

# *Exploratory precipitation metrics: spatiotemporal characteristics, process- oriented, and phenomena-based*

Article

Published Version

Creative Commons: Attribution 4.0 (CC-BY)

Open access

Leung, L. Ruby ORCID logoORCID: <https://orcid.org/0000-0002-3221-9467>, Boos, William R., Catto, Jennifer L., A. DeMott, Charlotte, Martin, Gill M., Neelin, J. David, O'Brien, Travis A., Xie, Shaocheng, Feng, Zhe, Klingaman, Nicholas P. ORCID logoORCID: <https://orcid.org/0000-0002-2927-9303>, Kuo, Yi-Hung, Lee, Robert W. ORCID logoORCID: <https://orcid.org/0000-0002-1946-5559>, Martinez-Villalobos, Cristian, Vishnu, S., Priestley, Matthew D. K., Tao, Cheng and Zhou, Yang (2022) Exploratory precipitation metrics: spatiotemporal characteristics, process-oriented, and phenomena-based. *Journal of Climate*, 35 (12). pp. 3659-3686. ISSN 0894-8755 doi: <https://doi.org/10.1175/JCLI-D-21-0590.1> Available at <https://centaur.reading.ac.uk/107157/>

It is advisable to refer to the publisher's version if you intend to cite from the work. See [Guidance on citing](#).

Published version at: <http://dx.doi.org/10.1175/JCLI-D-21-0590.1>

To link to this article DOI: <http://dx.doi.org/10.1175/JCLI-D-21-0590.1>

Publisher: American Meteorological Society

All outputs in CentAUR are protected by Intellectual Property Rights law, including copyright law. Copyright and IPR is retained by the creators or other copyright holders. Terms and conditions for use of this material are defined in the [End User Agreement](#).

[www.reading.ac.uk/centaur](http://www.reading.ac.uk/centaur)

## **CentAUR**

Central Archive at the University of Reading

Reading's research outputs online

## Exploratory Precipitation Metrics: Spatiotemporal Characteristics, Process-Oriented, and Phenomena-Based

L. RUBY LEUNG,<sup>a</sup> WILLIAM R. BOOS,<sup>b</sup> JENNIFER L. CATTO,<sup>c</sup> CHARLOTTE A. DEMOTT,<sup>d</sup> GILL M. MARTIN,<sup>e</sup> J. DAVID NEELIN,<sup>f</sup> TRAVIS A. O'BRIEN,<sup>g,h</sup> SHAOCHENG XIE,<sup>i</sup> ZHE FENG,<sup>a</sup> NICHOLAS P. KLINGAMAN,<sup>j,k</sup> YI-HUNG KUO,<sup>f</sup> ROBERT W. LEE,<sup>j,k</sup> CRISTIAN MARTINEZ-VILLALOBOS,<sup>f,l</sup> S. VISHNU,<sup>b</sup> MATTHEW D. K. PRIESTLEY,<sup>c</sup> CHENG TAO,<sup>i</sup> AND YANG ZHOU<sup>h</sup>

<sup>a</sup> Pacific Northwest National Laboratory, Richland, Washington

<sup>b</sup> University of California, Berkeley, Berkeley, California

<sup>c</sup> University of Exeter, Exeter, United Kingdom

<sup>d</sup> Colorado State University, Fort Collins, Colorado

<sup>e</sup> Met Office, Exeter, United Kingdom

<sup>f</sup> University of California, Los Angeles, Los Angeles, California

<sup>g</sup> Indiana University, Bloomington, Indiana

<sup>h</sup> Lawrence Berkeley National Laboratory, Berkeley, California

<sup>i</sup> Lawrence Livermore National Laboratory, Livermore, California

<sup>j</sup> National Centre for Atmospheric Science, Reading, United Kingdom

<sup>k</sup> University of Reading, Reading, United Kingdom

<sup>l</sup> Universidad Adolfo Ibáñez, Peñalolén, Santiago, Chile

(Manuscript received 4 August 2021, in final form 30 December 2021)

**ABSTRACT:** Precipitation sustains life and supports human activities, making its prediction one of the most societally relevant challenges in weather and climate modeling. Limitations in modeling precipitation underscore the need for diagnostics and metrics to evaluate precipitation in simulations and predictions. While routine use of basic metrics is important for documenting model skill, more sophisticated diagnostics and metrics aimed at connecting model biases to their sources and revealing precipitation characteristics relevant to how model precipitation is used are critical for improving models and their uses. This paper illustrates examples of exploratory diagnostics and metrics including 1) spatiotemporal characteristics metrics such as diurnal variability, probability of extremes, duration of dry spells, spectral characteristics, and spatiotemporal coherence of precipitation; 2) process-oriented metrics based on the rainfall–moisture coupling and temperature–water vapor environments of precipitation; and 3) phenomena-based metrics focusing on precipitation associated with weather phenomena including low pressure systems, mesoscale convective systems, frontal systems, and atmospheric rivers. Together, these diagnostics and metrics delineate the multifaceted and multiscale nature of precipitation, its relations with the environments, and its generation mechanisms. The metrics are applied to historical simulations from phases 5 and 6 of the Coupled Model Intercomparison Project. Models exhibit diverse skill as measured by the suite of metrics, with very few models consistently ranked as top or bottom performers compared to other models in multiple metrics. Analysis of model skill across metrics and models suggests possible relationships among subsets of metrics, motivating the need for more systematic analysis to understand model biases for informing model development.

**KEYWORDS:** Precipitation; Climate models; Diagnostics; Model evaluation/performance

### 1. Introduction

Precipitation is a key component of the water cycle connecting processes across the atmosphere, land, ocean, and cryosphere (Trenberth et al. 2007). Through decades of development, the current generation of climate models uses increasingly sophisticated, physically based subgrid parameterizations of convection

and cloud microphysics to simulate precipitation, although their horizontal resolutions are still typically much coarser than needed to explicitly resolve precipitation formation processes. When, where, how often, and how much precipitation falls has significant implications for the energy, water, and biogeochemical cycles of the Earth system. For example, biases in soil moisture can often be linked to biases in precipitation amount, frequency, and intensity, which influence the partitioning of precipitation into evapotranspiration, runoff, and soil moisture storage, with subsequent impact on surface temperature through evaporative cooling (Qian et al. 2006). Relatedly, biases in modeling the surface water and energy balance due to precipitation biases can

Denotes content that is immediately available upon publication as open access.

Supplemental information related to this paper is available at the Journals Online website: <https://doi.org/10.1175/JCLI-D-21-0590.s1>.

Corresponding author: L. Ruby Leung, [ruby.leung@pnnl.gov](mailto:ruby.leung@pnnl.gov)



This article is licensed under a Creative Commons Attribution 4.0 license (<http://creativecommons.org/licenses/by/4.0/>).

DOI: 10.1175/JCLI-D-21-0590.1

© 2022 American Meteorological Society. For information regarding reuse of this content and general copyright information, consult the AMS Copyright Policy ([www.ametsoc.org/PUBSReuseLicenses](http://www.ametsoc.org/PUBSReuseLicenses)).

influence clouds, convection, and precipitation through energetic constraints and land–atmosphere feedbacks. Because of the myriad Earth system interactions and feedbacks mediated by precipitation, skillful modeling of precipitation and understanding and attribution of precipitation biases are scientifically challenging (Dai 2006; Covey et al. 2016; Chen et al. 2021). As precipitation biases are among the most consequential in limiting the use of climate models for decision support, there is an urgent need to improve precipitation modeling across a wide range of spatial and temporal scales (Tapiador et al. 2019).

Quantifying and understanding model precipitation biases is an important step toward improving the overall quality of climate simulations and predictions. Metrics are objective measures for benchmarking model performance against observations and facilitating model intercomparison. Common metrics of precipitation have focused on aspects such as the spatial distribution of annual and seasonal mean precipitation, daily precipitation amount, frequency, and intensity, and the probability density function of precipitation rate (Deser et al. 2012; Chen and Dai 2018, 2019). Increasingly, metrics related to extremes such as annual maximum daily precipitation and consecutive dry days have also been used to evaluate precipitation characteristics connected more closely to societal impacts. These metrics have revealed multiple longstanding precipitation biases in climate models. For example, climate models tend to produce too frequent light daily precipitation, but not enough high-intensity daily precipitation compared to observations (Dai 2006; Stephens et al. 2010; Chen et al. 2021), while subdaily intensities can vary considerably between models (e.g., Klingaman et al. 2017). Most global climate models simulate a spurious intertropical convergence zone (ITCZ) in the southeastern Pacific and South Atlantic, resulting in a double-ITCZ bias that is most prominent during boreal winter (Mechoso et al. 1995; Lin 2007; Mapes and Neale 2011; Hwang and Frierson 2013; Oueslati and Bellon 2013; Hirota et al. 2014; Tian 2015; Tian and Dong 2020). Erroneous diurnal timing of precipitation over land is another common bias, which is most noticeable during boreal summer in regions such as the central United States featuring nocturnal peaks in precipitation (Dai et al. 1999; Tang et al. 2021). Precipitation biases have also been identified in regions with complex terrain such as the western United States (Mejia et al. 2018) and Europe (Mehran et al. 2014), in Amazonia (Yin et al. 2013), and in monsoon regions such as Asia (Sperber et al. 2013).

Although precipitation diagnostics and metrics have been incorporated in model evaluation and diagnostic packages such as ESMValTool (Eyring et al. 2020) and the PCMDI Metrics Package (PMP; Gleckler et al. 2016) used by climate modeling centers and the climate science community, they focus on limited aspects of precipitation for benchmarking global climate simulations. At the same time, over the past few years new precipitation diagnostics and metrics have been developed to deconvolve and better understand model precipitation biases. For example, Ma et al. (2013) proposed a set of metrics and diagnostics to evaluate and diagnose tropical precipitation biases and associated moist processes in climate models. Their proposed diagnostics include stratiform fraction of precipitation, probability density function of daily precipitation intensity, composites of column water vapor, column relative humidity, temperature, and specific

humidity profiles as a function of precipitation intensity, and composites of stratiform rainfall fraction as a function of column relative humidity. Klingaman et al. (2017) developed a set of diagnostics and metrics for analyzing precipitation intensity and coherence on a range of time and space scales.

This study represents a collaborative effort as an outgrowth of a workshop on “Benchmarking Simulated Precipitation in Earth System Models” (Pendergrass et al. 2020) to develop more advanced precipitation metrics and demonstrate their use in benchmarking diverse aspects of precipitation from climate simulations. Three types of precipitation diagnostics and metrics are presented: 1) spatiotemporal characteristics metrics, such as diurnal variability, probability of extremes, duration of dry spells, spectral characteristics, and spatiotemporal coherence of precipitation; 2) process-oriented metrics, based on the rainfall–moisture coupling and temperature–water vapor environments of precipitation; and 3) phenomena-based metrics, focusing on precipitation associated with weather phenomena such as low pressure systems, mesoscale convective systems, frontal systems, and atmospheric rivers. These diagnostics and metrics take advantage of analysis building on advances in understanding the thermodynamic environments of precipitation (e.g., Bretherton et al. 2004; Neelin et al. 2009; Kuo et al. 2018; Chen et al. 2020) and their role in modes of variability (e.g., Wolding et al. 2020), and in tracking weather features such as atmospheric rivers (e.g., Shields et al. 2018).

While examples of the above metrics have been reported in recent literature (e.g., Klingaman et al. 2017; Ahmed et al. 2020; Feng et al. 2021a), they are deemed exploratory partly because they have not been widely used or implemented in standard metrics and diagnostics packages and partly because they allow deeper exploration of precipitation characteristics and associated processes. Some of these diagnostics and metrics require variables besides precipitation to evaluate relationships with environmental conditions, or to track weather features, so their data requirements go beyond the baseline precipitation metrics already implemented in widely used metrics and diagnostics packages (Pendergrass et al. 2020). Furthermore, additional research may be needed on interpretations of results from use of these metrics, to standardize their use, or to address technical or computational issues. Here, we apply the exploratory metrics to a common set of climate simulations from phases 5 (CMIP5; Taylor et al. 2012) and 6 (CMIP6; Eyring et al. 2016) of the Coupled Model Intercomparison Project. While our aim is not to provide an exhaustive study of the ability of these models to represent precipitation, we illustrate how such diagnostics and metrics may be used to evaluate broader aspects of precipitation in climate simulations and to explore insights that may be gained through comparative analysis of multiple metrics. With increasing model resolutions to better resolve weather and large-scale environments (e.g., Haarsma et al. 2016), the exploratory diagnostics and metrics may be even more relevant not only for benchmarking models but also for understanding the causes of model precipitation biases. They also provide useful information to support the growing and more diverse uses of precipitation from climate models and improve communications of climate model performance by connecting precipitation to commonly understood weather phenomena. A collection of

TABLE 1. Observational and reanalysis data for benchmarking models. Variables  $P$ ,  $Q$ ,  $U$ ,  $V$ ,  $T$ , CVW, and IR Tb are precipitation, specific humidity, zonal wind, meridional wind, temperature, column water vapor, and infrared brightness temperature, respectively.

	Variables	Temporal resolution	Max spatial resolution	Period of coverage	Domain of coverage
GPCP	$P$	Monthly	$0.25^\circ \times 0.25^\circ$	1979–2020	Global
GPCP 1DD	$P$	Daily	$1^\circ$	1996–present	Global
CMORPH	$P$	30 min	8 km	1998–2017	$60^\circ\text{S}$ – $60^\circ\text{N}$
PERSIANN-CDR	$P$	Monthly	$0.25^\circ \times 0.25^\circ$	1983–2017	$60^\circ\text{S}$ – $60^\circ\text{N}$
TRMM 3B42	$P$	3-hourly	$0.25^\circ \times 0.25^\circ$	1998–2019	$50^\circ\text{S}$ – $50^\circ\text{N}$
TRMM-TMI	CVW	Twice-daily snapshot	$0.25^\circ \times 0.25^\circ$	2002–14	$40^\circ\text{S}$ – $40^\circ\text{N}$
TRMM PR 2A25	$P$	Twice-daily snapshot	5 km	2002–14	$40^\circ\text{S}$ – $40^\circ\text{N}$
GPM-IMERG	$P$	Hourly	$0.1^\circ \times 0.1^\circ$	2001–20	$60^\circ\text{S}$ – $60^\circ\text{N}$
ARMBE	$P$	Hourly	Single point	SGP: 1993–2018 MAO: 2014/15	Single point
VARANAL	$P$	SGP: hourly MAO: 3-hourly	$0.5^\circ \times 0.5^\circ$	SGP: 2004–18 MAO: 2014/15	SGP: $3^\circ \times 3^\circ$ MAO: $2^\circ \times 2^\circ$
NASA Global Merged IR V1	IR $T_b$	Hourly	Raw data at 4 km but coarsened to $0.1^\circ \times 0.1^\circ$	2000–19	$60^\circ\text{S}$ – $60^\circ\text{N}$
ERA-Interim	$Q$ , $U$ , $V$ , $T$	3-hourly	80 km	1979–2019	Global
ERA5	$Q$ , $U$ , $V$ , $T$	Hourly	30 km	1979–2019	Global
MERRA-2	$Q$ , $U$ , $V$ , $T$	3-hourly	50 km	1980–2019	Global
CFSR	$Q$	6-hourly	38 km	1979–2019	Global

such exploratory diagnostics and metrics is a valuable addition to the existing precipitation diagnostics and metrics packages that are used in the community.

We briefly summarize the observational data, climate model output, and the feature tracking methods in section 2. Key results are presented in sections 3, 4, and 5 for the spatiotemporal characteristics, process-oriented metrics, and phenomena-based metrics, respectively. Each area is presented as a module describing the diagnostics and metrics and the results of applying them to climate model outputs summarized in a multipanel figure. We conclude with discussion and summary in section 6.

## 2. Data and feature tracking methods

### a. Observational data and climate model outputs

Several observational precipitation data products are used for benchmarking precipitation from climate simulations. These include 1) the Tropical Rainfall Measurement Mission (TRMM) Multisatellite Precipitation Analysis (TMPA-RT) (3B42; Huffman et al. 2007); 2) the Remote Sensing Systems TRMM Microwave Imager (TMI) Daily Environmental Suite on  $0.25^\circ$  grid, version 7.1 (Wentz et al. 2015); 3) the TRMM Precipitation Radar (PR) Rainfall Rate and Profile L2 1.5 h V7 (2A25; TRMM 2011); 4) the monthly and daily Global Precipitation Climatology Project (GPCP) V3 combined precipitation dataset (Huffman et al. 2020); 5) CMORPH bias-corrected integrated satellite precipitation estimates (Joyce and Xie 2011); 6) Precipitation Estimation from Remotely Sensed Information (PERSIANN) (Ashouri et al. 2015); and 7) Global Precipitation Measurement (GPM) Multi-satellitE Retrievals (IMERG) precipitation data V06B (Tan et al. 2019). They represent a diverse set of precipitation data derived from satellite- and ground-based remote sensing retrievals. In addition, ground-based precipitation observations at the DOE Atmospheric

Radiation Measurement (ARM) Program's Southern Great Plains (SGP) and Manacapuru (MAO) sites are also used. The ARM data used in this study are from the ARM best estimate (ARMBE; Xie et al. 2010) data products and the ARM long-term continuous variational analysis (VARANAL; Xie et al. 2004). At these ARM sites, the available surface rain gauge measurements and/or radar retrievals provide additional information to validate satellite-based precipitation products. Table 1 summarizes the spatial and temporal resolution and domain coverage of these datasets. While the highest spatial resolution available for the dataset is given in Table 1, coarse-graining of the data for comparison to models is described with each metric. As different exploratory diagnostics and metrics have different requirements for precipitation data, we do not standardize the use of observational precipitation data in calculating the metrics, but recognize the need to address uncertainty in observed precipitation products in use and interpretation of metrics.

Besides precipitation data, several global reanalysis products are used to provide gridded data of the atmospheric environments needed for calculation of some process-oriented metrics and identification and tracking of weather features for the phenomena-based metrics: 1) ERA-Interim (Dee et al. 2011), 2) ERA5 (Hersbach et al. 2020; Hoffmann et al. 2019), 3) MERRA-2 (Gelaro et al. 2017), and 4) CFSR (Saha et al. 2010). Last, the NASA Global Merged IR V1 infrared brightness temperature ( $T_b$ ) data (Janowiak et al. 2017) are also used to track mesoscale convective systems (MCSs). The spatial and temporal resolutions of the reanalysis products and  $T_b$  data are also summarized in Table 1.

The exploratory metrics are applied to benchmark precipitation from the Coupled Model Intercomparison Project phase 5 (CMIP5; Taylor et al. 2012) and phase 6 (CMIP6; Eyring et al. 2016), with typical horizontal resolution of  $\sim 1^\circ$ . Two of the metrics on low pressure systems and mesoscale convective systems

TABLE 2. Variables and their temporal frequency used to calculate various precipitation metrics and the objectives of the metrics.

Metrics	Variables and temporal frequency	Objectives
Diurnal cycle of precipitation	3-hourly precipitation	Intercompare a large number of models with observations and with each other on the diurnal cycle of precipitation over different climate regimes
Extremes of daily precipitation and duration of dry spells	Daily precipitation	Use characteristic scales governing probabilities in the large-event regime for dry and wet precipitation extremes to capture the performance of models
Spectral analysis of precipitation	3-hourly and daily precipitation sampled over ~20 years of data, by season and annual.	Examine the ability of models to represent the range of precipitation intensities typically occurring at any location, on 3-hourly and daily time scales
Coherence analysis of precipitation	3-hourly and daily precipitation	Measure and compare the spatial and temporal scales of precipitation across observations and models
MJO east–west power ratio and Maritime Continent propagation	Daily precipitation	Evaluate the relationship between precipitation spatial coherence and MJO propagation across the Maritime Continent in models
Rainfall–moisture coupling	Daily precipitation and vertically integrated water vapor and saturation water vapor	Evaluate the coupling of tropical rainfall and moisture in models and how this coupling affects MJO simulation
Temperature–water vapor environment	3-hourly/hourly vertically integrated saturation humidity and snapshots of column water vapor (CWV) and precipitation	Quantify the thermodynamic environment that produces most precipitation at subdaily time scales
Low pressure systems	6-hourly 850-hPa values of zonal wind, meridional wind, temperature, and specific humidity; 6-hourly precipitation	Track LPS in observations and simulations and compare their depiction of number, structure, and rainfall
Mesoscale convective systems	Hourly outgoing longwave radiation and precipitation	Track MCSs in observations and simulations and compare their depiction of MCS number and rainfall
Frontal precipitation	6-hourly 850-hPa zonal and meridional wind components, specific humidity, temperature, and daily precipitation	Use fronts as a precipitation regime to decompose precipitation errors into frontal and nonfrontal, and to quantify the representation of the dynamical impact of fronts on precipitation intensity
Atmospheric rivers	3- or 6-hourly zonal and meridional wind components, specific humidity, surface pressure, and precipitation	Assess whether models simulate AR-related precipitation in the correct locations and with enough contrast between regions with high AR precipitation and low AR precipitation

are applied to precipitation from several high-resolution simulations from HighResMIP (Haarsma et al. 2016) as these weather features are better defined and more reasonably resolved at higher resolution. In HighResMIP, high-resolution simulations have nominal resolutions ranging from 0.25° to 0.5°, with their low-resolution counterparts ranging from 1.0° to 1.4°. Table 2 summarizes the variables and their temporal frequency used to calculate the various metrics.

