



25 1. Introduction

26 The onset of the South China Sea (SCS), also known as the East Sea in Vietnam, summer
27 monsoon (SCSSM) marks a rapid reorganization of the Asian monsoon system, characterized by a
28 transition from a dry, trade-wind-dominated spring regime to a convectively coupled summer
29 circulation (Wang and LinHo, 2002; Ding et al., 2004; Wang et al., 2004; Ding et al., 2015).
30 Climatologically occurring in mid-May, this transition is marked by a reversal of lower-tropospheric
31 zonal winds over the SCS, enhanced low-level moisture convergence, and the eastward retreat of the
32 western North Pacific (WNP) subtropical high (WNPSH), leading to the establishment of deep
33 convection over the basin (Zhou and Chan, 2007; Wang et al., 2009). As a precursor to the broader
34 East Asian summer monsoon, variability in the timing of SCSSM onset strongly regulates the
35 initiation of the regional rainy season over Southeast Asia and southern China. Previous studies have
36 further shown that this transition is sensitive to subseasonal disturbances and remote forcing,
37 including intraseasonal oscillations, tropical–extratropical wave interactions, and large-scale Indo–
38 Pacific circulation anomalies that modulate low-level westerlies and convective instability over the
39 SCS (Ding and Chan, 2005; Geen, 2021).

40 Although the SCSSM onset occurs every year, its interannual variability exerts a strong control
41 on regional climate through modulation of large-scale circulation, moisture transport, and convective
42 activity. Variations in onset date (OD) alter the establishment and persistence of low-level westerlies
43 over the SCS and the WNP, thereby redistributing summer rainfall across East and Southeast Asia
44 and shaping regional drought and flood risks. An early SCSSM OD is often followed by suppressed
45 summer rainfall over subtropical East Asia, including the middle and lower reaches of the Yangtze
46 River basin, whereas a delayed OD tends to prolong the pre-monsoon rainfall regime and enhance
47 flood potential (Jiang et al., 2018; He and Zhu, 2015). Over Southeast Asia, early OD favors
48 enhanced convection and an increased likelihood of extreme rainfall events, while delayed OD is
49 associated with reduced low-level moisture convergence and suppressed heavy rainfall in May (Hu et
50 al., 2022a, b). The SCSSM OD also modulates tropical cyclone activity over the WNP by influencing



51 background vorticity, vertical wind shear, and monsoon trough development, with earlier OD linked
52 to increased tropical cyclone genesis and a higher frequency of landfalling storms along the
53 southeastern coast of China (Chen et al., 2017; Huangfu et al., 2017; Wang and Chen, 2018).

54 Together, these mechanisms identify SCSSM OD variability as a key driver of seasonal-scale climate
55 variability with direct implications for agriculture, water resources, and disaster risk management.

56 Motivated by these impacts, substantial efforts have been devoted to predicting the SCSSM OD
57 using dynamical seasonal forecast systems. Despite the chaotic nature of atmosphere, previous
58 studies demonstrate that useful predictive skill for interannual variability in SCSSM onset exists at
59 seasonal lead times of approximately one to three months (e.g., Zhu and Li, 2017; Martin et al.,
60 2019; Chevuturi et al., 2019, 2021; Attada et al., 2022). This seasonal predictability is derived
61 primarily from slowly evolving boundary conditions, most notably the El Niño–Southern Oscillation
62 (ENSO), which modulates large-scale circulation over the WNP through the Indo-Pacific Ocean
63 Capacitor effect (Xie et al., 2016). As a result, operational forecast systems such as the UK Met
64 Office GloSea5 and the European Centre for Medium-Range Weather Forecasts (ECMWF) seasonal
65 forecasting system 5 (SEAS5) can skillfully discriminate between early and late onset years up to
66 three months in advance, particularly when onset is defined using robust circulation-based metrics,
67 such as the 850-hPa zonal wind (U850) index of Wang et al. (2004), rather than local precipitation
68 alone (Bombardi et al., 2017, 2020; Chevuturi et al., 2019; Martin et al., 2019).

69 However, extending this predictive capability to longer lead times of four to five months remains
70 a formidable challenge. State-of-the-art dynamical models encounter a pronounced spring
71 predictability barrier, during which the influence of ENSO conditions from the preceding winter
72 weakens and stochastic atmospheric variability associated with the seasonal transition becomes
73 dominant, often rendering forecasts initialized in winter (e.g., preceding December or January)
74 statistically insignificant for predicting a May onset (Martin et al., 2019). In addition to intrinsic
75 predictability limits, systematic model deficiencies persist. Coupled systems such as ECMWF
76 SEAS5 frequently exhibit a cold sea surface temperature (SST) bias over the SCS and struggle to



77 realistically simulate the northward propagation of intraseasonal oscillations, leading to a systematic
78 late-onset bias and degraded forecast skill at extended lead times (Chevuturi et al., 2021; Bui-Minh et
79 al., 2024). Together, these limitations suggest that existing linear indices and conventional dynamical
80 frameworks may not adequately capture the nonlinear precursors required to bridge the predictability
81 gap between seasonal (~three-month) and long-range (~five-month) forecasts.

82 In parallel with advances in dynamical seasonal prediction, recent studies have proposed
83 alternative circulation-based approaches to defining monsoon onset using clustering of synoptic-
84 scale atmospheric patterns (e.g., Borah et al., 2013; Dai et al., 2021; Bui-Minh et al., 2024). These
85 approaches typically employ self-organizing maps (SOM; Kohonen, 2001) to project high-
86 dimensional atmospheric fields onto a finite set of representative circulation patterns, which are
87 subsequently grouped using clustering algorithms such as K-means. The resulting clusters
88 objectively distinguish between pre-monsoon and monsoon circulation regimes, allowing onset to be
89 identified as the sustained transition from dry, pre-monsoon patterns to convectively active monsoon-
90 associated patterns. By construction, this framework captures the full spatial structure and temporal
91 evolution of the large-scale circulation, rather than relying on a single variable or threshold-based
92 index. Bui-Minh et al. (2024) shows that synoptic-pattern clustering provides a more physically
93 consistent and temporally coherent identification of monsoon onset than traditional index-based
94 definitions. Importantly, such pattern-based definitions are more closely aligned with the dynamical
95 structures resolved by numerical models, suggesting potential advantages for seasonal prediction.

96 Despite this promise, the extent to which synoptic-pattern-based definitions can improve the
97 predictability of SCSSM OD in operational seasonal forecast systems, particularly at lead times
98 exceeding three months, has not yet been systematically assessed. This study addresses this gap by
99 developing a synoptic clustering-based definition of SCSSM onset using the combination of SOM
100 and clustering techniques and evaluating its performance within a seasonal prediction framework.
101 Specifically, the proposed onset definition will be applied to retrospective forecasts from the
102 ECMWF SEAS5 system to assess whether it enhances prediction skill for SCSSM onset date at



103 extended lead times. By explicitly linking monsoon onset to large-scale circulation regimes rather
104 than single-variable indices, this work aims to provide a more robust and dynamically meaningful
105 framework for understanding and predicting the SCSSM onset. The remainder of this paper is
106 organized as follows. Section 2 describes the data and methodology. Section 3 presents the synoptic
107 clustering-based definition of the SCSSM OD. Section 4 evaluates its performance in seasonal
108 prediction. Finally, discussions and conclusions are given in Section 5.

109 **2. Datasets and methodology**

110 **2.1 Datasets**

111 **2.2.1 Observational data**

112 The fifth generation of the ECMWF Reanalysis dataset (ERA5; Hersbach et al., 2020) for the
113 period 1979 to the present, with a horizontal resolution of $0.25^\circ \times 0.25^\circ$, is used to analyze
114 atmospheric circulation. The analysis focuses on 850-hPa zonal wind (U850), meridional wind
115 (V850), and geopotential height (Z850). Pentad (5-day) means are computed from hourly ERA5
116 data. Monthly SST data are taken from the UK Met Office Hadley Centre Global Sea Ice and Sea
117 Surface Temperature dataset (HadISST; Rayner et al., 2003), with a resolution of $1^\circ \times 1^\circ$, available
118 from 1870 to the present. Anomalies are calculated by subtracting the 1981–2010 climatology from
119 detrended data over the corresponding period of record.

120 **2.2.2 Model data**

121 The ECMWF SEAS5 (Johnson et al., 2019) coupled forecasting system consists of the
122 Integrated Forecast System (IFS) atmospheric model, the HTESSEL land surface model, and the
123 Nucleus for European Modelling of the Ocean (NEMO). SEAS5 provides seasonal forecasts at an
124 O320 horizontal resolution (approximately 36 km) with 91 vertical levels in the atmosphere, and an
125 ORCA 0.25° grid (approximately 27 km) with 75 vertical levels in the ocean. The SEAS5 hindcasts
126 and forecasts are initialized on the first day of each month and integrated for seven months, using a



127 21-member ensemble for hindcasts over 1981–2016 and a 51-member ensemble for real-time
128 forecasts from 2017 onward.

129 **2.2 Methodology**

130 **2.2.1 Clustering**

131 The methodology applied in this study follows Nguyen-Le et al. (2017, 2019) and is designed to
132 classify synoptic-scale atmospheric circulation patterns over the SCS and surrounding regions before
133 and after the SCSSM onset using SOM. From a synoptic perspective, SCSSM onset is characterized
134 by the development of the low-level monsoon trough and cross-equatorial flow near 105°E, the
135 eastward retreat of the WNPSH, and the establishment of a coherent monsoonal meridional
136 circulation (Liu et al., 2016; Huangfu et al., 2017; Hu et al., 2018). Based on these dynamical
137 features, SOM input vectors are constructed from pentad-mean ERA5 reanalysis fields of Z850,
138 U850, and V850, which are concatenated for each pentad to form a single high-dimensional input
139 vector for SOM training. A total of 864 input vectors is generated, corresponding to 18 pentads per
140 year (from pentad 19, 1–5 April, to pentad 36, 25–29 June) over the 38-year period 1979–2016. The
141 spatial domain spans 5°S–25°N and 95°E–135°E, yielding 161 longitude points and 81 latitude
142 points. Consequently, each raw input vector contains $3 \times 161 \times 81 = 39123$ elements. Because these
143 elements are not statistically independent owing to strong spatial coherence and inter-variable
144 correlations, direct use of the full input vectors is neither efficient nor optimal. To reduce
145 dimensionality and remove redundant information, principal component analysis (PCA) is applied
146 jointly to the three original reanalysis fields. Prior to PCA, each variable is normalized to ensure
147 comparable variance contributions, as both PCA and SOMs are sensitive to variable scaling. Only
148 empirical orthogonal functions explaining 99% of the total variance are retained. The corresponding
149 principal components (PCs) are then used as input to the SOM in place of the original fields. This
150 procedure reduces the dimensionality of each input vector to $d = 144$, corresponding to a reduction



151 factor of approximately 270, thereby improving statistical independence among input elements and
152 substantially reducing computational cost.

