

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20

*Geophysical Research Letters*

Supporting Information for

**Improving low-cloud fraction prediction through machine learning**

Haipeng Zhang<sup>1,2</sup>, Youtong Zheng<sup>3,4</sup>, and Zhanqing Li<sup>1,2</sup>

<sup>1</sup>Department of Atmospheric and Oceanic Science, University of Maryland, College Park, MD, USA

<sup>2</sup>Earth System Science Interdisciplinary Center, University of Maryland, College Park, MD, USA

<sup>3</sup>Department of Atmospheric and Earth Science, University of Houston, Houston, TX, USA

<sup>4</sup>Institute of Climate and Atmospheric Science, University of Houston, Houston, TX, USA

**Contents of this file**

- Texts S1 to S2
- Tables S1 to S4
- Figures S1 to S4

21 **Texts**

22

23 **Text S1. Building XGB10 and XGB7**

24 To determine the best combination of the hyperparameters in XGB10 and XGB7, the  
25 optimization space of six key parameters (learning\_rate, n\_estimators, subsample,  
26 colsample\_bytree, and colsample\_bylevel) was explored using a Bayesian optimization technique  
27 (Snoek et al., 2012). This technique is highly efficient for parameter tuning, allowing finding the  
28 maximum value of a target function in as few iterations as possible based on Bayesian inference  
29 and the Gaussian process. To avoid overfitting on the training set, we used 5-fold cross-validation  
30 in each iteration during the optimization processes. The mean squared error (MSE) for the  
31 validation set reached the minimum, usually within 10 iterations. The investigated ranges of each  
32 parameter for hyperparameter optimization in XGB10 and XGB7 were summarized in Table S2.  
33 To address the uncertainties caused by random seeds in optimization processes, we performed  
34 Bayesian optimization ten times. This resulted in ten optimized parameter sets, creating ten  
35 ensemble members each for XGB10 and XGB7, as summarized in Tables S3 and S4, respectively.  
36 The ensemble-mean prediction results and absolute SHAP values from each model were used for  
37 evaluation in this study.

38

39

40

41

42

43

44

45

46

47 **Text S2. SHAP Explainability analysis**

48 SHAP (SHapley Additive exPlanations; Lundberg et al., 2018; Lundberg & Lee, 2017) is a  
49 high-fidelity and unified approach to exploring and interpreting the tree-based ML model (e.g.,  
50 XGBoost) behavior. It explains a model’s individual output as a sum of the contributions of each  
51 feature (or predictor) and the mean predicted value through an explanation model, which can be  
52 expressed as:

53 
$$y = \bar{y} + \sum_i \phi_i, \quad (1)$$

54 where  $y$  is the final prediction for one case,  $\bar{y}$  is the average prediction across all cases, and  $\phi_i$  is  
55 the contribution of the  $i$ -th feature to the prediction for this case (called SHAP values). Based on  
56 cooperative game theory, Lundberg & Lee, (2017) first proposed the KernelSHAP algorithm to  
57 calculate SHAP values by sampling the predictions of a machine learning model by replacing  
58 feature values with random values from the feature distribution. But KernelSHAP is  
59 computationally slow and ignores feature dependence. A more efficient and exact algorithm for  
60 tree ensemble models (TreeSHAP) was developed using the conditional expectation to estimate  
61 feature effects with no feature independence assumption required (Lundberg et al., 2018, 2020).  
62 The features with larger absolute SHAP values contribute more to a prediction than those with  
63 smaller values.

