface
920 state are informative and should be retained in a climate–surface–memory evaluation
921 framework. For this reason, lag-sensitive association analysis is treated here as part of the
922 prototype DesertMIP benchmark logic.

923 **3.7.4. Surface-memory signatures under extreme dust conditions**

924 The lag analysis also provides useful context for severe dust conditions. Extreme events were
925 defined as the upper 5 % of the AOD distribution. Surface conditions associated with these
926 events are summarized in Figure S3. Mean soil moisture during extreme events is about 0.026,
927 compared with about 0.094 under non-extreme conditions. Vegetation indicators are also lower
928 during extreme events, and the corresponding memory terms show similar reductions in the
929 months preceding those events.

930 These differences indicate that severe dust conditions are associated with drier and less
931 vegetated antecedent surface states. However, this result should not be interpreted as proof that
932 surface memory alone improves extreme-event reconstruction. Sect. 3.6 showed that the
933 addition of surface and memory predictors did not improve extreme-event detection metrics.
934 The present subsection therefore identifies lagged surface contrasts associated with severe
935 events, but it does not overturn the earlier conclusion that tail reconstruction remains limited.

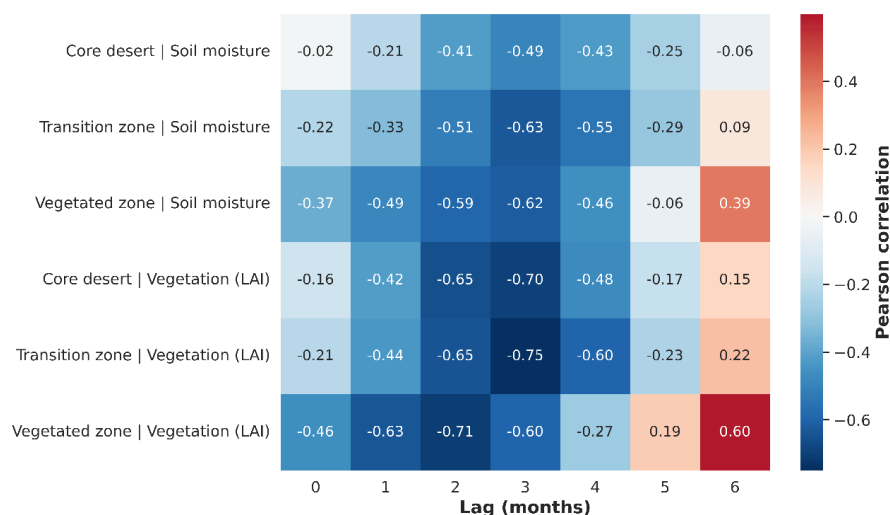
936 **3.7.5. Differences in lag response across benchmark regimes**

937 Lag-dependent associations also vary across benchmark regimes, as shown in Figure 12. In the
938 core-desert regime, the strongest correlations occur near a lag of three months, with values of
939 about -0.49 for soil moisture and -0.70 for LAI. Stronger lagged associations appear in
940 transition zones, where the corresponding values reach about -0.63 for soil moisture and -0.75
941 for LAI. Vegetated regions show a particularly strong LAI signal during the second and third
942 months, with correlations exceeding about -0.71 in some cases.

943 These results indicate that lag-sensitive benchmark behaviour is regime dependent. Transition
944 zones and more variable vegetated margins show stronger lag structure than the more
945 homogeneous desert core. This pattern is consistent with the broader finding that surface-state
946 and memory information are most informative where the land surface changes more strongly
947 through time. As in the domain-wide analysis, these regime-based correlations identify the lag
948 of strongest association, not necessarily the lag of best predictive skill. In summary, dust
949 variability in the present dataset is associated with both instantaneous and antecedent surface



950 conditions. These lag diagnostics therefore support the inclusion of memory-sensitive
 951 descriptors in future DesertMIP-style benchmark evaluations, while remaining distinct from
 952 claims of direct causal attribution.



953

954 **Figure 12.** Lag-dependent associations between land-surface variables and dust AOD across
 955 benchmark regimes, including the core desert, transition zone, and vegetated region. The
 956 strongest lag structure appears in transition and more variable vegetated environments. The
 957 figure summarizes regime-dependent association peaks rather than predictive-skill optima. The
 958 figure is intended to show regime-dependent benchmarking signals within the climate–surface–
 959 memory framework.

960 **3.8. Spatial Analysis of Dust Patterns and Model Error**

961 Observed, reconstructed, and error maps are treated here as mandatory benchmark products
 962 because they reveal whether model gains are regionally coherent or confined to limited parts of
 963 the domain. In this section, spatial diagnostics are evaluated for the full
 964 Climate+Surface+Memory configuration. Consistent with Sect. 3.5, bias is defined as predicted
 965 minus observed. Positive values therefore indicate overestimation, and negative values indicate
 966 underestimation.

967 **3.8.1. Spatial pattern of observed AOD**

968 The multi-year mean map of observed dust AOD is shown in Figure 13a. The highest dust
 969 burden occurs in the northern and central parts of the study area. A broad core band between
 970 about 15° N and 20° N shows the largest values. Mean AOD reaches about 0.4 in the most
 971 active source regions. Southward from this band, AOD decreases as vegetation cover increases
 972 and surface conditions become less favorable for sustained dust emission.