153 The SOM training outcome depends on both the lattice size and several training parameters,
154 including the learning rate, neighborhood radius, and training length. A trial-and-error approach is
155 therefore employed to identify an optimal SOM configuration. Multiple lattice sizes (4×4 , 5×5 , $6 \times$
156 6 , 7×7 , and 8×8) are tested and evaluated based on node occupancy, quantization error (QE), and
157 topographic error (TE). Specifically, an excessive number of samples assigned to individual nodes
158 (e.g., >10) indicates insufficient pattern separation and suggests the need for a larger lattice, whereas
159 the presence of empty nodes implies over-partitioning and an excessively large lattice. Based on
160 these criteria, a SOM with a 6×6 hexagonal lattice (36 nodes) is selected, using a learning rate of
161 0.2 and a neighborhood radius of three. To ensure numerical stability and convergence, the SOM is
162 trained for two million iterations. Given the relatively limited sample size, bootstrap resampling is
163 incorporated into the training procedure, whereby samples are randomly drawn with replacement
164 from the original dataset. A total of 1000 SOM realizations are generated using this configuration.
165 The final “master” SOM is selected as the realization exhibiting the lowest QE and TE (Fig. S1) and
166 a relatively smooth Sammon mapping (Sammon, 1969), indicating good preservation of the
167 topological relationships among patterns (Fig. S2).

168 Because the SOM is designed to preserve the topological structure of the original high-
169 dimensional input space in a two-dimensional lattice, nodes that are adjacent in the SOM represent
170 more similar circulation patterns than those that are farther apart (Vesanto and Alhoniemi, 2000).
171 This similarity can be quantified using the unified distance matrix (U-matrix), which measures the
172 average distance between the codebook vector of a given node and those of its nearest neighbors
173 (Ultsch and Siemon, 1990). To enhance the objectivity and robustness of pattern classification, a
174 two-step clustering strategy is adopted. First, the SOM organizes the atmospheric states according to
175 their topological similarity. Second, the resulting SOM nodes are further grouped into a smaller
176 number of circulation regimes based on the U-matrix structure using K-means clustering. A known



177 limitation of K-means is the need to predefine the number of clusters, K . To address this issue, K is
178 determined objectively by examining the spatial distribution of U-matrix values. Potential clusters
179 are identified around local minima of the U-matrix, while cluster boundaries correspond to nodes
180 with relatively large U-matrix values separating adjacent minima. Fig. S3 shows the U-matrix of the
181 master SOM, with values normalized to the unit interval. Based on the number, spatial separation,
182 and surrounding high-distance boundaries of the local minima, the SOM lattice is partitioned into six
183 distinct clusters ($K = 6$), representing six characteristic synoptic-scale circulation regimes during the
184 April–June period. Four clusters correspond to pronounced local minima located in the top-left,
185 bottom-left, bottom-middle, and bottom-right regions of the lattice. In addition, elevated U-matrix
186 values in the top-middle and top-right portions of the SOM indicate circulation patterns that are
187 markedly different from those elsewhere, supporting their identification as two additional,
188 independent regimes. Further methodological details regarding the combined use of SOMs and U-
189 matrix-guided clustering can be found in Nishiyama et al. (2007) and Nguyen-Le et al. (2017, 2019).

190 **2.2.1 Prediction**

191 In the prediction phase, the circulation clustering results are applied to prognostic atmospheric
192 fields from SEAS5. Pentad-mean forecasts of Z850, U850, and V850 fields from pentad 19 (1–5
193 April) to pentad 36 (25–29 June) for the period 2017–2024 are extracted over the same spatial
194 domain used in the SOM training. Forecasts initialized from December of the preceding year through
195 April, corresponding to lead times of approximately five to one months, are analyzed to assess lead-
196 time dependence. The forecast fields are normalized using the means and standard deviations derived
197 from the ERA5 training period and projected onto the same PCA space to construct a forecast vector
198 $\mathbf{p} = [p_1, p_2, \dots, p_d]$, where $d = 144$. The Euclidean distance between \mathbf{p} and the centroid of each
199 synoptic-scale circulation cluster (C1–C6) is then computed, and each forecast pentad is assigned to
200 the cluster with the minimum distance. The SCSSM onset date is subsequently identified using the
201 synoptic clustering-based definition introduced in Section 3.



202 **2.3 Measures of prediction skill**

203 Following Chevuturi et al. (2019, 2021), the prediction skill of the SCSSM OD is evaluated
204 separately for the dependent and independent periods by comparing onset dates diagnosed from
205 SEAS5 hindcasts (1981–2016) and forecasts (2017–2024) with those derived from ERA5 reanalysis.
206 OD are calculated for each individual ensemble member as well as for the ensemble mean, allowing
207 assessment of both deterministic and probabilistic forecast skill. Skill is further evaluated for
208 different initialization months to examine lead-time dependence.

209 **2.3.1 Deterministic prediction skill**

210 The ability of SEAS5 to reproduce the observed interannual variability of the SCSSM OD is
211 quantified using the Pearson correlation coefficient (r) between predicted and observed onset dates.
212 The statistical significance of the correlation is assessed using a two-tailed Student's t test, with $p <$
213 0.05 indicating significance.

214 Because correlation alone does not fully characterize predictability, potential predictability is
215 further evaluated using the ratio of predictable component (RPC) (Eade et al., 2014), defined as

$$216 \quad \text{RPC} = \frac{r}{\sqrt{\sigma_{\text{sig}}^2 / \sigma_{\text{tot}}^2}} \quad (1)$$

217 where r is the correlation between the ensemble-mean prediction and observations, σ_{sig}^2 is the
218 variance of the ensemble mean representing the predictable signal, and σ_{tot}^2 is the total variance
219 across all ensemble members. An RPC value close to unity indicates that the forecast system
220 effectively extracts the available predictable signal, whereas values substantially below unity suggest
221 underconfident or overly dispersed ensemble forecasts.

222 **2.3.2 Categorical skill of ensemble-mean forecasts**

223 To evaluate the skill of SEAS5 in predicting the relative timing of SCSSM onset, OD are
224 classified as early, normal, or late based on the forecast–observation difference. A forecast is



225 classified as early if the predicted OD is more than one pentad earlier than observed (< -1 pentad),
226 normal if the difference lies within ± 1 pentad, and late if the predicted OD is more than one pentad
227 later than observed (> 1 pentad). This categorical framework allows assessment of forecast skill in
228 capturing relative onset timing, even when the exact OD is not reproduced.

229 For deterministic forecasts based on ensemble-mean OD, categorical prediction skill is
230 quantified using Accuracy (ACC) and the Heidke Skill Score (HSS) (WCRP, 2015). ACC measures
231 the fraction of correct categorical forecasts relative to the total number of forecasts and is defined as

$$232 \quad ACC = \frac{1}{N} \sum_{i=1}^C n(F_i, O_i) \quad (2)$$

233 where $C(F_i, O_i)$ equals 1 if the forecast category F_i matches the observed category O_i (early–
234 normal–late), and 0 otherwise, and N is the total number of forecasts. ACC ranges from 0 (no skill)
235 to 1 (perfect accuracy).

236 The HSS evaluates forecast accuracy relative to random chance and is defined as

$$237 \quad HSS = \frac{\frac{1}{N} \sum_{i=1}^C n(F_i, O_i) - \frac{1}{N^2} \sum_{i=1}^C n(F_i) n(O_i)}{1 - \frac{1}{N^2} \sum_{i=1}^C n(F_i) n(O_i)} \quad (3)$$

238 where $C(F_i)$ and $n(O_i)$ denote the marginal frequencies of forecast and observed categories,
239 respectively. HSS ranges from negative values (worse than random) to 1 (perfect skill), with 0
240 indicating no skill relative to climatological chance.

241 **2.3.3 Probabilistic skill of ensemble forecasts**

242 Probabilistic prediction skill is evaluated using the full ensemble information. For each forecast,
243 probabilities for the three onset categories (early/normal/late) are defined as the fraction of ensemble
244 members predicting each category. These probabilistic forecasts are verified against observations
245 using the Brier Skill Score (BSS) and the Ranked Probability Skill Score (RPSS), which assess
246 forecast reliability, resolution, and discrimination relative to a climatological reference forecast. The
247 discrete BSS (dBSS) for each category is defined as

$$248 \quad dBSS = 1 - \frac{BS_c}{BS_{clim+D}} \quad (4)$$



249 where BS_c is the Brier Score for a given category,

$$250 \quad BS_c = \frac{1}{N} \sum_{i=1}^N (F_i - O_i)^2 \quad (5)$$

251 where F_i denoting the forecast probability, O_i the observed category indicator (1 for the
 252 observed category and 0 otherwise), and N the number of forecasts. The reference score BS_{clim} is
 253 computed assuming equal climatological probabilities (1/3). The dBSS ranges from negative values
 254 (worse than climatology) to 1 (perfect skill), with 0 indicating no skill.

255 The discrete RPSS (dRPSS) evaluates cumulative probability errors across ordered onset
 256 categories (early–normal–late) and is defined as

$$257 \quad dRPSS = 1 - \frac{RPS}{RPS_{clim} + D} \quad (6)$$

258 where the Ranked Probability Score (RPS) is given by

$$259 \quad RPS = \frac{1}{N} \sum_{i=1}^N \left[\frac{1}{(C-1)} \left(\sum_{j=1}^C (F_j - O_j)^2 \right) \right] \quad (7)$$

260 with $C = 3$ denoting the number of ordered categories. The reference score RPS_{clim} is
 261 calculated using climatological probabilities of 1/3 for each category. The dRPSS similarly ranges
 262 from negative values to 1, with 0 indicating no skill relative to climatology, and summarizes forecast
 263 errors arising from systematic biases while accounting for discrimination and resolution.

264 To account for finite ensemble size, a discrete correction term D is applied to both dBSS and
 265 dRPSS following Weigel et al. (2007),

$$266 \quad D = \frac{1}{M} \sum_{i=1}^K \left[p_i \left(1 - p_i - 2 \times \sum_{j=i+1}^K p_j \right) \right] \quad (8)$$

267 where M is the ensemble size and p_k is the climatological forecast probability for category k .