#### b. Feature identification and tracking methods

The phenomena-based metrics require identification and tracking of weather features in observations and simulated

precipitation. A brief description of methods used to track low pressure systems (LPS), mesoscale convective systems (MCS), frontal systems (FRT), and atmospheric rivers (AR) are provided below while more detailed descriptions are provided in the cited references.

##### 1) LOW PRESSURE SYSTEMS

The TempestExtremes feature tracking algorithm (Ullrich and Zarzycki 2017) is used to track tropical low pressure systems by identifying extrema in candidate tracking variables. A systematic assessment of multiple candidate variables, hundreds of quantitative tracking criteria, and several vertical levels led to selection of

the streamfunction of the 850-hPa horizontal wind (Vishnu et al. 2020) as the optimal tracking variable. Streamfunction minima were used to identify lower-tropospheric cyclonic vortices within 35° of the equator in the ERA5 reanalysis, four HighResMIP models, and the 0.25°-resolution E3SM model (Caldwell et al. 2019). The streamfunction was calculated from the horizontal wind for each dataset, with any wind velocities that were extrapolated below Earth's surface (e.g., in ERA5) set to zero before solving the Poisson problem for the streamfunction (Vishnu et al. 2020). The resulting track dataset for ERA5, together with tracks for four other reanalyses, are available in a Zenodo repository (<https://doi.org/10.5281/zenodo.3890646>).

## 2) MESOSCALE CONVECTIVE SYSTEMS

The FLEXTRKR algorithm is used to track MCSs in observations and model simulations. An MCS is defined as a convective system with 1) cold cloud shield (CCS)  $> 4 \times 10^4$  km<sup>2</sup> containing a precipitation feature (PF) with major axis length  $> 100$  km and 2) PF area, mean rain rate, rain rate skewness, and heavy rain volume ratio larger than corresponding lifetime dependent thresholds, with 3) both conditions 1 and 2 lasting continuously for longer than 4 h. As in Feng et al. (2021b), CCS is tracked using geostationary satellite Tb data and defined using a threshold of  $T_b < 241$  K. For model simulations, Tb is derived based on simulated outgoing longwave radiation following the empirical formulation provided by Yang and Slingo (2001). PF is tracked using the IMERG hourly precipitation data and PFs are defined as contiguous areas within the CCS with hourly rain rate  $> 2$  mm h<sup>-1</sup>.

## 3) FRONTAL PRECIPITATION

Fronts are identified using an automated method applied to 6-hourly gridded data at 2.5° resolution (Berry et al. 2011, Catto et al. 2015). This method calculates a thermal front parameter (TFP) as function of a thermal parameter:

$$\text{TFP}(\theta_w) = -\nabla|\nabla\theta_w| \cdot \left( \frac{\nabla\theta_w}{|\nabla\theta_w|} \right).$$

While many variables can be used to calculate the thermal front parameter (Thomas and Schultz 2019), we have used the wet bulb potential temperature ( $\theta_w$ ) as in Hewson (1998). After calculating the TFP, the field is masked where this is above a fixed negative threshold. Frontal points are then defined as the locations where the gradient of the TFP is equal to zero. These points are joined into contiguous lines and regridded as binary objects with an area of influence of plus and minus one grid box. Fronts can be separated into warm, cold, and quasi-stationary fronts, but here we have maintained simplicity by considering all fronts together.

## 4) ATMOSPHERIC RIVERS

With few exceptions, previous studies have utilized only a single AR detection tool (ARDT) in each study, whereas over 30 ARDTs currently exist (Shields et al. 2018; Rutz et al. 2019). Recent results from the Atmospheric River Tracking

Method Intercomparison Project (ARTMIP) have demonstrated that different ARDTs can produce different scientific results, which suggests that multiple ARDTs may need to be used when evaluating climate models in order to gain a complete picture of model skill in simulating ARs (O'Brien et al. 2020b). ARs are detected globally using six independently developed ARDTs, which we refer to by the following code names: ARCONNECT v2 (Sellars et al. 2017), GuanWaliser v2 (Guan et al. 2018), Lora v2 (Skinner et al. 2020), Mundhenk v3 (Mundhenk et al. 2016), TECA BARD v1.0 (O'Brien et al. 2020b), and Tempest LR (McClenney et al. (2020); O'Brien et al. 2022). These ARDTs were run on output from the MERRA-2 reanalysis as part of the ARTMIP Tier 1 experiment (Shields et al. 2018) and on output from the CMIP5 and CMIP6 multimodel ensembles as part of the ARTMIP Tier 2 CMIP5/6 experiment (O'Brien et al. 2022). The methods use a variety of heuristic rules to objectively identify atmospheric rivers from integrated vapor transport (IVT; the vertical integral of horizontal moisture transport) and/or integrated water content. For example, the widely used GuanWaliser v2 algorithm identifies ARs as continuous regions of integrated vapor transport exceeding the climatological 85th percentile, if the continuous regions meet specific geometric thresholds indicative of long and narrow regions of intense poleward moisture transport. We employ multiple ARDTs because recent literature indicates that different ARDTs may, in some instances, lead to qualitatively different answers to the same question (O'Brien et al. 2020a,b; Zhou et al. 2021).

## 3. Spatiotemporal characteristics metrics

Precipitation variability at different spatial and temporal scales is associated with specific processes such as convection driven by diurnal solar heating at the land surface, seasonal moisture convergence related to monsoon systems, disturbances related to convectively coupled equatorial waves, and large-scale atmosphere–ocean interactions. Diagnostics and metrics of spatiotemporal precipitation characteristics are therefore useful for relating model biases to specific mechanisms of precipitation generation at relevant ranges of spatiotemporal scales. These metrics are also useful for informing use of climate model precipitation data at appropriate spatiotemporal scales. Four metrics to benchmark the diurnal cycle of precipitation, daily precipitation and duration of dry spells, fractional contribution to the total mean rainfall from different intensities, and spatial and temporal coherence of precipitation are discussed in this section.

### a. Diurnal cycle of precipitation

The Fourier analysis has been widely applied to quantifying the diurnal cycle of precipitation in both observations and GCMs. However, the model-simulated rainfall is often quite noisy and therefore is poorly fit by low-order Fourier harmonics at single grid points. Covey et al. (2016) proposed a summary metric which illustrates the model-simulated Fourier amplitude and phase, averaged separately over all land and all ocean areas, in a single two-dimensional map. This metric enables

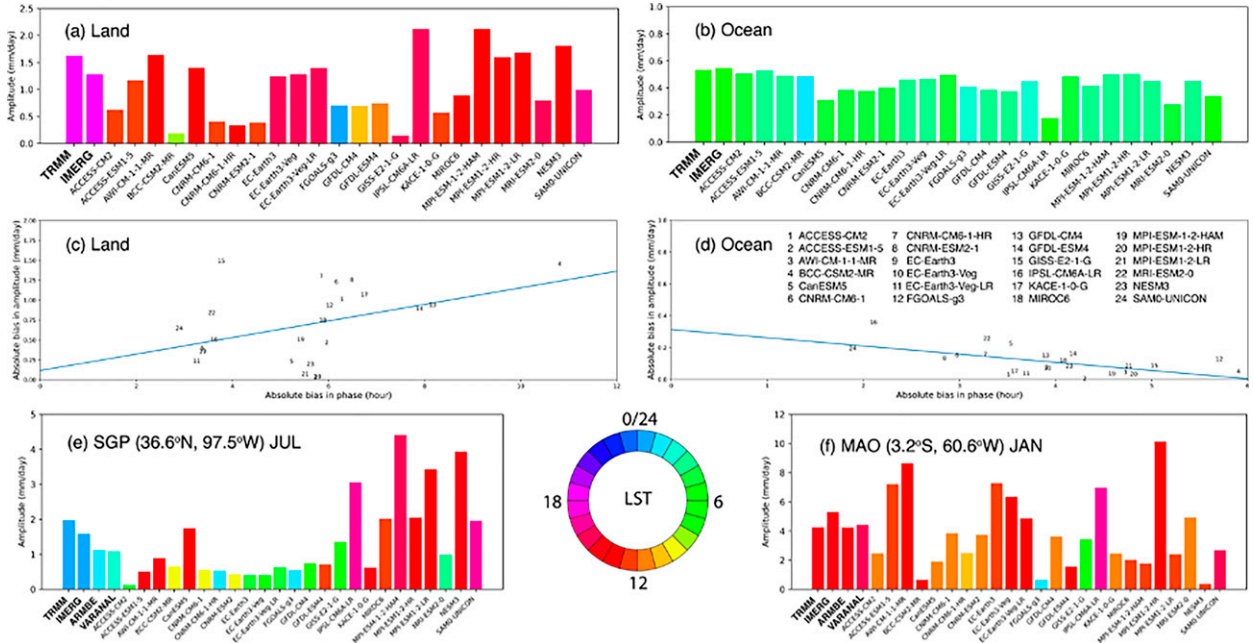


FIG. 1. (a) Bar plot of the composite mean diurnal harmonic amplitude (y axis) and phase in LST (colors) of summertime precipitation averaged over land. (b) As in (a), but over ocean. (c) Scatterplot of absolute bias in diurnal harmonic phase vs amplitude over land. (d) As in (c), but over ocean. (e),(f) As in (a), but for ARM SGP and MAO sites, respectively. Here, summertime precipitation refers to July for the Northern Hemisphere and January for the Southern Hemisphere. Model precipitation for 24 CMIP6 historical simulations is examined for the years 1996–2005.

intercomparison of climate models with observations and with each other over different climate regimes, but it becomes problematic when the number of models increases. Here we extend the procedure of Covey et al. (2016) and propose a metric that clearly displays the Fourier amplitude and phase of each individual model from a large number of groups in one bar plot (C. Tao et al. 2021, unpublished manuscript).

Figure 1a shows an example of the composite diurnal harmonic amplitude and phase (in LST) of summertime precipitation from 24 CMIP6 models versus observations over land. To generate the figure, we first produce a composite diurnal time series of precipitation, averaged over many years, for each grid point. We then apply Fourier analysis on the composite diurnal cycle of precipitation and focus on the first harmonic component, following Dai (2001). Here, the diurnal harmonic amplitude and phase are averaged over all land points between 50°S and 50°N using a vector averaging method, which automatically down weights the areas with a weak diurnal cycle (Covey and Gleckler 2014; Covey et al. 2016). Model precipitation is evaluated for the period of 1996–2005. Previous studies (e.g., Dai et al. 2007) have indicated that a stable diurnal cycle can be obtained with just a few years of data. As shown, the two satellite-based observations (TRMM 3B42 v7 and GPM-IMERG) agree quite well with each other in terms of both diurnal amplitude and phase. Over land, the major deficiency of the models is the too early diurnal precipitation peak, consistent with previous studies (e.g., Dai 2006; Xie et al. 2019). The majority of the models

show a diurnal harmonic phase peaking between 1200 and 1500 LST instead of early evening from the observations. The observed early morning diurnal harmonic phase over the ocean is generally captured by most of the CMIP6 models (Fig. 1b) while the corresponding diurnal harmonic amplitude is somewhat underestimated in all 24 CMIP6 models. To highlight the models with the best performance, Figs. 1c and 1d show the scatterplot of absolute model bias relative to TRMM observations in diurnal harmonic phase versus amplitude over land and over ocean, respectively. Over ocean, interestingly, models that perform better in the diurnal cycle phase tend to perform worse in amplitude (Fig. 1d). The relationship between model biases in precipitation diurnal phase and amplitude over land is less significant than that over ocean but there is a tendency for models with smaller bias in phase to have correspondingly smaller bias in amplitude (Fig. 1c). Particularly, EC-Earth3, EC-Earth3-Veg, and EC-Earth3-Veg-LR compare the best to the observations over land in terms of both diurnal amplitude and phase. Similar results are found by interpolating the data to a common grid (not shown). Generally, the impact of model resolution on the simulated diurnal cycle of precipitation is minimal.

The metric diagram can also be easily computed for smaller scales and at different locations where rich ground-based high-frequency observations are available. Figures 1e and 1f compare the simulated diurnal harmonic amplitude and phase to observations at the two ARM sites (SGP and MAO) where precipitation shows distinct diurnal variability with SGP



featuring a nocturnal peak whereas MAO has an afternoon peak. Despite some discrepancies, the satellite-based products agree fairly well with the ground-based rain gauge and/or radar measurements in general. As shown in Fig. 1e, there is a large model spread in both diurnal amplitude and phase at SGP, with most of the models (all but two) failing to capture the observed nocturnal peak around midnight in which half of the models actually show a diurnal precipitation peak in the afternoon. CNRM-CM6-1-HR and FGOALS-g3 simulate the diurnal phase much closer to that observed but both significantly underestimate the diurnal harmonic amplitude. The majority of models show a diurnal precipitation peak around noon at MAO, a few hours earlier than the observed (Fig. 1f). In general, the CMIP6 models show diverse results simulating the diurnal amplitude with some overestimating the observed value and some underestimating it, but they often show consistent biases in simulating the diurnal phase. Almost all the models peak too early during the day and miss the nocturnal diurnal peak at certain regions.

To summarize, the metric developed here provides a quick comparison with observations and among models, and reasonably summarizes the systematic model errors in reproducing the diurnal cycle of precipitation over both large areas and single point locations. Particularly, by displaying the diurnal harmonic amplitude and phase from the Fourier analysis in one bar plot, this metric enables the evaluations with a focus on individual model performance from a large number of models.

#### *b. Extremes: Daily precipitation and duration of dry spells*

Despite being seemingly contrasting variables, daily precipitation and the duration of dry spells share many features in the shape of their probability distributions. The probability density functions (PDFs) of both quantities are characterized by a power-law range, where the probability decreases slowly with each order of magnitude increase in precipitation rate or duration of dry spells, up to a cutoff scale (denoted  $P_L$  for daily precipitation and  $t_L$  for the duration of dry spells; see Figs. 2a,b) where the probability decreases roughly exponentially (Figs. 2a,b), ultimately controlling the size of extreme percentiles (Martinez-Villalobos and Neelin 2018, 2021; Chang et al. 2020). These quantities have connections with the moisture budget, with  $P_L$  (and hence also extreme percentiles) scaling with the amplitude of moisture convergence fluctuations within precipitating events (Neelin et al. 2017; Martinez-Villalobos and Neelin 2019), and  $t_L$  scaling with the balance between moisture convergence fluctuations at dry times and the mean moisture source tendency (Pierrehumbert et al. 2007; Stechmann and Neelin 2014).

Recently, Martinez-Villalobos and Neelin (2021) showed that the shape of the large daily precipitation probability tail and the spatial pattern of the cutoff scale are well simulated by GCMs but there is a bias in the magnitude of  $P_L$  compared to observational datasets (see also Fig. 2a). This suggests that two metrics can succinctly summarize the general model behavior of daily precipitation and dry-spell duration extremes. The first one is the spatial correlation coefficient over 50°S–50°N (the spatial extent of TRMM-3B42; see Table 1) between model simulated  $P_L$  and

$t_L$  patterns (see Figs. 2c,d for their CMIP6 multimodel mean and their observational counterparts (TRMM-3B42 in this case). The second metric is the scaling factor, defined as the model area weighted mean  $P_L$  or  $t_L$  divided by the TRMM-3B42 observational estimate of the same quantity. The first metric tests whether extremes are well simulated spatially regardless of magnitude (values can range between  $-1$  and  $1$ , with  $1$  denoting a model that simulates the spatial pattern of TRMM-3B42), and the second tests the overall magnitude of the pattern (values can range between  $0$  and infinity, with  $1$  being the best). To gauge model behavior, we also calculate the same metrics comparing GPCP versus TRMM-3B42 as a measure of observational uncertainty. The differences between observational precipitation products can be large, thus, model results may be sensitive to the choice of target observational product. This sensitivity is discussed in section 4 and in Fig. S2 in the online supplemental material. We note the caveat that part of the differences between models and observational products noted below may be the result of sampling different internal variability realizations (Deser et al. 2012) due to the relatively short span in which precipitation observational products have been available. However, different realizations from models of the same family (e.g., GFDL models, CNRM models) tend to perform similarly, which suggest that sampling variability has only a minor effect on the results. More details on these metrics and methodology are given in the online supplemental material.

Figures 2e and 2f show the results for 35 CMIP6 models and for the multimodel ensemble mean (MME) for  $P_L$  and  $t_L$  respectively. We first find that there is a substantial observational uncertainty for  $P_L$ . The overall magnitude of  $P_L$  in GPCP is about 70% (scaling factor of 0.68) of TRMM-3B42 magnitude and the correlation coefficient of the patterns is 0.81. There are several models that are closer to TRMM-3B42 than the observational uncertainty. Among these we highlight HadGEM3-GC31-MM as the model with the closest  $P_L$  spatial pattern ( $r = 0.89$ ) and GFDL-ESM4 as the model with the closest overall magnitude to TRMM-3B42 (scaling factor = 0.98). The MME benefits from the good performance of the best models in the spatial structure and averages the overall magnitude of  $P_L$  in the different models. This results in a multimodel mean that is closer than GPCP to TRMM-3B42 in both  $P_L$  spatial pattern and magnitude.

The model performance on the duration of dry spells is similarly encouraging. While all individual models and the MME simulate longer duration of dry spells than both TRMM-3B42 and GPCP (even after the models wet-day biases are greatly reduced; see the online supplemental material), the  $t_L$  pattern correlation in almost all models is comparable, although reduced, with the pattern correlation between TRMM-3B42 and GPCP. Even though the magnitude of  $P_L$  and  $t_L$  (hence also extreme percentiles) differs from TRMM-3B42 in almost all models, the fact that the patterns are well correlated helps boost confidence in model projections of relative (i.e., percent) changes of daily precipitation and dry-spell duration extremes.