973 This spatial structure is consistent with the benchmark regime framework used throughout the
974 study. It reflects the expected sequence from desert core to transition zone and then to more
975 vegetated areas. For this reason, the observed map provides the reference spatial benchmark
976 against which reconstructed patterns and model errors are evaluated.

977 **3.8.2. Spatial pattern of reconstructed AOD**

978 The reconstructed AOD field from the full Climate+Surface+Memory configuration is shown
979 in Figure 13b. The model reproduces the main regional structure of dust activity. In particular,
980 it captures the location of the main source band and the broad north–south gradient in mean
981 AOD. This agreement indicates that the model is able to recover the dominant spatial
982 organization of dust variability over North Africa.

983 At the same time, the comparison between observed and reconstructed means suggests that the
984 model tends to produce slightly higher values over part of the domain. This tendency is modest
985 in magnitude, but it is regionally widespread. The spatial error field provides a clearer view of
986 where this bias is concentrated.

987 **3.8.3. Spatial structure of model error**

988 The spatial error map is shown in Figure 13c and is defined as predicted minus observed. Under
989 this convention, positive values indicate overestimation and negative values indicate
990 underestimation. Most grid-cell differences fall within about -0.03 to $+0.03$. Positive values
991 are more widespread than negative values, which indicates that mild overestimation is the
992 dominant error mode in the full model configuration.

993 The strongest positive errors are concentrated along desert margins and in parts of the transition
994 zone. More localized negative errors occur in some central and eastern areas. This pattern
995 suggests that the main limitation of the model is not a failure to reproduce the overall spatial
996 structure of dust AOD. Instead, it is a tendency toward modest but spatially coherent bias in
997 specific benchmark regions.

998 **3.8.4. Regime-based summary of observed, reconstructed, and bias patterns**

999 The available regime-based error metrics indicate that the vegetated zone has the lowest
1000 reconstruction error, with $MAE = 0.047$ and $RMSE = 0.067$. The transition zone remains the
1001 most difficult regime, with the highest $RMSE = 0.078$. The desert core shows intermediate
1002 behavior, with $MAE \approx 0.053$. These results are consistent with the benchmark interpretation
1003 developed in Sect. 3.5. They suggest that model performance depends strongly on structural
1004 regime differences in land-surface state and heterogeneity (Table 7).

1005

1006



1007 **Table 7.** Regime-based spatial summary of observed, reconstructed, and bias fields for the
 1008 full Climate+Surface+Memory configuration

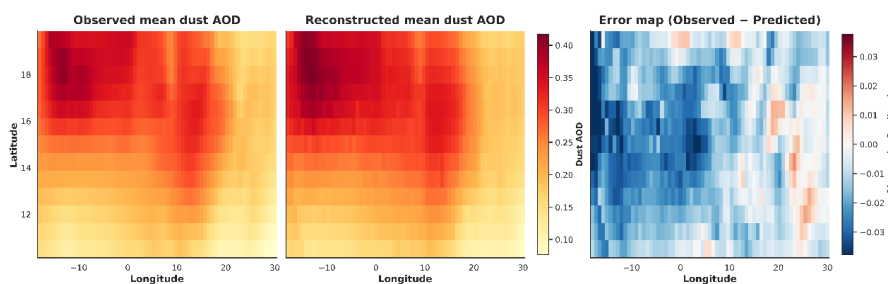
Surface regime	Mean observed AOD	Mean predicted AOD	Mean bias (Predicted - Observed)	MAE	RMSE
Core desert	0.275	0.289	0.014	0.054	0.076
Transition zone	0.256	0.271	0.015	0.055	0.078
Vegetated zone	0.178	0.192	0.013	0.047	0.067

1009

1010 **3.8.5. Relationship between spatial error and surface properties**

1011 The relationship between model error and surface characteristics is further examined in Figure
 1012 S4. The spatial comparison suggests two broad patterns. First, larger absolute errors tend to
 1013 occur in areas with stronger desert characteristics. Second, regions with stronger surface-
 1014 memory signals tend to exhibit lower absolute errors. Together, these patterns indicate that
 1015 model error is not randomly distributed, but is linked to spatial variations in aridity and
 1016 antecedent land-surface memory.

1017 These patterns are consistent with the broader benchmark interpretation of the study. Even in
 1018 the full Climate+Surface+Memory configuration, the most difficult areas remain those where
 1019 surface processes are highly threshold-dependent, spatially heterogeneous, or imperfectly
 1020 constrained by available observations. From a DesertMIP perspective, such spatial diagnostics
 1021 are useful because they identify where model gains are robust and where benchmark
 1022 performance remains limited.