64  
65  
66  
67

68 **Tables:**  
 69

70 **Table S1.** Summary of meteorological factors used as predictors in XGB10 and XGB7

Model	Predictor	Description
XGB10	$\theta_{1000}, \theta_{850}, \theta_{700}$	Potential temperature at 1000, 850, and 700 hPa (K)
	$RH_{1000}, RH_{850}, RH_{700}$	Relative humidity at 1000, 850, and 700 hPa (%)
	$U_{1000}$	Horizontal wind speed at 1000 hPa (m/s)
	$\omega_{700}$	Vertical velocity at 700 hPa (Pa/s)
	LHF	Latent heat flux ( $W/m^2$ )
	PWV	Column-integrated precipitable water vapor ( $kg/m^2$ )
XGB7	$RH_{700}, U_{1000}, \omega_{700}, LHF$	Same as those predictors used in XGB10
	LTS	Lower-tropospheric stability ( $\theta_{700} - \theta_{1000}$ ) (K)
	$\Delta q$	The moisture contrast between the boundary layer and free troposphere ( $q_{1000} - q_{700}$ ) (g/kg)
	$T_{adv}$	Horizontal temperature advection (K/day)

71  
 72  
 73

74

75

**Table S2.** Investigated range of hyperparameters for XGB10 and XGB7

Parameter	Investigated range
max_depth	3-7
n_estimators	100-1000
learning_rate	0.01-1.0
subsample	0.5-1.0
colsample_bytree	0.5-1.0
cosample_bylevel	0.5-1.0

76

77

78

79

**Table S3.** Optimized hyperparameters and training/test errors for ten XGB10 members

Member ID	colsample_bylevel	colsample_bytree	learning_rate	max_depth	n_estimators	subsample	MSE for training set	MSE for test set
<b>01</b>	0.6800	0.7890	0.6129	6	857	0.9990	0.0701	0.0739
<b>02</b>	0.7347	0.5853	0.5122	6	914	0.7131	0.0702	0.0740
<b>03</b>	0.8588	0.9309	0.6960	6	598	0.9372	0.0703	0.0741
<b>04</b>	0.7295	0.9263	0.3387	6	844	0.7752	0.0699	0.0737
<b>05</b>	0.9966	0.9833	0.2730	6	861	0.9342	0.0699	0.0737
<b>06</b>	0.7742	0.8292	0.4596	6	737	0.7562	0.0700	0.0738
<b>07</b>	0.8666	0.7718	0.2392	6	606	0.6046	0.0704	0.0741
<b>08</b>	0.8208	0.7708	0.3453	6	922	0.7314	0.0699	0.0737
<b>09</b>	0.7125	0.9662	0.3404	6	904	0.9243	0.0699	0.0737
<b>10</b>	0.5488	0.7314	0.5591	6	621	0.8558	0.0704	0.0742