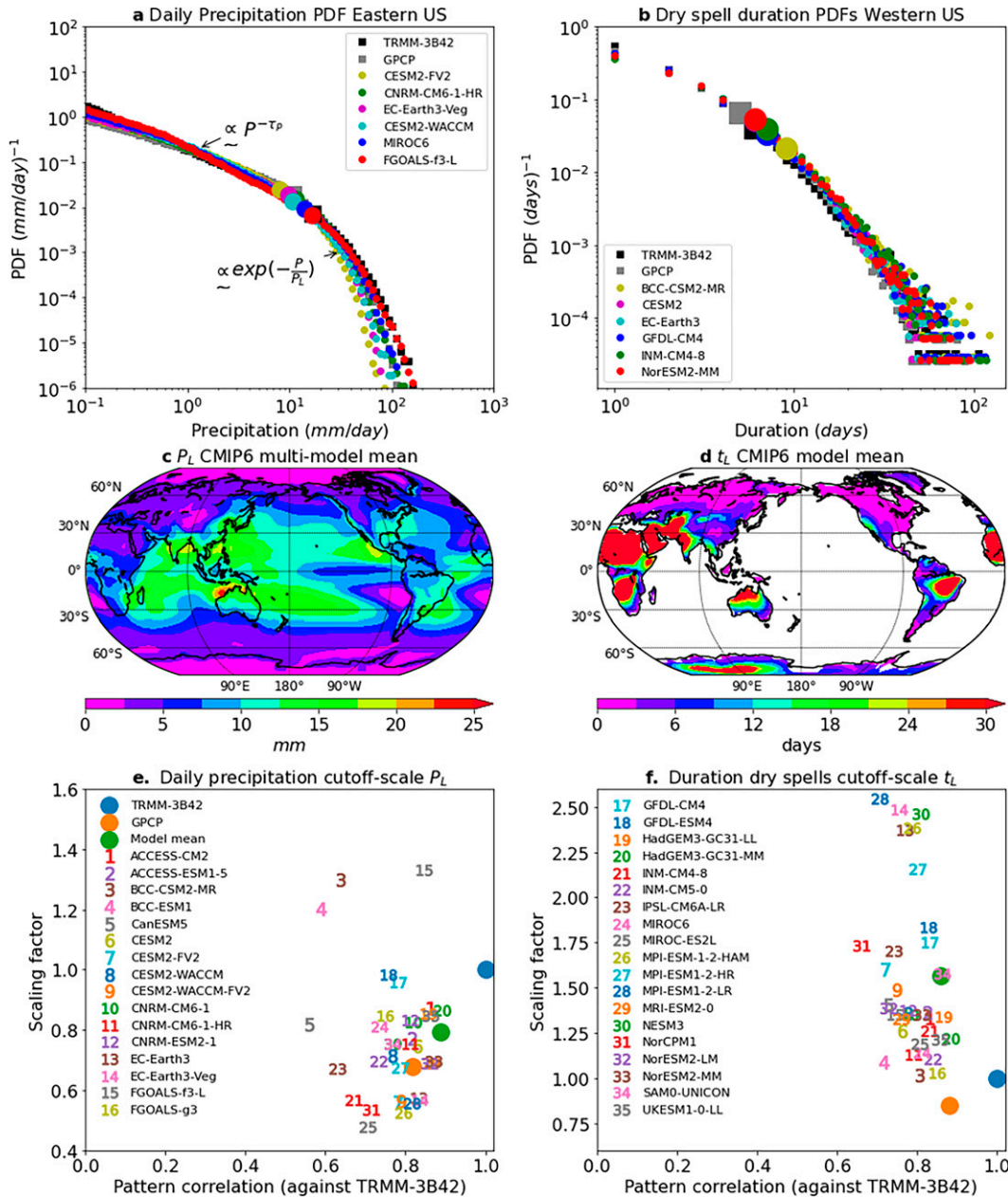


FIG. 2. Observational (GPCP and TRMM-3B42) and selected models: (a) daily precipitation PDFs in the eastern United States (25°–48°N, 257°–294°E) and (b) dry spell durations PDFs in the western United States (30°–48°N, 236°–257°E). In (a) and (b) the cutoff scales are shown by a large circle (for models) or large squares (for observational datasets). Note that the larger or longer the cutoff scale, the more extreme is the large event tail. (c). Multimodel mean (out of 35 models) of the daily precipitation cutoff-scale  $P_L$  pattern. (d). Multimodel mean of the dry spell duration cutoff-scale  $t_L$  pattern (with model-dependent dry-day precipitation threshold). (e). Scatterplot of the  $P_L$  scaling factor and pattern correlation coefficient against TRMM-3B42 for individual models [numbers; legend across (e) and (f) gives corresponding acronyms], multimodel mean (green dot), GPCP (orange dot), and TRMM-3B42 [blue dot at (1, 1) by definition]. (f). As in (e), but for  $t_L$  scaling factor and pattern correlation coefficient.

*c. Spectral analysis*

Following the method of Klingaman et al. (2017) implemented in Analyzing Scales of Precipitation (ASoP) version 1.0, we calculate the fractional contribution to the total mean rainfall

from different intensities, at 3-h and daily time scales, sorted into 100 bins of varying width ranging from 0.005 to 2360 mm day<sup>-1</sup>. This reveals the relative importance of precipitation events in a given intensity bin to the total precipitation. The

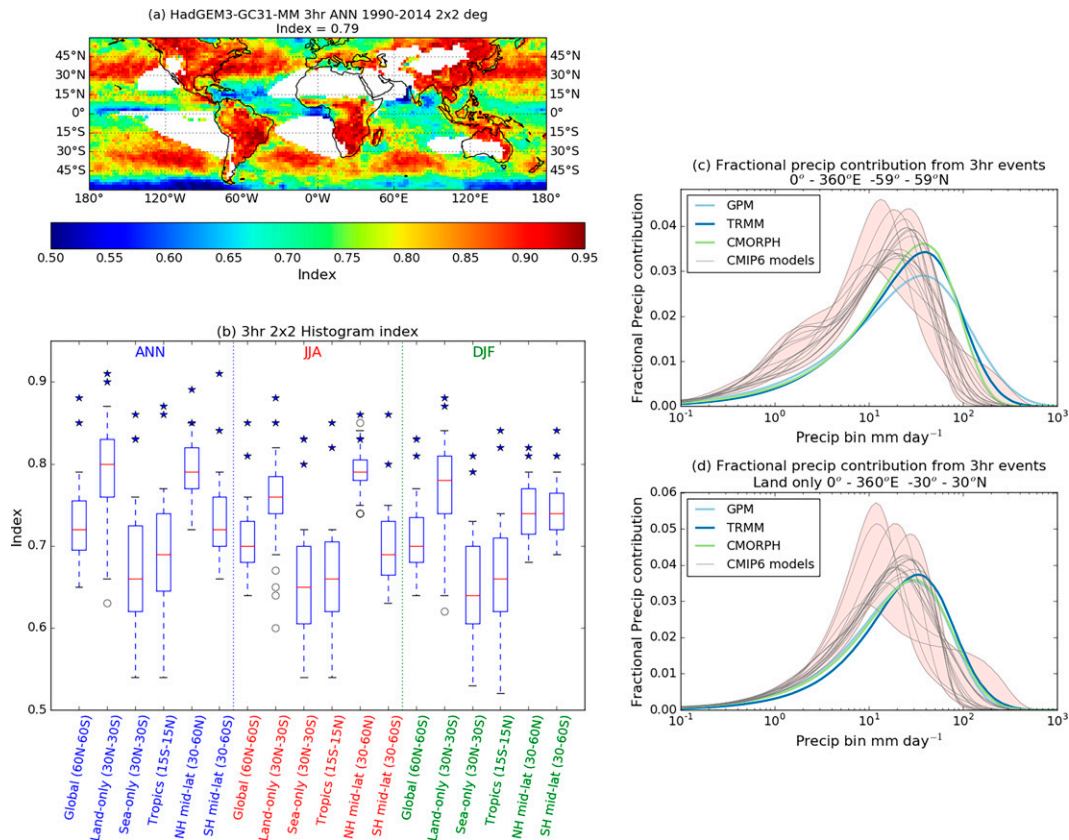


FIG. 3. (a) Example map of index from 3-h rainfall data from HadGEM3-GC31-MM vs GPM; (b) Summary metrics for different regions from time series of 3-h rainfall data from 23 CMIP6 models compared with GPM-IMERG observations. Boxes show the interquartile range while whiskers indicate the full range of model indices. Red line shows the median. Filled stars indicate other observational datasets (TRMM and CMORPH). (c),(d) Histograms of 3-h rainfall data from 23 CMIP6 models and 3 observational datasets for global (60°S–60°N) and land-only (30°S–30°N) domains, respectively. All model and observation data were averaged to a  $2^\circ \times 2^\circ$  grid, using conservative area-weighted averaging, before analysis.

calculation is performed at each grid box, using a horizontal resolution that is sufficiently coarse for at least some spatial averaging to be carried out for all of the models and the observations. To avoid removing important spatial detail, we limit this resolution to  $2^\circ \times 2^\circ$ , thereby requiring us to omit models whose resolution is similar to or coarser than this. Calculations are performed for the whole year (ANN) and for each season, over 25 years (1990–2014) of CMIP6 historical simulations.

To evaluate the models, we use a similarity index (Perkins et al. 2007) to compare the fractional histograms from each model with those obtained from 19 years of GPM-IMERG observations (2001–19) at each grid point between 60°S and 60°N. This measures the overlap between the model and observed histograms, with values closer to 1.0 indicating that the histograms match better and a value of 0.0 indicating they are entirely separated. Metrics are the spatial root-mean-square of these indices over selected regions. Any region could be chosen for metric evaluation; here we have used six regions: global (60°S–60°N), tropics (15°S–15°N), land-only (30°S–30°N), sea-only (30°S–30°N), Northern Hemisphere

(NH) midlatitudes (30°–60°N), and Southern Hemisphere (SH) midlatitudes (30°–60°S).

Figure 3a shows an example map of the indices from 3-h rainfall data from HadGEM3-GC31-MM versus GPM-IMERG. This suggests that performance is better over land than ocean, and over the midlatitudes than the tropics. Figure 3b shows the overall metric summary information for the 3-h time scale. This confirms that the pattern seen for HadGEM3-GC31-MM is similar for the other CMIP6 models and is consistent through the seasons. The stars indicate comparison of GPM-IMERG with other observation datasets, providing a measure of uncertainty. The metrics from the models nearly all lie outside this uncertainty range. Figure 3c provides additional information about the model–observation differences: the models are generally biased toward smaller rainfall accumulations, although there are a few models for which there is a greater than observed contribution from the largest rainfall accumulations. We find similar results for daily accumulations.

The metrics are a useful guide to the overall model performance, but the fact that the histograms are analyzed at each grid point, and that the calculations can be performed on any

temporal or spatial scale, means there is much more information available from these diagnostics to users and model developers that could be used to understand model errors on a range of time scales (see, e.g., [Martin et al. 2017](#)). There is also the potential for subsampling of rainfall associated with organized systems or phenomena (such as tropical cyclones, fronts, MJO) prior to the histogram analysis, which could increase our understanding of these systems as well as providing information on model errors. The metrics could also be used to examine the influence of model resolution, ocean–atmosphere coupling, and the inclusion of Earth system processes on the spread of rainfall intensities.

#### d. Coherence analysis

The “Analyzing Scales of Precipitation” (ASoP) diagnostics ([Klingaman et al. 2017](#)) can measure, and compare, the spatial and temporal scales of precipitation across observations and GCMs. The “ASoP Coherence” package was designed to produce a single diagnostic or metric for a chosen region. Here, we extend the package to operate on gridded data. We measure spatial and temporal coherence in 3-hourly and daily precipitation in GPM-IMERG observations and CMIP6 historical simulations. We perform these calculations on a common  $2^\circ \times 2^\circ$  grid, a horizontal resolution that is sufficiently coarse for at least some spatial averaging to be carried out for all of the models and the observations while also avoiding removing important spatial detail. This requires us to omit models whose resolution is similar to or coarser than this. The calculations are performed between  $60^\circ\text{S}$  and  $60^\circ\text{N}$ , neglecting any point with annual-mean rainfall  $< 1 \text{ mm day}^{-1}$  and, in the remaining points, any months in the dry season, defined as months that contribute, in the mean, less than  $1/24$  of the annual precipitation.

[Figures 4a–c](#) use 3-hourly data to show the temporal scale, defined as the first lag at which the temporal autocorrelation is  $< 0.2$ , for the CMIP6 historical multimodel mean ([Fig. 4a](#)), GPM-IMERG ([Fig. 4b](#)), and the multimodel mean bias ([Fig. 4c](#)). Throughout much of the tropical and subtropical oceanic regions, the CMIP6 multimodel mean precipitation is too persistent, highlighting an area for model improvement. [Figures 4d–f](#) use daily-mean data to show the spatial scale, which is computed from the temporal correlation of the precipitation between each grid point and its surrounding grid points, using intervals of radii given in the color bar beneath the panel. The scale is defined as the first search radius at which the spatial correlation is  $< 0.2$ . Daily precipitation spatial scales are larger in the CMIP6 multimodel mean ([Fig. 4d](#)) than in GPM ([Fig. 4e](#)), particularly in the eastern equatorial Pacific and Atlantic Oceans, and in near-equatorial regions of the Indian Ocean, as well as much of the subtropical oceans ([Fig. 4f](#)). Combined with the temporal scale results above, this suggests that CMIP6 models produce precipitation features that are too large and that last too long, particularly in the tropical oceans.

[Klingaman et al. \(2017\)](#) also defines spatial and temporal coherence metrics. The spatial metric is derived from the likelihood of coincidence of upper-quartile and lower-quartile precipitation at neighboring grid points; the temporal metric is derived from the likelihood of consecutive time steps of the

upper quartile and lower quartile at the same grid point. Quartiles are computed for each grid point and each month of the seasonal cycle. For the temporal coherence metric, we show the aggregated grid point metrics (computed  $60^\circ\text{S}$ – $60^\circ\text{N}$ ) as Taylor diagrams for global land ([Fig. 4g](#)), ocean ([Fig. 4h](#)), and all points ([Fig. 4i](#)). The CMIP6 models show higher centered RMS difference and lower correlations, against GPM-IMERG, for land points than for ocean points, indicating that persistence of land precipitation is another area for model improvement. The spatial standard deviation values of the coherence metrics shown in the Taylor diagrams can provide further insights for model improvements: models that have a smaller standard deviation than GPM-IMERG are typically too persistent across all grid points, as the mean bias is positive for nearly all models (not shown), while models that have a greater standard deviation are typically too persistent in some regions and too intermittent in others. These standard deviations show stronger negative biases over land than over ocean, indicating that models show little spatial variability in temporal coherence over land and hence cannot distinguish regions dominated by longer-lived rain-bearing systems from regions dominated by shorter-lived systems.

Next, we demonstrate the ability to compare the spatial scale of precipitation (now restricted to the tropical Indian Ocean:  $10^\circ\text{S}$ – $10^\circ\text{N}$ ,  $60^\circ$ – $90^\circ\text{E}$ ; using daily data; determined as correlations at a distance of 800 km) with two metrics of the MJO in CMIP6 models, two satellite observation datasets, and ERA5 ([Fig. 4j](#)). The satellite observations and ERA5 have an average precipitation spatial coherence of 0.06–0.09, and the CMIP6 models cover the range  $-0.03$  to  $+0.26$ . CMIP6 models have a relatively close relationship between precipitation spatial coherence and the MJO Maritime Continent propagation metric ([Ahn et al. 2020](#);  $R^2 = 0.489$ ). This suggests that those climate models with a higher spatial coherence of daily precipitation propagate the MJO more robustly east over the Maritime Continent. The relationship is weaker ( $R^2 = 0.114$ ) between precipitation spatial coherence and the MJO east/west power ratio, which measures MJO spatiotemporal structure (e.g., [Sperber and Kim 2012](#); [Ahn et al. 2017](#)). There is no relationship between precipitation temporal scale and either MJO metric. Comparisons between spatiotemporal characteristics metrics and process- or phenomena-based metrics may be able to lead to greater insights and understanding of the origins of biases and model errors.

#### 4. Process-oriented metrics

Although metrics of spatiotemporal characteristics are suggestive of the processes contributing to precipitation biases at different spatial and temporal scales, they do not by themselves represent processes related to precipitation. Here, process-oriented metrics are used to reveal relationships between precipitation and the thermodynamic environments, which provide important information on the ability of models to reproduce the observed relationships and the potential contributions of large-scale biases in the atmospheric environments to the precipitation biases. Here, we discuss two metrics highlighting the coupling of precipitation with the thermodynamic environments.

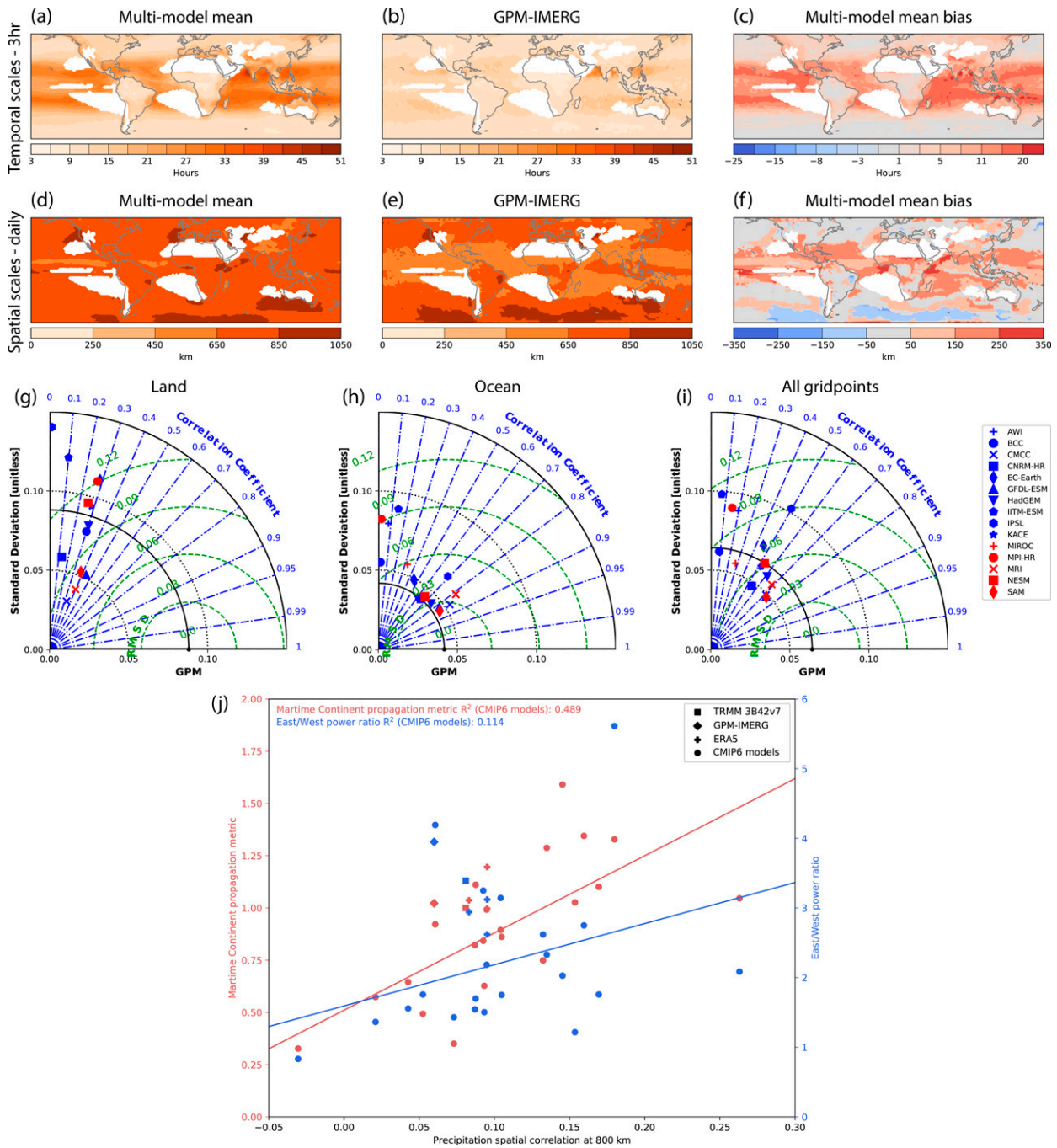


FIG. 4. (top) Temporal scale (h; defined as the first lag at which the autocorrelation of 3-h precipitation is  $< 0.2$ ; within  $60^{\circ}\text{S}$ – $60^{\circ}\text{N}$ ) (a) the CMIP6 historical multimodel mean, (b) GPM-IMERG, and (c) the multimodel mean bias. (second row) Spatial scale (km; defined as the first distance at which the correlation between a grid point and neighboring points within a distance bin is  $< 0.2$ , with bin edges given as divisions of the color bar) for (d) the CMIP6 historical multimodel mean, (e) GPM-IMERG, and (f) the multimodel mean bias. In (a)–(f), white shading denotes grid points with annual-mean precipitation  $< 1 \text{ mm day}^{-1}$ , which are not included in the analysis. (third row) Summary Taylor diagrams of temporal coherence, using 3-hourly data, over (g) land-only, (h) ocean-only, and (i) all grid points, within  $60^{\circ}\text{S}$ – $60^{\circ}\text{N}$ , for the CMIP6 models vs GPM-IMERG. (j) Tropical Indian Ocean ( $10^{\circ}\text{S}$ – $10^{\circ}\text{N}$ ,  $60^{\circ}$ – $90^{\circ}\text{E}$ ) precipitation spatial correlation at 800 km ( $4 \times$  the grid scale) vs MJO Maritime Continent propagation metric (Ahn et al. 2020; left-side axis; red) and MJO east–west power ratio (e.g., Sperber and Kim 2012, Ahn et al. 2017; right-side axis; blue) for two satellite observational datasets (TRMM and GPM-IMERG), ERA5 over three different periods, and the CMIP6 models. Quoted  $R^2$  values and lines of best fit are for CMIP6 models only.