1023

1024 **Figure 13.** Mandatory benchmark maps for the full Climate+Surface+Memory configuration.
 1025 (a) Multi-year mean observed dust AOD. (b) Multi-year mean reconstructed dust AOD. (c)



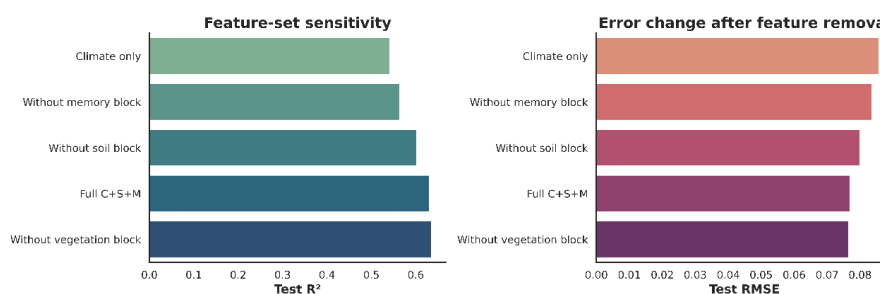
1026 Spatial error field, defined as predicted minus observed. Positive values indicate overestimation,
1027 and negative values indicate underestimation. Together, these maps show whether model
1028 performance is regionally coherent across the benchmark domain.

1029 3.9. Uncertainty and Stability Analysis

1030 In the context of a future DesertMIP protocol, sensitivity tests to variable blocks, temporal
1031 aggregation, threshold choice, and memory-window definition are recommended as minimum
1032 robustness checks for any benchmark implementation. The purpose of this section is therefore
1033 not to claim definitive robustness, but to evaluate whether the main findings remain
1034 qualitatively stable under several alternative analysis settings.

1035 3.9.1. Sensitivity to variable selection

1036 Feature-ablation experiments were used to test the dependence of model skill on major
1037 predictor blocks. The results shown in Figure 14 indicate three consistent patterns. First,
1038 removal of soil and vegetation variables reduces model performance. Second, removal of the
1039 memory block produces the largest decline in R^2 and the largest increase in RMSE. Third,
1040 omission of surface-memory terms reduces model skill to a level close to that of the Climate-
1041 only configuration. These results suggest that the gain in the full configuration does not arise
1042 solely from redundant predictors. Instead, it reflects additional information associated with
1043 dynamic surface state and antecedent surface history.



1044

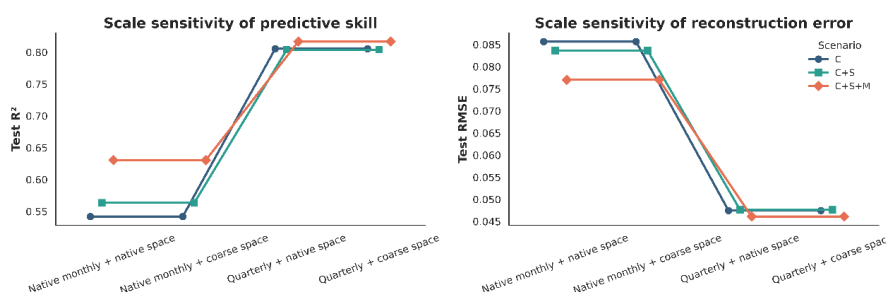
1045 **Figure 14.** Sensitivity of model performance to variable-block removal. The left panel shows
1046 changes in R^2 , and the right panel shows changes in RMSE after systematic removal of soil,
1047 vegetation, and memory predictor blocks. The figure highlights the contribution of surface-
1048 memory information to the full reconstruction framework.

1049 3.9.2. Sensitivity to spatial and temporal scale

1050 Sensitivity tests across temporal and spatial scales are shown in Figure 15. Temporal
1051 aggregation from monthly to seasonal resolution increases apparent skill, with R^2 rising from
1052 about 0.55 to values above 0.80 and RMSE decreasing accordingly. This pattern is consistent
1053 with the reduction of short-term temporal noise under seasonal averaging. The relative



1054 advantage of the Climate+Surface+Memory configuration remains most evident at the native
 1055 spatial scale. This result suggests that surface and memory information are especially relevant
 1056 for reconstructing finer-scale and shorter-term variability. At the same time, these comparisons
 1057 should be interpreted as sensitivity analyses rather than as evidence of fully independent model
 1058 generalization.



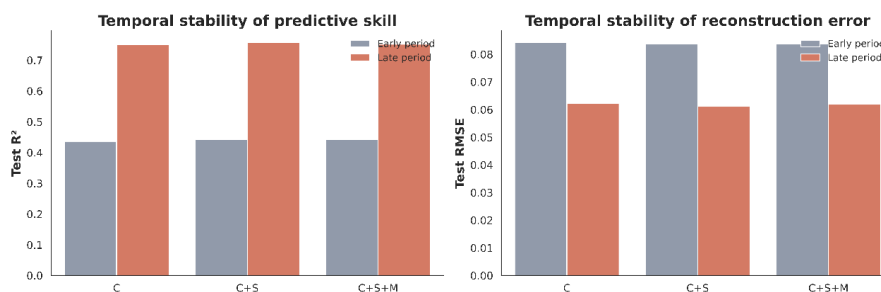
1059
 1060 **Figure 15.** Sensitivity of model performance to spatial and temporal scale. The left panel shows
 1061 changes in R^2 , and the right panel shows changes in RMSE across alternative aggregation
 1062 scales. Temporal aggregation increases apparent skill, while the advantage of the
 1063 Climate+Surface+Memory configuration remains most visible at the native scale.

