



1 **A physically guided deep learning reconstruction of terrestrial water storage anomalies at 0.1° across China**

2 Xueying Li ^{1*}, Yan Sun ^{2*}, Xihui Gu ³, Niko Wanders ⁴, Bridget R. Scanlon ⁵, and Louise J. Slater ¹

3 1. School of Geography and the Environment, University of Oxford, Oxford, UK

4 2. School of Computer Science, University of Sydney, Sydney, Australia

5 3. School of Geography and Information Engineering, China University of Geosciences, Wuhan,
6 China

7 4. Department of Physical Geography, Utrecht University, Utrecht, The Netherlands

8 5. Bureau of Economic Geology, Jackson School of Geosciences, The University of Texas at Austin,
9 Austin, USA

10 Correspondence: Xueying Li (xueying.li@ouce.ox.ac.uk); Yan Sun (ysun9899@uni.sydney.edu.au)

11 **Abstract**

12 Terrestrial water storage (TWS), comprising all surface and subsurface water components, is a key
13 indicator of water availability. The Gravity Recovery and Climate Experiment (GRACE) satellite
14 mission provides large-scale estimates of TWS anomalies (TWSA), but its coarse spatial resolution (3°,
15 approximately 300 km) limits the analysis of hydrologic processes at sub-regional scales. Using a
16 physically-guided deep learning framework, we downscale TWSA from the original 3° GRACE
17 mascons to 0.1° (approximately 10 km) across China, generating a standard version (2002–2019) with
18 comprehensive observations used for model constraints and independent evaluation and an extended
19 version (2020–2023) to support more recent hydrologic analyses. The downscaled TWSA preserves
20 large-scale GRACE signals at the 3° grid scale (median correlation coefficient (*CC*): 0.95; root-mean-
21 square error (*RMSE*): 1.38 cm) and basin scale (median *CC*: 0.94; *RMSE*: 1.72 cm), with a low median
22 uncertainty (0.88 cm) across China. Its reliability is supported by high consistency with physically
23 informed TWSA spatial patterns at the 0.1° resolution (median *CC*: 0.91) and internally consistent water
24 balance closure beyond the native GRACE resolution (median *CC*: 0.80; *RMSE*: 1.44 cm). Evaluation
25 against independent observations demonstrates that the downscaled TWSA agrees well with
26 groundwater variations in intensively irrigated regions (*CC*: 0.65 for irrigation intensity > 50%) and
27 annual glacier elevation change in cryospheric areas (*CC*: 0.97). The datasets improve fine-scale
28 characterization of TWS variability and associated hydrologic processes in China, and can be used as a
29 reference for evaluating performance of high-resolution hydrologic models. The two versions of the
30 dataset are available at <https://doi.org/10.5281/zenodo.19502906>.

31 **1 Introduction**

32 Terrestrial water storage (TWS), the sum of all surface and subsurface storage components, is a crucial
33 determinant of the global water budget (Immerzeel et al., 2020; Scanlon et al., 2023). TWS change
34 plays a key role in determining water availability (Rodell et al., 2018) and modulating water flux
35 interactions (Tapley et al., 2019), and is closely linked to hydro-climatic extremes such as droughts
36 (Long et al., 2013; Pokhrel et al., 2021) and floods (Reager et al., 2014; Sun et al., 2024). Over the past
37 two decades, TWS has changed substantially in response to climate change/variability (Guan et al.,
38 2023; Li et al., 2022; Rodell et al., 2018) and accelerating human pressures (Döll et al., 2014; Gu et al.,
39 2019; Guan et al., 2025). Monitoring TWS variations across different spatial scales is essential to
40 understand the response of the terrestrial water system to a changing environment, providing valuable
41 guidance for water resources management and climate change adaptation.

42 Despite its importance, TWS remains understudied relative to hydrologic fluxes, such as precipitation,
43 evapotranspiration (ET), and streamflow, primarily because TWS represents an integration of different
44 storage components and is difficult to measure directly in situ, particularly over large spatial domains.
45 Although hydrologic and land surface models simulate some water storage components (e.g., soil
46 moisture), they face challenges in explicitly resolving all TWS components and their interactions,
47 particularly for groundwater and glacier processes (Scanlon et al., 2018; Tiwari et al., 2025). Recent
48 advances in satellite gravimetry have provided an unprecedented opportunity to quantify TWS changes



49 at large spatial scales, with availability of 20+ years of TWS anomaly (TWSA) observations from the
50 Gravity Recovery and Climate Experiment (GRACE) and GRACE Follow-On (GRACE-FO) missions
51 (Tapley et al., 2019). Using GRACE and GRACE-FO retrievals, changes in TWS have been quantified
52 from regional to global scales (Li et al., 2022; Rodell et al., 2018; Scanlon et al., 2023), constraining
53 the estimation of individual storage components (Li et al., 2019a; Zhao et al., 2022), closure of the water
54 balance (Li et al., 2019b), and calibration of model parameters (Bai et al., 2018; Chen et al., 2017).
55 However, an inherent disadvantage of the GRACE data lies in its coarse spatial resolution (3° ;
56 approximately $110,000 \text{ km}^2$ at the Equator), linked to the design of the satellite orbit and accuracy of
57 the instruments (Rowlands et al., 2005; Wahr et al., 2006). Although post-processing (e.g.,
58 regularization algorithms) can be used to generate GRACE Level-3 products at nominal resolutions of
59 0.25° – 1° , the adjacent grid cells in these Level-3 products are not statistically independent and thus do
60 not represent an effective improvement in spatial resolution (Loomis et al., 2019; Save et al., 2016;
61 Wiese et al., 2016). The coarse resolution severely limits its application of satellite-based TWSA for
62 investigating hydrologic processes at sub-regional scales (Chen et al., 2022; Li et al., 2025), which are
63 often more sensitive to human intervention and climate variability than regional to continental scales.

