

1 **Title**

2 *Using supervised learning to develop BaRAD, a 40-year monthly bias-*  
3 *adjusted global gridded radiation dataset – supplementary information*

4

5 **Authors**

6 TC Chakraborty<sup>1</sup>, Xuhui Lee<sup>1</sup>

7

8 **Affiliations**

9 <sup>1</sup>School of the Environment, Yale University, New Haven, CT 06520, USA

10 corresponding author: TC Chakraborty ([tc.chakraborty@yale.edu](mailto:tc.chakraborty@yale.edu))

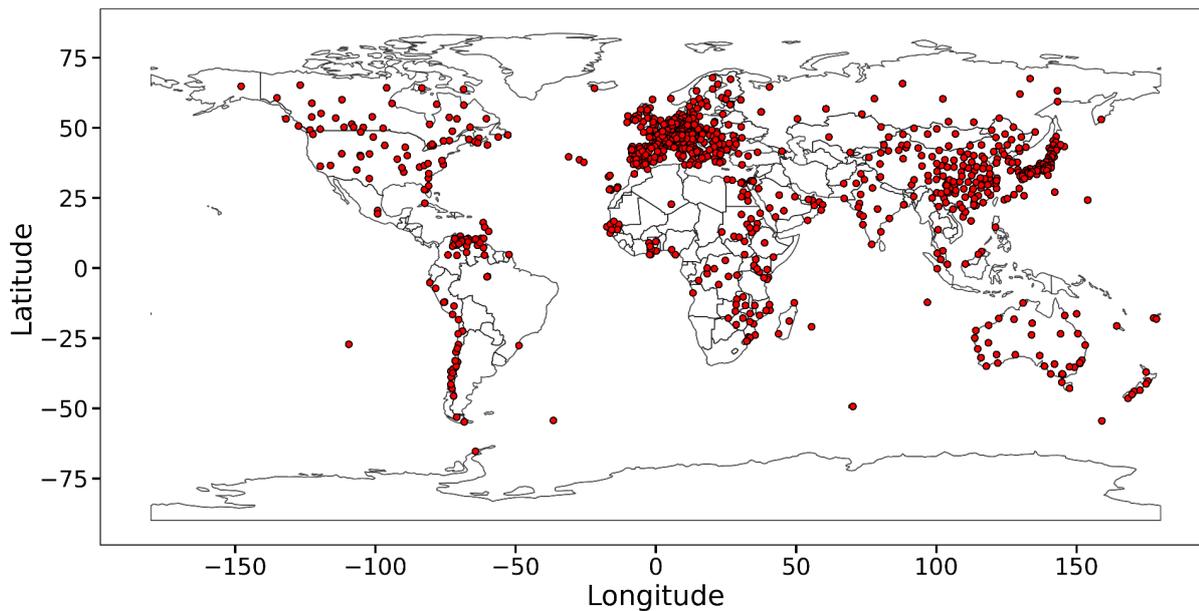
11

12 **Table of Contents:**

13 Figures S1 to S7 from pages 2 to 8

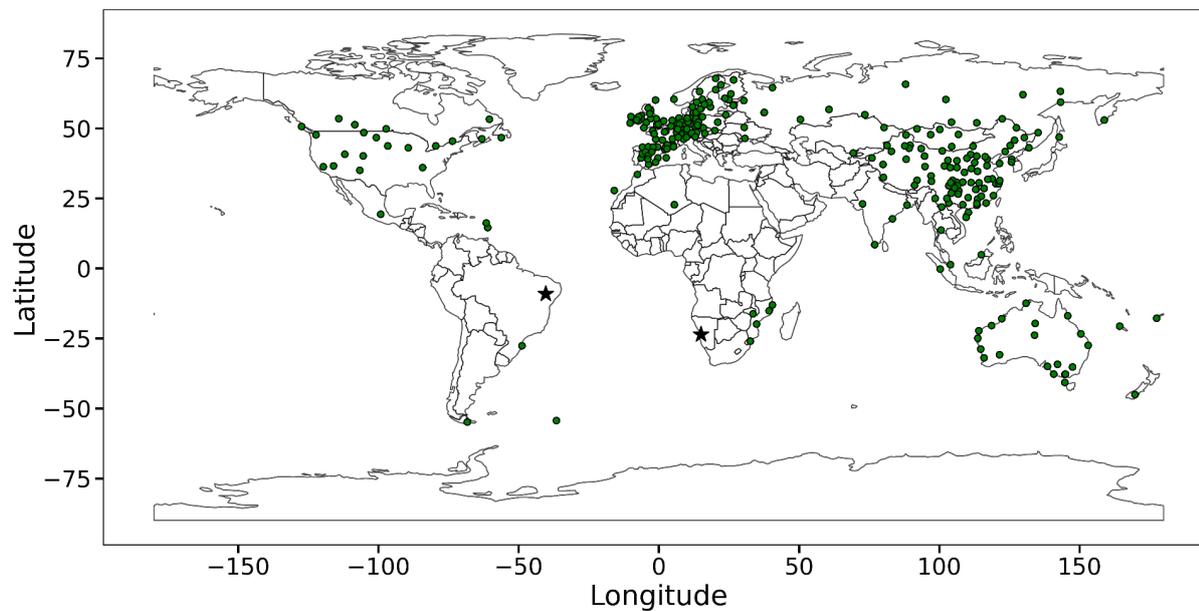
14 Tables S1 and S2 on pages 9 and 10

**a**



15

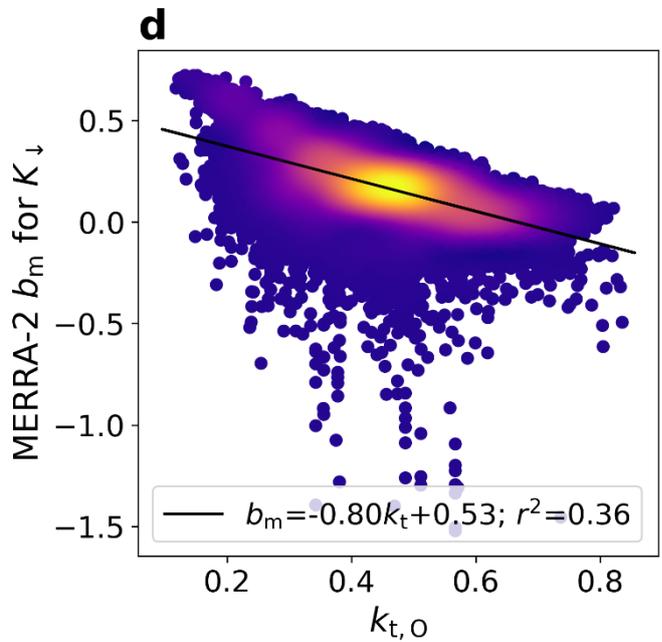
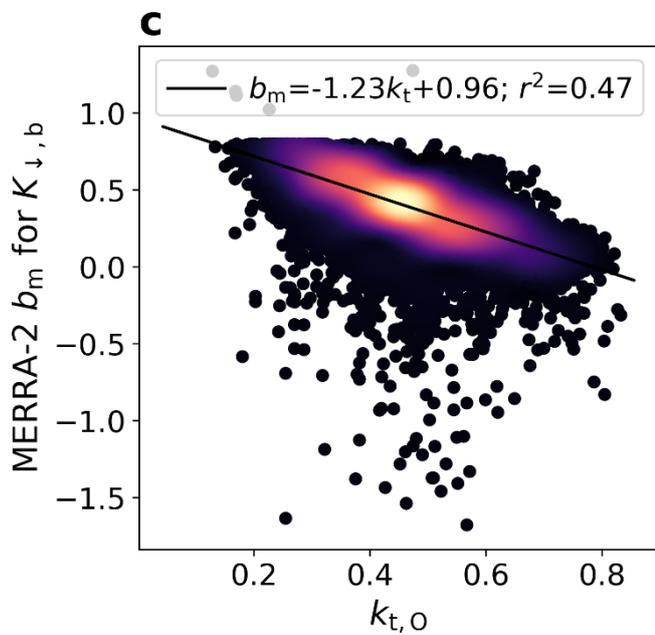
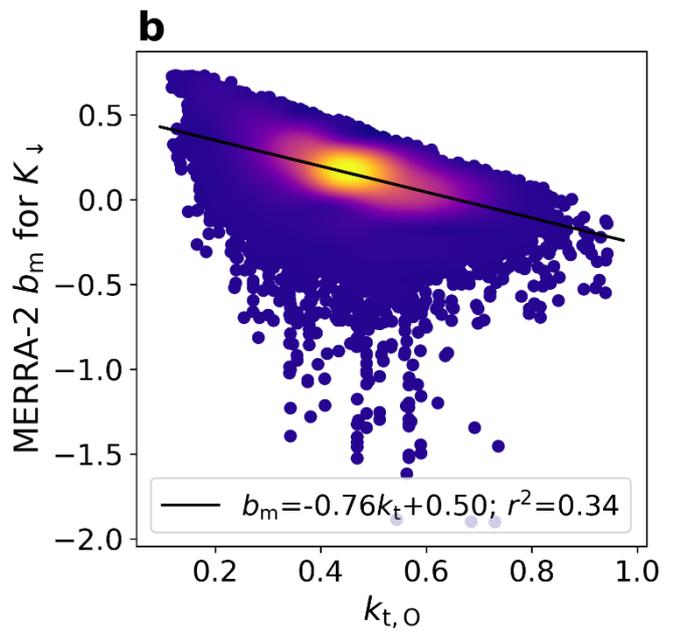
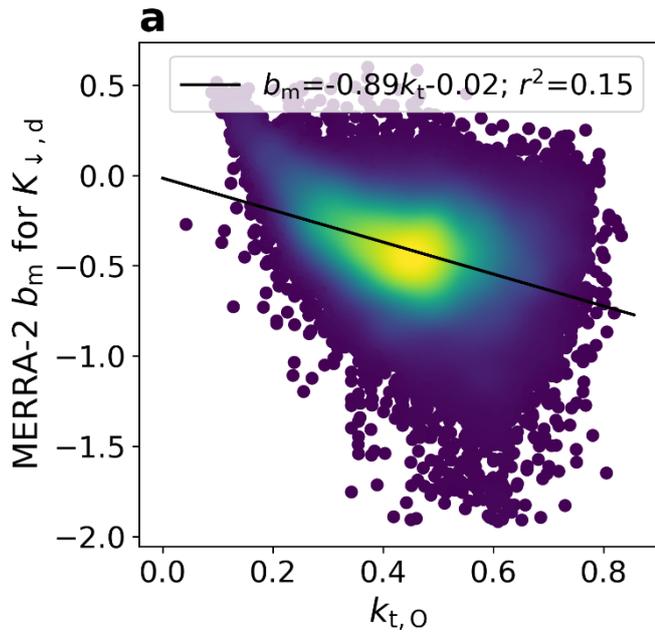
**b**