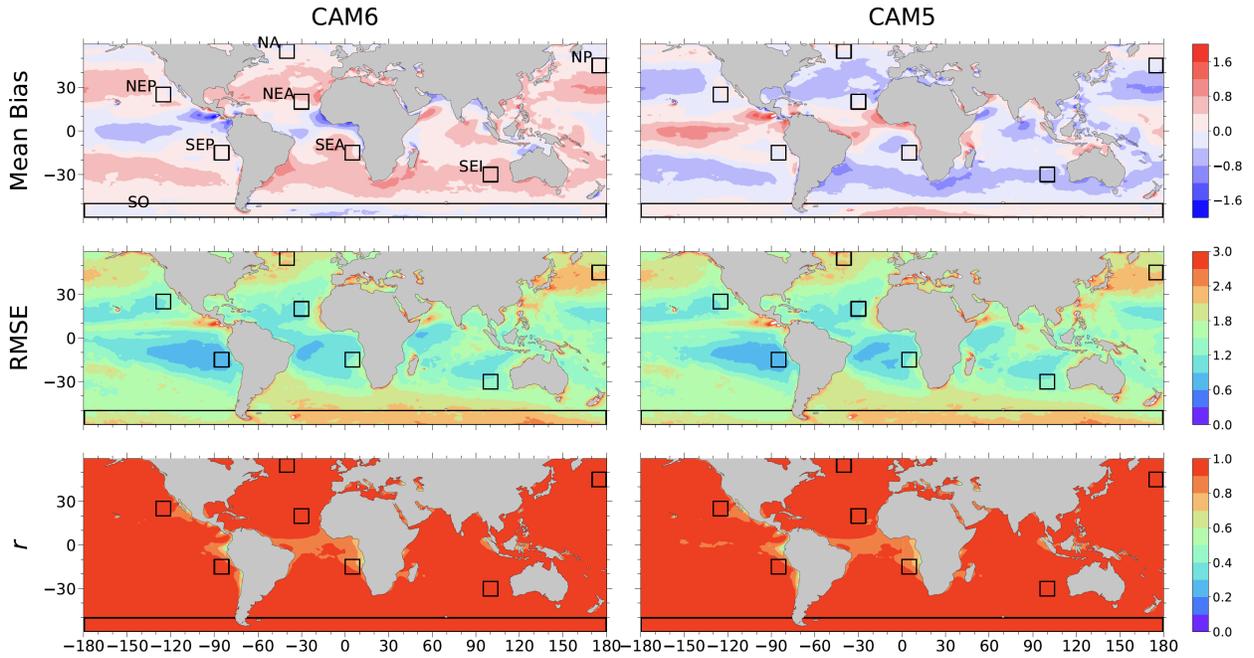
83  
84  
85

**Table S4.** Optimized hyperparameters and training/test errors for ten XGB7 members

Member ID	colsample_bylevel	colsample_bytree	learning_rate	max_depth	n_estimators	subsample	MSE for training set	MSE for test set
<b>01</b>	0.8515	0.7095	0.4283	6	475	0.7146	0.0816	0.0842
<b>02</b>	0.6609	0.9254	0.4318	6	837	0.6013	0.0813	0.0840
<b>03</b>	0.5205	0.9203	0.4727	6	999	0.8017	0.0813	0.0840
<b>04</b>	0.8626	0.7625	0.6250	6	637	0.9427	0.0814	0.0841
<b>05</b>	0.8614	0.9740	0.3785	6	781	0.5637	0.0813	0.0840
<b>06</b>	0.9062	0.5898	0.3028	6	619	0.9218	0.0815	0.0841
<b>07</b>	0.9456	0.7859	0.6579	5	871	0.8049	0.0815	0.0841
<b>08</b>	0.9559	0.5823	0.4667	6	994	0.6967	0.0814	0.0840
<b>09</b>	0.7334	0.9058	0.7428	6	843	0.8529	0.0813	0.0841
<b>10</b>	0.9634	0.7916	0.8672	4	937	0.7346	0.0818	0.0843