64 Frameworks of downscaling GRACE/-FO observations are mainly categorized into dynamic (model-
65 based) and statistical (data-based) approaches (Pascal et al., 2022). Dynamic downscaling, such as data
66 assimilation methods, integrates GRACE observations into hydrologic or land surface models (Eicker
67 et al., 2014; Li et al., 2019a). Despite its physical interpretability, this approach entails substantial
68 computational demands. More critically, dynamic downscaling largely relies on the structure,
69 parameters, and process description of physical models, potentially distorting the original signals of
70 GRACE measurements (Gerdener et al., 2023). Statistical downscaling frameworks are more widely
71 adopted, using data-driven approaches such as machine learning (ML) to establish relationships
72 between GRACE/-FO TWSA observations and selected predictors. To date, downscaled TWSA datasets
73 have primarily been generated by regression-based ML algorithms, including multiple linear regression
74 (Pascal et al., 2022), partial least squares regression (Vishwakarma et al., 2021), random forest (Arshad
75 et al., 2024; Wang et al., 2024), support vector machine (Kalu et al., 2024; Yazdian et al., 2023), and
76 long-short term memory (Li and Kusche, 2026). These approaches link coarse-resolution TWSA with
77 aggregated predictors and subsequently infer fine-scale variability, assuming a consistent relationship
78 between TWSA and predictors across different spatial scales. One limitation of this assumption is that
79 spatially heterogeneous interactions among variables across scales are not explicitly represented, which
80 limits model generalization in regions with complex hydrologic processes. A recent application of
81 convolutional neural network (CNN) to TWSA downscaling (Gou and Soja, 2024) captures spatial
82 structural dependencies and multi-scale feature interactions through end-to-end learning, incorporating
83 both GRACE signals and physically informed TWSA patterns to constrain training processes.

84 Despite these efforts, current TWSA downscaling remains insufficient for sub-regional investigations,
85 particularly in regions subject to strong climate variability, intensive human interventions, and
86 heterogeneity in the land surface, such as China. Available downscaled TWSA products in China,
87 derived from both global and regional ML models, are provided at resolutions of 0.25° (Li and Kusche,
88 2026) and 0.5° (Gou and Soja, 2024; Vishwakarma et al., 2021; Xiong et al., 2025). The spatial
89 resolution of these products remains coarse for fine-scale analysis, particularly when compared with
90 commonly used hydrologic models and reanalysis datasets, such as PCR-GLOBWB 2 (Sutanudjaja et
91 al., 2018) and ERA5-Land (Muñoz-Sabater et al., 2021), both available at a finer resolution of 0.1°
92 (approximately 10 km). In addition, glacier dynamics are not considered in most datasets, although
93 glacier mass balance is a sensitive indicator of climate warming and a major contributor to TWS
94 changes in high mountain regions of Southwest and Northwest China (Hugonnet et al., 2021; Li et al.,
95 2022). Moreover, current downscaled TWSAs are primarily produced using climatic and hydrologic
96 predictors, which may limit their representation in regions with intensive human-water interactions,
97 such as highly irrigated areas relying on ground and surface water abstractions. Therefore, the relatively
98 coarse spatial resolution, omission of glacier dynamics, and lack of anthropogenic predictors call for
99 further improvements of TWSA downscaling in China, especially for sub-regional analysis.

100 Here we use a deep learning framework to downscale TWSA from the original 3° GRACE/-FO mascons
101 ($\text{TWSA}_{\text{GRACE}}$) to a much higher spatial resolution of 0.1° (TWSA_{DS}), incorporating physically informed



102 TWSA patterns ($TWSA_{PI}$) jointly derived from hydrologic simulations and satellite altimetry. A CNN-
103 based U-Net model was selected for multi-scale feature fusion, via its encoder-decoder architecture that
104 enhances generalization by isolating high-resolution details. Twelve predictors were included to
105 comprehensively characterize TWS variations, associated with the hydrologic cycle, land processes,
106 and human-water interactions. We provide a standard dataset of monthly 0.1° TWSA across China for
107 2002–2019, a period with comprehensive observations available for robust quality assessment. The gap
108 period between GRACE and GRACE-FO (July 2017–May 2018) is reconstructed using the U-Net deep
109 learning model. An extended dataset (2020–2023) is also provided to support more recent hydrologic
110 analyses, but due to limited data availability, glacier mass balance is not constrained by independent
111 satellite observations during 2020–2023 in the current version (Li and Sun, 2026a). The downscaled
112 0.1° TWSA was evaluated through: (1) large-scale comparisons with original GRACE/-FO observations;
113 (2) assessments against physically informed spatial patterns; (3) water balance closure across small to
114 medium catchments beyond the native GRACE resolution; (4) validation across irrigated and
115 glacierized regions using independent observations; and (5) comparisons with representative
116 downscaling studies (Gou and Soja, 2024; Xiong et al., 2025). The new datasets are valuable not only
117 for investigating sub-regional TWS changes and related hydrologic analyses in China, but also for
118 improving performance of high-resolution hydrologic models and providing insights into fine-scale
119 water resources management.

120 2 Study area and data processing

121 2.1 Study area

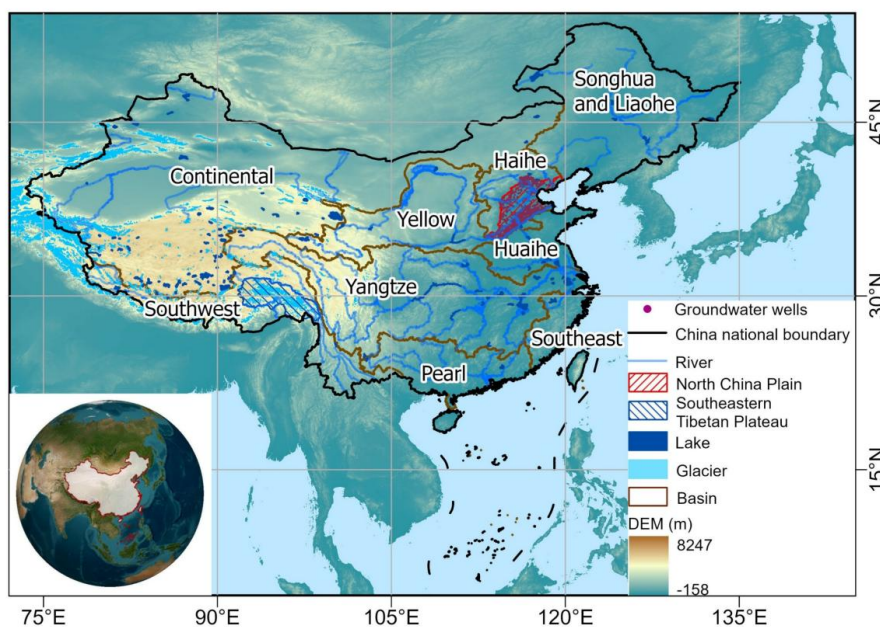
122 We generate high-resolution TWSA across China and surrounding regions (15° – 57° N, 69° – 138° E; Fig.
123 1). Covering approximately 9.6 million km^2 , China encompasses strong gradients in climate,
124 topography, and human activities (Liu et al., 2017; Piao et al., 2010). Nine major river basins are
125 delineated, including the Songhua-Liaohe, Haihe, Yellow, Huaihe, Yangtze, Pearl, Southeast, Southwest,
126 and Continental basins. Northern basins such as the Haihe and Yellow River basins are characterized
127 by limited water resources and intensive groundwater abstraction (Feng et al., 2013), whereas southern
128 basins including the Yangtze and Pearl are dominated by monsoon-driven precipitation and frequent
129 flooding (Zhang et al., 2006). In addition, the Tibetan Plateau and surrounding high mountain regions
130 act as major water sources for downstream basins, with glacier and snowpack contributing significantly
131 to regional water storage variability (Li et al., 2022; Li et al., 2025; Yao et al., 2022).