1064 **3.9.3. Stability across independent time periods**

1065 Model performance was also compared across two non-overlapping periods, referred to here as
 1066 Early and Late. As shown in Figure 16, performance is higher in the Late period, with $R^2 \approx$
 1067 0.75, than in the Early period, with $R^2 \approx 0.44$. This difference may reflect changes in data
 1068 quality, evolving surface conditions, or temporal variation in predictor–response relationships.
 1069 Despite this shift in absolute performance, the ordering of the three explanatory configurations
 1070 remains unchanged:

$$\text{Climate-only} < \text{Climate+Surface} < \text{Climate+Surface+Memory}$$

1071 This consistency suggests that the relative benefit of surface and memory information is not
 1072 confined to a single part of the record. However, the comparison is temporal rather than spatial,
 1073 and it therefore does not remove the possibility that some optimism remains due to spatial
 1074 autocorrelation.

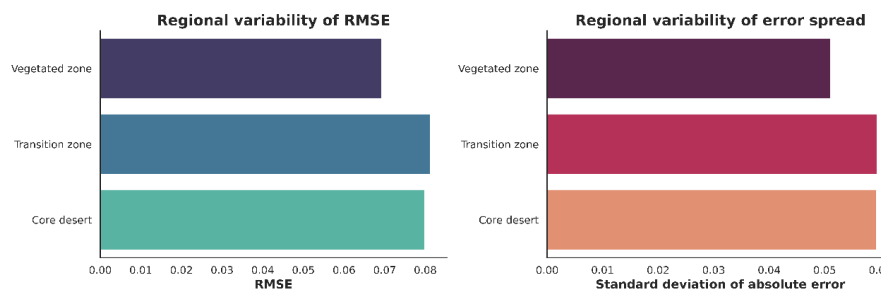


1076

1077 **Figure 16.** Stability of scenario ranking across two independent time periods. Performance is
 1078 shown using R^2 and RMSE for the Early and Late periods. Although absolute skill differs
 1079 between the two periods, the ordering of the three explanatory configurations remains
 1080 unchanged.

1081 **3.9.4. Inter-subregion variability**

1082 Figure 17 summarizes model behavior across the major surface regimes. The transition zone
 1083 shows the highest RMSE, while the core desert follows closely behind. The vegetated region
 1084 shows the lowest reconstruction error. The figure also indicates that the spread of absolute error
 1085 is largest in the transition zone and remains similarly high in the core desert, whereas it is clearly
 1086 smaller in the vegetated regime. This pattern is consistent with the regime-based results
 1087 reported in Section 3.5 and Table 7, which identify transition environments as the most difficult
 1088 benchmark regions and vegetated areas as the most stable ones. Together, these regime-level
 1089 comparisons provide descriptive evidence of spatial heterogeneity in model performance. They
 1090 should not, however, be interpreted as a substitute for stricter spatial validation such as leave-
 1091 region-out testing.



1092

1093 **Figure 17.** Regime-level summary of model error characteristics across the vegetated zone,
 1094 transition zone, and core desert. The transition zone exhibits the highest RMSE, whereas the
 1095 core desert shows a similarly large error spread. The vegetated zone has the lowest overall error
 1096 levels. These contrasts highlight the spatial heterogeneity of model performance and confirm
 1097 that transition environments remain the most challenging benchmark regime for dust-AOD

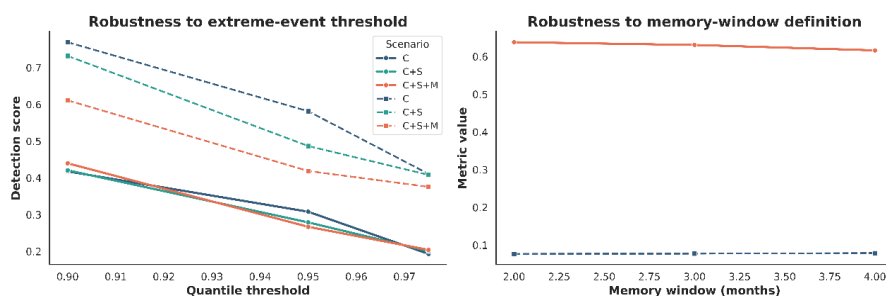


1098 reconstruction.

1099 3.9.5. Robustness across threshold and memory-window choices

1100 Figure 18 examines two additional sensitivity dimensions: the choice of extreme-event
1101 threshold and the definition of the memory window. The left panel shows that stricter thresholds,
1102 from quantile 0.90 to 0.97, reduce event-detection scores. This pattern is expected because more
1103 severe tail events are rarer and harder to identify. The right panel shows that changing the
1104 memory window from two to four months produces only minor changes in model skill. Together,
1105 these results suggest that the main conclusions about the usefulness of surface and memory
1106 information are not highly sensitive to moderate changes in tuning choices.

1107 Importantly, these robustness results should not be interpreted as evidence that the
1108 Climate+Surface+Memory configuration improves all aspects of performance under all
1109 thresholds. Section 3.6 showed that overall improvements in reconstruction skill do not
1110 automatically translate into better extreme-event detection, and that tail metrics such as POD
1111 and CSI can decline even when aggregate skill improves. The present sensitivity tests therefore
1112 support the stability of the general climate–surface–memory framework, but they do not
1113 overturn the earlier conclusion that severe-event reconstruction remains a separate challenge.