16

17 **Figure S1: Spatial distribution of ground observations.** Distribution of GEBA sites used for evaluating and training bias-  
18 correction algorithms in the present study for (a) shortwave radiation and (b) diffuse radiation. Sub-figure (b) also shows the  
19 locations (as black stars) of the two BSRN stations used to independently validate the BaRAD product.

20

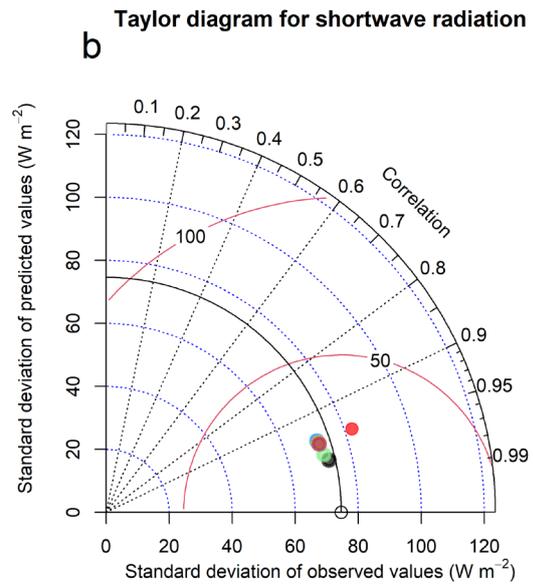
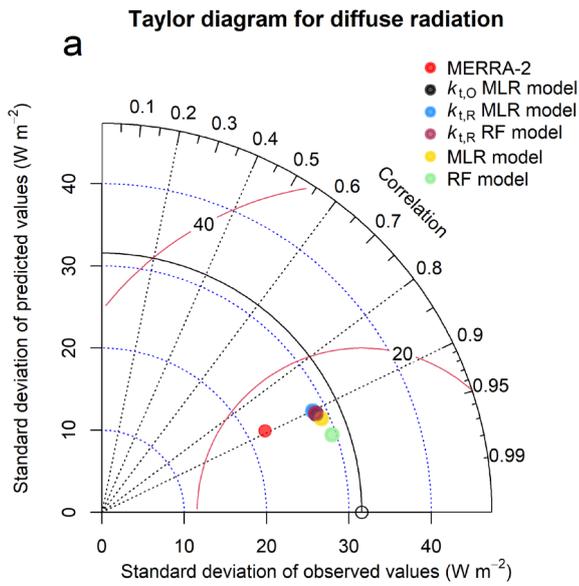