132 Recent climate change/variability and accelerating human pressures have jointly altered TWS across
133 China. Severe threats to water security include warming-driven glacier mass losses from high mountains
134 (Li et al., 2022; Li et al., 2025; Yao et al., 2022) and large-scale groundwater depletion from intensive
135 irrigation (Feng et al., 2013; Long et al., 2025; Long et al., 2020). These processes pose two major
136 challenges for hydrologic analysis: cryospheric dynamics and intensive human-water interactions,
137 which remain inadequately resolved. In this study, we selected the Southeastern Tibetan Plateau (SETP)
138 and North China Plain (NCP) as two representative subregions for assessing downscaled TWSA based
139 on independent observations. The SETP (28° – 32° N, 92° – 98° E) covers the Eastern Nyainqentanglha
140 Ranges, Eastern Himalayas, and Western Hengduan Mountains, and features the largest maritime
141 glaciers on the Tibetan Plateau with a mean elevation exceeding 4 km (Fig. 1). Glaciers over the SETP
142 are mainly distributed in the Parlung Tsangpo River basin and adjacent areas, where glacier meltwater
143 accounts for more than 50% of the total runoff in this region (Zhao et al., 2022). Under a warming
144 climate, glaciers of the SETP have been shrinking for decades, leading to significant losses in TWS
145 over this region.

146 In contrast to climate-driven TWS declines in the SETP, the NCP has emerged as a global hotspot of
147 groundwater depletion, largely driven by extensive irrigation supported by groundwater abstraction
148 (Scanlon et al., 2023). Here the NCP refers to the region bordered by the Yanshan Mountains in the
149 north, the Taihang Mountains in the west, the Yellow River in the south, and the Bohai Gulf in the
150 northeast (35° – 40° N, 113° – 120° E) (Yang et al., 2022). Reported groundwater levels in the NCP
151 declined at a rate of 1–2 m/year in the late 20th century (Yang et al., 2021), leading to an estimated
152 cumulative depletion of about 60 km^3 from the 1960s to 2008 (Cao et al., 2013). China's South-to-North
153 Water Diversion Project, operating since 2014, has reduced the rate of groundwater depletion in the



154 NCP in recent years (Long et al., 2025; Long et al., 2020; Yang et al., 2022), but this region remains a
 155 critical hotspot for analysing TWS variability due to its long-term groundwater stress.



156

157 **Fig. 1** China's national boundary, and the locations and elevations of major rivers, lakes, glaciers, and nine large
 158 river basins. Two representative subregions, the North China Plain and the Southeastern Tibetan Plateau, are
 159 marked, and locations of in-situ groundwater wells in the North China Plain are also shown.

160 2.2 Data processing

161 Datasets used in this study include labels and predictors of the U-Net model and ancillary data for model
 162 validation (Table 1). We combine the state-of-the-art datasets from multiple sources, including satellite
 163 remote sensing, hydrologic simulations, and reanalysis data to better represent the climatic and
 164 anthropogenic drivers underlying TWS variability.

165 Table 1 Information of data sets used in this study.

Variable	Product name / data accessibility	Original spatial resolution	Temporal resolution	Temporal coverage
<i>Coarse-resolution label of TWSA observation</i>				
TWSA _{GRACE}	JPL-M	3°	Monthly	Apr 2002–Jun 2017; Jun 2018–present
<i>High-resolution label of TWSA spatial pattern</i>				
TWSA _{PI}	PCR GLOBWB 2	5 arcmin	Monthly	1981–2024
	Hugonnet et al. (2021)	100 m	Monthly	2000–2019
<i>Dynamic predictors</i>				
Precipitation	ERA5L	0.1°	Monthly	1950–present
Evapotranspiration	ERA5L	0.1°	Monthly	1950–present
Runoff	ERA5L	0.1°	Monthly	1950–present
Temperature	ERA5L	0.1°	Monthly	1950–present
NDVI	MOD13 A3	1 km	Monthly	2000–present
<i>Static predictors</i>				
Irrigation intensity	FAO	5 arcmin	–	2005



Groundwater irrigation area	FAO	5 arcmin	–	2005
Population	GPWv4	5 arcmin	–	2000, 2005, 2010, 2015, 2020
Elevation	SRTM	1 arc-second	–	2000

Ancillary data sets

Groundwater level	National Earth System Science Data Center	–	Monthly	2005–2018
Glacier mass balance	Zhao et al. (2022)	Region average	Annual	2003–2019

166 *Note:* Glacier mass balance estimates with broad spatial coverage are available up to 2019 and are therefore
 167 included only in the standard version of downscaled TWSA (2002–2019). In the extended period (2020–2023),
 168 the reconstructed TWSA primarily reflects variations in other water storage components (e.g., lakes, reservoirs,
 169 snow, soil moisture, and groundwater) simulated by PCR-GLOBWB 2 (Li and Sun, 2026a).

170 **2.2.1 TWSA from GRACE and GRACE-FO mascon solution**

171 We used the GRACE/-FO mascon solutions provided by NASA Jet Propulsion Laboratory (JPL) (Wiese
 172 et al., 2016) based on the latest Release Number 06 (RL06) version. Compared to previous spherical
 173 harmonic solutions, mascon solutions eliminate the need for post-processing and suffer less from
 174 leakage errors, and thus better preserve GRACE and GRACE-FO signals (Landerer and Swenson, 2012;
 175 Wiese et al., 2022). The JPL-M (mascon) provides monthly TWSA as anomalies relative to the 2004–
 176 2009 time-mean baseline. Missing monthly TWSA during the GRACE period (Jan 2003–Jun 2017)
 177 caused by instrument failure was linearly interpolated using the nearest two monthly estimates, whereas
 178 the missing data in the gap period between the GRACE and GRACE-FO missions (Jul 2017–May 2018)
 179 were reconstructed through the U-Net model. Several post-processing algorithms have been applied to
 180 the JPL-M to reprocess the original signals to provide user-friendly products; that is, applying the
 181 Coastal Resolution Improvement (CRI) filter to reduce land/ocean leakage errors and using gridded
 182 gain factors provided by the Community Land Model (CLM) to redistribute mass at a resolution of 0.5°
 183 within each 3° mascon element (Wiese et al., 2016). In this study, we use the JPL-M data sets without
 184 the land-grid scaling factors so that the satellite retrievals are entirely independent of hydrologic models.
 185 In addition, we used the block mean method to upgrade the JPL-M data from the nominal resolution of
 186 0.5° into its native resolution 3°, aiming to reveal the effective resolution of GRACE/-FO satellites.