1114

1115 **Figure 18.** Robustness checks for the prototype benchmark framework. The left panel shows
1116 the sensitivity of extreme-event detection scores to the choice of quantile threshold. The right
1117 panel shows the sensitivity of overall model skill to the definition of the memory window from
1118 two to four months. The figure is intended as a protocol-level stability check rather than as
1119 evidence that all diagnostics improve uniformly under the full configuration.

1120 3.9.6. Summary of robustness checks

1121 Taken together, the sensitivity analyses support three main conclusions. First, the direction of
1122 the main findings remains broadly stable across variable-removal tests, scale changes, temporal
1123 partitions, and memory-window choices. Second, the addition of surface and memory
1124 information retains a consistent advantage in overall explanatory skill across these settings.
1125 Third, the present evidence should be interpreted as supportive rather than definitive because



1126 the validation design is still primarily temporal and not fully spatially independent. In this sense,
1127 the robustness analysis provides a useful protocol-level check for the proposed benchmark
1128 framework, while also identifying the need for stricter spatial validation in future DesertMIP-
1129 style implementations.

1130 **4 Discussion**

1131 **4.1. What is robustly supported by the present results**

1132 The present results support four main conclusions. First, climate variables alone are not
1133 sufficient to explain the observed variability of dust AOD over North Africa. The addition of
1134 dynamic land-surface variables improves reconstruction skill in a consistent way. This result
1135 agrees with earlier studies that identified wind as a necessary but not sufficient control on dust
1136 emission. Soil moisture, vegetation cover, surface roughness, and source characteristics also
1137 shape erosion thresholds and emission efficiency (Mousavi et al., 2023; Wu et al., 2016). The
1138 present study does not introduce entirely new controls. Its contribution is to evaluate them
1139 within a single hierarchical framework and to quantify their incremental value step by step.

1140 Second, lagged surface predictors add a further, although modest, improvement beyond
1141 instantaneous surface variables. This result is consistent with the idea that dust variability is not
1142 fully captured by a purely synchronous representation of land-surface state. Instead, part of the
1143 signal appears to be associated with antecedent conditions such as progressive drying,
1144 vegetation decline, and multi-month changes in surface resistance. This interpretation is
1145 broadly consistent with studies from the Sahel and the Sahara, which show that rainfall, soil
1146 moisture, and vegetation can influence dust activity over both short and extended timescales
1147 (Bergametti et al., 2017). At the same time, the lag analysis in the present study should be
1148 interpreted cautiously. It supports the usefulness of memory-sensitive predictors and
1149 diagnostics, but it does not by itself isolate a unique causal mechanism.

1150 Third, regime dependence is strong. Model error is not distributed uniformly across the domain.
1151 Transition zones show the highest uncertainty and the strongest sensitivity to surface and
1152 memory information. This pattern is physically plausible. Transitional environments combine
1153 sharper gradients in vegetation, moisture, and surface roughness than the more homogeneous
1154 desert core. Recent work has also emphasized the importance of spatial heterogeneity and
1155 resolution in dust modeling (Chappell et al., 2024). In that context, the present results support
1156 the broader view that fixed source masks and static boundary descriptions can be least adequate
1157 where ecological boundaries are dynamic (LeGrand et al., 2023).

1158 Fourth, severe events remain difficult to reconstruct. The model reproduces large-scale spatial
1159 structure and mean behaviour more successfully than the upper tail of the AOD distribution.
1160 This limitation is consistent with earlier evaluations of operational and research dust models,



1161 which often reproduce broad patterns more successfully than peak magnitude or event
1162 frequency (L. Zhang et al., 2022). In the present case, the addition of surface and memory
1163 predictors improves overall reconstruction skill, but it does not improve extreme-event
1164 detection metrics. This distinction is important and should be retained throughout the
1165 interpretation of the study.

1166 **4.2. What do these findings mean for DesertMIP?**

1167 The present findings support the need for a benchmark structure in which dust-model
1168 performance is evaluated not only against overall error metrics, but also with respect to surface-
1169 state sensitivity, regime dependence, antecedent-memory effects, and tail behaviour. In this
1170 sense, the current study provides a pilot diagnostic blueprint toward DesertMIP.

1171 This implication arises directly from the structure of the results. The hierarchical comparison
1172 shows that the step from climate-only predictors to climate plus surface predictors is
1173 informative. The further step from instantaneous surface state to surface memory also remains
1174 informative, even though the quantitative gains are modest. Regime-based analyses show that
1175 these gains are not spatially uniform. Extreme-event analyses show that improvements in mean-
1176 state performance do not automatically transfer to severe outbreaks. Taken together, these
1177 results point to a benchmark logic that is more structured than a single global score.

1178 From this perspective, DesertMIP should not be understood simply as another intercomparison
1179 label. It should be understood as a framework for organizing dust-model evaluation around the
1180 climate–surface–memory problem. That problem includes at least four linked dimensions:
1181 overall skill, regime-specific behaviour, lag-sensitive surface influence, and tail performance.
1182 The present study suggests that all four dimensions are necessary if dust-model behaviour is to
1183 be evaluated in a way that is both physically interpretable and diagnostically useful.

1184 **4.3. What should a DesertMIP-style benchmark contain?**

1185 A useful DesertMIP-style benchmark should include a small number of clearly defined
1186 elements. First, it should include hierarchical experiments. At minimum, these should
1187 distinguish climate-only predictors, climate plus instantaneous surface state, and climate plus
1188 surface plus memory. Such a structure allows the added value of each information block to be
1189 tested directly.