21

22

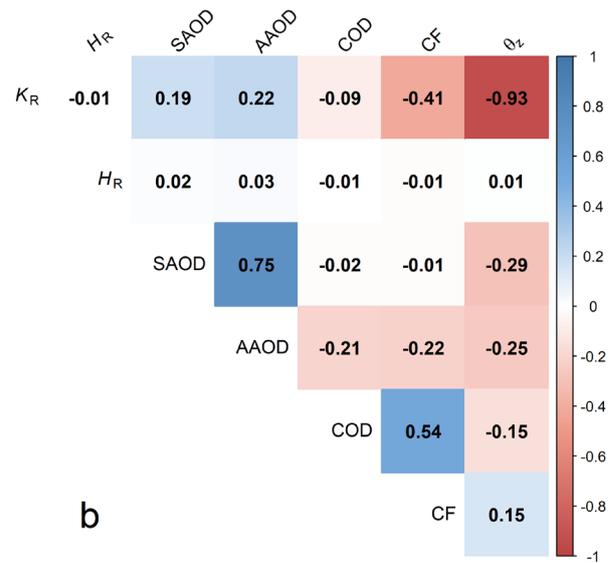
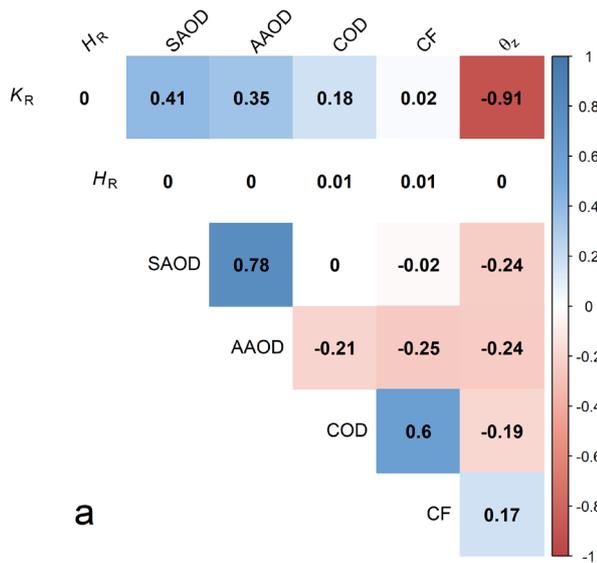
23 **Figure S2: Control of clearness index on biases in the MERRA-2 dataset.** (a) Bias in diffuse radiation ( $K_{\downarrow,d}$ ), (b) bias in  
 24 total shortwave radiation ( $K_{\downarrow}$ ), (c) bias in direct beam radiation ( $K_{\downarrow,b}$ ), and (d) bias total shortwave radiation ( $K_{\downarrow}$ ) for the sites  
 25 that also have  $K_{\downarrow,b}$  measurements. Statistical summaries of the associations are noted. Color indicates data density.

26



27  
 28 **Figure S3: Taylor diagrams of bias-correction models.** The Taylor diagrams represent the observed radiation values and  
 29 predicted values from MERRA-2, the  $k_{t,O}$  models, the  $k_{t,R}$  models, the MLR models, and the RF models for the consolidated  
 30 validation data for (a)  $K_{\downarrow,d}$  and (b)  $K_{\downarrow}$ .