187 **2.2.2 High-resolution spatial pattern of physically informed TWSA**

188 Physically informed TWSA ($TWSA_{PI}$) was the sum of individual water storage components, including
 189 simulated storage anomalies derived from PCR-GLOBWB 2 (Sutanudjaja et al., 2018) and glacier
 190 anomalies derived from satellite altimetry (Hugonnet et al., 2021). Note that $TWSA_{PI}$ is used as a
 191 reference for the U-Net model to learn the pixel-scale structure (spatial distribution of high and low
 192 values), instead of constraining model outputs as ground truth.

193 PCR-GLOBWB 2 (Sutanudjaja et al., 2018) is a global hydrologic and water resources model that
 194 accounts for anthropogenic water use and management (e.g., irrigation, groundwater pumping, and
 195 reservoir regulation), and simulates surface (rivers, lakes, reservoirs, canopy, and snow) and subsurface
 196 (soil and groundwater) water storage. As one of the model outputs, monthly TWS simulated by PCR-
 197 GLOBWB 2 is the sum of these storage components, with high spatial resolution of 5 arcmin
 198 (approximately 0.0833°). We interpolated the 5-arcmin data to 0.1° resolution, and the model-simulated
 199 TWS was calculated as TWSA by subtracting the 2004–2009 time-mean baseline, to be consistent with
 200 GRACE data processing. Importantly, because PCR-GLOBWB 2 does not include data assimilation of
 201 GRACE measurements, the modelled TWSAs and GRACE TWSAs are entirely independent.

202 Due to the lack of glacier processes in PCR-GLOBWB 2 (Janzing et al., 2025), we incorporate
 203 multisource satellite altimetry retrievals to add glacier changes, obtained from the latest estimation of
 204 global glacier mass balance by Hugonnet et al. (2021). This glacier database provides monthly changes
 205 in glacier volume relative to the year 2000 at a spatial resolution of 100 m during 2000–2019. It was
 206 generated based on glacier inventories from Randolph Glacier Inventory 6.0 (RGI 6.0) (Pfeffer et al.,
 207 2014) and elevation changes from a large number of satellite archives, such as ASTER stereo images
 208 and ArcticDEM (Hugonnet et al., 2021). We first obtained monthly time series of cumulative changes
 209 in glacier volume relative to the year 2000, and divided them by glacier areas determined by RGI 6.0



210 to calculate cumulative changes in glacier thickness (Eq.1). Second, changes in glacier thickness were
 211 converted into liquid water equivalent (LWE) by multiplying the ratio of glacier density to water density
 212 (Eq.2). Finally, the water-equivalent contribution from glacier changes is area-weighted based on the
 213 area of each 0.1° grid cell, to obtain changes in glacier storage (expressed in water equivalent) at the
 214 mean of 0.1° grid cell (Eqs.3–4). For consistency with GRACE/-FO data, remote sensing-based glacier
 215 changes were calculated as anomalies by subtracting the 2004–2009 time-mean baseline.

$$216 \quad h_g = \frac{V_g}{A_g} \quad (1)$$

$$217 \quad LWE = h_g \cdot \frac{\rho_{glacier}}{\rho_{water}} \quad (2)$$

$$218 \quad LWE' = LWE \cdot \frac{A_g}{A_{grid}} \quad (3)$$

$$219 \quad A_{grid} = R^2 \cdot (\omega_2 - \omega_1)(\sin \varphi_2 - \sin \varphi_1) \quad (4)$$

220 where h_g , V_g , and A_g represent thickness, volume, and area of glaciers, respectively; LWE indicates
 221 liquid water equivalent corresponding to ice thickness, while LWE' denotes area-weighted water
 222 equivalent of each 0.1° grid cell; ρ is average density, where average density of glaciers ($\rho_{glacier}$) is 850
 223 kg/m^3 (Huss, 2013) and average water density (ρ_{water}) is 1000 kg/m^3 ; A_{grid} denotes area of each grid cell,
 224 computed based on a spherical Earth approximation; R is mean radius of the Earth (6,371 km); ω_1 and
 225 ω_2 are the western and eastern longitudinal boundaries of the grid cell (in radians); and φ_1 and φ_2 denote
 226 the southern and northern latitudinal boundaries (in radians).

227 By summing storage anomalies derived from PCR-GLOBWB 2 and glacier anomalies derived from
 228 satellite altimetry, our physically informed TWSA ($TWSA_{PI}$) explicitly accounts for multiple
 229 components of TWS change to ensure a comprehensive representation. The normalization relative to
 230 the same baseline period with GRACE processing (2004–2009) reduces potential inconsistency
 231 associated with different data sources and processing schemes. Given that physical models may
 232 introduce substantial noise (Gou and Soja, 2024), we applied an outlier filtering approach to $TWSA_{PI}$
 233 to improve the training robustness. As training was conducted at a monthly time step (Section 3.2), the
 234 lower and upper outlier bounds were determined as the 1st and 99th percentiles of the data distribution
 235 within each monthly map, and excluded from the computation of loss functions.

236 2.2.3 Input predictors

237 We considered a total of twelve predictors from four categories to capture the combined effects of
 238 climatic and anthropogenic factors on TWS variations. First, precipitation, ET, and runoff are three
 239 critical hydrologic fluxes in the water balance and highly related to TWS variations. We obtained
 240 monthly estimates of these variables at a spatial resolution of 0.1° from the latest land component of
 241 the fifth generation of European ReAnalysis (ERA5) data, ERA5-Land (simplified as ERA5L) (Muñoz-
 242 Sabater et al., 2021), provided by the European Centre for Medium Range Forecasts (ECMWF).