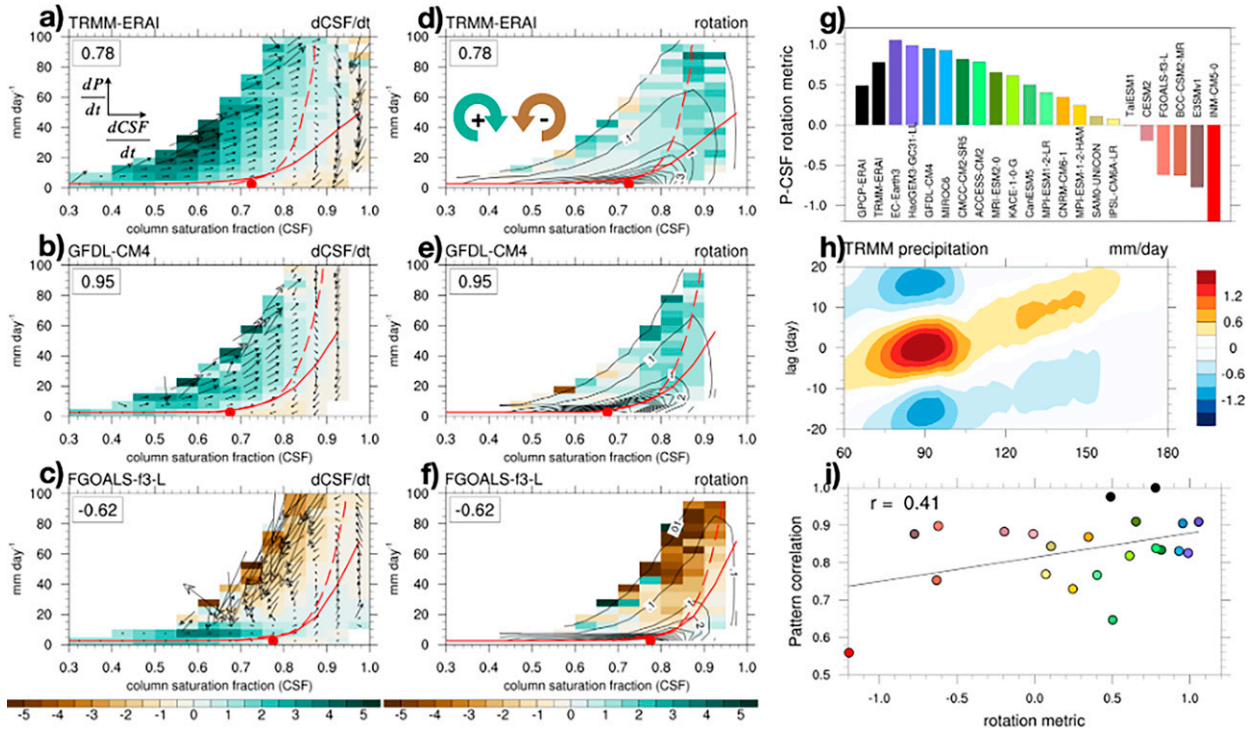


FIG. 5. (left) Daily mean column saturation fraction (CSF) tendency ( $d\text{CSF}/dt$ ; in  $\% \text{ day}^{-1}$ ; shading) and the daily mean joint CSF rainfall-rate tendencies (vectors) as a function of CSF rainfall rate ( $P$ ) for the Indo-Pacific warm pool (ocean-only grid points in  $20^{\circ}\text{S}$ – $20^{\circ}\text{N}$ ,  $30^{\circ}$ – $180^{\circ}\text{E}$ ) from ERA-Interim and (top) TRMM 3B42, (middle) GFDL-CM4 (the median high-performing model), and (bottom) FGOALS-f3-L (the median low-performing model). Red filled circle is mode of observations; red solid and dashed lines are mean rainfall rate and CSF, respectively, for a given CSF or rainfall rate bin. (center) Nondimensional rotation ( $d[d\text{CSF}/dt]/d\text{CSF} - d[dP/dt]/dP$ ; shading with clockwise rotation shaded green; value shown in upper left of each panel) and CSF- $P$  probability distribution function (contours). (right) Shown are (top) frequency-weighted mean rotation for ERA-Interim–TRMM and CMIP6 models, i.e., the “rotation metric”; (middle) the lagged regression of TRMM 3B42 tropical precipitation ( $10^{\circ}\text{S}$ – $10^{\circ}\text{N}$ -averaged) anomalies onto the 20–100-day filtered eastern Indian Ocean area-averaged ( $5^{\circ}\text{S}$ – $5^{\circ}\text{N}$ ,  $85^{\circ}\text{E}$ – $95^{\circ}\text{E}$ ) rainfall; and (bottom) scatterplot of individual model convection–moisture rotation metric against the Jiang et al. (2015) MJO propagation metric (see text for details), where colors of dots match bar colors above. The correlation of the two metrics is  $r = 0.49$ .

### a. Rainfall–moisture coupling

Latent heating from tropical rainfall formation forces large-scale circulation anomalies that affect weather patterns globally through the tropical–extratropical teleconnection response (Stan et al. 2017). The onset of tropical heavy rainfall is critically dependent upon the relative saturation of the atmosphere (Bretherton et al. 2004; Neelin et al. 2009), while the teleconnection response is sensitive to the spatial and temporal scale of the heating anomaly (Yadav and Straus 2017; Wang et al. 2020). The MJO is a prominent example of a large-scale tropical disturbance that is strongly governed by column moisture (Adames and Kim 2016) and is also a major driver of tropical–extratropical teleconnections (e.g., Henderson et al. 2017). With this section, we aim to understand how tropical rainfall and moisture are coupled and how this coupling affects MJO simulation in CMIP6 models.

Following Wolding et al. (2020), daily tendencies of precipitation ( $P$ ) and column saturation fraction (CSF; i.e., vertically integrated column water vapor divided by vertically integrated saturation column water vapor) over the Indo-Pacific warm

pool are averaged within conditionally sampled CSF and  $P$  bins. All data are first remapped onto a common  $2.5^{\circ} \times 2.5^{\circ}$  grid. In Figs. 5a–c, joint CSF and  $P$  (CSF– $P$  for short) tendencies are shown with vectors, which indicate if CSF– $P$  departures above or below the mean CSF– $P$  line lead to column moistening or drying. In observations and in most CMIP6 models, the vectors rotate clockwise about the mode (red circles in Figs. 5a–f) that corresponds to the quasi-equilibrium state (Neelin et al. 2008; Wolding et al. 2020). This clockwise rotation indicates that anomalously high precipitation for a given CSF is associated with column moistening, while anomalously low precipitation is associated with column drying. The strength of this rotation in each CSF– $P$  bin can be diagnosed using a vorticity-like metric based on nondimensionalized CSF and  $P$  tendencies where positive values denote clockwise rotation (Fig. 5b). A scalar rotation metric  $R$  is then computed as the frequency-weighted rotation in CSF– $P$  space.

For models with  $R > 0$ , positive moistening and rainfall tendencies are largest during the dry-to-moist transition when  $P$  is much greater than its mean value for a given CSF (solid red line

in Figs. 5a–c). Analysis of radar data collected over the tropical Indian Ocean indicate that this state is associated with a transition from trade wind cumulus to cumulus congestus (Wolding et al. 2020). Negative moistening and rainfall tendencies are largest when CSF is greater than its average value for a given  $P$  (red dashed line in Figs. 5a–c), a state associated with widespread stratiform rainfall with embedded convection. For models with  $R < 0$ , higher-than-average rainfall at intermediate CSF is associated with strong drying; positive  $P$  tendencies are only observed at high CSF. Rainfall–moisture coupling in  $R < 0$  models suggests that exaggerated depletion of column water vapor by rainfall leads to excessive drying at intermediate CSF, thus reducing the likelihood of subsequent heavy precipitation. Heavy precipitation in these models is only observed at high CSF, where the environment cannot be rapidly dried by rainfall.

Correlations between the  $R$  metric and several MJO propagation “pattern correlation” metrics for a subset of CMIP6 models suggest that tropical rainfall–moisture coupling plays an important role in regulating MJO periodicity. Various MJO pattern correlation metrics have been used to assess MJO propagation in models by correlating simulated and observed rainfall lagged regressions over the warm pool. Jiang et al. (2015) computed pattern correlations of regression coefficient using the composite propagation plotted in Fig. 5h (i.e., the “full” metric). Wang et al. (2018) and DeMott et al. (2019) reduced the influence of MJO period on the pattern correlation by masking coefficients within  $\pm 15^\circ$  longitude of the rainfall basepoint (the “masked” metric), while Ahn et al. (2020) completely removed periodicity effects by only considering positive coefficients in a small portion of the domain east the Maritime Continent (the “MC-crossing metric”). Correlations between the  $R$ -metric and the full, masked, and MC-crossing propagation metrics are 0.47, 0.23, and 0.11, respectively. The correlation is only statistically significant for the full pattern correlation metric, which measures the combined effects of MJO propagation and period.

### b. Temperature–water vapor environment

The aim of this module is to create metrics that capture the typical range of moisture and temperature over which precipitation is produced by condensing information from prior diagnostics (which also provides information on sensitivity to sampling and resolution; Kuo et al. 2018, 2020). Here we use a thermodynamic space in which temperature is measured by the vertically integrated saturation humidity,  $q_{\text{sat}}$ , and moisture is measured by column relative humidity,  $\text{CRH} = \text{CWV}/q_{\text{sat}}$ , where CWV is column water vapor, for each  $q_{\text{sat}}$ . Figure 6a shows, for  $q_{\text{sat}} = 65.5$  mm over tropical oceans, the conditional mean precipitation rate (circles) and precipitation contribution (lines) from observations and one model instance. For observations, we use precipitation from the TRMM PR, column water vapor from the TRMM Microwave Imager (TMI), and ERA5 temperature for computing  $q_{\text{sat}}$  (for an alternative combination of observations, we use MERRA-2 temperature in Figs. 6c,d). The PR is coarse-grained to  $0.25^\circ \times 0.25^\circ$ , compatible with the CWV resolution; results are insensitive to resolution up to  $1.5^\circ$  (Kuo et al. 2018). The observed precipitation rate sharply picks

up as CRH increases above a certain threshold. The precipitation contribution peaks near this value because the system spends less time at the high precipitation values and the many occurrences of low CRH contribute little to precipitation. The MIROC-E2SL model exhibits qualitatively similar behavior, although the precipitation pickup is too weak and begins at lower CRH than observed, as seen more clearly in the peak of the precipitation contribution. To characterize the moisture range over which precipitation is produced, we identify the CRH values associated with the 25th and 75th percentiles of precipitation contribution for each  $q_{\text{sat}}$ . These CRH values for  $q_{\text{sat}}$  (tropospheric temperature environment) between the 25th and 75th percentiles of  $q_{\text{sat}}$  (blue lines) are shown in Fig. 6b, together with the precipitation contribution as a function of CRH and  $q_{\text{sat}}$  (color contours). A notable feature is that the CRH values associated with the 25th and 75th percentiles as well as peak of precipitation contribution decrease as  $q_{\text{sat}}$  increases (i.e., precipitation is produced at lower CRH in a warmer environment).

The values associated with these percentiles provide a good summary of the observed thermodynamic range associated with precipitation, shown by the blue trapezoid in Fig. 6b. We choose a visual reference range (gray box) and repeat it in Figs. 6c and 6d. Figure 6c presents typical thermodynamic ranges associated with precipitation from a subset of CMIP6 historical simulations and two observational combinations. Deviations of the trapezoids from the observed along the  $q_{\text{sat}}$  axis indicate cold/warm biases in the simulation, and deviations along the CRH axis indicate that models tend to produce precipitation outside the observed CRH range. Figure 6d exhibits the thermodynamic ranges as in Fig. 6c, but for the 17 available CMIP6 models, ranked by the precipitation contribution error defined as the  $L^2$  difference between the observed and model-simulated precipitation contribution (i.e., the mean square of the dotted area in Fig. 6a), averaged over the four most probable  $q_{\text{sat}}$  bins. This scalar metric focuses on relative humidity rather than temperature bias. It is encouraging to see that some of the models can produce most of their precipitation in a thermodynamic environment close to the observed range both by the scalar metric and rhomboid location. Other models fare poorly by these measures. Most models capture the decrease in the CRH for the 75th-percentile precipitation contribution with increasing temperature, but only about half capture this feature for the 25th percentile.

## 5. Phenomena-based metrics

Phenomena-based metrics emphasize weather features such as synoptic systems and different types of storms that generate precipitation. While synoptic systems such as fronts may be broadly resolved by GCMs at typical  $1^\circ$  resolution, storms such as tropical cyclones, LPS, and MCS require higher-resolution modeling. Models’ ability to simulate these storms is critical as they are key contributors to extreme precipitation in many regions. Feature tracking (briefly summarized in section 2b) is used to identify and track the weather features, allowing precipitation associated with these features be isolated and evaluated using different metrics that measure model–observation

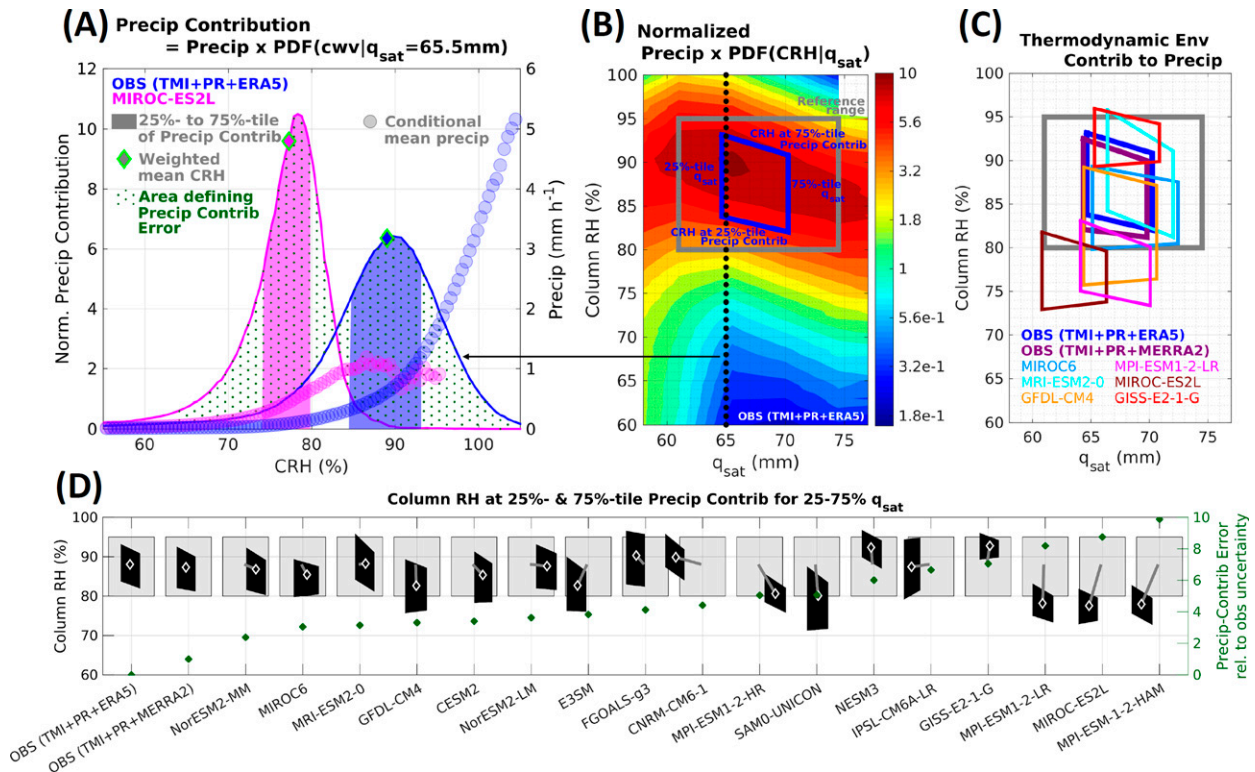


FIG. 6. (a) Observed (blue; TMI + PR + ERA5; see text) and an example model (magenta) conditional mean precipitation rate (circles) and precipitation contribution (lines) as a function of column relative humidity (CRH) for column-integrated saturation humidity  $q_{sat} = 65.5$  mm (bin width: 4.5 mm) for tropical oceans within  $20^{\circ}\text{S}$ – $20^{\circ}\text{N}$ . Note that precipitation contributions here are normalized so the area under each curve is one. The 25th–75th percentiles of precipitation contributions are indicated by shaded areas, and the mean CRH weighted by 25%–75% precipitation contribution by diamonds. The precipitation contribution error [in (d)] is defined by considering the  $L^2$  difference between the observed and model-simulated precipitation contributions, i.e., the mean square of the dotted area. (b) Color contours indicate precipitation contributions normalized for each  $q_{sat}$ . Blue trapezoid is the CRH at the 25th and 75th percentiles of the precipitation contribution between the 25th and 75th percentiles of  $q_{sat}$ . The gray box indicates a visual reference range, which remains invariant in (c) and (d). (c) Trapezoids are as in (b), but from a set of CMIP6 historical simulations; the observed trapezoid from (b) is repeated (blue) and an additional observational combination (TMI + PR + MERRA2; see text) shown in purple. (d) Black trapezoids and gray boxes are as in (c). The white diamonds indicate the 50th-percentile  $q_{sat}$  and the mean CRH weighted by the precipitation contribution within the rhomboid range. The precipitation-contribution error (dark green) is defined as the  $L^2$  difference between the observed and model-simulated precipitation contribution, averaged over the four most probable  $q_{sat}$  bins. Here the difference between the two observational combinations provides a simple measure of observational uncertainty and is used to normalize the precipitation contribution error.

differences. Here, four examples of weather features and associated precipitation are discussed.

#### a. Low pressure systems

A wide variety of synoptic-scale disturbances that consist of balanced flow around a pressure minimum produce precipitation in Earth's tropics and extratropics. Classic examples are midlatitude baroclinic waves, which often produce intense precipitation through semigeostrophic uplift in their frontal zones, and tropical cyclones, which produce precipitation through the radial, frictionally balanced component of their circulation. Understanding the mechanisms by which such systems amplify and generate precipitation requires tracking the systems from initial genesis; this can be a difficult task, requiring data of sufficiently fine resolution and algorithms of adequate robustness to unambiguously represent the weak and sometimes horizontally small low pressure center. Here we illustrate how a strategic

choice of variables allows for improved tracking of low pressure systems (LPS) in the South Asian monsoon, which produce a large fraction of that region's annual mean rainfall as well as many extreme precipitation events. This tracking exercise allows the relationship of circulation with precipitation to be characterized in observations and model ensembles.

Tropical LPS are most commonly identified and tracked using lower-tropospheric vorticity or sea level pressure. Even for strong tropical cyclones, ambiguities in the criteria used in the tracking algorithm can lead to large uncertainties in the number of storms identified in observationally constrained gridded data (e.g., Murakami 2014). This issue is even more problematic for weak LPS, where the noisiness of the vorticity field produces irregular, broken tracks for systems that seem to move smoothly when tracked subjectively using a standard suite of meteorological data (Fig. 7a). Sea level pressure, which is less noisy, is sometimes used to track LPS instead but



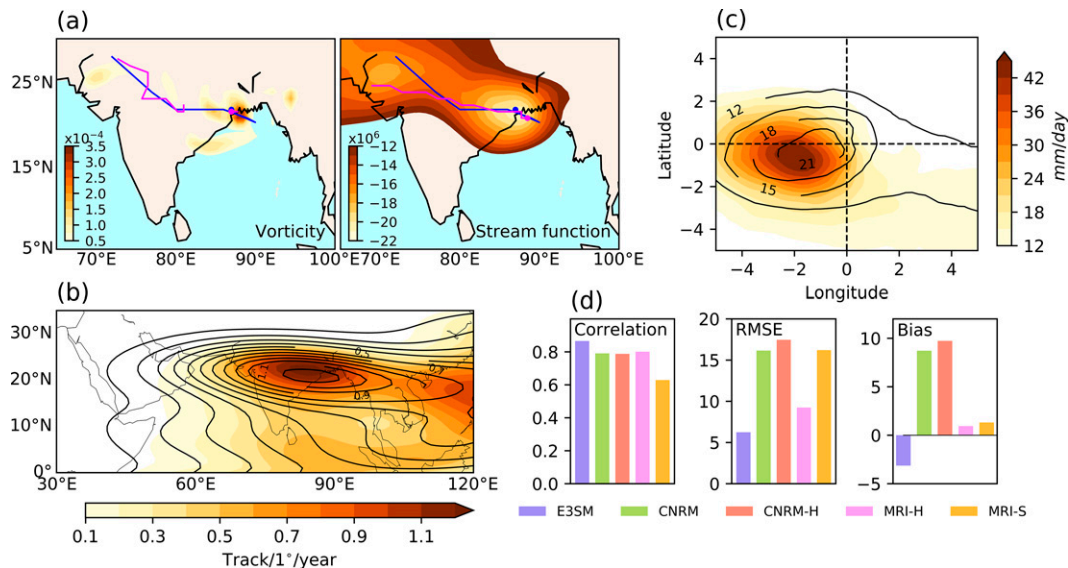


FIG. 7. (a) Example of the influence of variable choice on tracking skill: compared to the 850-hPa relative vorticity, which is commonly used to track tropical disturbances, a more continuous track that better matches the subjectively analyzed reference track is obtained using the streamfunction of the 850-hPa horizontal wind (magenta lines show tracks obtained from an automated algorithm applied to ERA5, while blue lines show the reference track). (b),(c) Comparison of model (black contours, for E3SM) and observed (shading) representations of (b) climatological mean track density and (c) vortex-centered composite rain rate for South Asian monsoon low pressure systems. E3SM simulates a reasonable track density but produces disturbances that rain too little with peak rainfall biased slightly toward the vortex center. Observed tracks are from ERA5; observed precipitation from TRMM. (d) Metrics showing skill of E3SM and four High-ResMIP models in simulating the spatial structure of rainfall in South Asian low pressure systems. For vortex-centered composite rain rates [as in (c)], we show the correlation coefficient, root-mean-square error (in  $\text{mm day}^{-1}$ ), and horizontal mean bias (in  $\text{mm day}^{-1}$ , averaged over a  $10^\circ \times 10^\circ$  box around the composite vortex center) compared to TRMM. Note that E3SM has the highest correlation and the lowest RMSE but a larger magnitude bias in horizontal-mean precipitation than the MRI models; the MRI model skill degrades at finer resolution (MRI-S is finer resolution than MRI-H), while CNRM model skill has little sensitivity to resolution.

is ill suited for South Asian LPS, which typically have winds that peak around 3 km above the surface; geopotential height near the level of maximum wind also does not capture the full rotational flow given the low latitude and high Rossby number of these storms. Physical reasoning, as well as systematic assessment of multiple candidate variables with hundreds of combinations of quantitative tracking criteria, showed that the streamfunction of the horizontal 850-hPa wind is an optimal variable to use for tracking these LPS (Fig. 7b; Vishnu et al. 2020). This streamfunction represents the full nondivergent wind, even when geostrophic balance does not hold, yet retains the smoothness of the geopotential or sea level pressure fields; it was inverted using a method to avoid contamination by any wind data extrapolated below Earth's surface (Vishnu et al. 2020).