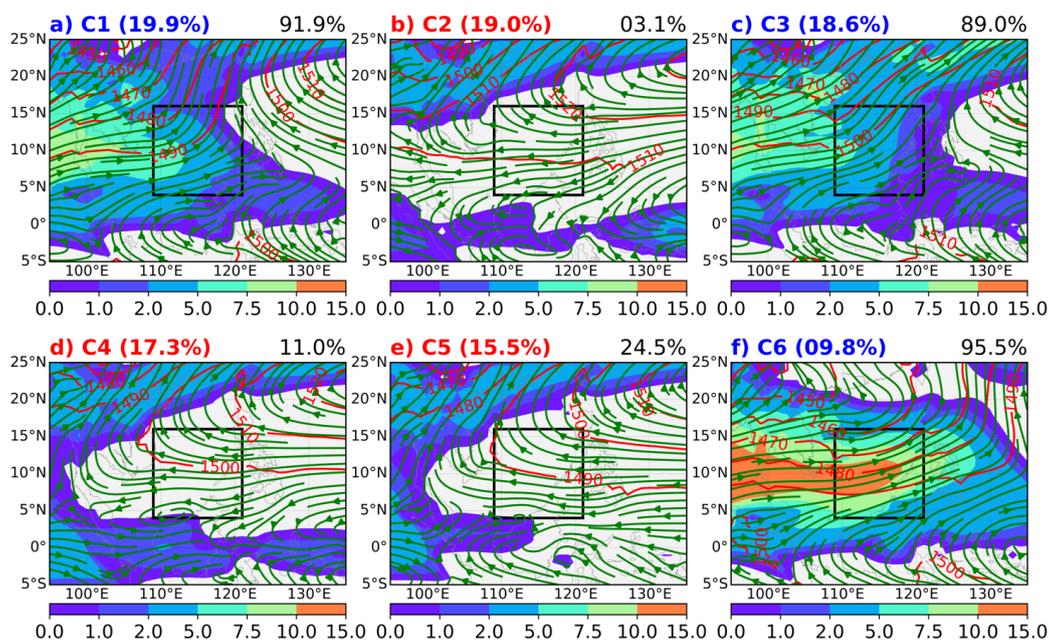
268 The correction term is averaged across all years to account for variations in ensemble size and
 269 becomes negligible for large ensembles. This correction is particularly important for small
 270 ensembles, where limited probability spread can otherwise lead to artificially inflated skill estimates.



271 **3. Synoptic clustering–based definition of the South China Sea summer monsoon onset**

272 **3.1 Definition of SCSSM onset based on synoptic circulation regimes**

850hPa Streamline, Zonal Wind and Geopotential Height



273
274 **Figure 1:** Synoptic circulation regimes associated with the SCSSM identified from SOM clustering. Shown
275 are composites of 850-hPa winds (vectors), zonal wind (shades) and geopotential height (contours) for six
276 clusters (C1–C6). The black box denotes the SCS region used for onset definition based on the U850 index of
277 Wang et al. (2004). Cluster labels and their relative frequencies of occurrence (%) are shown in each panel,
278 while numbers in the upper-right corner indicate the percentage of occurrences after the SCSSM OD defined
279 by Wang et al. (2004). Clusters (a) C1, (c) C3, and (f) C6 correspond to monsoon regimes, whereas (b) C2, (d)
280 C4, and (e) C5 represent pre-monsoon regimes.

281 Figure 1 presents the six synoptic circulation regimes (C1–C6) identified from the combination
282 of SOM and K-means clustering, illustrated by composites of 850-hPa winds and geopotential
283 height. Together, these clusters represent the dominant low-level circulation regimes spanning the
284 seasonal transition from pre-monsoon to monsoon conditions. Specifically, clusters C2, C4, and C5
285 are representative of pre-monsoon regimes. These patterns are characterized by easterly or
286 northeasterly low-level flow over the SCS region, associated with the westward extension or



287 persistence of the WNPSH. Geopotential height contours are largely zonal and tightly packed north
288 of the SCS, indicating weak low-level convergence and unfavorable conditions for sustained deep
289 convection over the basin. Consistent with this interpretation, these clusters occur almost exclusively
290 before the SCSSM onset defined by the U850 index of Wang et al. (2004), as reflected by their low
291 post-onset occurrence frequencies (3.1% for C2, 11.0% for C4, and 24.0% for C5).

292 In contrast, clusters C1, C3, and C6 correspond to monsoon regimes, characterized by the
293 establishment of low-level westerlies over the SCS and enhanced cyclonic curvature of the flow.
294 These regimes are associated with an eastward retreat and weakening of the WNPSH, strengthened
295 geopotential height gradients across the SCS, and enhanced low-level convergence conducive to
296 deep convection. Among the monsoon regimes, C6 represents the most mature and canonical
297 monsoon state, characterized by strong and spatially coherent southwesterly flow extending across
298 the SCS into the WNP and a pronounced eastward retreat of the WNPSH. This regime occurs almost
299 exclusively after the onset defined by Wang et al. (2004), with a post-onset occurrence frequency of
300 95.5%, indicating a fully established monsoon circulation. C3 also corresponds to a mature monsoon
301 regime, exhibiting stronger low-level westerlies over the SCS and a further eastward withdrawal of
302 the WNPSH compared to C1, reflecting a more developed monsoon structure. While C1 represents a
303 relatively weaker or early-stage monsoon circulation, the high post-onset occurrence frequencies of
304 C3 (89.0%) and C1 (91.9%) confirm that all three regimes are robustly associated with post-onset
305 monsoon conditions and collectively capture variability in monsoon strength and maturity.

306 The clear dynamical separation between pre-monsoon and monsoon circulation patterns,
307 together with their contrasting post-onset occurrence frequencies, demonstrates that the combination
308 of SOM and K-means clustering effectively captures the large-scale circulation transition underlying
309 SCSSM onset. Importantly, the presence of multiple monsoon regimes indicates that SCSSM onset is
310 not associated with a single circulation pattern, but rather with a regime transition toward persistent
311 low-level westerlies and enhanced convergence marked by the eastward retreat of the WNPSH over



312 the SCS. Based on this regime framework, the OD of the SCSSM can be defined objectively using
313 synoptic circulation transitions rather than a single-variable threshold.

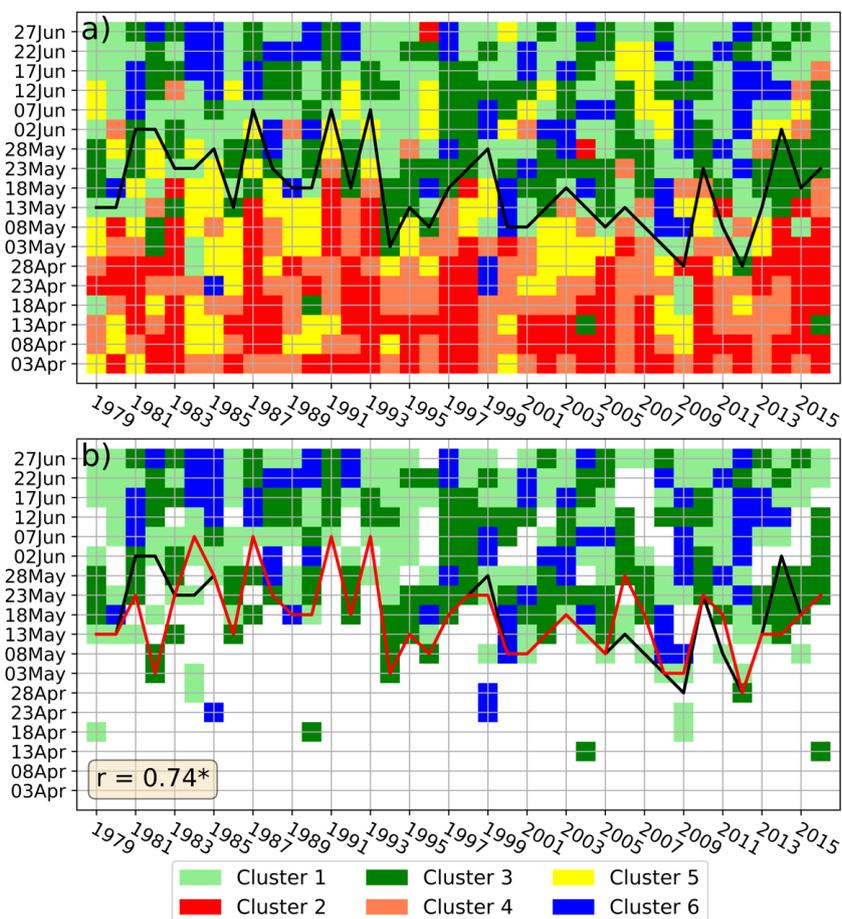
314 *The OD of the SCSSM is defined as the first pentad from pentad 24 (26–30 April onward) that*
315 *satisfies the following three criteria: a) **Immediate regime transition:** The synoptic circulation in the*
316 *onset pentad and the subsequent pentad must both belong to the monsoon regimes (C1, C3, or C6);*
317 *b) **Persistence criterion:** Within the subsequent four pentads (including the onset pentad), the*
318 *synoptic circulation must belong to the monsoon regimes in at least three pentads; and c) **Mature***
319 *monsoon condition: Within the same four-pentad window (including the onset pentad), cluster C3 or*
320 *C6 must occur at least once.*

321 The first criterion ensures that the identified onset corresponds to a genuine transition into
322 monsoon-type circulation rather than a transient synoptic fluctuation, thereby excluding “bogus”
323 onsets associated with brief westerly intrusions (Wang et al., 2004). The second criterion enforces
324 circulation persistence, filtering out short-lived or weak monsoon-like disturbances commonly
325 observed during the pre-monsoon period. The third criterion further requires the emergence of C3
326 and C6, which represents a more developed monsoon circulation state, thereby ensuring that the
327 onset is associated with a sufficiently mature and dynamically robust monsoon configuration.
328 Together, these criteria explicitly link SCSSM onset to a sustained regime transition characterized by
329 persistent low-level westerlies over the SCS and the retreat of the WNPSH, providing a dynamically
330 grounded alternative to traditional single-index onset definitions.

331 Figure 2 compares the SCSSM onset dates identified using the synoptic clustering-based
332 definition (red line) with those derived from the U850 index of Wang et al. (2004) (black line). The
333 two onset definitions show strong agreement, with a statistically significant correlation of 0.74 ($p <$
334 0.01), indicating that the new definition captures the dominant interannual variability of SCSSM
335 onset timing identified by the conventional circulation-based index. This overall consistency
336 demonstrates that the synoptic clustering-based framework remains aligned with established onset
337 metrics while being explicitly grounded in large-scale circulation regimes. Despite this overall



338 coherence, differences occur in individual years. The clustering-based onset is sometimes identified
339 later than the Wang et al. (2004) onset when brief or weak westerlies are not followed by sustained
340 monsoon circulation, reflecting the imposed persistence and maturity criteria. In contrast, earlier
341 onset identification occurs in years with an early and coherent transition to monsoon-type circulation
342 preceding the U850 threshold. These differences highlight the sensitivity of index-based definitions
343 to short-lived zonal wind reversals and underscore the advantage of a regime-based framework that
344 explicitly accounts for circulation persistence. Overall, the synoptic clustering-based definition
345 preserves the large-scale timing of SCSSM onset while providing a more dynamically constrained
346 and physically interpretable onset identification, supporting its application to seasonal prediction.