243 Second, we considered cumulative anomalies of precipitation, ET, and runoff to characterize the
 244 memory effect of hydrologic fluxes on TWS variations. On the basis of the water balance principle
 245 (input fluxes – output fluxes = TWS changes), TWS is affected by accumulated flux variables in a given
 246 period (Eq.5). Cumulative anomalies of hydrologic fluxes represent long-term water surplus or deficit
 247 and have been shown to improve the performance of deep learning models in reconstructing TWS in
 248 previous studies (Li et al., 2022; Wang et al., 2024). In this study, flux variables (precipitation, ET, and
 249 runoff) were calculated as monthly anomalies relative to the 2004–2009 baseline, consistent with
 250 GRACE data processing, and then cumulatively integrated from a fixed starting point of April 2002.

$$252 \quad TWS_n = TWS_0 + \Delta TWS_1 + \Delta TWS_2 + \dots + \Delta TWS_n$$

$$253 \quad = TWS_0 + (flux_{In,1} - flux_{Out,1}) + (flux_{In,2} - flux_{Out,2}) + \dots + (flux_{In,n} - flux_{Out,n})$$

$$251 \quad = TWS_0 + (flux_{In,1} + flux_{In,2} + \dots + flux_{In,n}) - (flux_{Out,1} + flux_{Out,2} + \dots + flux_{Out,n}) \quad (5)$$

254 Third, temperature, vegetation dynamics, and topography were included to describe land surface
 255 processes. Temperature largely affects the energy budget and is a dominant driver of TWS changes over
 256 cryospheric regions (Li et al., 2022). We obtained 2 m air temperature from ERA5L datasets with 0.1°



257 resolution. Vegetation dynamics were represented by Normalized Difference Vegetation Index (NDVI),
258 which was obtained from the Moderate Resolution Imaging Spectroradiometer (MODIS) vegetation
259 index product (MOD13 A3) at the monthly scale and 1 km resolution. NDVI reflects the vegetation
260 greenness and canopy density, closely linked to transpiration and soil moisture uptake. Topography was
261 represented by the Digital Elevation Model (DEM) from the Shuttle Radar Topography Mission (SRTM)
262 1 Arc-Second Global data set (approximately 30 m resolution). The elevation plays a key role in
263 regulating runoff pathways and groundwater recharge, further influencing the spatial variability of TWS.

264 Fourth, we incorporated predictors related to human interventions in the water cycle, including
265 irrigation intensity (fraction of area equipped for irrigation), fraction of irrigated area supplied by
266 groundwater, and population density. Irrigation information was derived from the latest version (v5.0)
267 of the Global Map of Irrigation Areas from the Food and Agriculture Organization (FAO) at 5 arcmin
268 (0.0833°). Population density serves as a proxy for domestic and municipal water demand, provided by
269 the Gridded Population of the World dataset, Version 4 (GPWv4), with a resolution of 5 arcmin
270 (0.0833°). This dataset is primarily based on national censuses and population registers, without
271 extensive post-processing or ancillary data-driven refinement.

272 The native spatial resolutions of these twelve predictors are equal to or finer than the target downscaling
273 resolution of 0.1°, thereby providing sufficient spatial details to learn high-resolution TWS patterns.
274 Eight predictors are dynamic, with monthly values related to the water balance (precipitation,
275 evapotranspiration, runoff, and their accumulations), energy budget (temperature), and vegetation
276 (NDVI). In contrast, four predictors describing irrigation, population, and topography are treated as
277 quasi-static factors due to their relatively slow temporal variability. Subject to data availability, we input
278 the average state (annual average population during 2000–2020) or representative state (the irrigation
279 map obtained in 2005 and DEM map in 2000) of these static predictors into the U-Net model as time-
280 invariant input channels. All features were linearly normalized on the basis of their 1st percentiles and
281 99th percentiles to reduce the impacts of outliers.

282 2.2.4 Ancillary datasets

283 Ancillary datasets include in-situ observations of groundwater levels across the NCP and independent
284 estimation of cumulative glacier surface elevation over the SETP. We collected observed groundwater
285 level during 2005–2018 of 559 groundwater wells across the NCP, provided by the National Earth
286 System Science Data Center, National Science and Technology Infrastructure of China. These
287 groundwater wells provide monthly groundwater level variations with at least 12 months of records,
288 which are compared with TWSA at the corresponding grid cells. Because we aim to evaluate the
289 temporal variability of TWSA rather than absolute magnitude (Section 4.3), observed groundwater level
290 variations are considered representative without the need for conversion to groundwater storage. Such
291 conversion requires the estimation of specific yield, which is highly uncertain over the NCP and may
292 introduce additional errors.

293 In addition, glacier estimates over the SETP were obtained from Zhao et al. (2022), based on satellite
294 altimetry observations from ICESat and CryoSat-2, which are independent of those included in TWSA_{PI}.
295 This dataset provides regionally averaged cumulative surface elevation over the SETP at seasonal and
296 annual scales, and is used to evaluate the temporal variability of TWSA in this glacierized area.

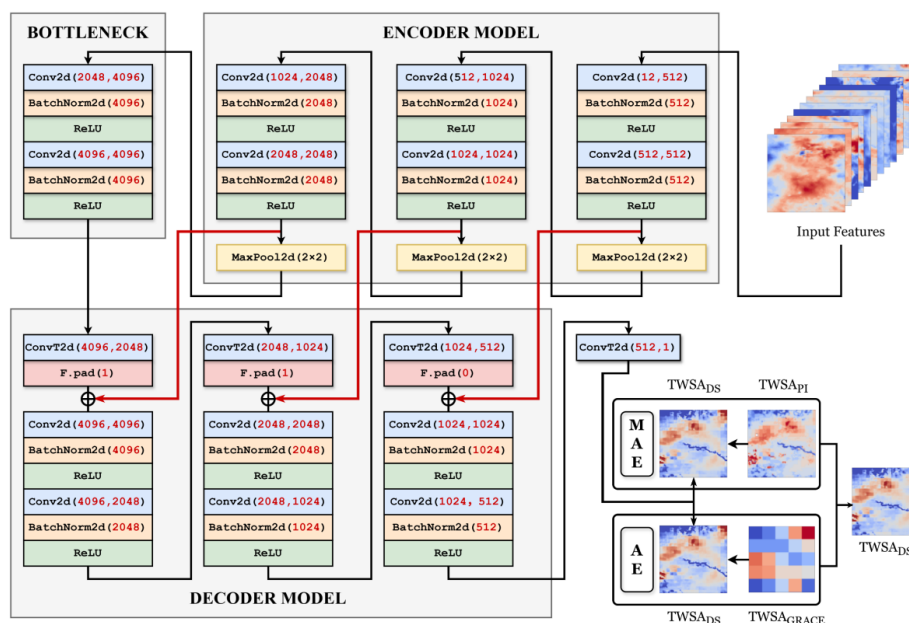
297 3 Methodology

298 3.1 U-Net model structure

299 The input to the U-Net is a four-dimensional tensor of shape $B \times C \times 30 \times 30$, where B denotes the
300 global batch size and C represents the number of features channels derived from input predictors. The
301 global batch size was set to 512 ($B = 512$), with a total of 12 feature channels, including 8 dynamic and
302 4 static features ($C = 12$). All predictors were processed to 0.1° and grouped into 30×30 spatial patches,
303 corresponding to the spatial extent of a single 3° GRACE grid cell (i.e., $3^\circ = 30 \times 0.1^\circ$). The network
304 maps multi-channel input patches to a single-channel output at 0.1° resolution, representing the
305 predicted TWSA at each grid cell. The U-Net model adopts an encoder-decoder architecture with
306 symmetric skip connections to capture both large-scale contextual information and fine-scale spatial
307 details (Fig. 2). The encoder progressively extracts hierarchical features through convolutional blocks



308 with downsampling, while the decoder restores spatial resolution via upsampling. The number of feature
 309 channels increases in the encoder and decreases in the decoder. Skip connections are applied at each
 310 stage to preserve spatial information. Detailed layer configurations, including channel dimensions and
 311 feature map sizes, are illustrated in Fig. 2. Finally, a 1×1 convolution layer maps the decoded features
 312 to a single-channel output, representing the predicted TWSA at 0.1° resolution.