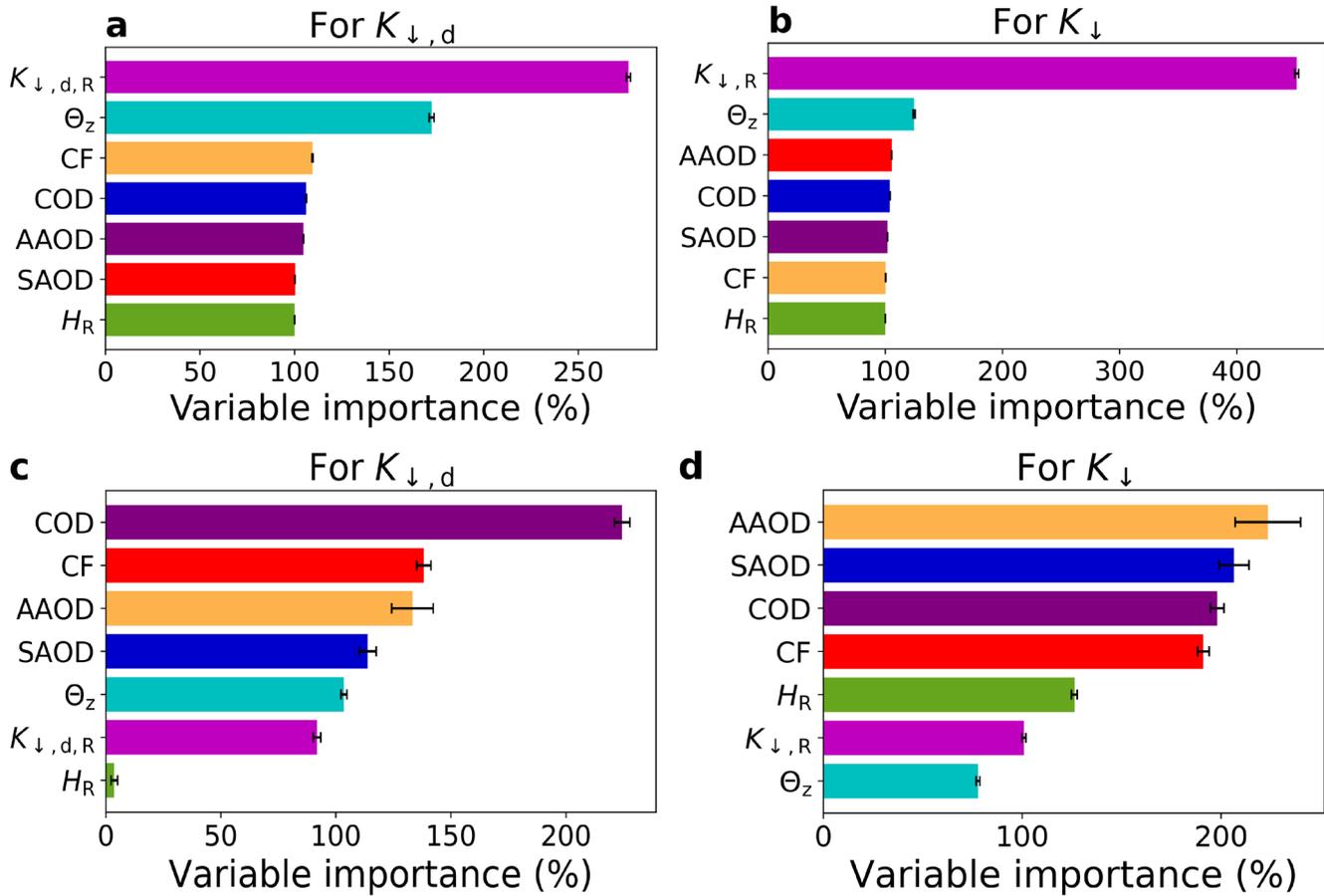
31  
 32  
 33



34

35 **Figure S4: Correlation matrices of features.** The correlation matrices of the features selected for training the supervised  
 36 machine learning models for (a)  $K_{l,d}$  and (b)  $K_l$ .