Precipitation in South Asian monsoon LPS is known to fall southwest of the storm center, where the interaction of the storm's rotational flow with the background vertical shear produces quasigeostrophic uplift (Rao and Rajamani 1970; Sanders 1984). This placement of peak precipitation is well captured when compositing TRMM precipitation relative to ERA5 LPS tracks (Fig. 7c). ERA5 also accurately represents the well-known distribution of track

density, with storm frequency peaking strongly over the northwest Bay of Bengal (Fig. 7b). Recent work has shown that LPS frequency likely peaks in that small region because the large-scale, low-level monsoon winds are barotropically unstable there (Diaz and Boos 2019) and vapor pressures are large with strong horizontal gradients (Ditchek et al. 2016; Adames and Ming 2018). Wind-enhanced evaporation from the Bay of Bengal may also enhance LPS intensity there (Murthy and Boos 2020; Fujinami et al. 2020; Diaz and Boos 2021).

By tracking LPS in ensembles of GCMs, we can create composites that allow model precipitation bias to be assessed in a phenomenon-based system rather than in a space- or time-based system that averages many types of atmospheric disturbances. One high-resolution GCM (E3SM integrated at  $0.25^\circ$  resolution) represents the track density of South Asian monsoon LPS well, in addition to the spatial structure of precipitation relative to the vortex center (Figs. 7b,c). This is notable given the poor ability of some coarse-resolution GCMs to simulate these LPS (Praveen et al. 2015). However, the E3SM model simulates monsoon LPS rainfall that is too weak, with the peak storm-centered composite precipitation being about half that observed (Fig. 7c). Other models exhibit

a variety of biases in their representation of monsoon LPS precipitation with differing sensitivities to model resolution. Storm-centered composites in the CNRM models have overly strong precipitation with little sensitivity to model resolution, while the MRI models produce roughly the right amount of precipitation over the entire storm but with a spatial pattern that, unexpectedly, degrades at finer model resolution (Fig. 7d). These biases are large for some models, exceeding 50% of the system-averaged TRMM rain rate of  $15 \text{ mm day}^{-1}$ ; interannual variations in LPS activity and storm-centric rain rates are substantially more modest (e.g., Sikka 1980; Krishnamurthy and Ajayamohan 2010; Vishnu et al. 2020).

Such assessment of model skill in representing the synoptic systems that produce extreme rainfall, such as monsoon LPS, is an important step in producing reliable projections of future extreme rainfall. The LPS dataset used here, which is available for five modern reanalysis products, provides LPS tracks throughout the global tropics that can be used to better understand a variety of synoptic-scale phenomena, including the weak progenitors of tropical cyclones.

### b. Mesoscale convective systems

Mesoscale convective systems (MCSs) are ubiquitous over the tropics year-round and in the midlatitudes during the warm season. Besides contributing to over 50% of the annual precipitation in most regions of the tropics and selected regions in the midlatitudes (Nesbitt et al. 2006; Feng et al. 2021b), MCSs are also key contributors to extreme precipitation, partly because of their larger size and longer lifetime compared to individual convective storms (Stevenson and Schumacher 2014). Because of the distinctive nocturnal timing of MCS, erroneous diurnal timing of summer precipitation produced by models has been used to infer their failure in simulating MCSs. Recent efforts in developing algorithms to identify and track MCSs in observations (Feng et al. 2018) and model simulations (Feng et al. 2021a) have provided unprecedented opportunities to directly evaluate MCSs and their characteristics in weather and climate models using MCS-specific metrics.

Using FLEXTRKR, an algorithm developed to track MCSs using both infrared brightness temperature ( $T_b$ ) and precipitation features (PFs) (Feng et al. 2018, 2019), a global ( $60^\circ\text{S}$ – $60^\circ\text{N}$ ) MCS tracking database has been developed at  $\sim 10 \text{ km}$  and hourly resolution (Feng et al. 2021b). Combining the track locations and precipitation, this database can be used to derive information of the MCS number, MCS precipitation and its fractional contribution to the total precipitation, MCS maximum precipitation rate, MCS lifetime, and MCS translation speed and direction. As MCSs are not well defined at coarser spatial resolution, we develop MCS metrics mainly for use in evaluating high-resolution weather and climate simulations with grid spacing  $< 50 \text{ km}$ . Instead of coarse graining the observations and model outputs, which correspond to a range of grid spacing, to a common resolution, we use specific PF criteria derived for a given resolution for MCS tracking to facilitate comparison across datasets of different resolutions (Feng et al. 2021a).

Figures 8a and 8b compare the MCS number tracked using two algorithms, a more commonly used method that tracks MCSs using  $T_b$  only versus FLEXTRKR that tracks MCSs using both  $T_b$  and PF. These two methods produce similar observed total MCS number and spatial distribution in the tropics, but larger differences are noticeable in the midlatitudes. Including PF in MCS tracking noticeably reduces the number of MCSs in the midlatitudes by disqualifying large cold cloud systems (e.g., synoptically forced) with small area and/or low rainfall intensity PF as MCSs. Using only IR  $T_b$ , the model (E3SM) simulates too many MCSs (blue contours) except in a few locations. In contrast, using both IR  $T_b$  and PF, E3SM simulates too few MCSs (magenta contours) except in a few locations. These results show that large cold cloud systems are produced by the model too frequently but many of them fail to meet the PF thresholds. This is supported by the composited MCS rain rates shown in Fig. 8c for northeast moving MCSs in the central United States during spring (MAM) and summer (JJA). The simulated and observed rain rate composites have similar size, but the model produces much lower peak rain rates. A higher fraction (65%) of MCSs in the model have a northeast propagation than observed (44%).

Figure 8d summarizes the MCS precipitation metrics for four models in HighResMIP. The pattern correlation, root-mean-square error (RMSE), and bias are calculated based on comparison of the observed and simulated composited MCS rain rates over the central United States. Since hourly  $T_b$  is not available from the HighResMIP models except E3SM, MCSs are tracked using an algorithm that depends only on PF, trained using MCSs tracked using both  $T_b$  and PF (Feng et al. 2016). Note that E3SM is a free-running fully coupled simulation with constant 1950 forcing while other simulations are atmosphere-only simulations driven by observed sea surface temperature and sea ice distribution. The models exhibit a range of biases from larger negative (E3SM) to larger positive (NICAM) and the skills are generally lower during summer than spring. The seasonal difference is particularly large for NICAM. Unlike the other models that parameterized deep convection, no deep convection scheme was used in NICAM at 56-km grid spacing. Last, it is worth noting that metrics based on composited MCS precipitation can only reveal differences in PF qualified as MCS. All models evaluated here display significant dry bias in the summer, consistent with the ubiquitous warm, dry bias noted in CMIP5 (Lin et al. 2017), as the models simulate much lower numbers of MCSs compared to observations. Therefore, we emphasize the importance of using multiple metrics for comprehensive evaluation of precipitation in models.

### c. Frontal precipitation

Fronts have been identified using the method described in section 2, applied to ERA-Interim and five CMIP6 models, giving gridded front objects on a  $2.5^\circ$  grid. The fronts are linked to daily precipitation, using GPCP 1DD as an observational precipitation estimate. The precipitation data are regridded to the same resolution as the fronts in order to

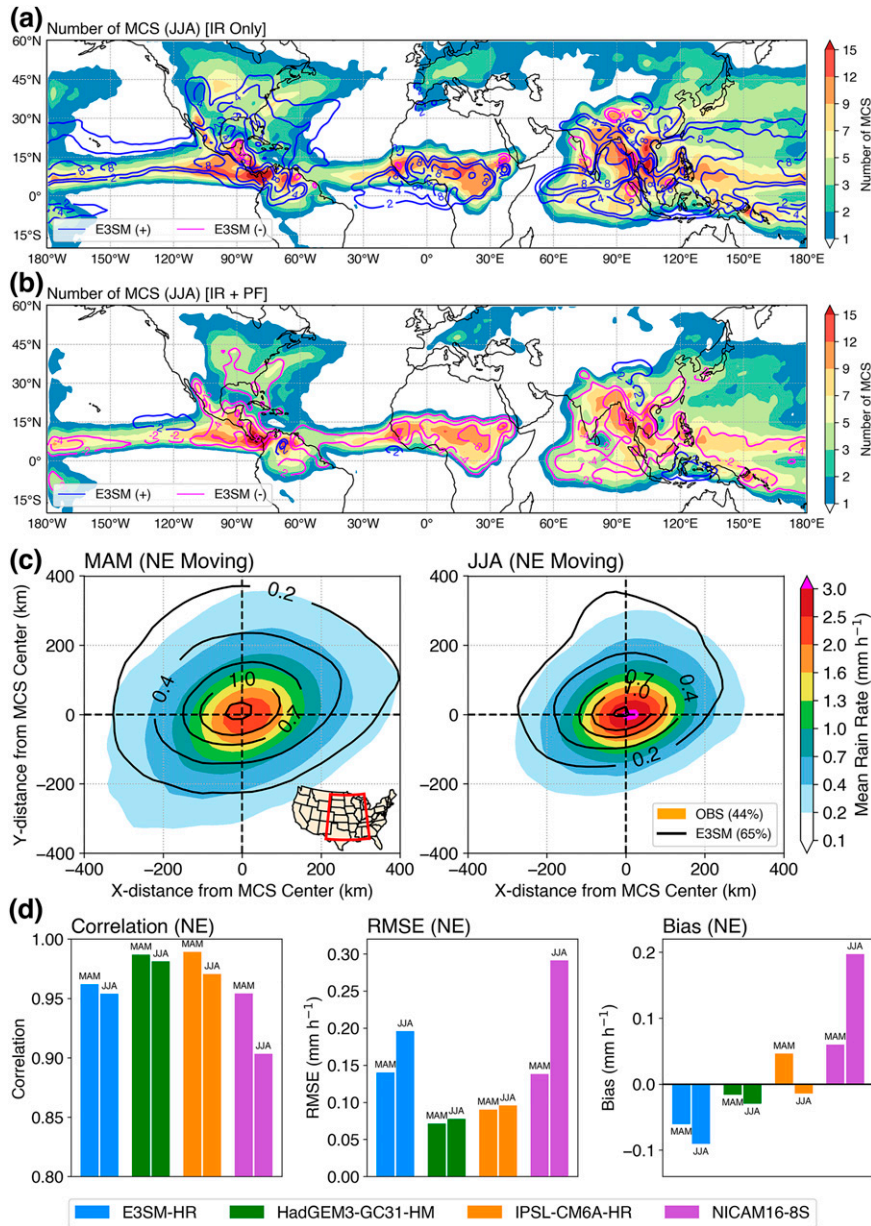


FIG. 8. Influence of variable choice on MCS tracking. MCSs are tracked using (a) only infrared brightness temperature ( $IR T_b$ ) and (b) both  $IR T_b$  and precipitation feature (PF). The observed number of MCS is shown in color shading and the model bias is shown in color contours (blue and magenta) for positive and negative) bias, respectively. (c) Comparison of simulated (black contours, for E3SM) and observed (colored shading) MCS rain rates ( $mm h^{-1}$ ) composited with a center collocated with the geometric centroid of the MCS PF. Composites are shown for (left) spring (MAM) and (right) summer (JJA) for northeast-moving MCSs inside the central United States (red region) shown in the inset in the left panel. Trained on the MCS statistics tracked using both  $IR T_b$  and PF, the MCSs used in these composites are tracked using only PF to facilitate comparison with other models for which hourly precipitation but not hourly outgoing longwave radiation is available. (d) Metrics showing skill of E3SM and three HighResMIP models in simulating the spatial structure of MCS rainfall in the central United States. Based on the rain rate composites [as shown in (c) for E3SM], three metrics—correlation coefficient, root-mean-square error ( $mm h^{-1}$ ), and mean bias ( $mm h^{-1}$ )—are used to evaluate different aspects of the model MCS rainfall.

make the linking simpler. We consider precipitation only if it is above a threshold of 1 mm, which is the minimum 24-h precipitation a gauge can measure, and this eliminates some of the “drizzle problem” that models tend to have (Stephens et al. 2010). The precipitation is associated with a front if it lies within the front area of influence (which is equivalent to being in the same grid box or the surrounding eight grid boxes) during any of the four 6-hourly reanalysis times in the 24-h precipitation period. From this association of fronts and precipitation, we can produce the diagnostics of frontal (and nonfrontal) precipitation frequency ( $F_f$ ,  $F_{nf}$ ), frontal (and nonfrontal) precipitation intensity ( $I_f$ ,  $I_{nf}$ ), frontal amplification factor ( $A_f = I_f/I_{nf}$ ), and fraction of total precipitation from fronts ( $P_f$ ) [see Catto and Pfahl (2013) and Catto et al. (2015) for full details]. Comparing the model diagnostics to the observational estimates from ERA-Interim and GPCP, we can produce a number of metrics, including the correlation, RMSE, and bias of these values.

Since precipitation biases ( $E_p$ ) in the models depend on the frequency of fronts, the frequency of precipitation, and the intensity of the precipitation, we can also decompose the bias of each model into components associated with these characteristics as follows:

$$E_p = \Delta F_f I_{f,o} + F_{f,o} \Delta I_f + \Delta F_f \Delta I_f + \Delta F_{nf} I_{nf,o} + F_{nf,o} \Delta I_{nf} + \Delta F_{nf} \Delta I_{nf},$$

where subscript  $o$  represents the observational estimate, and  $\Delta$  represents the difference between model and observational estimate. The cross terms (3 and 6) are generally very small and are not shown.

Maps (Fig. 9a) of the error decomposition for term 1 (contribution from frequency of frontal precipitation) show that there are large regions of positive bias contribution. Errors are largely confined to the regions of maximum storm track activity and in the NH the largest positive bias contributions can be seen over the Kuroshio, over western Europe and parts of the North Atlantic, and at the end of the Pacific storm track into North America. In the SH the largest positive contributions are in a band between 30° and 40°S, particularly around the south coast of Australia. Term 2 errors (contribution from intensity of frontal precipitation) are generally largest in the same regions and indicate negative contributions to the total bias, with this being particularly notable over the North Atlantic region. The maps indicate a compensation of biases between terms 1 and 2, which is confirmed for each of the models in Fig. 9b and is consistent with the CMIP5 models (Catto et al. 2015). In the midlatitudes the contribution to the total precipitation error from the nonfrontal precipitation terms is small (Fig. 9b), as expected due to the high frequency of fronts.

The models all overestimate  $A_f$  due to larger negative biases in the nonfrontal precipitation intensity than the negative biases in frontal precipitation intensity (not shown). These biases are large compared to the GPCP  $A_f$  of 1.28 in NH DJF and 1.35 in SH JJA and are strongly correlated with the model biases in the intensity of the frontal precipitation (not shown). The spatial correlation is between 0.4 and 0.6 in

the NH and between 0.3 and 0.4 in the SH, indicating a better representation in the NH.

The proportion of total precipitation associated with fronts in the winter seasons is 0.50 in the NH and 0.54 in the SH for GPCP and ERA-Interim. The biases in this quantity range between 0.02 and 0.27 (Fig. 9d), with most models showing a better representation in the SH. The models that perform better for the proportion do not necessarily show better performance in the  $A_f$  metric, indicating the utility of looking at more than one metric.

Analyzing the ranks of the models using the various calculated metrics, we can see that some models that perform well in metrics that quantify magnitude differences (e.g., the decomposition terms and biases) also perform poorly in their spatial correlation, (e.g., IPSL-CM6A-LR). Again, this points to the importance of considering a number of different metrics to investigate the model performance.

#### d. Atmospheric rivers

Atmospheric rivers (ARs) are long narrow bands of poleward vapor transport often associated with the warm sector in advance of midlatitude cyclone cold fronts (Ralph et al. 2018). They account for a large fraction of wet-season precipitation in a number of regions (Dettinger 2011; Rutz et al. 2014; Guan and Waliser 2015), and they account for a majority of the poleward moisture transport (Gimeno et al. 2014). Previous studies examining ARs in climate model simulations have assessed the ability of models to adequately simulate relevant characteristics of ARs, including global and landfalling frequency, intensity, precipitation, duration, life cycle, and so on (Dettinger 2011; Payne and Magnusdottir 2015; Shields and Kiehl 2016; Goldenson et al. 2018). In this module, we present two metrics aimed at answering the following questions: 1) Do models simulate AR-related precipitation in the correct locations? 2) Do models simulate enough contrast between regions with high AR precipitation and low AR precipitation? 3) Does the diversity of AR detection and tracking (ARDTs) affect the above conclusions?

We utilize output from six global ARDTs that participated in the ARTMIP Tier 1 experiment and Tier 2 CMIP5/6 experiment (see section 2b); these ARDTs identified ARs in MERRA-2 and in historical simulations from nine members of the CMIP5 and CMIP6 multimodel ensembles. We quantitatively define “AR-precipitation” for each ARDT as precipitation occurring when AR conditions are identified by a given ARDT. We calculate AR-precipitation for MERRA-2 (using the precipitation field from MERRA-2) and for the CMIP5 and CMIP6 simulations. We calculate 30-yr averages of these quantities and regrid all to a common 2° × 2° grid to facilitate direct comparison of the fields between the simulations and the reanalysis. Additionally, we calculate AR-precipitation for ERA 20C (1900–2010) to provide a combined estimate of observational uncertainty and natural variability (since we use a different time period than with MERRA-2). Figures 10a and 10b show the bias in AR-precipitation between two CMIP6 models, with one model’s bias field indicating some regional biases in AR-precipitation (Fig. 10a) and another model’s bias field indicating systematically too little AR-precipitation (Fig. 10b).