313
 314 **Fig. 2** The structure of U-Net model for TWSA downscaling. Multi-channel predictors at 0.1° resolution are
 315 mapped to high-resolution TWSA through an encoder-decoder framework with skip connections. The encoder
 316 and decoder models are both built upon a unified modular structure consisting of two convolutional layers. The
 317 red lines denote skip connections, which are used to fuse low-level features with high-level features. The symbol
 318 \oplus denotes the concatenation of two tensors with the same shape along the channel dimension. Model training is
 319 constrained by a dual-loss design, enforcing agreement with GRACE TWSA at the patch scale and with physically
 320 informed patterns at the pixel scale.

321 3.2 Training strategy and uncertainty estimation

322 For model training, the study domain was represented as a 420×690 matrix of 0.1° pixels,
 323 corresponding to a 14×23 matrix of 3° grid cells (i.e., 30×30 pixels per grid cell), covering China and
 324 its surrounding regions (15° – 57° N, 69° – 138° E). Coastal GRACE grid cells with a land fraction larger
 325 than 20% were retained to ensure adequate terrestrial signal. The spatial domain was partitioned into 7
 326 \times 11 non-overlapping blocks. Each block consists of 60×60 pixels and corresponds to four effective
 327 GRACE grid cells, improving the stability of the patch-scale constraint. The easternmost 30 columns
 328 were excluded from block partitioning, as they contain very few valid terrestrial signals within China.
 329 At each monthly time step, the non-overlapping blocks were randomly partitioned into training (70%;
 330 54 blocks) and validation (30%; 23 blocks) subsets. The split was conducted within each month rather
 331 than across time, as the objective was to learn spatial relationships rather than temporal forecasting. To
 332 increase the number of training samples and improve spatial continuity, a sliding window with a stride
 333 of 10 is applied to extract overlapping patches within the training and validation sets, respectively.
 334 Parameters are optimized using the Adam optimizer with a learning rate schedule to ensure stable
 335 convergence during training. The model is trained for 800 epochs, with the learning rate gradually
 336 reduced over time.

337 The training objective combines two complementary losses, which balance amplitude constraints from
 338 GRACE/-FO observations at large scales and spatial patterns from physically informed TWSA at fine



339 resolutions. We computed the loss of Absolute Error (AE) between the averaged GRACE TWSAs and
340 averaged predicted TWSAs over each patch (corresponding to the GRACE original resolution of 3°),
341 and then averaged all patches within a batch to obtain the aggregation error at the GRACE scale (Eq.6).
342 To capture high-resolution spatial patterns, a loss of Mean Absolute Error (MAE) was additionally
343 introduced at the pixel level of 0.1° . As for the MAE loss constrained by physically informed TWSA,
344 we exclude invalid pixels, including missing values in the ocean and outliers of hydrologic simulations
345 (1^{st} and 99^{th} percentiles of the data distribution within each monthly map). Because the number of valid
346 pixels varies across patches, the MAE loss was computed as the average over all valid pixels across the
347 batch (Eq.7).

$$348 \quad L_{AE} = \frac{1}{B} \sum_{b=1}^B \left| \frac{1}{N} \sum_{i=1}^N \hat{y}_i - \frac{1}{N} \sum_{i=1}^N y_{i,G} \right| \quad (6)$$

$$349 \quad L_{MAE} = \frac{1}{\sum_{b=1}^B |\Omega_b|} \sum_{b=1}^B \sum_{i \in \Omega_b} |\hat{y}_i - y_{i,PI}| \quad (7)$$

350 where B denotes the batch size, which is set to 512 in this study; N represents the number of pixels
351 within each patch; \hat{y}_i indicates the predicted TWSA at pixel i ; y_i represents the reference value at each
352 pixel, corresponding to GRACE TWSA for the AE loss ($y_{i,G}$ in Eq.6) and physically informed TWSA
353 for the MAE loss ($y_{i,PI}$ in Eq.7); and Ω represents the valid pixels in each patch with $|\Omega|$ denoting the
354 number of valid values in this patch.

355 The AE loss constrains the spatially averaged TWSA within each 30×30 patch, thereby preserving the
356 coarse-resolution signals from GRACE measurements. In contrast, the MAE loss is computed at the
357 pixel level, and improves fine-scale spatial details by minimizing the absolute differences between
358 predicted TWSA and the high-resolution TWSA_{PI}. The final loss function was defined as a weighted
359 combination of the two losses constrained by GRACE TWSA at the patch level and physically informed
360 TWSA at the pixel level (Eq.8).

$$361 \quad L = \lambda_{AE} L_{AE} + \lambda_{MAE} L_{MAE} \quad (8)$$

362 Two parameters of λ_{AE} and λ_{MAE} control the trade-off between large-scale GRACE consistency and
363 fine-scale spatial details. Increasing the ratio of λ_{AE} to λ_{MAE} forces the model to better match the
364 GRACE statistics but may suppress fine-scale spatial variability. Conversely, decreasing this ratio
365 enhances the model's ability to resolve local structural details but may weaken the consistency with
366 GRACE-derived TWS variations if the fine-scale constraint becomes overly dominant.

367 In this study, we first calculated the average uncertainty of GRACE JPL-M across China (1.96 cm).
368 Using this as a benchmark, we designed the loss-weighting strategy to ensure that the AE loss for both
369 the training and validation sets converged to values below the average GRACE uncertainty across China.
370 This criterion guarantees that the downscaled TWSA has comparable accuracy relative to GRACE
371 observation at large scales. Within this constraint, we gradually increased the ratio of λ_{MAE} to λ_{AE}
372 to improve the model's ability to resolve fine-scale spatial variability while preserving large-scale
373 consistency. Through grid search over the combinations of $[2.0, 1.5, 1.0] \times [2.0, 1.5, 1.0]$, the optimal
374 weighting was determined as $\lambda_{AE} = 1.0$ and $\lambda_{MAE} = 1.0$.