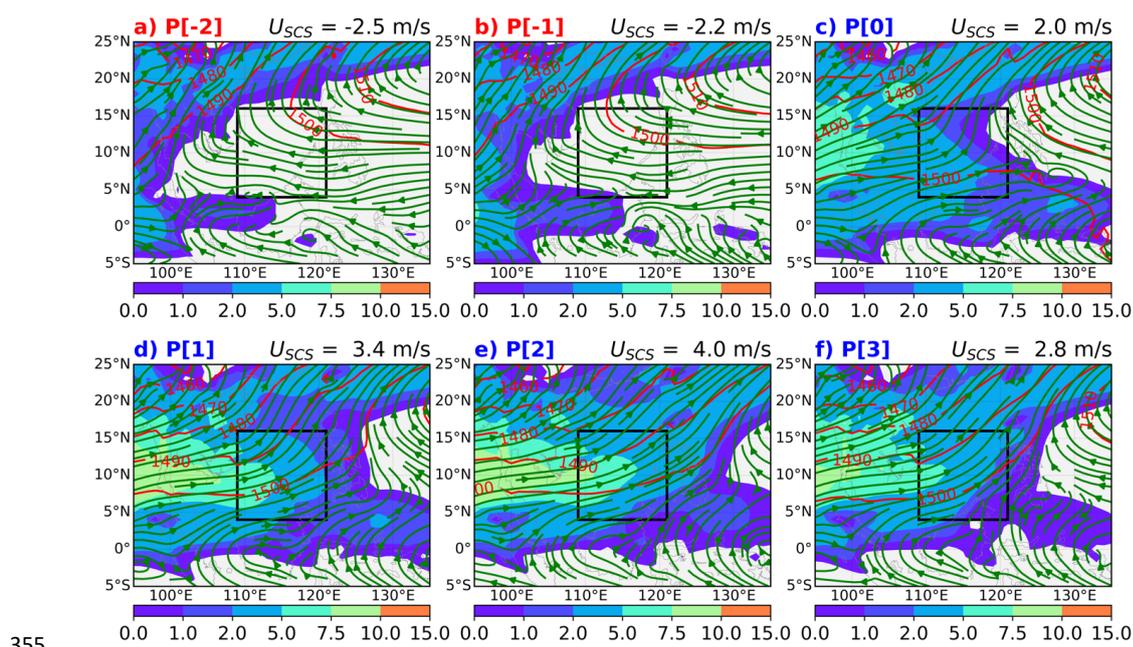
347



348 **Figure 2:** Interannual evolution of synoptic circulation regimes and SCSSM OD. Colored shading indicates
 349 the dominant synoptic circulation regime (C1–C6) for each pentad from early April to late June for individual
 350 years during 1979–2016. (a) OD identified using U850 index of Wang et al. (2004) (red line). (b) Comparison
 351 between OD identified by the synoptic clustering–based definition (red line) and the U850-based definition of
 352 Wang et al. (2004) (black line). The correlation between the two onset time series is statistically significant (r
 353 $= 0.74$, $p < 0.01$). Only monsoon regimes (C1, C2, and C6) are shown in (b).

354 3.2 Climatology and interannual variability of the SCSSM onset

850hPa Streamline, Zonal Wind and Geopotential Height



355

356 **Figure 3:** Climatological evolution of 850-hPa winds (vectors), zonal wind (shades) and geopotential height
 357 (contours) from two pentads before to three pentads after the SCSSM onset, based on the synoptic clustering–
 358 based definition. Panels show composites for P[–2] and P[–1] (pre-onset), P[0] (onset), and P[1]–P[3] (post-
 359 onset). The black box denotes the SCS region used to compute the area-averaged zonal wind (U_{SCS}), whose
 360 values (m s^{-1}) are indicated in the upper-right corner of each panel.

361 Figure 3 illustrates the climatological evolution of low-level circulation from two pentads prior
 362 to SCSSM onset to three pentads after onset, based on the synoptic clustering–based definition. The
 363 composites show 850-hPa winds and geopotential height, with the black box indicating the South
 364 China Sea (SCS) region used to calculate area-averaged zonal wind (U_{SCS}). During the pre-onset

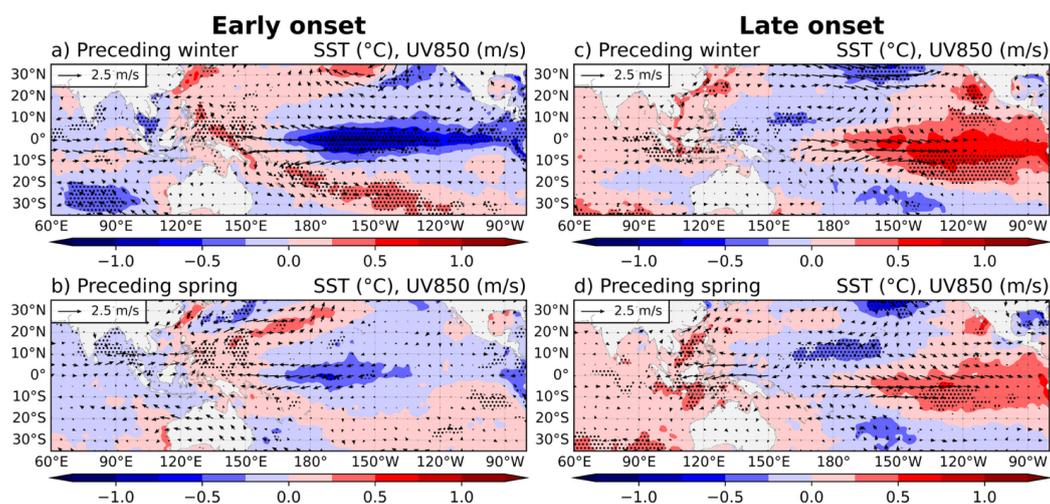


365 phase P[-2] and P[-1] (Figs. 3a–b), the SCS is dominated by easterly to northeasterly low-level
366 flow, reflected by negative U_{SCS} values (-2.5 and -2.2 m s⁻¹, respectively). This circulation is
367 associated with a westward-extending WNPSH, as indicated by the orientation and packing of
368 geopotential height contours north of the SCS. Low-level convergence over the basin is weak,
369 consistent with suppressed convective conditions characteristic of the pre-monsoon regime.

370 At the onset pentad P[0] (Fig. 3c), an abrupt circulation transition occurs. Low-level winds over
371 the SCS reverse to westerly, with U_{SCS} becoming positive (2.0 m s⁻¹), marking the establishment of
372 monsoon flow. This transition is accompanied by an eastward retreat of the WNPSH and enhanced
373 meridional geopotential height gradients across the SCS, indicating strengthened low-level
374 convergence favorable for deep convection.

375 In the post-onset period P[1]–P[3] (Figs. 3d–f), westerly flow intensifies and becomes more
376 spatially coherent across the SCS, with U_{SCS} increasing to 3.4 – 4.0 m s⁻¹ before slightly weakening by
377 P[3]. The circulation exhibits a well-developed monsoon structure, characterized by persistent
378 southwesterlies extending from the equatorial Indian Ocean into the SCS and WNP. The sustained
379 positive U_{SCS} and stable geopotential height configuration demonstrate that the onset identified by
380 the clustering-based definition corresponds to a robust and persistent transition into the summer
381 monsoon regime rather than a transient wind reversal.

382 The climatological evolution shown in Figure 3 is consistent with the combined influence
383 of slowly varying boundary forcing and subseasonal atmospheric variability on SCSSM onset. On
384 seasonal timescales, the transition from pre-onset easterlies to sustained post-onset westerlies reflects
385 the gradual weakening and eastward retreat of the WNPSH, a process that is strongly modulated by
386 ENSO through the Indo-Pacific Ocean Capacitor mechanism (Xie et al., 2009, 2016). This ENSO-
387 related background modulation is further illustrated in Figure 4, which shows composite anomalies
388 of SST and 850-hPa winds for early and late SCSSM onset years during the preceding winter
389 (December–January–February) and spring (March–April–May).



390

391 **Figure 4:** Composite anomalies of SST (shading, °C) and 850-hPa winds (vectors, m/s) associated with early
392 and late SCSSM onset years. (a, b) Early-onset composites for the preceding winter and spring, respectively.
393 (c, d) Late-onset composites for the preceding winter (December–January–February) and spring (March–
394 April–May). Stippling denotes regions where SST anomalies are statistically significant at the 95% confidence
395 level. Wind vectors indicate anomalous low-level circulation, with the reference vector shown in the upper-
396 left corner of each panel.

397 For early onset years (Figs. 4a–b), the preceding winter is characterized by cold SST anomalies
398 over the central–eastern equatorial Pacific, warm anomalies in the WNP and cold anomalies in the
399 Indian Ocean, consistent with La Niña–like conditions. Associated low-level wind anomalies exhibit
400 weakened easterlies over the western Pacific and enhanced westerly tendencies extending toward the
401 SCS. During the preceding spring, these anomalies persist and intensify, promoting a background
402 circulation favorable for the early establishment of low-level westerlies over the SCS and facilitating
403 an earlier monsoon transition.

404 In contrast, late onset years (Figs. 4c–d) are associated with warm and cold SST anomalies over
405 the central–eastern equatorial Pacific and the WNP, respectively, during the preceding winter,
406 indicative of El Niño conditions, together with warming in the tropical Indian Ocean. The
407 corresponding low-level wind anomalies show strengthened easterlies over the western Pacific,
408 reinforcing the WNPSH and suppressing westerly development over the SCS. These conditions



409 persist into spring, delaying the seasonal weakening of the WNPSH and inhibiting the establishment
410 of monsoon-type circulation over the SCS, thereby favoring a later onset.

411 Superimposed on this ENSO-modulated seasonal preconditioning, the abrupt reversal of low-
412 level winds at the onset pentad and the rapid strengthening of westerlies in the subsequent pentads
413 (Fig. 3) highlight the critical role of subseasonal disturbances in triggering SCSSM onset once the
414 large-scale background state becomes favorable. Intraseasonal oscillations, particularly those
415 associated with the Madden–Julian Oscillation, have been shown to modulate low-level westerlies,
416 convection, and moisture convergence over the SCS, providing the immediate dynamical trigger for
417 monsoon onset (Wang et al., 2009). The coexistence of gradual ENSO-related background evolution
418 and abrupt synoptic-scale transitions underscores that SCSSM onset emerges from the interaction
419 between low-frequency boundary forcing and high-frequency atmospheric variability. It suggests that
420 the synoptic clustering–based onset definition naturally captures this multi-timescale behavior by
421 linking onset to both a persistent circulation regime transition and the large-scale background state.
422 This makes the proposed definition suitable for applications to seasonal prediction.

423 **4. Application to seasonal prediction of SCSSM onset**

424 This section evaluates the applicability of the synoptic clustering–based SCSSM onset definition
425 for seasonal prediction using SEAS5 forecasts. Prediction skill is assessed separately for a dependent
426 training period (1981–2016) and an independent forecast period (2017–2024) to examine both in-
427 sample performance and out-of-sample robustness. Results obtained using the clustering-based
428 definition of the present study (NL26) are systematically compared with those based on the
429 conventional U850 index of Wang et al. (2004; W04).