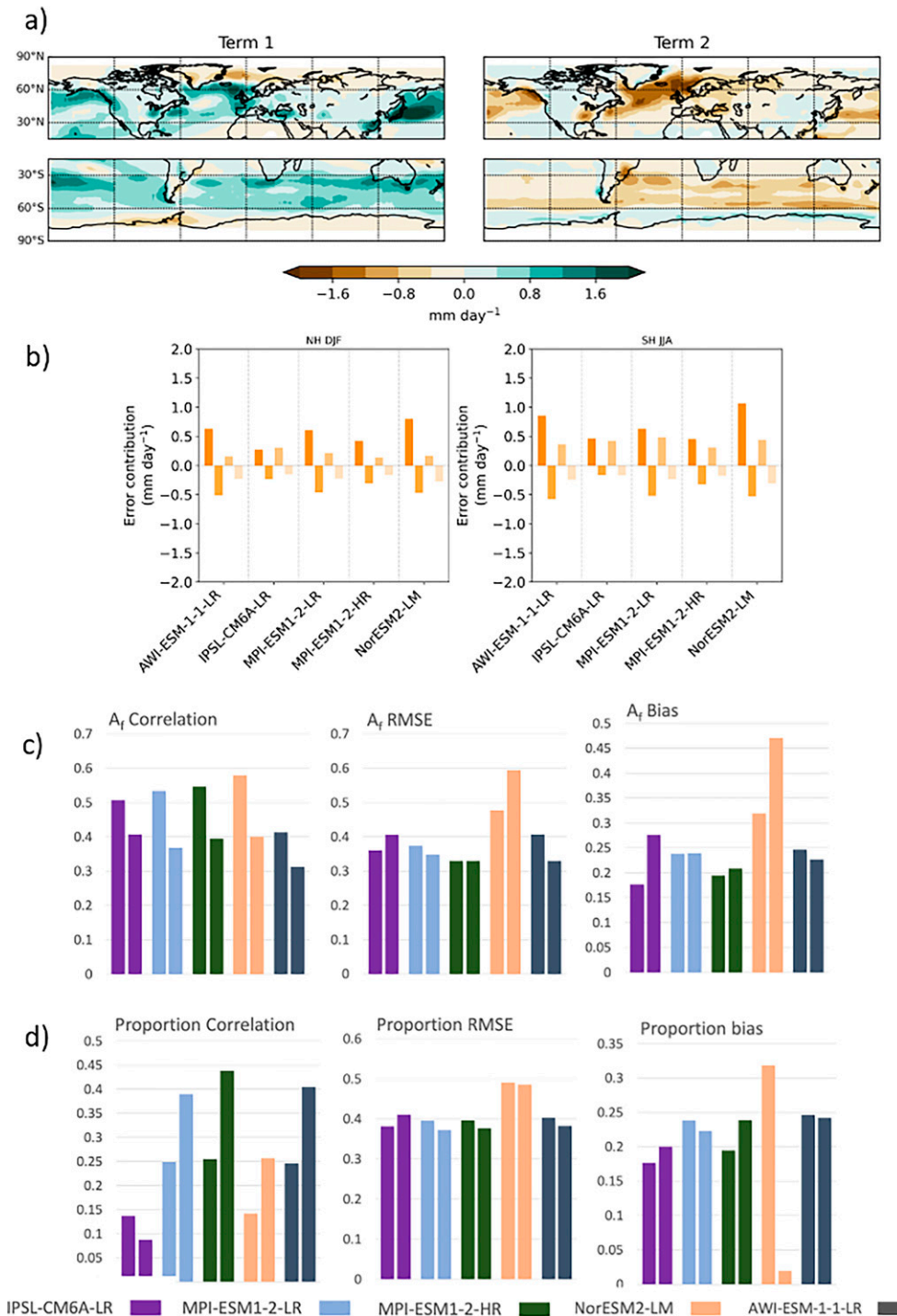


FIG. 9. Representation of frontal precipitation in five CMIP6 models (1980–2014) compared to ERA-Interim fronts with GPCP daily precipitation (1997–2017) for winter (DJF in the NH and JJA in the SH). (a) Multimodel mean of the first and second terms of the decomposition in  $\text{mm day}^{-1}$ . (b) Area-averaged decomposition terms (terms 1, 2, 4, and 5) for each of the models in the NH and SH extratropics (15°–90°). (c) The (left) correlation, (center) RMSE, and (right) bias for the frontal amplification factor  $A_f = P_f/P_{\text{nf}}$ . For each model the left bar is the NH extratropics in DJF and the right bar is the SH extratropics in JJA. (d) The (left) correlation, (center) RMSE, and (right) bias for the proportion of precipitation associated with fronts.

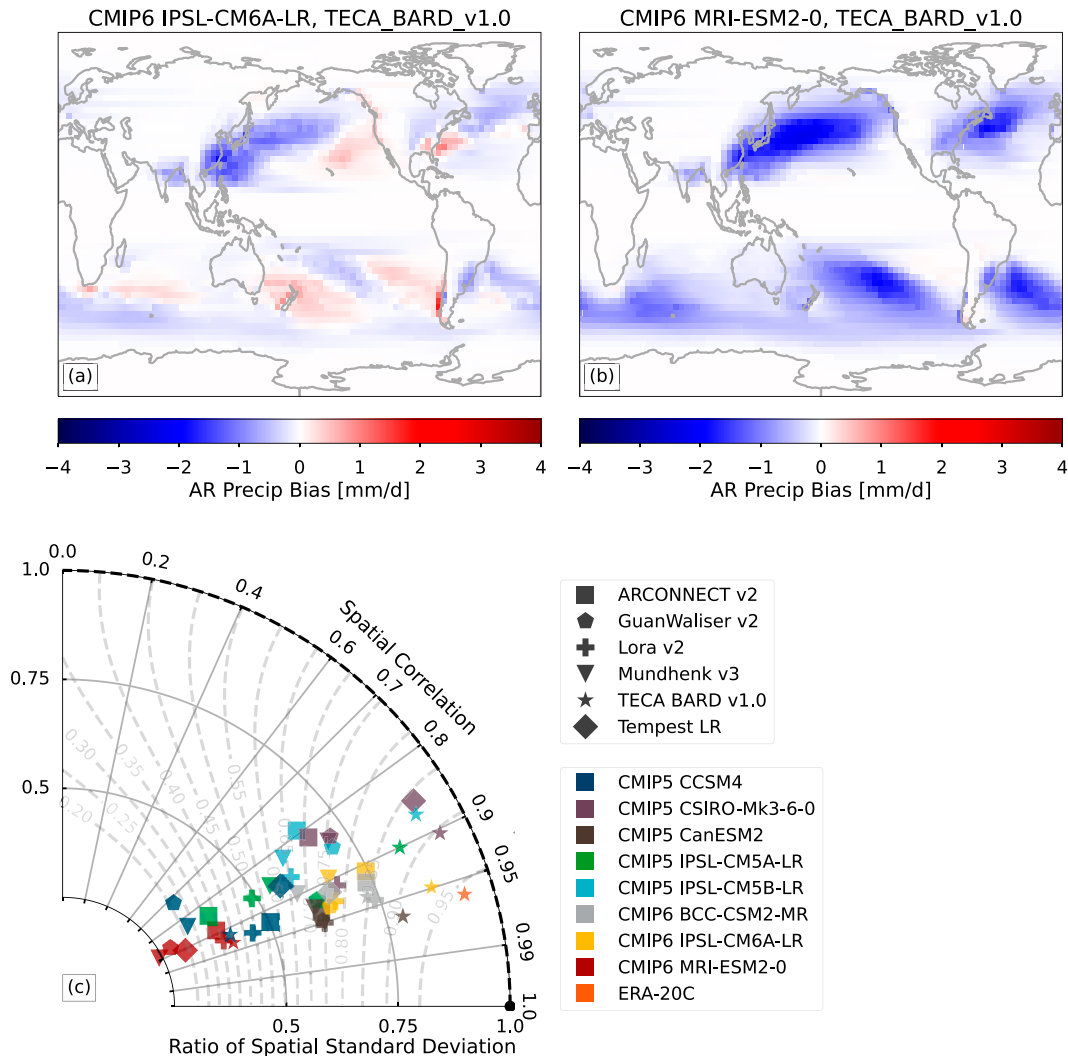


FIG. 10. AR precipitation metrics considering AR detection diversity. Bias in mean annual precipitation ( $\text{mm day}^{-1}$ ) associated with ARs detected using the TECA BARD v1.0 ARDT for (a) the CMIP6 historical simulation from the IPSL-CM6A-LR model (1950–86) and MERRA-2 (1980–2016) and (b) the CMIP6 historical simulation from the MRI-ESM2-0 model (1950–86) and MERRA-2 (1980–2016). (c) A Taylor diagram comparing the spatial correlations and spatial standard deviations of AR-precipitation between simulations and MERRA-2, using multiple ARDTs. Colors are associated with models, and markers are associated with ARDTs. The dashed gray curves in (c) show contours of constant Taylor skill metric.

The spatial correlation coefficient of AR-precipitation between each model simulation and MERRA-2 is used to answer question 1 above, and the ratio of the spatial standard deviation of AR-precipitation between each model and MERRA-2 is used to assess question 2. These quantities are calculated for all available model–ARDT pairs in order to assess question 3. Figure 10c shows a Taylor diagram constructed by plotting the spatial correlations on the azimuthal axis and the ratio of the standard deviations on the radial axis.

It appears that models generally produce AR-precipitation in the correct regions, but they do not have enough spatial variability in AR-precipitation. The models have relatively high spatial correlation coefficients (regardless of which ARDT is used), with most models having coefficients between 0.8 and

0.95. It is notable, however, that the value of the spatial correlation coefficient can depend strongly on which ARDT is used. Consider results from the CMIP5 CCSM4 simulation (navy blue markers), which range from about 0.7 when evaluated with the GuanWaliser v2 ARDT to over 0.9 with the ARCONNECT v2, Lora v2, and TECA BARD v1.0 ARDTs. In contrast to the spatial correlation, all models have less variability than the MERRA-2 simulation, and models exhibit a wide range of skill in this metric.

Across the ARDTs used, some models form distinct clusters in the Taylor diagram, with the CMIP6 MRI-ESM2-0 and CMIP5 CCSM4 simulations having systematically low Taylor skill values and the CMIP5 CanESM2 simulation having systematically high Taylor skill values. These distinct clusters

indicate consensus among the ARDTs about the model skill. In contrast, some models span the Taylor diagram; for example, the skill of the CMIP5 IPSL-CM5A-LR simulation depends strongly on which ARDT is used, with the TECA-BARD v1.0 giving a Taylor skill score of approximately 0.87 and ARCONNECT v2 giving a skill score of only about 0.32. Comparing between generations, the CMIP6-CM6A-LR simulation has systematically higher Taylor skill scores than either of the CMIP5 IPSL simulations. Further, the CMIP6-CM6A-LR simulation forms a distinct cluster in the Taylor diagram, suggesting a consensus among ARDTs that the CMIP6 version of the IPSL model is superior to the CMIP5 versions.

The ARDTs exhibit distinctive differences in model evaluation. Metrics calculated with the TECA-BARD v1.0 ARDT (star markers in Fig. 10c) are systematically higher than any other ARDT, and most models evaluated by TECA-BARD v1.0 appear skillful at simulating AR-precipitation. The notable exceptions are the CMIP6 MRI-ESM2-0 and CMIP5 CCSM4 simulations, which—as noted previously—have low metric scores no matter which ARDT is used, which is due to a systematic low bias in AR-precipitation in the simulations (e.g., Fig. 10b). Other ARDTs, such as ARCONNECT v2, have a wide spread in the AR-precipitation metrics.

These differences among ARDTs are partly related to their designs. ARCONNECT v2 utilizes an absolute threshold in IVT when identifying ARs, which would make the ARDT much more sensitive to biases in model humidity and/or winds. If a simulation has a systematic low bias in IVT, for example, then the ARCONNECT v2 ARDT will detect systematically fewer ARs in that simulation. Other ARDTs, such as Lora v2 and TECA-BARD v1.0, utilize relative IVT thresholds, which may be less sensitive to model bias.

## 6. Discussion and summary

With a primary goal of introducing a suite of exploratory precipitation metrics and demonstrating their use in evaluating precipitation in climate models, we minimized the hurdle by allowing different groups to apply their diagnostics and metrics to readily available model outputs using their preferred or readily available benchmark datasets. Although most of the metrics were applied to CMIP6 simulations including HighResMIP, the number of models evaluated ranges between 4 and 35. Because feature tracking generally requires more variables and higher temporal frequency data, the LPS, MCS, FRT, and AR metrics were demonstrated using only 4–9 simulations. Although all other metrics were applied to a much larger number of CMIP6 simulations (17–35), differences in the specific simulations used and whether a single or multiple members of a model family were used make comparison across models and metrics difficult.

Despite the difficulty in drawing broad conclusions, some general observations can be made for each metric and by comparing across models and metrics. For precipitation diurnal cycle, models generally perform much better over ocean than over land, as models have a tendency to produce peak precipitation in the afternoon over land while the observed peak precipitation occurs in the late afternoon/early evening. There

is a relatively strong negative intermodel correlation between biases in the diurnal amplitude and phase over ocean but such correlation is positive and weaker over land. Almost all the examined models fail to capture the nocturnal peak observed at the ARM SGP site. For precipitation and dry spells, models perform well in simulating the spatial pattern of both daily precipitation and duration of dry-spell cutoff scales, which means that models would also do well in simulating the spatial distribution of extremes. However, there is a larger spread in terms of scaling factor (i.e., the overall magnitude of the patterns), with the daily precipitation cutoff scale closer to observations than the dry spell duration cutoff scale. Pattern correlation and scaling factor are largely independent metrics as their intermodel correlations are relatively low. In contrast with the precipitation diurnal cycle, spectral analysis shows that models perform better over land than ocean (between 30°S and 30°N) and better over the NH midlatitudes (30°–60°N) than the tropics (15°S–15°N). The majority of the models analyzed have their spectra overlapping with observations by more than 60% in all of the regions and seasons, but the metrics from the models nearly all lie outside the spread of the observation datasets used. The temporal and spatial coherence analysis highlights that the CMIP6 models generally produce precipitation features that are too large and that last too long, particularly in the tropical oceans. Despite these general tendencies, models have a wide range of abilities, with some producing good spatial and temporal variability while others perform poorly at both. There are stronger negative biases over land than over ocean, indicating that models show little spatial variability in temporal coherence over land and hence cannot distinguish regions dominated by longer-lived rain-bearing systems from regions dominated by shorter-lived systems. In the tropical Indian Ocean there are some relationships between the precipitation coherence and MJO metrics (Maritime Continent propagation).

For the process-oriented metrics, coupling of rainfall tendencies and CSF tendencies over the Indo-Pacific warm pool (Fig. 5) is well simulated in 5 of the 20 models analyzed for that metric, and poorly simulated in 8 models; the remaining models with neutral skill may either overestimate or underestimate the rainfall-moisture “rotation” metric derived from this diagnostic. While the rotation metric is modestly correlated with the MJO pattern correlation metric ( $r = 0.41$ ), several models may perform well in one metric, but poorly in another, indicating that rainfall-moisture coupling alone is not a good predictor of a model’s ability to simulate the MJO. For the temperature–water vapor environment, almost half of the models produce most of their precipitation over tropical oceans in a temperature–moisture environment that is reasonably close to the observed range (using twice the distance between the two observational estimates as the reference range). This reflects that the deep-convective parameterizations in these models have included a substantial dependence of convective updrafts on lower free-tropospheric humidity (Kuo et al. 2017). Such a precipitation–temperature–water vapor relationship, however, is not perfectly aligned with other metrics related to precipitation and atmospheric moisture, as will be discussed further below.

In the category of phenomena-based metrics, all HighResMIP models examined here simulated synoptic-scale vortices (i.e., LPS) over South Asia with the qualitatively correct spatial structure of rainfall, with no improvement in model skill at finer horizontal resolution in the two models for which low- and high-resolution versions were examined. This contrasts with prior studies that found LPS were simulated more accurately at finer resolutions; different result may be due to use of a range of coarser resolutions than examined here (Praveen et al. 2015) or the use of only one model (Sabin et al. 2013). In contrast to the general skill in simulating the spatial structure of precipitation within LPS, models exhibited a wide range of biases in representing the amplitude of LPS precipitation, with the three main models examined showing large negative bias, large positive bias, and low bias, with the bias magnitude changing little or, unexpectedly, even degrading at finer resolution. For MCS metrics, the four HighResMIP models evaluated show varying skill in reproducing the observed composited MCS rainfall in the central United States, with model ability to simulate intense convective precipitation a distinguishing factor. Skill scores are worse in summer than spring in all models, consistent with the more dominating frontal large-scale environments of MCS in spring, which are more skillfully simulated by global models (Song et al. 2019). The precipitation error decomposition into frontal precipitation frequency and intensity indicates that all the models evaluated have compensating biases. They produce frontal precipitation too frequently, with intensity that is too low. This is consistent with the results from CMIP5 in Catto et al. (2015), although the CMIP6 models so far seem to have smaller errors. The total precipitation coming from fronts is well represented in the models, including the spatial patterns, indicating good representation of fronts themselves. For the AR precipitation metric, ERA-20C has a Taylor skill score of 0.96 relative to MERRA-2 when assessed using the TECA\_BARD\_ARDT AR tracking method, which provides a measure of observational uncertainty in the metric. Considering the inter-ARDT spread in the Taylor skill score, no models perform well in simulating AR precipitation as none is within one standard deviation of ERA 20C score.

As our diagnostic analysis has been summarized succinctly using scalar metrics, meta-analysis of model skill can be facilitated by developing a matrix of skill scores for models versus metrics to reveal possible relationships among metrics and models. Comparing across metrics and models, it is clear that model skill varies substantially. To help reveal potential relationships among metrics and models, we identified the top-5 and bottom-5 simulations evaluated by each category of metrics (e.g., diurnal precipitation) and its subcategories (e.g., amplitude and phase of diurnal precipitation). The results of this relative model ranking are not shown, as we focus on insights that can be gained from the comparative analysis rather than highlighting the performance of specific models. Consistent with the diverse model skill exhibited across metrics and models, only two model families are in the top-5 group for more than three different categories of metrics and are not in the bottom-5 group in any metrics. Similarly, only one model family is in the bottom-5 group for more than three categories of metrics and is not in the top-5 group in any metrics. Many models perform well in some metrics but

poorly in other metrics. There is a general tendency for simulations produced by the same model family but using different resolutions, model versions, or model configurations, to perform similarly, although some exceptions can also be found.

Focusing on the actual model skill for each metric, we also identified the good and poor performing models in an absolute sense to determine how well models perform for each metric, and subsequently ranked the metrics according to those in which most models performed well or poorly. This absolute skill and ranking was determined by the developers of each metric based on their own judgement, which generally involved comparing model skill relative to some uncertainty related to observation data, and for ARs, uncertainty in tracking methods is also considered. A few metrics that stand out with more models performing well and poorly are highlighted here. Notably, more than 50% of the simulations evaluated based on the diurnal amplitude and phase of precipitation over ocean are considered skillful, while the same is true for the evaluation of spectral characteristics over land and the NH midlatitudes, and for the scaling factor of daily precipitation and the pattern correlation of the cutoff scale between the simulated and observed duration of dry spells. In contrast, two metrics stand out as more challenging for models, with more than 50% of the simulations considered to be performing poorly. These are correlation coefficients of the Taylor skill score for spatial coherence over both land and ocean and the AR precipitation Taylor skill score. Last, more than 50% of the simulations are considered neutral (neither skillful nor poor) with respect to several metrics including diurnal amplitude and phase over land; spectral analysis over ocean, tropics, and SH midlatitudes; and MJO pattern correlation. For other metrics, models are more mixed in how well they represent the specific precipitation characteristics evaluated.

Based on the relative and absolute ranking, additional insights can be gained with regard to the potential relationships among the metrics by calculating the correlation coefficients between the model ranking based on different metrics for the overlapping models, although not all metrics should be connected (e.g., due to geographical differences). For illustrative purposes, we calculated the correlation coefficients between the model ranking based on the temperature–water vapor environment and the model ranking based on other metrics for the overlapping models. We found relatively strong correlations ( $r > 0.5$ ) of model skill in temperature–water vapor environment with model skill in precipitation cutoff scale (both pattern correlation and scaling factor), spectral analysis, temporal and spatial coherence, and MJO propagation based on the relative ranking. On the other hand, model skill in temperature–water vapor environment has very low ( $r < 0.2$ ) or negative correlations with model skill in diurnal precipitation over land (both amplitude and phase), diurnal precipitation over ocean (amplitude only), and dry spell cutoff scale (pattern correlation). Notably, correlations with the rotation metric and MJO east/west power ratio metric are also rather low ( $< 0.3$ ).

The above analysis is suggestive of some predictive power of the model skill in the temperature–water vapor environment on the model skill in several other precipitation characteristics.



This motivates future work to understand these relationships by performing additional diagnostic analysis, and also to apply the exploratory metrics more systematically to the same set of model simulations using comparable benchmark datasets in order to support quantitative analysis of skill across models and metrics. This may reveal less obvious relationships among metrics and models, reflecting relationships among processes and/or weather phenomena highlighted by the metrics, or relationships among models due to commonality such as parameterization schemes. Such information is useful for guiding model development and model tuning. Machine learning approaches may be used to develop predictive models of the relationships among the different metrics presented here, or between those metrics and others such as metrics for the modes of climate variability (e.g., MJO, ENSO), circulation indices (e.g., monsoon), sea surface temperature pattern, etc. Such mapping of model skill scores across metrics may help focus efforts on improving model prediction skill given the important role of modes of variability in predicting precipitation at various time scales.

Going beyond baseline metrics that evaluate basic precipitation features, data requirements are an obstacle for systematic application of the metric suite because high-temporal-frequency data and certain variables (e.g., outgoing longwave radiation) are not commonly available from the CMIP data archive. Table 2 is a good starting point for expansion in the future when more exploratory metrics are added. Communicating the data requirements to community efforts such as CMIP and demonstrating the usefulness of the exploratory metrics are both important for increasing awareness of, and advocacy for, the data needs of model evaluation and diagnostics to support the broad use of climate model output.