375 Uncertainty of downscaled TWSA was estimated from three independent training runs with different
376 random seeds, capturing variability from initialization, data shuffling, and stochastic optimization.
377 Despite identical training settings, stochasticity in the learning process introduces variability in model
378 parameters and predictions. For each independent run, the trained model produced a complete set of
379 downscaled TWSA fields, from which the evaluation metrics are computed. The final downscaled
380 TWSA was generated by averaging TWSA outputs across the three runs, while the corresponding
381 standard deviation was estimated as uncertainty. For each monthly map, those pixels with estimated
382 uncertainty exceeding the 99^{th} percentile of all pixels across China were identified as low-confidence
383 predictions. These values were subsequently replaced using nearest-neighbour interpolation from
384 surrounding valid grid cells to ensure spatial continuity.



385 3.3 Closure of water balance and evaluation indices

386 The water balance serves as the water mass-conservation constraint at the catchment scale,
387 demonstrating the relationship between TWS change (TWSC) and the flux integration of precipitation
388 (P), evapotranspiration (ET), and runoff (R) (Eq.9).

$$389 \quad TWSC = P - ET - R \quad (9)$$

390 To reduce potential high-frequency artifacts from the finite differences, both TWSC and hydrologic
391 fluxes were smoothed, using Eq.10 and Eq.11, respectively (Gou and Soja, 2024; Landerer et al., 2010;
392 Xiong et al., 2025).

$$393 \quad TWSC(t) = \frac{TWSA(t+1) - TWSA(t-1)}{2\Delta t} \quad (10)$$

$$394 \quad \bar{X} = \frac{1}{4}X(t-1) + \frac{1}{2}X(t) + \frac{1}{4}X(t+1) \quad (11)$$

395 where t denotes the target month; Δt is the time interval taken as a month; and X represents the
396 hydrologic flux of P , ET , or R . Water balance closure was used to assess whether the downscaled TWSA
397 preserves mass conservation beyond the native GRACE resolution and captures physically consistent
398 relationships among hydrologic fluxes. Precipitation, ET , and runoff were derived from ERA5L,
399 consistent with the input predictors, to ensure methodological consistency. This analysis is intended to
400 evaluate internal model consistency rather than provide an independent validation, and therefore,
401 consistent datasets avoid additional uncertainties introduced by heterogeneous data sources. We
402 compared TWS change derived from the storage and flux term (Eq.9) across 163 small to medium
403 catchments in China. These basins are defined based on HydroBASINS Level 5 database, representing
404 a medium spatial scale with the drainage area on the order of 10,000–100,000 km². We further excluded
405 small basins (area < 10,000 km²), desert regions (long-term annual mean precipitation < 200 mm), and
406 basins strongly affected by intensive irrigation (irrigated area fraction > 5%) or cryospheric processes
407 (glacier area fraction > 5%). This exclusion aims to minimize the influence of processes that are not
408 explicitly accounted for in the water balance equation (e.g., additional storage loss due to human water
409 use), thereby ensuring a more reliable assessment.

410 We selected two commonly used indices to evaluate the performance of downscaled TWSA, including
411 the Pearson correlation coefficient (CC) (Eq.12) and root mean square error ($RMSE$) (Eq.13).

$$412 \quad CC = \frac{\sum_{j=1}^n (x_j - \bar{x})(y_j - \bar{y})}{\sqrt{\sum_{j=1}^n (x_j - \bar{x})^2 \sum_{j=1}^n (y_j - \bar{y})^2}} \quad (12)$$

$$413 \quad RMSE = \sqrt{\frac{\sum_{j=1}^n (x_j - y_j)^2}{n}} \quad (13)$$

414 where x_j represents the variable to be evaluated; y_j represents the corresponding variable taken as the
415 baseline; an overbar denotes mean; and n is number of data pairs for evaluation.

416 Here CC was used to indicate the degree of coherence between the evaluated and reference datasets,
417 while $RMSE$ was employed to quantify the magnitude of discrepancies. The other index of Kling-Gupta
418 efficiency (KGE) (Gupta et al., 2009; Kling et al., 2012) was not adopted, because it is less suitable for
419 anomaly-based TWSA evaluation when the mean TWSA could approach zero. Specifically, KGE is a
420 comprehensive indicator including correlation, ratio of means, and ratio of standard deviations of the
421 two paired datasets. The value of ratio of means can be extremely high if its denominator (the mean of
422 TWSA) is approaching to zero, further resulting in a very large negative value of KGE over the regions
423 with relatively stable TWS changes.

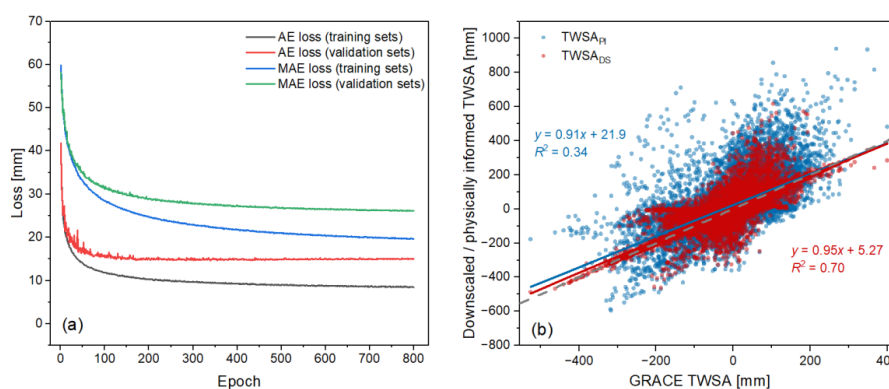
424 4 Results

425 4.1 Performance of downscaled TWSA at large scales

426 The AE loss for both training and validation sets converged after approximately 400 epochs (Fig. 3a),
427 representing the GRACE-based constraint. By the end of training (800 epochs), the AE loss converged
428 to 0.84 cm for the training set and 1.50 cm for the validation set, both lower than the spatially averaged