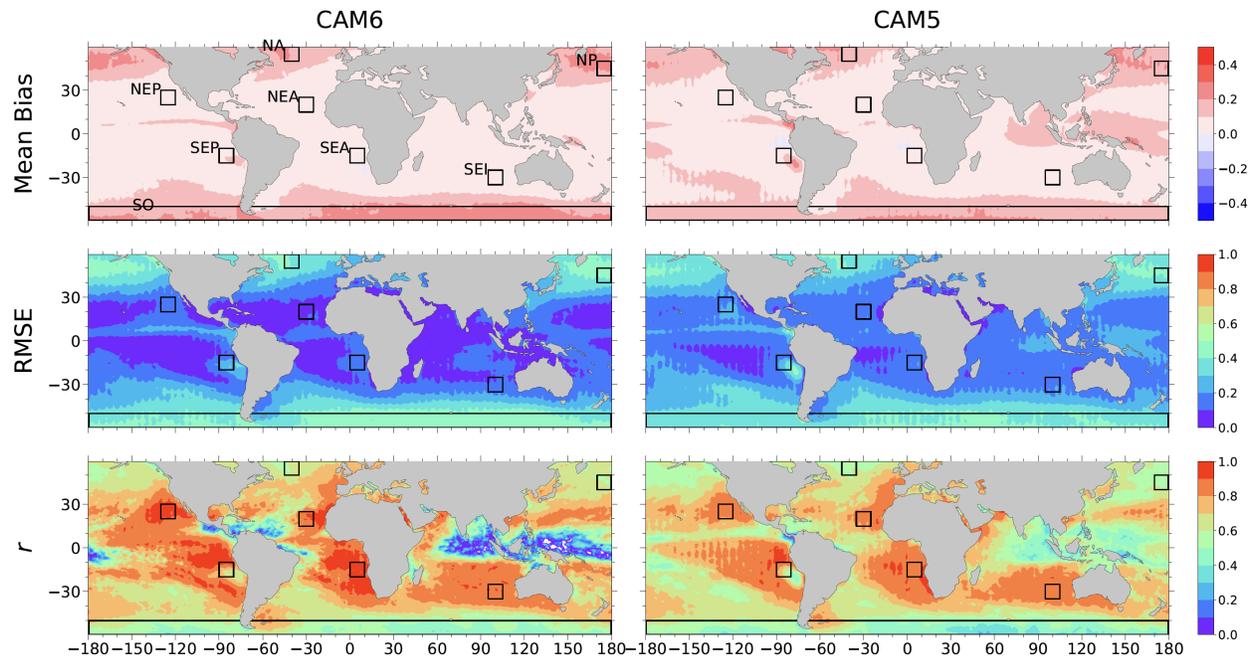
86  
87  
88  
89  
90

**Figures:**



91  
92  
93  
94  
95  
96  
97  
98

**Figure S1.** Comparison of near-surface zonal wind speeds between model nudging experiment outputs and ERA5 regarding three metrics: the mean bias, root mean squared error (RMSE), and correlation coefficients ( $r$ ) (from top to bottom rows, respectively). The first column shows the comparison for CAM6 outputs, with the second column for CAM5 outputs.



99

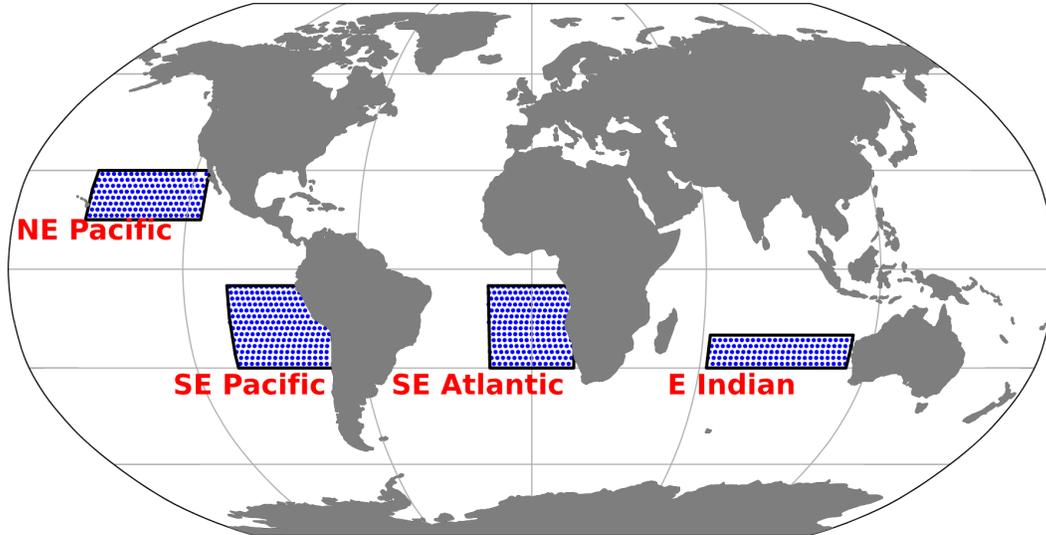
100 **Figure S2.** Same as Figure S1 but for comparisons between non-COSP low-cloud fraction (LCF)  
 101 and COSP-enabled LCF in CAM6 and CAM5.