While the metrics described in this study are useful individually, combining or connecting them may potentially provide more powerful metrics to benchmark models as well as revealing the underlying reasons or sources of the model biases. At the same time, decomposition of the metrics into independent components is useful for attributing model biases to multiple factors. Future work to standardize the metrics, addressing uncertainties in observation data and tracking methods, and improving interpretations of the metrics, may facilitate more robust use of the exploratory metrics. There may also be a need to reconcile features attributed to different phenomena simultaneously. For example, precipitation from MCSs embedded within frontal systems could potentially be attributed to both MCS and frontal precipitation (e.g., Dowdy and Catto 2017; Catto and Dowdy 2021). In regions such as the Bay of Bengal where LPS and MCS are both prominent, it is not clear if certain precipitation events could be attributed simultaneously to LPS and MCS and what implications this may have on metrics built upon these phenomena. Coding and software aspects may also require some attention in the future to facilitate implementation of the exploratory metrics in community packages for broader adoption and use.

Through this study, we have developed methodologies and analysis codes to calculate metrics and track weather phenomena. Applying them to the CMIP6 output and observation data has generated intermediate quantities and datasets such as tracks of LPS, MCS, fronts, and AR

and associated precipitation and large-scale environments. These datasets are useful not only for model evaluation but also for scientific investigations. For example, datasets derived from the historical simulations could be combined with similar datasets for simulations of the future climate to investigate the response of various precipitation metrics to radiative forcing. Different metrics may also be combined to understand the connections between different weather phenomena and storm types and their connections to the temperature–water vapor environments and modes of variability. Among the metrics described in this study, the spectral and coherence metrics have already been included in ASoP, and some atmospheric river tracking algorithms are available from Coordinated Model Evaluation Capabilities (CMEC). Efforts are ongoing to coordinate the development and implementation of metrics to be incorporated in community diagnostic packages to facilitate broader use to improve quantification and understanding of precipitation biases in weather and climate models.

*Acknowledgments.* This study represents a collaborative effort as an outgrowth of a workshop on “Benchmarking Simulated Precipitation in Earth System Models” sponsored by the Office of Science of the U.S. Department of Energy (DOE) Biological and Environmental Research through the Regional and Global Model Analysis (RGMA) program area. RGMA also supported Leung and Feng under the WACCEM scientific focus area, O’Brien and Zhou under the CASCADE scientific focus area, Boos and Vishnu under Award DE-SC0019367, DeMott under Award DE-SC0020092, and Klingaman and Lee under Award DE-SC0020324. O’Brien’s efforts were also partially supported by the Environmental Resilience Institute, funded by Indiana University’s Prepared for Environmental Change Grand Challenge initiative. Neelin, Kuo, and Martinez-Villalobos were supported by National Science Foundation Grant AGS-1936810 and National Oceanic and Atmospheric Administration Grants NA18OAR4310280 and NA21OAR 4310354. Martinez-Villalobos was also supported by Proyecto Corfo Ingeniería 2030 código 14ENI2-26865. Catto and Priestley were supported by the UK Natural Environment Research Council Grant NE/S004645/1. Work at LLNL was supported by the DOE Office of Science Biological and Environmental Research through the Earth System Model Development program area and the Atmospheric Radiation Measurement program, and performed under the auspices of the U.S. DOE by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344. Pacific Northwest National Laboratory is operated for the Department of Energy by Battelle Memorial Institute under Contract DE-AC05-76RL01830. Martin was supported by the U.K.–China Research and Innovation Partnership Fund through the Met Office Climate Science for Service Partnership (CSSP) China, as part of the Newton Fund, and by the Weather and Climate Science for Service Partnership (WCSSP) India, a collaborative initiative between the Met Office, supported by the U.K. Government’s Newton Fund, and the Indian Ministry of Earth Sciences (MoES). This research used resources of the National Energy Research

Scientific Computing Center (NERSC), also supported by the Office of Science of the U.S. Department of Energy under Contract DE-AC02-05CH11231. We acknowledge the World Climate Research Programme, which, through its Working Group on Coupled Modelling, coordinated and promoted CMIP6. We thank the climate modeling groups for producing and making available their model output, the Earth System Grid Federation (ESGF) for archiving the data and providing access, and the multiple funding agencies who support CMIP6 and ESGF. We thank DOE's RGMA program area, the Data Management program, and NERSC for making this coordinated CMIP6 analysis activity possible.

## REFERENCES

- Adames, Á. F., and D. Kim, 2016: The MJO as a dispersive, convectively coupled moisture wave: Theory and observations. *J. Atmos. Sci.*, **73**, 913–941, <https://doi.org/10.1175/JAS-D-15-0170.1>.
- , and Y. Ming, 2018: Interactions between water vapor and potential vorticity in synoptic-scale monsoonal disturbances: Moisture vortex instability. *J. Atmos. Sci.*, **75**, 2083–2106, <https://doi.org/10.1175/JAS-D-17-0310.1>.
- Ahmed, F., Á. F. Adames, and J. D. Neelin, 2020: Deep convective adjustment of temperature and moisture. *J. Atmos. Sci.*, **77**, 2163–2186, <https://doi.org/10.1175/JAS-D-19-0227.1>.
- Ahn, M.-S., and Coauthors, 2017: MJO simulation in CMIP5 climate models: MJO skill metrics and process-oriented diagnosis. *Climate Dyn.*, **49**, 4023–4045, <https://doi.org/10.1007/s00382-017-3558-4>.
- , D. Kim, D. Kang, J. Lee, K. R. Sperber, P. J. Gleckler, X. Jiang, H. Yoo-Geun, and H. Kim, 2020: MJO propagation across the Maritime Continent: Are CMIP6 models better than CMIP5 models? *Geophys. Res. Lett.*, **47**, e2020GL087250, <https://doi.org/10.1029/2020GL087250>.
- Ashouri, H., K.-L. Hsu, S. Sorooshian, D. K. Braithwaite, K. R. Knapp, L. D. Cecil, B. R. Nelson, and O. P. Prat, 2015: PERSIANN-CDR: Daily precipitation climate data record from multisatellite observations for hydrological and climate studies. *Bull. Amer. Meteor. Soc.*, **96**, 69–83, <https://doi.org/10.1175/BAMS-D-13-00068.1>.
- Berry, G., M. J. Reeder, and C. Jakob, 2011: A global climatology of atmospheric fronts. *Geophys. Res. Lett.*, **38**, L04809, <https://doi.org/10.1029/2010GL046451>.
- Bretherton, C. S., M. E. Peters, and L. E. Back, 2004: Relationships between water vapor path and precipitation over the tropical oceans. *J. Climate*, **17**, 1517–1528, [https://doi.org/10.1175/1520-0442\(2004\)017<1517:RBWVPA>2.0.CO;2](https://doi.org/10.1175/1520-0442(2004)017<1517:RBWVPA>2.0.CO;2).
- Caldwell, P. M., and Coauthors, 2019: The DOE E3SM coupled model version 1: Description and results at high resolution. *J. Adv. Model. Earth Syst.*, **11**, 4095–4146, <https://doi.org/10.1029/2019MS001870>.
- Catto, J. L., and S. Pfahl, 2013: The importance of fronts for extreme precipitation. *J. Geophys. Res. Atmos.*, **118**, 10791, <https://doi.org/10.1002/jgrd.50852>.
- , C. Jakob, and N. Nicholls, 2015: Can the CMIP5 models represent winter frontal precipitation? *Geophys. Res. Lett.*, **42**, 8596–8604, <https://doi.org/10.1002/2015GL066015>.
- , and A. J. Dowdy, 2021: Understanding compound hazards from a weather system perspective. *Wea. Climate Extremes*, **32**, 100313, <https://doi.org/10.1016/j.wace.2021.100313>.
- Chang, M., B. Liu, C. Martinez-Villalobos, G. Ren, S. Li, and T. Zhou, 2020: Changes in extreme precipitation accumulations during the warm season over continental China. *J. Climate*, **33**, 10799–10811, <https://doi.org/10.1175/JCLI-D-20-0616.1>.
- Chen, D., and A. Dai, 2018: Dependence of estimated precipitation frequency and intensity on data resolution. *Climate Dyn.*, **50**, 3625–3647, <https://doi.org/10.1007/s00382-017-3830-7>.
- , and —, 2019: Precipitation characteristics in the Community Atmosphere Model and their dependence on model physics and resolution. *J. Adv. Model. Earth Syst.*, **11**, 2352–2374, <https://doi.org/10.1029/2018MS001536>.
- , —, and A. Hall, 2021: Precipitation partitioning and the “drizzling” bias in CMIP5 models. *J. Geophys. Res. Atmos.*, **126**, e2020JD034198, <https://doi.org/10.1029/2020JD034198>.
- Chen, J., A. Dai, and Y. Zhang, 2020: Linkage between projected precipitation and atmospheric thermodynamic changes. *J. Climate*, **33**, 7155–7178, <https://doi.org/10.1175/JCLI-D-19-0785.1>.
- Covey, C., and P. Gleckler, 2014: Standard diagnostics for the diurnal cycle of precipitation. Lawrence Livermore National Laboratory Tech. Rep. LLNL-TR-659685, 11 pp., <https://www.osti.gov/servlets/purl/1165787>.
- , —, C. Doutriaux, D. N. Williams, A. Dai, J. Fasullo, K. Trenberth, and A. Berg, 2016: Metrics for the diurnal cycle of precipitation: Toward routine benchmarks for climate models. *J. Climate*, **29**, 4461–4471, <https://doi.org/10.1175/JCLI-D-15-0664.1>.
- Dai, A., 2001: Global precipitation and thunderstorm frequencies. Part II: Diurnal variations. *J. Climate*, **14**, 1112–1128, [https://doi.org/10.1175/1520-0442\(2001\)014<1112:GPATFP>2.0.CO;2](https://doi.org/10.1175/1520-0442(2001)014<1112:GPATFP>2.0.CO;2).
- , 2006: Precipitation characteristics in eighteen coupled climate models. *J. Climate*, **19**, 4605–4630, <https://doi.org/10.1175/JCLI3884.1>.
- , F. Giorgi, and K. E. Trenberth, 1999: Observed and model simulated diurnal cycles of precipitation over the contiguous United States. *J. Geophys. Res.*, **104**, 6377–6402, <https://doi.org/10.1029/98JD02720>.
- , X. Lin, and K.-L. Hsu, 2007: The frequency, intensity, and diurnal cycle of precipitation in surface and satellite observations over low- and mid-latitudes. *Climate Dyn.*, **29**, 727–744, <https://doi.org/10.1007/s00382-007-0260-y>.
- Dee, D. P., and Coauthors, 2011: The ERA-Interim reanalysis: Configuration and performance of the data assimilation system. *Quart. J. Roy. Meteor. Soc.*, **137**, 553–597, <https://doi.org/10.1002/qj.828>.
- DeMott, C. A., N. P. Klingaman, W. L. Tseng, M. A. Burt, Y. Gao, and D. A. Randall, 2019: The convection connection: How ocean feedbacks affect tropical mean moisture and MJO propagation. *J. Geophys. Res. Atmos.*, **124**, 11910–11931, <https://doi.org/10.1029/2019JD031015>.
- Deser, C., A. Phillips, V. Bourdette, and H. Teng, 2012: Uncertainty in climate change projections: The role of internal variability. *Climate Dyn.*, **38**, 527–546, <https://doi.org/10.1007/s00382-010-0977-x>.
- Dettinger, M., 2011: Climate change, atmospheric rivers, and floods in California—A multimodel analysis of storm frequency and magnitude changes. *J. Amer. Water Resour. Assoc.*, **47**, 514–523, <https://doi.org/10.1111/j.1752-1688.2011.00546.x>.
- Diaz, M., and W. R. Boos, 2019: Monsoon depression amplification by moist barotropic instability in a vertically sheared

- environment. *Quart. J. Roy. Meteor. Soc.*, **145**, 2666–2684, <https://doi.org/10.1002/qj.3585>.
- , and —, 2021: The influence of surface heat fluxes on the growth of idealized monsoon depressions. *J. Atmos. Sci.*, **78**, 2013–2027, <https://doi.org/10.1175/JAS-D-20-0359.1>.
- Ditchek, S. D., W. R. Boos, S. J. Camargo, and M. K. Tippett, 2016: A genesis index for monsoon disturbances. *J. Climate*, **29**, 5189–5203, <https://doi.org/10.1175/JCLI-D-15-0704.1>.
- Dowdy, A., and J. L. Catto, 2017: Extreme weather caused by concurrent cyclone, front and thunderstorm occurrences. *Sci. Rep.*, **7**, 40359, <https://doi.org/10.1038/srep40359>.
- Eyring, V., S. Bony, G. A. Meehl, C. A. Senior, B. Stevens, R. J. Stouffer, and K. E. Taylor, 2016: Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design and organization. *Geosci. Model Dev.*, **9**, 1937–1958, <https://doi.org/10.5194/gmd-9-1937-2016>.
- , and Coauthors, 2020: Earth System Model Evaluation Tool (ESMValTool) v2.0—An extended set of large-scale diagnostics for quasi-operational and comprehensive evaluation of Earth system models in CMIP. *Geosci. Model Dev.*, **13**, 3383–3438, <https://doi.org/10.5194/gmd-13-3383-2020>.
- Feng, Z., L. R. Leung, S. Hagos, R. A. Houze, C. D. Burleyson, and K. Balaguru, 2016: More frequent intense and long-lived storms dominate the trend in central U.S. rainfall. *Nat. Commun.*, **7**, 13429, <https://doi.org/10.1038/ncomms13429>.
- , —, R. A. Houze Jr., S. Hagos, J. Hardin, Q. Yang, B. Han, and J. Fan, 2018: Structure and evolution of mesoscale convective systems: Sensitivity to cloud microphysics in convection-permitting simulations over the United States. *J. Adv. Model. Earth Syst.*, **10**, 1470–1494, <https://doi.org/10.1029/2018MS001305>.
- , R. A. Houze Jr., L. R. Leung, F. Song, J. Hardin, J. Wang, W. Gustafson Jr., and C. Homeyer, 2019: Spatiotemporal characteristics and large-scale environment of mesoscale convective systems east of the Rocky Mountains. *J. Climate*, **32**, 7303–7328, <https://doi.org/10.1175/JCLI-D-19-0137.1>.
- , F. Song, K. Sakaguchi, and L. R. Leung, 2021a: Evaluation of mesoscale convective systems in climate simulations: Methodological development and results from MPAS-CAM over the United States. *J. Climate*, **34**, 2611–2633, <https://doi.org/10.1175/JCLI-D-20-0136.1>.
- , L. R. Leung, N. Liu, J. Wang, R. A. Houze Jr., J. Li, J. C. Hardin, and J. Guo, 2021b: A global high-resolution mesoscale convective system database using satellite-derived cloud tops, surface precipitation, and tracking. *J. Geophys. Res. Atmos.*, **126**, <https://doi.org/10.1029/2020JD034202>.
- Fujinami, H., H. Hirata, M. Kato, and K. Tsuboki, 2020: Mesoscale precipitation systems and their role in the rapid development of a monsoon depression over the Bay of Bengal. *Quart. J. Roy. Meteor. Soc.*, **146**, 267–283, <https://doi.org/10.1002/qj.3672>.
- Gelaro, R., and Coauthors, 2017: The Modern-Era Retrospective Analysis for Research and Applications, version 2 (MERRA-2). *J. Climate*, **30**, 5419–5454, <https://doi.org/10.1175/JCLI-D-16-0758.1>.
- Gimeno, L., R. Nieto, M. Vázquez, and D. A. Lavers, 2014: Atmospheric rivers: A mini-review. *Front. Earth Sci.*, **2**, 1–6, <https://doi.org/10.3389/feart.2014.00002>.
- Gleckler, P. J., C. Doutriaux, P. J. Durack, K. E. Taylor, Y. Zhang, D. N. Williams, E. Mason, and J. Servonnat, 2016: A more powerful reality test for climate models. *Eos*, **97**, <https://doi.org/10.1029/2016EO051663>.
- Goldenson, N., L. R. Leung, C. M. Bitz, and E. Blanchard-Wrigglesworth, 2018: Influence of atmospheric rivers on mountain snowpack in the western United States. *J. Climate*, **31**, 9921–9940, <https://doi.org/10.1175/JCLI-D-18-0268.1>.
- Guan, B., and D. E. Waliser, 2015: Detection of atmospheric rivers: Evaluation and application of an algorithm for global studies. *J. Geophys. Res. Atmos.*, **120**, 12514–12535, <https://doi.org/10.1002/2015JD024257>.
- , —, and F. M. Ralph, 2018: An intercomparison between reanalysis and dropsonde observations of the total water vapor transport in individual atmospheric rivers. *J. Hydrometeorol.*, **19**, 321–337, <https://doi.org/10.1175/JHM-D-17-0114.1>.
- Haarsma, R. J., and Coauthors, 2016: High Resolution Model Intercomparison Project (HighResMIP v1.0) for CMIP6. *Geosci. Model Dev.*, **9**, 4185–4208, <https://doi.org/10.5194/gmd-9-4185-2016>.
- Henderson, S. A., E. D. Maloney, and S. W. Son, 2017: Madden-Julian oscillation Pacific teleconnections: The impact of the basic state and MJO representation in general circulation models. *J. Climate*, **30**, 4567–4587, <https://doi.org/10.1175/JCLI-D-16-0789.1>.
- Hersbach, H., and Coauthors, 2020: The ERA5 global reanalysis. *Quart. J. Roy. Meteor. Soc.*, **146**, 1999–2049, <https://doi.org/10.1002/qj.3803>.
- Hewson, T. D., 1998: Objective fronts. *Meteor. Appl.*, **5**, 37–65, <https://doi.org/10.1017/S1350482798000553>.
- Hirota, H., Y. N. Takayabu, M. Watanabe, M. Kimoto, and M. Chikira, 2014: Role of convective entrainment in spatial distributions of and temporal variations in precipitation over tropical oceans. *J. Climate*, **27**, 8707–8723, <https://doi.org/10.1175/JCLI-D-13-00701.1>.
- Hoffmann, L., and Coauthors, 2019: From ERA-Interim to ERA5: The considerable impact of ECMWF's next-generation reanalysis on Lagrangian transport simulations. *Atmos. Chem. Phys.*, **19**, 3097–3124, <https://doi.org/10.5194/acp-19-3097-2019>.
- Huffman, G. J., and Coauthors, 2007: The TRMM Multisatellite Precipitation Analysis (TMPA): Quasi-global, multiyear, combined-sensor precipitation estimates at fine scales. *J. Hydrometeorol.*, **8**, 38–55, <https://doi.org/10.1175/JHM560.1>.
- , R. F. Adler, A. Behrangi, D. T. Bolvin, E. Nelkin, and Y. Song, 2020: Algorithm Theoretical Basis Document (ATBD) for Global Precipitation Climatology Project version 3.1 precipitation data, 32 pp., [https://docserv.gesdisc.eosdis.nasa.gov/public/project/MEaSURES/GPCP/GPCP\\_ATBD\\_V3.1.pdf](https://docserv.gesdisc.eosdis.nasa.gov/public/project/MEaSURES/GPCP/GPCP_ATBD_V3.1.pdf).
- Hwang, Y.-T., and D. M. W. Frierson, 2013: Link between the double-intertropical convergence zone problem and cloud biases over the Southern Ocean. *Proc. Natl. Acad. Sci. USA*, **110**, 4935–4940, <https://doi.org/10.1073/pnas.1213302110>.
- Janowiak, J., B. Joyce, and P. Xie, 2017: NCEP/CPC L3 half hourly 4 km global (60°S–60°N) merged IR V1, accessed 1 March 2021, <https://doi.org/10.5067/P4HZB9N27EKU>.
- Jiang, X., and Coauthors, 2015: Vertical structure and physical processes of the Madden-Julian oscillation: Exploring key model physics in climate simulations. *J. Geophys. Res. Atmos.*, **120**, 4718–4748, <https://doi.org/10.1002/2014JD022375>.
- Joyce, R. J., and P. Xie, 2011: Kalman filter-based CMORPH. *J. Hydrometeorol.*, **12**, 1547–1563, <https://doi.org/10.1175/JHM-D-11-022.1>.
- Klingaman, N. P., G. M. Martin, and A. F. Moise, 2017: ASoP (v1.0): A set of methods for analyzing scales of precipitation in general circulation models. *Geosci. Model Dev.*, **10**, 57–83, <https://doi.org/10.5194/gmd-10-57-2017>.