1190 Second, it should use public and reproducible benchmark inputs. This is important for
1191 transferability and for independent repetition of the benchmark across models and groups. Third,
1192 it should include regime masks that separate desert cores, transition zones, and more vegetated
1193 margins. The present results show that these zones behave differently and should therefore be
1194 treated as mandatory reporting strata rather than optional regional summaries.



1195 Fourth, the framework should include tail diagnostics in addition to standard regression metrics.
1196 The present study shows clearly that aggregate skill and extreme-event skill are not
1197 interchangeable. A model may improve in R^2 and RMSE while still failing to improve POD,
1198 CSI, or extreme-event bias. Fifth, the framework should retain lag diagnostics because they
1199 provide a useful description of memory-sensitive behaviour, even if they do not by themselves
1200 establish causality. Finally, a DesertMIP-style benchmark should include reproducible scripts,
1201 experiment definitions, and standardized output products. Without these elements, the
1202 framework remains conceptual rather than operational.

1203 **4.4. What is not yet demonstrated**

1204 The present study does not constitute a finalized intercomparison protocol and does not directly
1205 evaluate full process-based ESM dust schemes. Instead, it provides an observation-constrained
1206 pilot framework that can help define where such models should be tested more systematically.

1207 Several limitations follow from this scope. First, the benchmark is based on an interpretable
1208 machine-learning reconstruction framework, not on direct evaluation of process-based dust
1209 schemes within Earth system models. The present results can therefore identify useful
1210 diagnostic dimensions, but they cannot by themselves demonstrate the exact source of structural
1211 error in full ESM dust formulations. Second, the validation design is primarily temporal rather
1212 than fully spatially independent. This means that reported skill may still be influenced by spatial
1213 autocorrelation. Third, the lagged relationships reported here may partly reflect seasonality or
1214 temporal autocorrelation in addition to physically meaningful surface persistence. For this
1215 reason, the lag diagnostics should be interpreted as benchmark descriptors rather than as stand-
1216 alone evidence of mechanism.

1217 A further limitation concerns severe dust events. The present study shows that surface and
1218 memory predictors improve overall reconstruction but do not improve extreme-event detection.
1219 This means that a climate–surface–memory framework is necessary for a more complete
1220 benchmark, but it is not sufficient to solve the tail problem. Future progress will likely require
1221 tail-sensitive methods, stricter validation, and data structures that better resolve submonthly and
1222 mesoscale processes.

1223 Despite these limitations, the main value of the present study is clear. It shows that dust-model
1224 evaluation can be organized more effectively when climate forcing, instantaneous surface state,
1225 antecedent surface memory, regime dependence, and tail behaviour are treated as linked but
1226 distinct diagnostic dimensions. In that sense, the study provides a practical and reproducible
1227 first step toward DesertMIP.

1228 **5 Conclusion**

1229 This study examined whether dynamic land-surface information and antecedent surface-



1230 memory terms improve the reconstruction of dust AOD over North Africa. Three hierarchical
1231 configurations were evaluated: Climate-only, Climate plus Surface, and Climate plus Surface
1232 plus Memory. The results showed that climate variables alone did not capture the full variability
1233 of dust AOD. Adding instantaneous surface predictors improved reconstruction skill, and
1234 adding lagged surface information provided a further, although modest, gain.

1235 The memory-sensitive results should be interpreted with care. They show that antecedent
1236 surface conditions contain useful information for dust reconstruction, but they do not prove a
1237 single causal memory mechanism or a uniquely optimal lag. Instead, they support the use of
1238 lag-sensitive diagnostics as part of a broader climate–surface–memory benchmark structure.

1239 The results also showed that model performance was strongly regime dependent. Transition
1240 zones were the most difficult benchmark regions, with higher uncertainty and stronger
1241 sensitivity to surface and memory information. More homogeneous desert interiors showed
1242 smaller gains, while vegetated and semi-vegetated areas benefited more clearly from dynamic
1243 surface variables. This finding highlights the need for regime-aware evaluation in dust-model
1244 assessment.

1245 Severe dust outbreaks remained difficult to reconstruct. Although surface and memory
1246 predictors improved overall reconstruction skill, they did not improve extreme-event detection
1247 metrics. Tail behavior therefore remains a separate challenge and should be evaluated with
1248 dedicated diagnostics rather than inferred from average model performance.

1249 Taken together, the findings do not establish a finalized community intercomparison protocol.
1250 They do, however, provide a reproducible pilot architecture with clearly defined experiments,
1251 diagnostics, and benchmark products. The main contribution of the study is the definition of a
1252 climate–surface–memory benchmark logic that can guide future evaluation of dust models in
1253 dynamic desert environments.

1254 **Code and data availability**

1255 The DesertMIP project page is available at <https://earthsystemsci.github.io/DesertMIP/>. The
1256 model code used in this study is publicly available through Zenodo at
1257 <https://doi.org/10.5281/zenodo.20744633> (Mokarram, 2026a). The processed datasets and
1258 benchmark masks used to generate the results reported in this manuscript are publicly available
1259 through Zenodo at <https://doi.org/10.5281/zenodo.19727406> (Mokarram, 2026b). The code
1260 archive provides the publicly accessible DesertMIP-CSM v1.0 model code used in this study,
1261 and the data archive provides the processed datasets, benchmark masks, and derived outputs
1262 required to reproduce the reported results.