37



38

39

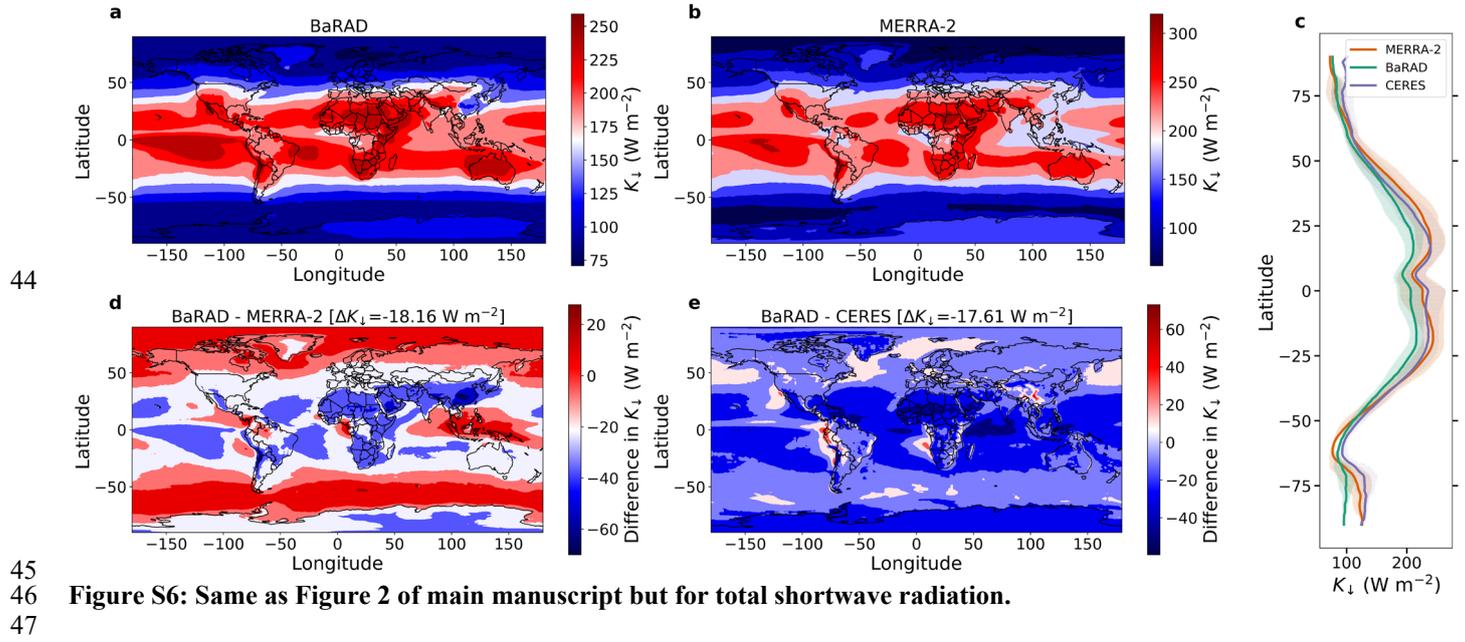
40

41

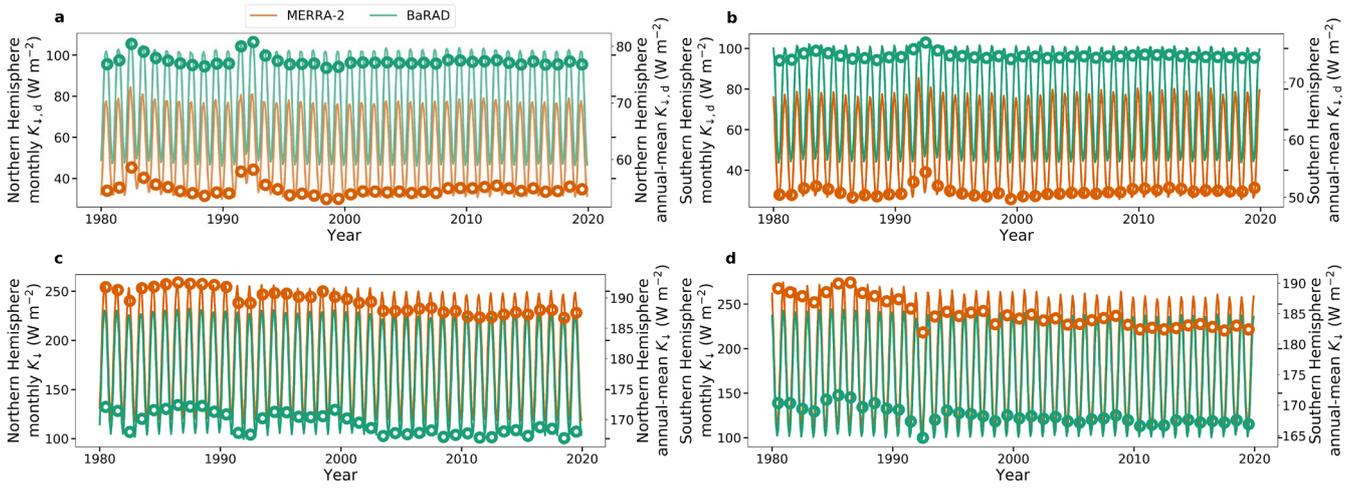
42

43

**Figure S5: Rank of variable importance for bias-correction.** (a) and (b): Permutation importance scores for the MLR model; (c) and (d): permutation importance scores for the RF model. The error bars show the standard deviation across the ten folds.



48



49

50 **Figure S7: Long-term trends at regional scale.** Sub-figures (a), (b), (c), and (d) show the long-term regional trends in  $K_{j,d}$   
 51 and  $K_j$  for northern and southern hemispheres, respectively. The monthly values are plotted on the left y-axes as lines and  
 52 the annual averages (plotted as circles) are on the right y-axes.  
 53

54 **Table S1: Summary of features.** Summary of the features, including their symbols and data source, used in the MLR and  
 55 RF bias-correction algorithms.

Feature Name	Feature Symbol	Description	Data Source
Incoming radiation at surface	$K_{\downarrow, R}$	Monthly grid-averaged value of incoming radiation at the surface. Can be either the total shortwave radiation reaching the surface ( $K_{\downarrow}$ ), or its diffuse component ( $K_{\downarrow, d}$ ), which is the portion after the light is scattered.	MERRA-2 reanalysis
Scattering Aerosol Optical Depth	SAOD	Monthly grid-averaged optical depth of scattering aerosols in the atmospheric column.	MERRA-2 reanalysis
Absorbing Aerosol Optical Depth	AAOD	Monthly grid-averaged optical depth of absorbing aerosols in the atmospheric column.	MERRA-2 reanalysis
Cloud Optical Depth	COD	Monthly grid-averaged optical depth of all clouds in the atmospheric column.	MERRA-2 reanalysis