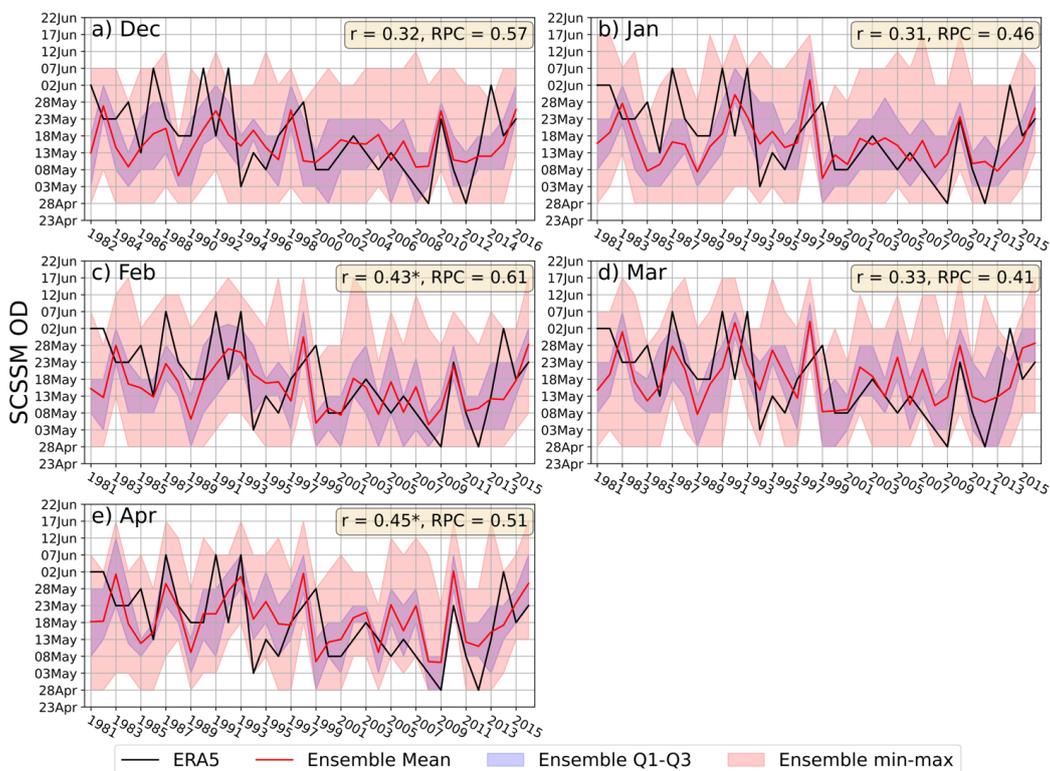
430 **4.1 Prediction skill during the training (hindcast) period (1981–2016)**

431 First, performance of the ECMWF SEAS5 seasonal prediction system in forecasting SCSSM
432 OD using both the W04 and NL26 definitions is evaluated during the dependent training period
433 (1981–2016). Forecasts initialized from December (~five-month lead time) to April (one-month lead



434 time) are examined to assess lead-time dependence, and both deterministic and probabilistic skill
435 metrics are considered.

436 4.1.1 Deterministic prediction skill

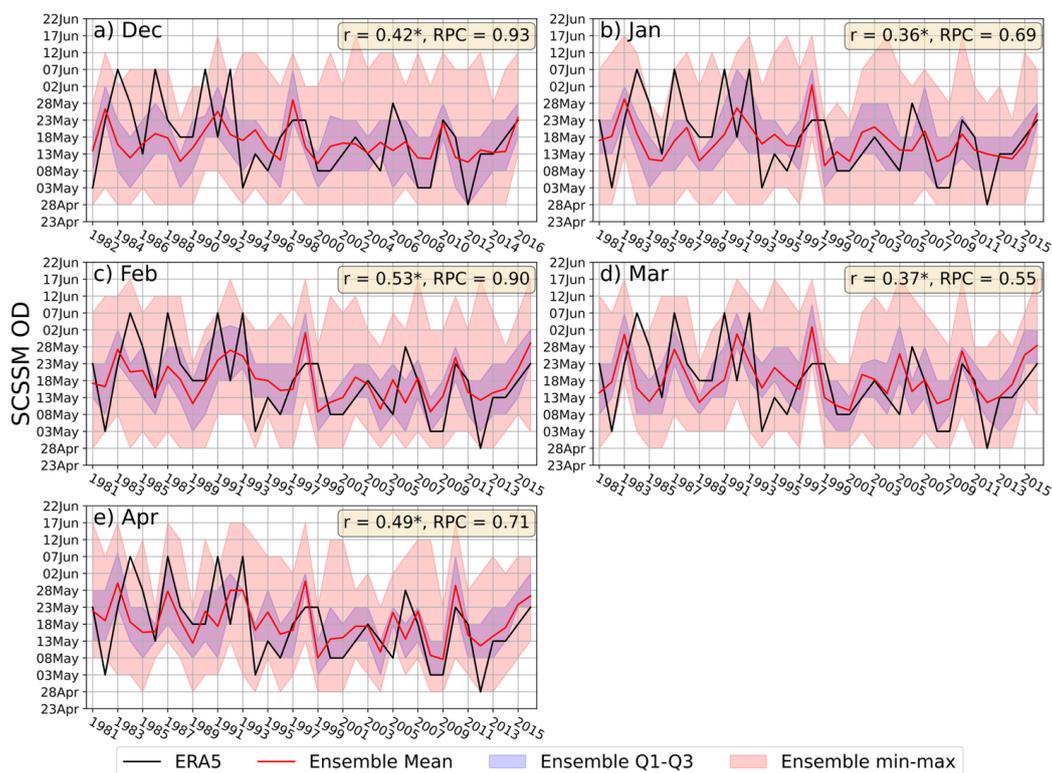


437
438 **Figure 5:** Hindcast performance of the ECMWF SEAS5 system in predicting SCSSM onset during the
439 dependent training period (1981–2016) using the U850-based definition of Wang et al. (2004). Time series
440 show observed onset dates from ERA5 (black) and ensemble-mean predictions (red) for forecasts initialized
441 from December to April. Blue shading denotes the interquartile range (Q1–Q3), while red shading indicates
442 the full ensemble spread. The Pearson correlation coefficient (r) and ratio of predictable component (RPC) are
443 shown for each initialization month. The * denotes r is statistically significant at the 95% confidence level.

444 Figure 5 shows the deterministic prediction skill of SEAS5 when SCSSM onset is defined using
445 the W04. Correlation coefficients (r) between ensemble-mean predictions and observed onset dates
446 range from 0.31 to 0.45, with modest improvements toward shorter lead times, and statistically



447 significant correlations obtained only for forecasts initialized in February and April. The ratio of
448 predictable component (RPC) generally remains below 0.65, indicating that a substantial fraction of
449 ensemble spread is not associated with the predictable signal. This behavior suggests that OD defined
450 purely by a local zonal-wind threshold are sensitive to short-lived wind reversals and synoptic noise,
451 which limits deterministic predictability even during the dependent period.



452

453 **Figure 6:** Same as Fig. 5, but for SCSSM OD defined using the synoptic clustering–based definition (NL26).

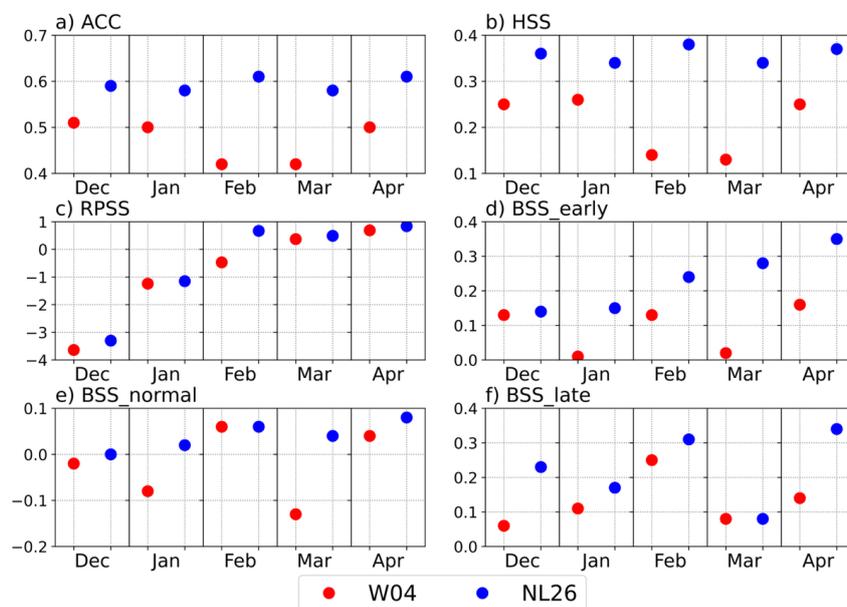
454 In contrast, Figure 6 demonstrates a clear improvement in deterministic forecast skill when the
455 NL26 definition is applied. Correlation coefficients increase to 0.36–0.53 for all initialization
456 months, with statistically significant correlations throughout and maximum skill for February
457 initialization ($r = 0.53$) and April initialization ($r = 0.49$). Importantly, skill improvements are already
458 evident at long lead times: for December initialization, the correlation increases from 0.32 (W04) to



459 0.42 (NL26), indicating enhanced sensitivity to slowly varying boundary forcing when onset is
 460 defined in terms of circulation regimes rather than a single-point wind index.

461 Even more striking is the improvement in potential predictability. RPC values increase
 462 substantially under the NL26 definition, reaching 0.90–0.93 for December and February
 463 initializations and remaining above 0.55 for all start months. These near-unity RPC values indicate
 464 that the ensemble mean captures most of the predictable signal and that ensemble dispersion is more
 465 consistent with signal amplitude. This result suggests that the clustering-based definition effectively
 466 filters out unpredictable synoptic-scale fluctuations and aligns the onset metric more closely with
 467 large-scale circulation features that are inherently predictable on seasonal timescales.

468 **4.1.2 Categorical and probabilistic prediction skill**



469
 470 **Figure 7:** Categorical and probabilistic prediction skill of SEAS5 for SCSSM OD during the dependent
 471 training period (1981–2016), comparing the U850-based definition of Wang et al. (2004; W04) and the
 472 synoptic clustering–based definition (NL26). Shown are (a) Accuracy (ACC), (b) Heidke Skill Score (HSS),
 473 (c) Ranked Probability Skill Score (RPSS), and Brier Skill Scores (BSS) for (d) early, (e) normal, and (f) late
 474 onset categories, for forecasts initialized from December to April.



475 Figure 7 further compares forecast skill between the W04 and NL26 definitions using
476 categorical and probabilistic metrics, including Accuracy (ACC), Heidke Skill Score (HSS), Ranked
477 Probability Skill Score (RPSS), and Brier Skill Scores (BSS) for early, normal, and late onset
478 categories. For deterministic categorical forecasts (Figs. 7a–b), NL26 consistently outperforms W04
479 across all initialization months. ACC increases from approximately 0.42–0.51 under W04 to 0.58–
480 0.61 under NL26, while HSS improves from values near 0.13–0.26 to 0.34–0.38. These
481 improvements indicate both higher hit rates and substantially greater skill relative to random chance,
482 confirming that the regime-based onset classification yields more robust categorical predictions.