1263 **Acknowledgements**



1264 The authors acknowledge the institutional and research support provided by Anhui Xinhua
1265 University, Shiraz University, and Nanjing University. The authors also acknowledge Nanjing
1266 University for supporting the article processing charge.

1267 **Author contributions**

1268 **Marzieh Mokarram** designed the study, developed the machine-learning framework,
1269 conducted the analysis, and wrote the first draft of the manuscript. **Mohammad Jafar**
1270 **Mokarram** contributed to data preparation, methodological discussion, and manuscript
1271 revision. **Huayu Lu** contributed to the conceptual design of the research, supervised the
1272 scientific methodology, and reviewed the manuscript. All authors reviewed the results and
1273 approved the final version of the manuscript.

1274 **Competing interests**

1275 The authors declare that they have no competing interests.

1276 **Disclaimer**

1277 The publisher remains neutral with regard to jurisdictional claims in published maps,
1278 institutional affiliations, and geographical representations included in this article. The
1279 responsibility for the content of the article rests entirely with the authors.

1280 **Financial support**

1281 This work was supported by the National Natural Science Foundation of China under Grant No.
1282 2025zrzdi07.

1283 **References**

- 1284 Adebisi, A.A., Huang, Y., Samset, B.H., Kok, J.F., 2023. Observations suggest that North
1285 African dust absorbs less solar radiation than models estimate. *Commun. Earth Environ.*
1286 4, 168. <https://doi.org/10.1029/2021JF006073>
- 1287 Aryal, Y.N., Evans, S., 2021. Global dust variability explained by drought sensitivity in CMIP6
1288 models. *J. Geophys. Res. Earth Surf.* 126, e2021JF00607.
1289 <https://doi.org/10.1038/s43247-023-00825-2>
- 1290 Bergametti, G., Marticorena, B., Rajot, J.-L., Chatenet, B., Féron, A., Gaimoz, C., Siour, G.,
1291 Coulibaly, M., Koné, I., Maman, A., 2017. Dust uplift potential in the central Sahel: An
1292 analysis based on 10 years of meteorological measurements at high temporal resolution.
1293 *J. Geophys. Res. Atmos.* 122, 12–433. <https://doi.org/10.1002/2017JD027471>
- 1294 Chappell, A., Hennen, M., Schepanski, K., Dhital, S., Tong, D., 2024. Reducing resolution
1295 dependency of dust emission modeling using albedo-based wind friction. *Geophys. Res.*
1296 *Lett.* 51, e2023GL106540. <https://doi.org/10.1029/2023GL106540>
- 1297 Eyring, V., Bony, S., Meehl, G.A., Senior, C.A., Stevens, B., Stouffer, R.J., Taylor, K.E., 2016.



- 1298 Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental
1299 design and organization. *Geosci. Model Dev.* 9, 1937–1958. [https://doi.org/10.5194/gmd-](https://doi.org/10.5194/gmd-9-1937-2016)
1300 9-1937-2016
- 1301 Ginoux, P., Prospero, J.M., Gill, T.E., Hsu, N.C., Zhao, M., 2012. Global-scale attribution of
1302 anthropogenic and natural dust sources and their emission rates based on MODIS Deep
1303 Blue aerosol products. *Rev. Geophys.* 50.
- 1304 Green, J.K., Berry, J., Ciais, P., Zhang, Y., Gentine, P., 2020. Amazon rainforest photosynthesis
1305 increases in response to atmospheric dryness. *Sci. Adv.* 6, eabb7232.
- 1306 Huneus, N., Schulz, M., Balkanski, Y., Griesfeller, J., Prospero, J., Kinne, S., Bauer, S.,
1307 Boucher, O., Chin, M., Dentener, F., 2011. Global dust model intercomparison in
1308 AeroCom phase I. *Atmos. Chem. Phys.* 11, 7781–7816.
- 1309 Kim, D., Chin, M., Schuster, G., Yu, H., Takemura, T., Tuccella, P., Ginoux, P., Liu, X., Shi,
1310 Y., Matsui, H., 2024. Where dust comes from: Global assessment of dust source
1311 attributions with AeroCom models. *J. Geophys. Res. Atmos.* 129, e2024JD041377.
- 1312 Kok, J.F., Adebisi, A.A., Albani, S., Balkanski, Y., Checa-Garcia, R., Chin, M., Colarco, P.R.,
1313 Hamilton, D.S., Huang, Y., Ito, A., 2021. Contribution of the world’s main dust source
1314 regions to the global cycle of desert dust. *Atmos. Chem. Phys.* 21, 8169–8193.
- 1315 Kok, J.F., Storelvmo, T., Karydis, V.A., Adebisi, A.A., Mahowald, N.M., Evan, A.T., He, C.,
1316 Leung, D.M., 2023. Mineral dust aerosol impacts on global climate and climate change.
1317 *Nat. Rev. Earth Environ.* 4, 71–86.
- 1318 LeGrand, S.L., Letcher, T.W., Okin, G.S., Webb, N.P., Gallagher, A.R., Dhital, S., Hodgdon,
1319 T.S., Ziegler, N.P., Michaels, M.L., 2023. Application of a satellite-retrieved sheltering
1320 parameterization (v1.0) for dust event simulation with WRF-Chem v4.1. *Geosci. Model*
1321 *Dev.* 16, 1009–1038.
- 1322 Mahowald, N., Albani, S., Kok, J.F., Engelstaeder, S., Scanza, R., Ward, D.S., Flanner, M.G.,
1323 2014. The size distribution of desert dust aerosols and its impact on the Earth system.
1324 *Aeolian Res.* 15, 53–71.
- 1325 Mokarram, M., 2026a. DesertMIP-CSM v1.0 model code. Zenodo.
1326 <https://doi.org/10.5281/zenodo.20744633>.
- 1327 Mokarram, M., 2026b. DesertMIP processed datasets and benchmark masks, version 1 [data
1328 set]. Zenodo. <https://doi.org/10.5281/zenodo.19727406>
- 1329 Molnar, C., 2020. *Interpretable machine learning: A guide for making black box models*
1330 *explainable*. Lulu.com. <https://christophm.github.io/interpretable-ml-book/>
- 1331 Mousavi, S.V., Karami, K., Tilmes, S., Muri, H., Xia, L., Rezaei, A., 2023. Future dust
1332 concentration over the Middle East and North Africa region under global warming and
1333 stratospheric aerosol intervention scenarios. *Atmos. Chem. Phys.* 23, 10677–10695.
- 1334 Pu, B., Ginoux, P., 2018. How reliable are CMIP5 models in simulating dust optical depth?