- Krishnamurthy, V., and R. S. Ajayamohan, 2010: Composite structure of monsoon low pressure systems and its relation to Indian rainfall. *J. Climate*, **23**, 4285–4305, <https://doi.org/10.1175/2010JCLI2953.1>.
- Kuo, Y.-H., J. D. Neelin, and C. R. Mechoso, 2017: Tropical convective transition statistics and causality in the water vapor–precipitation relation. *J. Atmos. Sci.*, **74**, 915–931, <https://doi.org/10.1175/JAS-D-16-0182.1>.
- , K. A. Schiro, and J. D. Neelin, 2018: Convective transition statistics over tropical oceans for climate model diagnostics: Observational baseline. *J. Atmos. Sci.*, **75**, 1553–1570, <https://doi.org/10.1175/JAS-D-17-0287.1>.
- , and Coauthors, 2020: Convective transition statistics over tropical oceans for climate model diagnostics: GCM evaluation. *J. Atmos. Sci.*, **77**, 379–403, <https://doi.org/10.1175/JAS-D-19-0132.1>.
- Lin, J.-L., 2007: The double-ITCZ problem in IPCC AR4 coupled GCMs: Ocean–atmosphere feedback analysis. *J. Climate*, **20**, 4497–4525, <https://doi.org/10.1175/JCLI4272.1>.
- Lin, Y., W. Dong, M. Zhang, Y. Xie, W. Xue, J. Huang, and Y. Luo, 2017: Causes of model dry and warm bias over central U.S. and impact on climate projections. *Nat. Commun.*, **8**, 881, <https://doi.org/10.1038/s41467-017-01040-2>.
- Ma, H.-Y., S. Xie, J. S. Boyle, S. A. Klein, and Y. Zhang, 2013: Metrics and diagnostics for precipitation-related processes in climate model short-range hindcasts. *J. Climate*, **26**, 1516–1534, <https://doi.org/10.1175/JCLI-D-12-00235.1>.
- Mapes, B., and R. Neale, 2011: Parameterizing convective organization to escape the entrainment dilemma. *J. Adv. Model. Earth Syst.*, **3**, M06004, <https://doi.org/10.1029/2011MS000042>.
- Martin, G. M., N. P. Klingaman, and A. F. Moise, 2017: Connecting spatial and temporal scales of tropical precipitation in observations and the MetUM-GA6. *Geosci. Model Dev.*, **10**, 105–126, <https://doi.org/10.5194/gmd-10-105-2017>.
- Martinez-Villalobos, C., and J. D. Neelin, 2018: Shifts in precipitation accumulation extremes during the warm season over the United States. *Geophys. Res. Lett.*, **45**, 8586–8595, <https://doi.org/10.1029/2018GL078465>.
- , and —, 2019: Why do precipitation intensities tend to follow gamma distributions? *J. Atmos. Sci.*, **76**, 3611–3631, <https://doi.org/10.1175/JAS-D-18-0343.1>.
- , and —, 2021: Climate models capture key features of extreme precipitation probabilities across regions. *Environ. Res. Lett.*, **16**, 024017, <https://doi.org/10.1088/1748-9326/abd351>.
- McClenny, E. E., P. A. Ullrich, and R. Grotjahn, 2020: Sensitivity of atmospheric river vapor transport and precipitation to uniform sea-surface temperature increases. *J. Geophys. Res. Atmos.*, **21**, e2020JD033421, <https://doi.org/10.1029/2020JD033421>.
- Mechoso, C. R., and Coauthors, 1995: The seasonal cycle over the tropical Pacific in coupled ocean–atmosphere general circulation models. *Mon. Wea. Rev.*, **123**, 2825–2838, [https://doi.org/10.1175/1520-0493\(1995\)123<2825:TSCOTT>2.0.CO;2](https://doi.org/10.1175/1520-0493(1995)123<2825:TSCOTT>2.0.CO;2).
- Mehran, A., A. AghaKouchak, and T. J. Phillips, 2014: Evaluation of CMIP5 continental precipitation simulations relative to satellite-based gauge-adjusted observations. *J. Geophys. Res. Atmos.*, **119**, 1695–1707, <https://doi.org/10.1002/2013JD021152>.
- Mejia, J. F., D. Koraćin, and E. M. Wilcox, 2018: Effect of coupled global climate models sea surface temperature biases on simulated climate of the western United States. *Int. J. Climatol.*, **38**, 5386–5404, <https://doi.org/10.1002/joc.5817>.
- Mundhenk, B. D., E. A. Barnes, and E. D. Maloney, 2016: All-season climatology and variability of atmospheric river frequencies over the North Pacific. *J. Climate*, **29**, 4885–4903, <https://doi.org/10.1175/JCLI-D-15-0655.1>.
- Murakami, H., 2014: Tropical cyclones in reanalysis data sets. *Geophys. Res. Lett.*, **41**, 2133–2141, <https://doi.org/10.1002/2014GL059519>.
- Murthy, V. S., and W. R. Boos, 2020: Quasigeostrophic controls on precipitating ascent in monsoon depressions. *J. Atmos. Sci.*, **77**, 1213–1232, <https://doi.org/10.1175/JAS-D-19-0202.1>.
- Neelin, J. D., O. Peters, J. W. B. Lin, K. Hales, and C. E. Holloway, 2008: Rethinking convective quasi-equilibrium: Observational constraints for stochastic convective schemes in climate models. *Philos. Trans. Roy. Soc.*, **366A**, 2579–2602, <https://doi.org/10.1098/rsta.2008.0056>.
- , —, and K. Hales, 2009: The transition to strong convection. *J. Atmos. Sci.*, **66**, 2367–2384, <https://doi.org/10.1175/2009JAS2962.1>.
- , S. Sahany, S. N. Stechmann, and D. N. Bernstein, 2017: Global warming precipitation accumulation increases above the current-climate cutoff scale. *Proc. Natl. Acad. Sci. USA*, **114**, 1258–1263, <https://doi.org/10.1073/pnas.1615333114>.
- Nesbitt, S. W., R. Cifelli, and S. A. Rutledge, 2006: Storm morphology and rainfall characteristics of TRMM precipitation features. *Mon. Wea. Rev.*, **134**, 2702–2721, <https://doi.org/10.1175/MWR3200.1>.
- O'Brien, T. A., and Coauthors, 2020a: Detection uncertainty matters for understanding atmospheric rivers. *Bull. Amer. Meteor. Soc.*, **101** (6), E790–E796, <https://doi.org/10.1175/BAMS-D-19-0348.1>.
- , and Coauthors, 2020b: Detection of atmospheric rivers with inline uncertainty quantification: TECA-BARD v1.0.1. *Geosci. Model Dev.*, **13**, 6131–6148, <https://doi.org/10.5194/gmd-13-6131-2020>.
- , and Coauthors, 2022: Increases in future AR count and size: Overview of the ARTMIP Tier 2 CMIP5/6 experiment. *J. Geophys. Res. Atmos.*, **127**, e2021JD036013, <https://doi.org/10.1029/2021JD036013>.
- Oueslati, B., and G. Bellon, 2013: Convective entrainment and large-scale organization of tropical precipitation: Sensitivity of the CNRM-CM5 hierarchy of models. *J. Climate*, **26**, 2931–2946, <https://doi.org/10.1175/JCLI-D-12-00314.1>.
- Payne, A. E., and G. Magnusdottir, 2015: An evaluation of atmospheric rivers over the North Pacific in CMIP5 and their response to warming under RCP 8.5. *J. Geophys. Res. Atmos.*, **120**, 11 173–11 190, <https://doi.org/10.1002/2015JD023586>.
- Pendergrass, A. G., P. J. Gleckler, L. R. Leung, and C. Jakob, 2020: Benchmarking simulated precipitation in Earth system models. *Bull. Amer. Meteor. Soc.*, **101** (6), E814–E816, <https://doi.org/10.1175/BAMS-D-19-0318.1>.
- Perkins, S. E., A. J. Pitman, N. J. Holbrook, and J. McAneney, 2007: Evaluation of the AR4 climate models' simulated daily maximum temperature, minimum temperature, and precipitation over Australia using probability density functions. *J. Climate*, **20**, 4356–4376, <https://doi.org/10.1175/JCLI4253.1>.
- Pierrehumbert, R. T., H. Brogniez, and R. Roca, 2007: On the relative humidity of the Earth's atmosphere. *The Global Circulation of the Atmosphere: Phenomena, Theory, Challenges*, T. Schneider and A. H. Sobel, Eds., Princeton University Press, 143–185.
- Praveen, V., S. Sandeep, and R. S. Ajayamohan, 2015: On the relationship between mean monsoon precipitation and low pressure systems in climate model simulations. *J. Climate*, **28**, 5305–5324, <https://doi.org/10.1175/JCLI-D-14-00415.1>.

- Qian, T., A. Dai, K. E. Trenberth, and K. W. Oleson, 2006: Simulation of global land surface conditions from 1948–2004. Part I: Forcing data and evaluation. *J. Hydrometeorol.*, **7**, 953–975, <https://doi.org/10.1175/JHM540.1>.
- Ralph, F. M., M. D. Dettinger, M. M. Cairns, T. J. Galarneau, and J. Eylander, 2018: Defining “atmospheric river”: How the *Glossary of Meteorology* helped resolve a debate. *Bull. Amer. Meteor. Soc.*, **99**, 837–839, <https://doi.org/10.1175/BAMS-D-17-0157.1>.
- Rao, K. V., and S. Rajamani, 1970: Diagnostic study of a monsoon depression by geostrophic baroclinic model. *MAUSAM*, **21**, 187–194, <https://doi.org/10.54302/mausam.v21i2.5366>.
- Rutz, J. J., W. J. Steenburgh, and F. Martin Ralph, 2014: Climatological characteristics of atmospheric rivers and their inland penetration over the western United States. *Mon. Wea. Rev.*, **142**, 905–921, <https://doi.org/10.1175/MWR-D-13-00168.1>.
- , and Coauthors, 2019: The Atmospheric River Tracking Method Intercomparison Project (ARTMIP): Quantifying uncertainties in atmospheric river climatology. *J. Geophys. Res. Atmos.*, **124**, 13 777–13 802, <https://doi.org/10.1029/2019JD030936>.
- Sabin, P., T. Krishnan, R. Ghattas, S. Denvil, J. L. Dufresne, F. Hourdin, and T. Pascal, 2013: High resolution simulation of the South Asian monsoon using a variable resolution global climate model. *Climate Dyn.*, **41**, 173–194, <https://doi.org/10.1007/s00382-012-1658-8>.
- Saha, S., and Coauthors, 2010: The NCEP Climate Forecast System Reanalysis. *Bull. Amer. Meteor. Soc.*, **91**, 1015–1058, <https://doi.org/10.1175/2010BAMS3001.1>.
- Sanders, F., 1984: Quasi-geostrophic diagnosis of the monsoon depression of 5–8 July 1979. *J. Atmos. Sci.*, **41**, 538–552, [https://doi.org/10.1175/1520-0469\(1984\)041<0538:OGDOTM>2.0.CO;2](https://doi.org/10.1175/1520-0469(1984)041<0538:OGDOTM>2.0.CO;2).
- Sellers, S. L., B. Kawzenuk, P. Nguyen, F. M. Ralph, and S. Sorooshian, 2017: Genesis, pathways, and terminations of intense global water vapor transport in association with large-scale climate patterns. *Geophys. Res. Lett.*, **44**, 12 465–12 475, <https://doi.org/10.1002/2017GL075495>.
- Shields, C. A., and J. T. Kiehl, 2016: Simulating the Pineapple Express in the half degree Community Climate System Model, CCSM4. *Geophys. Res. Lett.*, **43**, 7767–7773, <https://doi.org/10.1002/2016GL069476>.
- , and Coauthors, 2018: Atmospheric River Tracking Method Intercomparison Project (ARTMIP): Project goals and experimental design. *Geosci. Model Dev.*, **11**, 2455–2474, <https://doi.org/10.5194/gmd-11-2455-2018>.
- Sikka, D. R., 1980: Some aspects of the large scale fluctuations of summer monsoon rainfall over India in relation to fluctuations in the planetary and regional scale circulation parameters. *Proc. Indian Acad. Sci. Earth Planet. Sci.*, **89**, 179–195, <https://doi.org/10.1007/BF02913749>.
- Skinner, C. B., J. M. Lora, A. E. Payne, and C. J. Poulsen, 2020: Atmospheric river changes shaped mid-latitude hydroclimate since the mid-Holocene. *Earth Planet. Sci. Lett.*, **541**, 116293, <https://doi.org/10.1016/j.epsl.2020.116293>.
- Song, F., Z. Feng, L. R. Leung, R. A. Houze Jr., J. Wang, J. Hardin, and C. Homeyer, 2019: Contrasting the spring and summer large-scale environments associated with mesoscale convective systems over the U.S. Great Plains. *J. Climate*, **32**, 6749–6767, <https://doi.org/10.1175/JCLI-D-18-0839.1>.
- Sperber, K. R., and D. Kim, 2012: Simplified metrics for the identification of the Madden-Julian oscillation in models. *Atmos. Sci. Lett.*, **13**, 187–193, <https://doi.org/10.1002/asl.378>.
- , H. Annamalai, I.-S. Kang, A. Kitoh, A. Moise, A. Turner, B. Wang, and T. Zhou, 2013: The Asian summer monsoon: An intercomparison of CMIP5 vs. CMIP3 simulations of the late 20th century. *Climate Dyn.*, **41**, 2711–2744, <https://doi.org/10.1007/s00382-012-1607-6>.
- Stan, C., D. M. Straus, J. S. Frederiksen, H. Lin, E. D. Maloney, and C. Schumacher, 2017: Review of tropical–extratropical teleconnections on intraseasonal time scales. *Rev. Geophys.*, **55**, 902–937, <https://doi.org/10.1002/2016RG000538>.
- Stechmann, S. N., and J. D. Neelin, 2014: First-passage-time prototypes for precipitation statistics. *J. Atmos. Sci.*, **71**, 3269–3291, <https://doi.org/10.1175/JAS-D-13-0268.1>.
- Stephens, G. L., and Coauthors, 2010: Dreary state of precipitation in global models. *J. Geophys. Res.*, **115**, D24211, <https://doi.org/10.1029/2010JD014532>.
- Stevenson, S. N., and R. S. Schumacher, 2014: A 10-year survey of extreme rainfall events in the central and eastern United States using gridded multisensor precipitation analyses. *Mon. Wea. Rev.*, **142**, 3147–3162, <https://doi.org/10.1175/MWR-D-13-00345.1>.
- Tan, J., G. J. Huffman, D. T. Bolvin, and E. J. Nelkin, 2019: Diurnal cycle of IMERG V06 precipitation. *Geophys. Res. Lett.*, **46**, 13 584–13 592, <https://doi.org/10.1029/2019GL085395>.
- Tang, S., P. Gleckler, S. Xie, J. Lee, M. S. Ahn, C. Covey, and C. Zhang, 2021: Evaluating the diurnal and semidiurnal cycle of precipitation in CMIP6 models using satellite- and ground-based observations. *J. Climate*, **34**, 3189–3210, <https://doi.org/10.1175/JCLI-D-20-0639.1>.
- Tapiador, F. J., R. Roca, A. Del Genio, B. Dewitte, W. Petersen, and F. Zhang, 2019: Is precipitation a good metric for model performance? *Bull. Amer. Meteor. Soc.*, **100**, 223–233, <https://doi.org/10.1175/BAMS-D-17-0218.1>.
- Taylor, K. E., R. J. Stouffer, and G. A. Meehl, 2012: An overview of CMIP5 and the experiment design. *Bull. Amer. Meteor. Soc.*, **93**, 485–498, <https://doi.org/10.1175/BAMS-D-11-00094.1>.
- Thomas, C. M., and D. M. Schultz, 2019: What are the best thermodynamic quantity and function to define a front in gridded model output? *Bull. Amer. Meteor. Soc.*, **100**, 873–895, <https://doi.org/10.1175/BAMS-D-18-0137.1>.
- Tian, B., 2015: Spread of model climate sensitivity linked to double-Intertropical Convergence Zone bias. *Geophys. Res. Lett.*, **42**, 4133–4141, <https://doi.org/10.1002/2015GL064119>.
- , and X. Dong, 2020: The double-ITCZ bias in CMIP3, CMIP5, and CMIP6 models based on annual mean precipitation. *Geophys. Res. Lett.*, **47**, e2020GL087232, <https://doi.org/10.1029/2020GL087232>.
- Trenberth, K. E., L. Smith, T. Qian, A. Dai, and J. Fasullo, 2007: Estimates of the global water budget and its annual cycle using observational and model data. *J. Hydrometeorol.*, **8**, 758–769, <https://doi.org/10.1175/JHM600.1>.
- TRMM, 2011: TRMM Precipitation Radar rainfall rate and profile L2 1.5 hours V7. Goddard Earth Sciences Data and Information Services Center, accessed 19 August 2016, [https://disc.gsfc.nasa.gov/datacollection/TRMM\\_2A25\\_7.html](https://disc.gsfc.nasa.gov/datacollection/TRMM_2A25_7.html).
- Ullrich, P. A., and C. M. Zarzycki, 2017: TempestExtremes: A framework for scale-insensitive pointwise feature tracking on unstructured grids. *Geosci. Model Dev.*, **10**, 1069–1090, <https://doi.org/10.5194/gmd-10-1069-2017>.
- Vishnu, S., W. R. Boos, P. A. Ullrich, and T. A. O’Brien, 2020: Assessing historical variability of South Asian monsoon lows and depressions with an optimized tracking algorithm. *J. Geophys. Res. Atmos.*, **125**, e2020JD032977, <https://doi.org/10.1029/2020JD032977>.

- Wang, B., and Coauthors, 2018: Dynamics-oriented diagnostics for the Madden–Julian oscillation. *J. Climate*, **31**, 3117–3135, <https://doi.org/10.1175/JCLI-D-17-0332.1>.
- Wang, J., H. Kim, D. Kim, S. A. Henderson, C. Stan, and E. D. Maloney, 2020: MJO teleconnections over the PNA region in climate models. Part II: Impacts of the MJO and basic state. *J. Climate*, **33**, 5081–5101, <https://doi.org/10.1175/JCLI-D-19-0865.1>.
- Wentz, F. J., C. Gentemann, and K. A. Hilburn, 2015: Remote Sensing Systems TRMM TMI Daily Environmental Suite on 0.25 deg grid, version 7.1. Remote Sensing Systems, accessed 8 July 2016, [www.remss.com/missions/tmi](http://www.remss.com/missions/tmi).
- Wolding, B., J. Dias, G. Kiladis, F. Ahmed, S. W. Powell, E. Maloney, and M. Branson, 2020: Interactions between moisture and tropical convection. Part I: The coevolution of moisture and convection. *J. Atmos. Sci.*, **77**, 1783–1799, <https://doi.org/10.1175/JAS-D-19-0225.1>.
- Xie, S., and Coauthors, 2010: Clouds and more: ARM climate modeling best estimate data. *Bull. Amer. Meteor. Soc.*, **91**, 13–20, <https://doi.org/10.1175/2009BAMS2891.1>.
- , and Coauthors, 2019: Improved diurnal cycle of precipitation in E3SM with a revised convective triggering function. *J. Adv. Model. Earth Syst.*, **11**, 2290–2310, <https://doi.org/10.1029/2019MS001702>.
- , R. T. Cederwall, and M. Zhang, 2004: Developing long-term single-column model/cloud system–resolving model forcing data using numerical weather prediction products constrained by surface and top of the atmosphere observations. *J. Geophys. Res.*, **109**, D01104, <https://doi.org/10.1029/2003JD004045>.
- Yadav, P., and D. M. Straus, 2017: Circulation response to fast and slow MJO episodes. *Mon. Wea. Rev.*, **145**, 1577–1596, <https://doi.org/10.1175/MWR-D-16-0352.1>.
- Yang, G. Y., and J. Slingo, 2001: The diurnal cycle in the tropics. *Mon. Wea. Rev.*, **129**, 784–801, [https://doi.org/10.1175/1520-0493\(2001\)129<0784:TDCITT>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<0784:TDCITT>2.0.CO;2).
- Yin, L., R. Fu, E. Shevliakova, and R. E. Dickinson, 2013: How well can CMIP5 simulate precipitation and its controlling processes over tropical South America? *Climate Dyn.*, **41**, 3127–3143, <https://doi.org/10.1007/s00382-012-1582-y>.
- Zhou, Y., T. A. O'Brien, P. A. Ullrich, W. D. Collins, C. M., Patricola, and A. M. Rhoades, 2021: Uncertainties in atmospheric river lifecycles by detection algorithms: Climatology and variability. *J. Geophys. Res. Atmos.*, **126**, e2020JD033711, <https://doi.org/10.1029/2020JD033711>.