483 Differences are even more pronounced for probabilistic skill metrics. Under the W04 definition,
484 RPSS values are negative or near zero for December–January initializations, indicating little to no
485 skill relative to climatology (Fig. 7c). In contrast, RPSS becomes positive for all initialization months
486 under NL26, with values approaching unity for February to April initializations, reflecting a marked
487 improvement in the reliability and discrimination of ensemble-based probabilistic forecasts when
488 onset is defined in terms of persistent circulation regimes. Consistent improvements are also evident
489 in category-specific BSS (Figs. 7d–f). For early and late onset categories, NL26 yields systematically
490 higher BSS values across all lead times, indicating enhanced probabilistic resolution for extreme
491 onset years. Skill for the normal category remains modest for both definitions, reflecting the intrinsic
492 difficulty of predicting near-average onset timing; nevertheless, NL26 generally maintains neutral to
493 positive BSS, whereas W04 often exhibits negative skill.

494 Overall, results from the dependent training period demonstrate that defining SCSSM onset
495 using the NL26 enhances deterministic, categorical, and probabilistic prediction skill relative to the
496 WL04. By explicitly incorporating circulation persistence and regime transitions, the clustering-
497 based definition aligns onset timing more closely with large-scale, predictable circulation features
498 represented in SEAS5, leading to higher correlations, improved ensemble reliability, and more
499 skillful probabilistic forecasts. While these results establish the potential of the proposed definition
500 under in-sample conditions, it is essential to assess whether the skill improvements persist under true

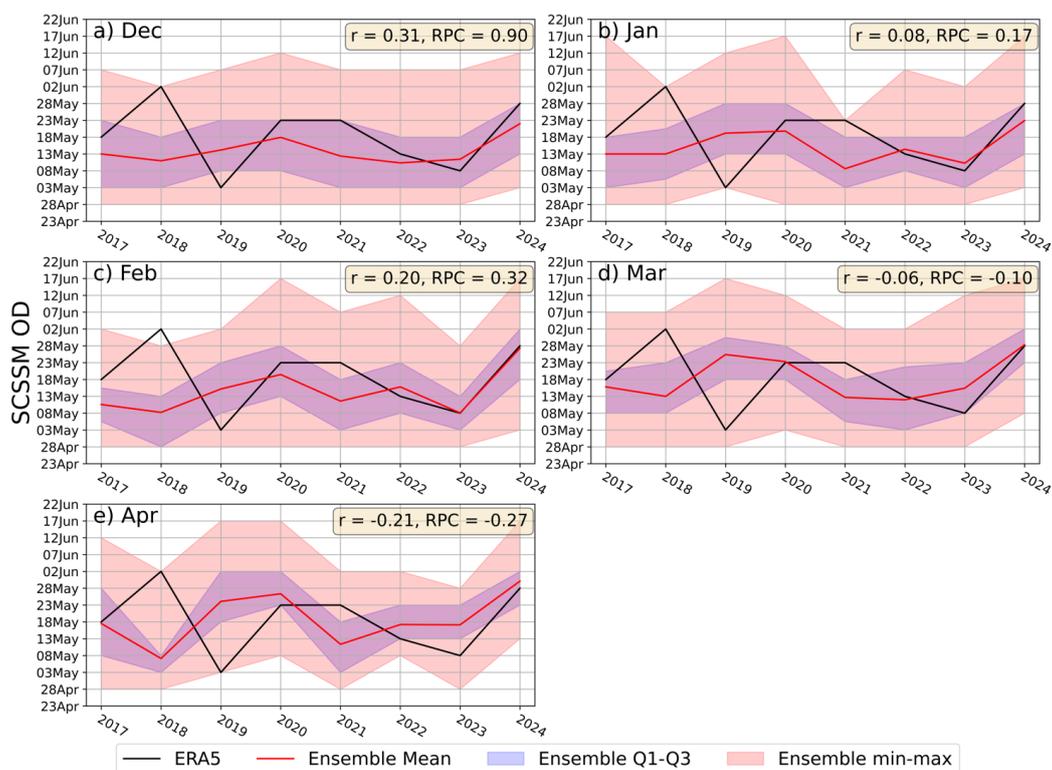


501 out-of-sample prediction. This is examined in the following subsection using independent SEAS5
502 forecasts for the period 2017–2024.

503 4.2 Prediction skill during the independent forecast period (2017–2024)

504 This subsection evaluates the out-of-sample performance of the synoptic clustering–based
505 SCSSM onset definition using independent SEAS5 forecasts for 2017–2024, which were not
506 included in the training of the SOM circulation regimes.

507 4.2.1 Deterministic prediction skill

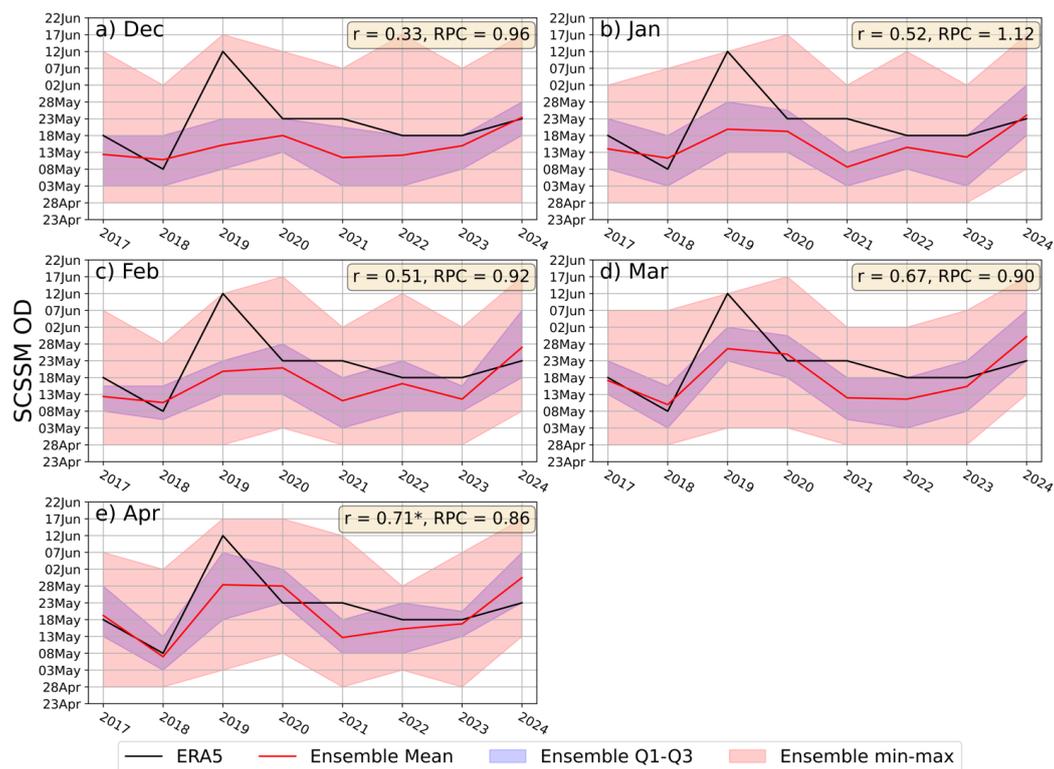


508

509 **Figure 8:** Forecast performance of the ECMWF SEAS5 system in predicting SCSSM onset during the
510 independent forecast period (2017–2024) using the U850-based definition of Wang et al. (2004). Time series
511 show ERA5 onset dates (black) and ensemble-mean forecasts (red), with blue shading denotes the
512 interquartile range (Q1–Q3) and red shading indicates the full ensemble spread. The Pearson correlation
513 coefficient (r) and ratio of predictable component (RPC) are shown for each initialization month. The *
514 denotes r is statistically significant at the 95% confidence level.



515 As expected for an independent forecast period with a limited sample size, overall deterministic
516 prediction skill is reduced compared to the dependent period (Fig. 8). When SCSSM OD is defined
517 using W04, correlations between SEAS5 ensemble-mean and ERA5 are generally weak and often
518 not statistically significant, particularly for forecasts initialized from January onward. In several
519 cases, both the correlation coefficient and the RPC are close to zero or negative, indicating little
520 usable predictability. Meanwhile, the NL26 yields consistently higher r and markedly improved RPC
521 values across most initialization months (Fig. 9). Forecasts initialized from January to April show the
522 largest gains, with r exceeding 0.5 in several months and RPC values approaching or exceeding
523 unity, indicating that a substantial fraction of the predictable signal is captured by the ensemble
524 mean. These improvements suggest that the clustering-based definition better aligns forecast
525 variability with predictable large-scale circulation features, even under out-of-sample conditions.

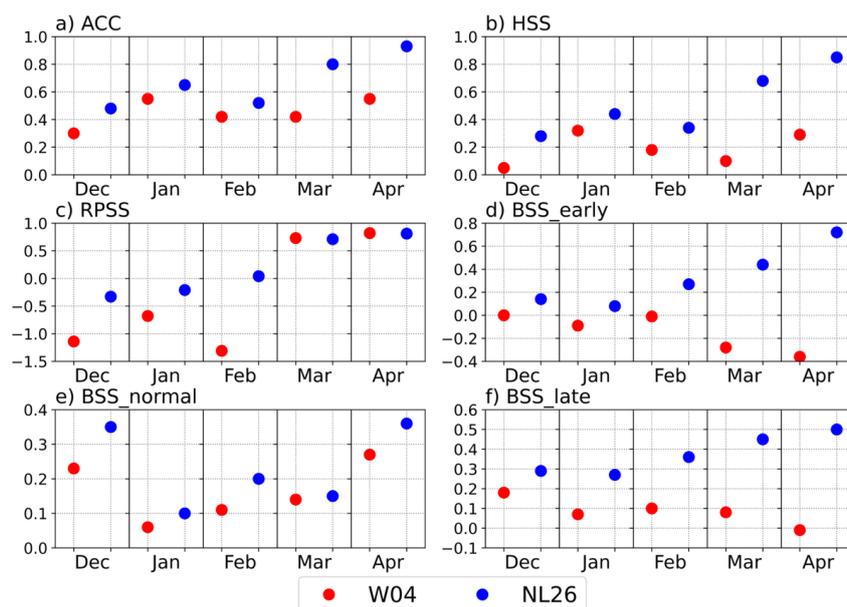


526

527 **Figure 9:** Same as Fig. 8, but for SCSSM OD defined using the synoptic clustering-based definition (NL26).



528 **4.2.2 Categorical and probabilistic prediction skill**



529

530 **Figure 10:** Categorical and probabilistic prediction skill of SEAS5 for SCSSM onset during the independent
531 forecast period (2017–2024), comparing the U850-based definition of Wang et al. (2004; W04) and the
532 synoptic clustering–based definition (NL26). Shown are (a) Accuracy (ACC), (b) Heidke Skill Score (HSS),
533 (c) Ranked Probability Skill Score (RPSS), and Brier Skill Scores (BSS) for (d) early, (e) normal, and (f) late
534 onset categories, for forecasts initialized from December to April.