Cloud Fraction	CF	Monthly grid-averaged cloud fraction.	MERRA-2 reanalysis
Zenith Angle	$\theta_z$	Monthly grid-averaged zenith angle, for the angle between the sun and the vertical direction.	Calculated
Altitude	$H_R$	Average altitude of the grid.	MERRA-2 reanalysis

56  
57

58 **Table S2: Summary of data products.** List of data products included in the present study, along with their temporal and  
 59 spatial resolution, and a few advantages and disadvantages  
 60

Dataset	Reference	Spatial Resolution	Finest temporal Resolution	Years of data availability	Advantages	Disadvantages
MERRA-2	24	0.5° x 0.625°	Hourly	1980 – Present	Physical model; Constrained by assimilated observations Simplified radiative transfer model;	Model parameterizations; Large biases in surface radiation Model parameterizations;
CERES	35	1° x 1°	Hourly	2001 – Present	Constrained by satellite observations	Large biases in surface radiation; Note available before 2001
DSCOVER/EPIC	21	0.1° x 0.1°	Hourly	June 2015 – June 2019	Data-driven model; Constrained by <i>in situ</i> and satellite observations	Limited period of availability; Not continuous at hourly scale; Data-driven
GEBA	28	Point	Monthly	Various; site-specific	Observations	Uneven geographic and temporal distribution; Sensor errors
BSRN	37	Point	Every Minute	Various; site-specific	Observations	Uneven geographic and temporal distribution; Sensor errors
BaRAD	Present study	0.5° x 0.625°	Monthly	1980 – 2019	Data-driven model; Constrained by <i>in situ</i> observations and MERRA-2 fields	Monthly scale; Data-driven; Influenced by sampling bias in training data