102

103

104

105



106

107

108

109

110

**Figure S3.** Starting points of trajectories sampled over the four selected regions where the stratocumulus-to-cumulus transition dominates, following Eastman & Wood (2018): Northeast Pacific (-155 to -115°E, 15 to 30°N), Southeast Pacific (-105 to -70°E, -30 to -5°N), Southeast Atlantic (-15 to 15°E, -30 to -5°N), and East Indian (62.5 to 112.5°E, -30 to -20°N).

111

112

113

114

115

116

117

118

119

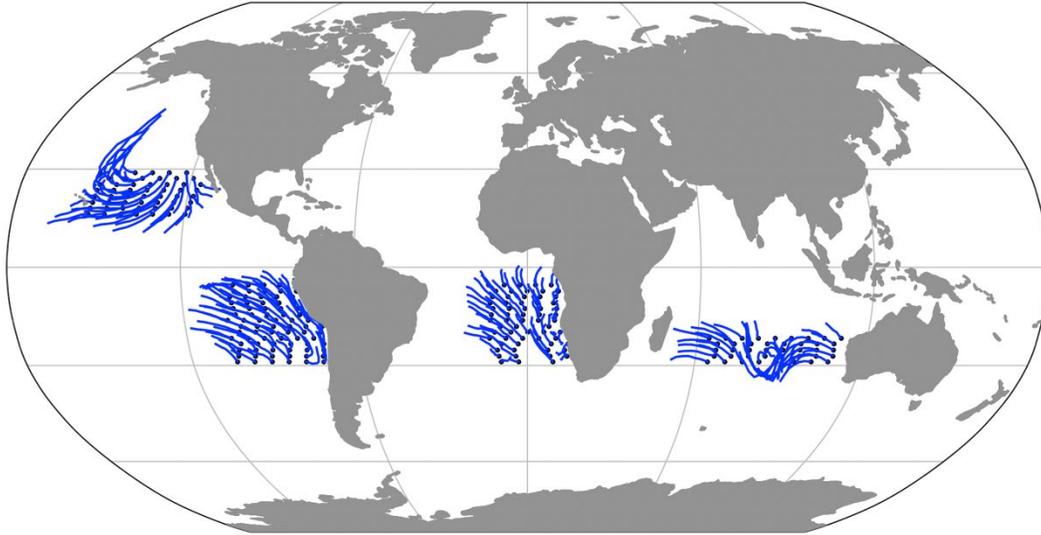
120

121

122

123

124



125

126 **Figure S4.** Subsets of forward trajectories (36 hours) over the four subtropical regions starting at  
127 6:00 p.m. on March 31, 2004. The black points denote the starting points.

128

129

130

131 **Reference**

132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151

Eastman, R., & Wood, R. (2018). The competing effects of stability and humidity on subtropical stratocumulus entrainment and cloud evolution from a Lagrangian perspective. *Journal of the Atmospheric Sciences*, 75(8), 2563–2578. <https://doi.org/10.1175/JAS-D-18-0030.1>

Lundberg, S. M., & Lee, S.-I. I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems, 2017-Decem*(Section 2), 4766–4775. Retrieved from <http://arxiv.org/abs/1705.07874>

Lundberg, S. M., Erion, G. G., & Lee, S.-I. (2018). Consistent individualized feature attribution for tree ensembles. Retrieved from <http://arxiv.org/abs/1802.03888>

Lundberg, S. M., Erion, G., Chen, H., DeGrave, A., Prutkin, J. M., Nair, B., et al. (2020). From local explanations to global understanding with explainable AI for trees. *Nature Machine Intelligence*, 2(1), 56–67. <https://doi.org/10.1038/s42256-019-0138-9>

Snoek, J., Larochelle, H., Adams, R. P., & Jeffery, C. (2012). Practical Bayesian optimization of machine learning algorithms. *Religion and the Arts*, 17(1–2), 57–73. <https://doi.org/10.48550/ARXIV.1206.2944>