535 Categorical skill metrics further highlight the advantages of the synoptic clustering–based
536 definition during the independent forecast period (Fig. 10). For deterministic ensemble-mean
537 forecasts, both ACC and HSS obtained with NL26 consistently exceed those derived from W04
538 across all initialization months, with particularly pronounced improvements for forecasts initialized
539 in March and April. HSS values under W04 are generally small and occasionally near zero,
540 indicating limited skill relative to climatology, whereas NL26 yields positive and substantially larger
541 HSS values, demonstrating improved discrimination among early, normal, and late onset categories.

542 Differences are even more pronounced for probabilistic prediction skill. RPSS values based on
543 W04 are frequently negative or close to zero, especially for early-season initializations, indicating



544 that probabilistic forecasts often perform no better than climatology. In contrast, RPSS values under
545 NL26 are positive for most initialization months and increase systematically toward shorter lead
546 times, reflecting enhanced reliability and resolution of ensemble-based probabilities. Category-
547 specific BSS further confirm that the clustering-based definition improves probabilistic skill across
548 all onset categories. The largest gains occur for early and late onset categories, which are most
549 strongly linked to persistent circulation regime transitions and therefore benefit directly from the
550 regime-based onset criteria. Skill for the “normal” category is also generally higher under NL26,
551 although improvements are more modest, consistent with the broader spread and inherently lower
552 predictability of near-climatological onset timing.

553 Overall, results from the independent forecast period indicate that the synoptic clustering-based
554 definition of SCSSM onset consistently outperforms the conventional U850-based index across
555 deterministic, categorical, and probabilistic metrics, despite the expected reduction in absolute skill
556 under true out-of-sample conditions. The clustering-based definition yields higher correlations,
557 improved ensemble reliability, and more skillful categorical and probabilistic forecasts,
558 demonstrating its robustness and transferability beyond the training period.

559 **5. Discussions and Conclusions**

560 **5.1 Discussions**

561 This study demonstrates that defining SCSSM OD in terms of persistent synoptic circulation
562 regimes substantially improves seasonal prediction skill compared with traditional index-based
563 definitions. The superior performance of the NL26 approach can be understood by considering both
564 the physical processes governing monsoon onset and the intrinsic limits of seasonal predictability.
565 Conventional onset definitions, such as the WL04, rely on local threshold exceedance of zonal winds
566 and are therefore highly sensitive to short-lived wind reversals associated with transient synoptic
567 disturbances. Such events may satisfy index criteria without representing a genuine transition into the
568 summer monsoon circulation, introducing noise into diagnosed OD and degrading forecast



569 verification. This sensitivity is reflected in the relatively modest deterministic correlations and low
570 RPC obtained using the W04 definition, particularly at longer lead times. In contrast, the synoptic
571 clustering-based definition explicitly links SCSSM onset to a regime transition in the large-scale
572 circulation, characterized by the establishment and persistence of low-level westerlies over SCS and
573 the eastward retreat of the WNPSH. By requiring both continuity and persistence of monsoon-type
574 regimes, the definition filters out spurious or short-lived events and emphasizes circulation structures
575 that evolve on longer timescales and are therefore more predictable. This physically grounded
576 framing explains the systematic improvements in deterministic, categorical, and probabilistic skill
577 metrics obtained with the clustering-based definition.

578 The results further underscore the multi-timescale nature of SCSSM onset. On interannual
579 timescales, slowly varying boundary forcing, most notably ENSO-related SST anomalies, modulates
580 the large-scale Indo-Pacific circulation and preconditions the timing of low-level westerly
581 establishment over the SCS. This source of predictability is effectively captured by SEAS5 when
582 onset is defined in terms of circulation regimes. Superimposed on this background evolution,
583 subseasonal disturbances, including intraseasonal oscillations, provide the immediate dynamical
584 trigger for onset once the large-scale environment becomes favorable. The clustering-based
585 definition naturally accommodates this interaction by identifying onset as a sustained regime shift
586 rather than a single instantaneous threshold crossing.

587 **5.2 Conclusions**

588 This study introduces a synoptic clustering-based definition of SCSSM onset that conceptualizes
589 onset as a persistent transition between large-scale circulation regimes, rather than a local or
590 instantaneous wind threshold. By identifying onset through sustained occupation of monsoon-type
591 circulation patterns derived from self-organizing maps and clustering of low-level fields, the
592 proposed definition reduces sensitivity to short-lived synoptic fluctuations and more faithfully
593 represents the dynamical evolution of the monsoon system. Application to ECMWF SEAS5 forecasts



594 demonstrates that the clustering-based definition yields systematic and robust improvements in
595 seasonal prediction skill relative to the conventional U850-based index of Wang et al. (2004). Skill
596 gains are evident across deterministic, potential predictability, categorical, and probabilistic metrics
597 during the dependent training period (1981–2016) and persist during the independent forecast period
598 (2017–2024), indicating that the approach generalizes well beyond the training sample.

599 The improved forecast performance is physically consistent with the dominant controls on
600 SCSSM onset variability, which arise from the interaction between slowly varying boundary forcing
601 and higher-frequency atmospheric variability. By emphasizing circulation persistence and structural
602 maturity, the regime-based framework more effectively isolates the predictable component of onset
603 variability and improves ensemble reliability, particularly for early and late onset events. These
604 results indicate that predictability-oriented, circulation-regime-based definitions offer a physically
605 meaningful and practically useful alternative to traditional index-based approaches for diagnosing
606 and forecasting monsoon onset. Although this study focuses on the SCSSM, the methodology is
607 readily applicable to other monsoon systems and climate transition processes in which circulation
608 persistence plays a key role in predictability. Future work should also explore the integration of
609 advanced pattern-recognition approaches, including deep learning, to further enhance regime
610 identification and forecast skill.

611 **Supplementary Information** The manuscript contains supplementary figures.

612 **Acknowledgments** The author gratefully acknowledges Dr. Le Duc for valuable guidance and
613 assistance with the implementation of the SOM and K-means clustering techniques.

614 **Fundings** This research did not receive any specific grant from funding agencies in the public,
615 commercial, or not-for-profit sectors.

616 **Data availability** All datasets analyzed in this study are publicly available from the following
617 sources: ERA5 (ECMWF; <https://www.ecmwf.int/en/forecasts/dataset/ecmwf-reanalysis-v5>);



618 HadISST (UK Met Office Hadley Centre; <https://www.metoffice.gov.uk/hadobs/hadisst>); SEAS5
619 (ECMWF; <https://cds.climate.copernicus.eu/datasets>).

620 **Code availability** The code used in this study is available from the corresponding author upon
621 reasonable request.

622 **Author contributions** The author conceived the study, performed the analysis, and wrote the
623 manuscript.

624 **Declarations**

625 **Competing interests** The author has no conflicts of interest to declare.

626 **Ethical approval** Ethical approval is not applicable to this article.

627 **References**

628 Attada R, Ehsan MA, Pillai PA (2022) Evaluation of Potential Predictability of Indian Summer
629 Monsoon Rainfall in ECMWF's Fifth-Generation Seasonal Forecast System (SEAS5). Pure Appl
630 Geophys 179:4639–4655. <https://doi.org/10.1007/s00024-022-03184-9>

631 Bombardi, R. J., Pegion, K. V., Kinter, J. L., Cash, B. A., and Adams, J. M.: Sub-seasonal
632 Predictability of the Onset and Demise of the Rainy Season over Monsoonal Regions, Front. Earth
633 Sci., 5, <https://doi.org/10.3389/feart.2017.00014>, 2017.

634 Bombardi, R. J., Moron, V., and Goodnight, J. S.: Detection, variability, and predictability of
635 monsoon onset and withdrawal dates: A review, Intl Journal of Climatology, 40, 641–667,
636 <https://doi.org/10.1002/joc.6264>, 2020.

637 Borah, N., Sahai, A. K., Chattopadhyay, R., Joseph, S., Abhilash, S., and Goswami, B. N.: A self-
638 organizing map–based ensemble forecast system for extended range prediction of active/break cycles



- 639 of Indian summer monsoon, *JGR Atmospheres*, 118, 9022–9034, <https://doi.org/10.1002/jgrd.50688>,
640 2013.
- 641 Bui-Minh, T., Doan, Q.-V., Nguyen, K.-C., Phan-Van, T., Trinh-Tuan, L., Cong, T., and Trinh-
642 Minh, N.: Determining the onset of summer rainfall over Vietnam using self-organizing maps, *Clim*
643 *Dyn*, 62, 9189–9206, <https://doi.org/10.1007/s00382-024-07385-x>, 2024.
- 644 Chen, T.-C., Tsay, J.-D., Matsumoto, J., and Alpert, J.: Impact of the Summer Monsoon Westerlies
645 on the South China Sea Tropical Cyclone Genesis in May, *Weather and Forecasting*, 32, 925–947,
646 <https://doi.org/10.1175/WAF-D-16-0189.1>, 2017.
- 647 Chevuturi, A., Turner, A. G., Woolnough, S. J., Martin, G. M., and MacLachlan, C.: Indian summer
648 monsoon onset forecast skill in the UK Met Office initialized coupled seasonal forecasting system
649 (GloSea5-GC2), *Clim Dyn*, 52, 6599–6617, <https://doi.org/10.1007/s00382-018-4536-1>, 2019.
- 650 Chevuturi, A., Turner, A. G., Johnson, S., Weisheimer, A., Shonk, J. K. P., Stockdale, T. N., and
651 Senan, R.: Forecast skill of the Indian monsoon and its onset in the ECMWF seasonal forecasting
652 system 5 (SEAS5), *Clim Dyn*, 56, 2941–2957, <https://doi.org/10.1007/s00382-020-05624-5>, 2021.
- 653 Dai, L., Cheng, T. F., and Lu, M.: Define East Asian Monsoon Annual Cycle via a Self-Organizing
654 Map-Based Approach, *Geophysical Research Letters*, 48, e2020GL089542,
655 <https://doi.org/10.1029/2020GL089542>, 2021.
- 656 Ding, Y.: SEASONAL MARCH OF THE EAST-ASIAN SUMMER MONSOON, in: *World*
657 *Scientific Series on Asia-Pacific Weather and Climate*, vol. 02, WORLD SCIENTIFIC, 3–53,
658 https://doi.org/10.1142/9789812701411_0001, 2004.
- 659 Ding, Y., Li, C., and Liu, Y.: Overview of the South China sea monsoon experiment, *Adv. Atmos.*
660 *Sci.*, 21, 343–360, <https://doi.org/10.1007/BF02915563>, 2004.
- 661 Ding, Y. and Chan, J. C. L.: The East Asian summer monsoon: an overview, *Meteorol. Atmos.*
662 *Phys.*, 89, 117–142, <https://doi.org/10.1007/s00703-005-0125-z>, 2005.