- 1335 Atmos. Chem. Phys. 18, 12491–12510.
- 1336 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., Prabhat, F.,
1337 2019. Deep learning and process understanding for data-driven Earth system science.
1338 Nature 566, 195–204.
- 1339 Ridley, D.A., Heald, C.L., Prospero, J.M., 2014. What controls the recent changes in African
1340 mineral dust aerosol across the Atlantic? Atmos. Chem. Phys. 14, 5735–5747.
1341 <https://doi.org/10.5194/acp-14-5735-2014>
- 1342 Rohrmann, A., Heermance, R., Kapp, P., Cai, F., 2013. Wind as the primary driver of erosion
1343 in the Qaidam Basin, China. Earth Planet. Sci. Lett. 374, 1–10.
- 1344 Soci, C., Hersbach, H., Simmons, A., Poli, P., Bell, B., Berrisford, P., Horányi, A., Muñoz-
1345 Sabater, J., Nicolas, J., Radu, R., 2024. The ERA5 global reanalysis from 1940 to 2022.
1346 Q. J. R. Meteorol. Soc. 150, 4014–4048.
- 1347 Van De Kerchove, R., Zanaga, D., Keersmaecker, W., Souverijns, N., Wevers, J., Brockmann,
1348 C., Grosu, A., Paccini, A., Cartus, O., Santoro, M., 2021. ESA WorldCover: Global land
1349 cover mapping at 10 m resolution for 2020 based on Sentinel-1 and 2 data., in: AGU Fall
1350 Meeting Abstracts. pp. GC451-0915.
- 1351 Wang, W., Evan, A.T., Flamant, C., Lavaysse, C., 2015. On the decadal scale correlation
1352 between African dust and Sahel rainfall: The role of Saharan heat low–forced winds. Sci.
1353 Adv. 1, e1500646.
- 1354 Wen, J., Tagliabue, G., Rossini, M., Fava, F. Pietro, Panigada, C., Merbold, L., Leitner, S., Sun,
1355 Y., 2025. Detection of fast-changing intra-seasonal vegetation dynamics of drylands using
1356 solar-induced chlorophyll fluorescence (SIF). Biogeosciences 22, 2049–2067.
- 1357 Wen, Q., Li, Y., Du, M., Song, W., Wei, L., Wang, Z., Li, X., 2025. Interdecadal shift in the
1358 impact of winter land–sea thermal contrasts on following spring transcontinental dust
1359 transport pathways in North Africa. Atmos. Chem. Phys. 25, 10853–10867.
- 1360 Wu, C., Lin, Z., He, J., Zhang, M., Liu, X., Zhang, R., Brown, H., 2016. A process-oriented
1361 evaluation of dust emission parameterizations in CESM: Simulation of a typical severe
1362 dust storm in East Asia. J. Adv. Model. Earth Syst. 8, 1432–1452.
- 1363 Zhang, B., Xiong, D., Tang, Y., Liu, L., 2022. Land surface roughness impacted by typical
1364 vegetation restoration projects on aeolian sandy lands in the Yarlung Zangbo River valley,
1365 southern Tibetan plateau. Int. Soil Water Conserv. Res. 10, 109–118.
- 1366 Zhang, L., Montuoro, R., McKeen, S.A., Baker, B., Bhattacharjee, P.S., Grell, G.A., Henderson,
1367 J., Pan, L., Frost, G.J., McQueen, J., 2022. Development and evaluation of the aerosol
1368 forecast member in the National Center for Environment Prediction (NCEP)’s global
1369 ensemble forecast system (GEFS-Aerosols v1). Geosci. Model Dev. 15, 5337–5369.