- 663 Ding, Y., Liu, Y., Song, Y., and Zhang, J.: From MONEX to the global monsoon: A review of
664 monsoon system research, *Adv. Atmos. Sci.*, 32, 10–31, <https://doi.org/10.1007/s00376-014-0008-7>,
665 2015.
- 666 Eade, R., Smith, D., Scaife, A., Wallace, E., Dunstone, N., Hermanson, L., and Robinson, N.: Do
667 seasonal-to-decadal climate predictions underestimate the predictability of the real world?,
668 *Geophysical Research Letters*, 41, 5620–5628, <https://doi.org/10.1002/2014GL061146>, 2014.
- 669 Geen, R.: Forecasting South China Sea Monsoon Onset Using Insight From Theory, *Geophysical*
670 *Research Letters*, 48, e2020GL091444, <https://doi.org/10.1029/2020GL091444>, 2021.
- 671 He, J. and Zhu, Z.: The relation of South China Sea monsoon onset with the subsequent rainfall over
672 the subtropical East Asia, *Intl Journal of Climatology*, 35, 4547–4556,
673 <https://doi.org/10.1002/joc.4305>, 2015.
- 674 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J.,
675 Peubey, C., Radu, R., Schepers, D., Simmons, A., Soci, C., Abdalla, S., Abellan, X., Balsamo, G.,
676 Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., De Chiara, G., Dahlgren, P., Dee, D., Diamantakis,
677 M., Dragani, R., Flemming, J., Forbes, R., Fuentes, M., Geer, A., Haimberger, L., Healy, S., Hogan,
678 R. J., Hólm, E., Janisková, M., Keeley, S., Laloyaux, P., Lopez, P., Lupu, C., Radnoti, G., De
679 Rosnay, P., Rozum, I., Vamborg, F., Villaume, S., and Thépaut, J.: The ERA5 global reanalysis,
680 *Quart J Royal Meteor Soc*, 146, 1999–2049, <https://doi.org/10.1002/qj.3803>, 2020.
- 681 Hu, P., Chen, W., Huang, R., and Nath, D.: On the weakening relationship between the South China
682 Sea summer monsoon onset and cross-equatorial flow after the late 1990s, *Intl Journal of*
683 *Climatology*, 38, 3202–3208, <https://doi.org/10.1002/joc.5472>, 2018.
- 684 Hu, P., Chen, W., Chen, S., Liu, Y., Wang, L., and Huang, R.: The Leading Mode and Factors for
685 Coherent Variations among the Subsystems of Tropical Asian Summer Monsoon Onset, *Journal of*
686 *Climate*, 35, 1597–1612, <https://doi.org/10.1175/JCLI-D-21-0101.1>, 2022a.



- 687 Hu, P., Chen, W., Li, Z., Chen, S., Wang, L., and Liu, Y.: Close Linkage of the South China Sea
688 Summer Monsoon Onset and Extreme Rainfall in May over Southeast Asia: Role of the Synoptic-
689 Scale Systems, *Journal of Climate*, 35, 4347–4362, <https://doi.org/10.1175/JCLI-D-21-0740.1>,
690 2022b.
- 691 Huangfu, J., Huang, R., and Chen, W.: Relationship between the South China Sea summer monsoon
692 onset and tropical cyclone genesis over the western North Pacific, *Intl Journal of Climatology*, 37,
693 5206–5210, <https://doi.org/10.1002/joc.5141>, 2017.
- 694 Jiang, X., Wang, Z., and Li, Z.: Signature of the South China Sea summer monsoon onset on spring-
695 to-summer transition of rainfall in the middle and lower reaches of the Yangtze River basin, *Clim
696 Dyn*, 51, 3785–3796, <https://doi.org/10.1007/s00382-018-4110-x>, 2018.
- 697 Johnson, S. J., Stockdale, T. N., Ferranti, L., Balmaseda, M. A., Molteni, F., Magnusson, L.,
698 Tietsche, S., Decremmer, D., Weisheimer, A., Balsamo, G., Keeley, S. P. E., Mogensen, K., Zuo, H.,
699 and Monge-Sanz, B. M.: SEAS5: the new ECMWF seasonal forecast system, *Geoscientific Model
700 Development*, 12, 1087–1117, <https://doi.org/10.5194/gmd-12-1087-2019>, 2019.
- 701 Li, X., Wang, L., Zhong, S., and Liu, L.: Comparative analysis of indices in capturing the onset and
702 withdrawal of the South Asian Summer Monsoon, *Environ. Res. Commun.*, 6, 031007,
703 <https://doi.org/10.1088/2515-7620/ad352b>, 2024.
- 704 Liu, B., Zhu, C., Yuan, Y., and Xu, K.: Two Types of Interannual Variability of South China Sea
705 Summer Monsoon Onset Related to the SST Anomalies before and after 1993/94, *Journal of
706 Climate*, 29, 6957–6971, <https://doi.org/10.1175/JCLI-D-16-0065.1>, 2016.
- 707 Martin, G. M., Chevuturi, A., Comer, R. E., Dunstone, N. J., Scaife, A. A., and Zhang, D.:
708 Predictability of South China Sea Summer Monsoon Onset, *Adv. Atmos. Sci.*, 36, 253–260,
709 <https://doi.org/10.1007/s00376-018-8100-z>, 2019.



- 710 Nguyen-Le, D. and Yamada, T. J.: Using Weather Pattern Recognition to Classify and Predict
711 Summertime Heavy Rainfall Occurrence over the Upper Nan River Basin, Northwestern Thailand,
712 *Weather and Forecasting*, 34, 345–360, <https://doi.org/10.1175/WAF-D-18-0122.1>, 2019.
- 713 Nguyen-Le, D., Yamada, T. J., and Tran-Anh, D.: Classification and forecast of heavy rainfall in
714 northern Kyushu during Baiu season using weather pattern recognition, *Atmospheric Science Letters*,
715 18, 324–329, <https://doi.org/10.1002/asl.759>, 2017.
- 716 Nishiyama, K., Endo, S., Jinno, K., Bertacchi Uvo, C., Olsson, J., and Berndtsson, R.: Identification
717 of typical synoptic patterns causing heavy rainfall in the rainy season in Japan by a Self-Organizing
718 Map, *Atmospheric Research*, 83, 185–200, <https://doi.org/10.1016/j.atmosres.2005.10.015>, 2007.
- 719 Rayner, N. A., Parker, D. E., Horton, E. B., Folland, C. K., Alexander, L. V., Rowell, D. P., Kent, E.
720 C., and Kaplan, A.: Global analyses of sea surface temperature, sea ice, and night marine air
721 temperature since the late nineteenth century, *J. Geophys. Res.*, 108, 2002JD002670,
722 <https://doi.org/10.1029/2002JD002670>, 2003.
- 723 Sammon, J. W.: A Nonlinear Mapping for Data Structure Analysis, *IEEE Transactions on*
724 *Computers*, C–18, 401–409, <https://doi.org/10.1109/T-C.1969.222678>, 1969.
- 725 Ultsch, A. and Siemon, H. P: Kohonen’s Self Organizing Feature Maps for Exploratory Data
726 Analysis. in *Proceedings of the International Neural Network Conference (INNC-90)*, Paris, France,
727 July 9–13, 1990 1. Dordrecht, Netherlands (eds. Widrow, B. & Angeniol, B.) 1, 305–308 (Kluwer
728 Academic Press, 1990).
- 729 Vesanto, J. and Alhoniemi, E.: Clustering of the self-organizing map, *IEEE Transactions on Neural*
730 *Networks*, 11, 586–600, <https://doi.org/10.1109/72.846731>, 2000.
- 731 Wang, B., Huang, F., Wu, Z., Yang, J., Fu, X., and Kikuchi, K.: Multi-scale climate variability of the
732 South China Sea monsoon: A review, *Dynamics of Atmospheres and Oceans*, 47, 15–37,
733 <https://doi.org/10.1016/j.dynatmoce.2008.09.004>, 2009.



- 734 Wang, B. and LinHo: Rainy Season of the Asian–Pacific Summer Monsoon*, *J. Climate*, 15, 386–
735 398, [https://doi.org/10.1175/1520-0442\(2002\)015%3C0386:RSOTAP%3E2.0.CO;2](https://doi.org/10.1175/1520-0442(2002)015%3C0386:RSOTAP%3E2.0.CO;2), 2002.
- 736 Wang, B., LinHo, Zhang, Y., and Lu, M.-M.: Definition of South China Sea Monsoon Onset and
737 Commencement of the East Asia Summer Monsoon*, *Journal of Climate*, 17, 699–710,
738 <https://doi.org/10.1175/2932.1>, 2004.
- 739 Wang, B., Ding, Q., Fu, X., Kang, I., Jin, K., Shukla, J., and Doblas-Reyes, F.: Fundamental
740 challenge in simulation and prediction of summer monsoon rainfall, *Geophysical Research Letters*,
741 32, 2005GL022734, <https://doi.org/10.1029/2005GL022734>, 2005.
- 742 Wang, L. and Chen, G.: Relationship between South China Sea summer monsoon onset and
743 landfalling tropical cyclone frequency in China, *Intl Journal of Climatology*, 38, 3209–3214,
744 <https://doi.org/10.1002/joc.5485>, 2018.
- 745 World Climate Research Programme (WCRP): Forecast verification—issues, methods, and FAQ,
746 <http://www.cawcr.gov.au/projects/verification> (last access: 25 January 2026), 2015.
- 747 Weigel, A. P., Liniger, M. A., and Appenzeller, C.: The Discrete Brier and Ranked Probability Skill
748 Scores, *Monthly Weather Review*, 135, 118–124, <https://doi.org/10.1175/MWR3280.1>, 2007.
- 749 Xie, S.-P., Hu, K., Hafner, J., Tokinaga, H., Du, Y., Huang, G., and Sampe, T.: Indian Ocean
750 Capacitor Effect on Indo–Western Pacific Climate during the Summer following El Niño, *Journal of*
751 *Climate*, 22, 730–747, <https://doi.org/10.1175/2008JCLI2544.1>, 2009.
- 752 Xie, S.-P., Kosaka, Y., Du, Y., Hu, K., Chowdary, J. S., and Huang, G.: Indo-western Pacific ocean
753 capacitor and coherent climate anomalies in post-ENSO summer: A review, *Adv. Atmos. Sci.*, 33,
754 411–432, <https://doi.org/10.1007/s00376-015-5192-6>, 2016.
- 755 Zhou, W. and Chan, J. C. L.: ENSO and the South China Sea summer monsoon onset, *Intl Journal of*
756 *Climatology*, 27, 157–167, <https://doi.org/10.1002/joc.1380>, 2007.



- 757 Zhu, Z. and Li, T.: Empirical prediction of the onset dates of South China Sea summer monsoon,
758 Clim Dyn, 48, 1633–1645, <https://doi.org/10.1007/s00382-016-3164-x>, 2017.