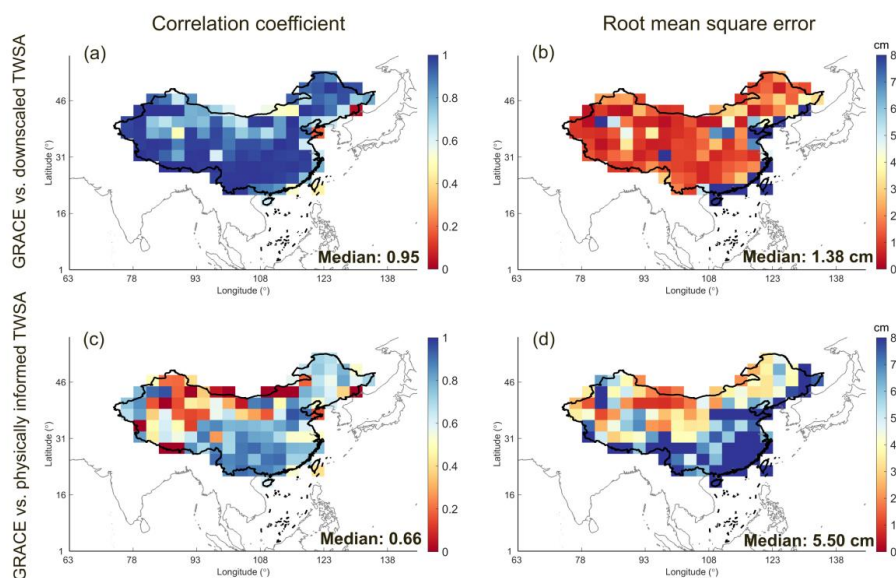
429 uncertainty of GRACE JPL-M across China (1.96 cm). The MAE loss, associated with physically
430 informed spatial patterns, converged more gradually and reached a stable level after approximately 600
431 epochs, with final values of 1.97 cm for the training set and 2.61 cm for the validation set, respectively.
432 For both AE and MAE losses, the validation curves consistently follow the trends of the training curves
433 without divergence, and the gap between them remains relatively stable after convergence, indicating
434 stable model training with no evident overfitting. We upscaled the high-resolution TWSA from 0.1° to
435 3° using the block mean method, and find that compared to physically informed TWSA, the downscaled
436 product substantially improved the agreement with GRACE/-FO observations (Fig. 3b). Across all 3°
437 grid cells in China during 2002–2019, the coefficient of determination (R^2) increased from 0.34
438 ($TWSA_{PI}$ vs. $TWSA_{GRACE}$) to 0.70 ($TWSA_{DS}$ vs. $TWSA_{GRACE}$), indicating enhanced large-scale
439 consistency after incorporating GRACE constraints.



440

441 **Fig. 3 (a)** Loss curves of the multi-constraint U-Net model. The curves are shown separately for the AE loss and
442 MAE loss, each including both training and validation sets. **(b)** Scatterplots between $TWSA_{GRACE}$ vs. $TWSA_{PI}$
443 (blue dots) and $TWSA_{GRACE}$ vs. $TWSA_{DS}$ (red dots) for all 3° grid cells in China during 2002–2019. Solid lines
444 indicate linear regression fits, with the corresponding regression equations and coefficients of determination (R^2).
445 Dots, regression lines, and corresponding statistics share consistent colors for each comparison. The dashed grey
446 line indicates the 1:1 line.

447 We mapped the grid-level evaluation metrics across China to evaluate the ability of downscaled TWSA
448 to preserve GRACE signals at the 3° resolution during the 2002–2019 period (Fig. 4). The close
449 agreement between $TWSA_{DS}$ and $TWSA_{GRACE}$ is reflected by a high correlation (median CC : 0.95; Fig.
450 4a) and low discrepancy (median $RMSE$: 1.38 cm; Fig. 4b). Relatively large discrepancies of $TWSA_{DS}$
451 are mainly found in coastal areas, which may be attributed to both land-ocean signal leakage in GRACE
452 observations and missing values over adjacent ocean areas in physical simulations. Nevertheless, over
453 the majority of inland China, $TWSA_{DS}$ demonstrates consistently high agreement with GRACE/-FO
454 observations. In comparison, the physically informed TWSA ($TWSA_{PI}$) exhibits much weaker
455 agreement with GRACE, characterized by a lower median CC (0.66) (Fig. 4c) and larger $RMSE$ (5.50
456 cm) (Fig. 4d). Low correlations between $TWSA_{PI}$ and $TWSA_{GRACE}$ are located across the arid and semi-
457 arid regions in Northwest China, whereas large discrepancies are found in the humid Southeast China.
458 These results show the limitation of model simulations to reveal TWSA in both terms of variability and
459 magnitude at large scales. The improvements in both CC and $RMSE$ highlight the effectiveness of
460 incorporating GRACE-based constraints in enhancing large-scale statistical consistency with satellite
461 observations.

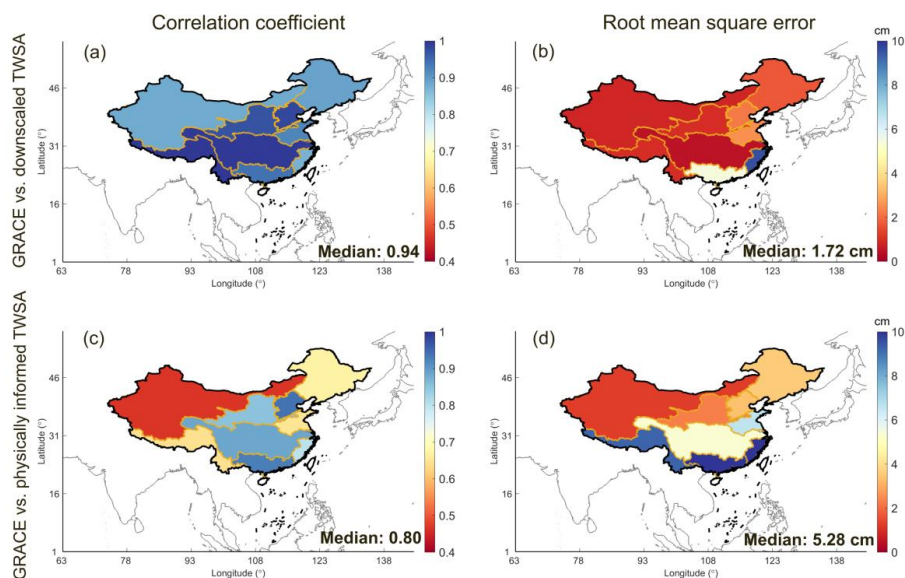


462

463 **Fig. 4** Spatial distribution of evaluation indices of (a)–(b) downscaled TWSA and (c)–(d) physically informed
464 TWSA relative to GRACE TWSA at each 3° grid cell, including correlation coefficient (left column) and
465 root mean square error (right column). Median values of the evaluation index over all grid cells are labelled in bottom-
466 right of each panel.

467 In addition to the grid-cell scale, TWS variability was assessed over the nine large basins in China to
468 characterize regional hydrologic changes (Fig. 5). Downscaled and GRACE TWSA show good
469 agreement for these basins, with a median CC of 0.94 (Fig. 5a) and RMSE of 1.72 cm (Fig. 5b).
470 Performance over the Southwest, Yangtze, Yellow, Haihe, and Huaihe river basins is particularly
471 reliable (CC: 0.92–1.00; RMSE: 0.46–2.57 cm), whereas relatively large discrepancies are found for the
472 Southeast and Pearl river basins in the humid Southeast China (RMSE: 9.60 cm of the Southeast basin
473 and 5.43 cm of the Pearl basin), which may be attributed to uncertainties from land-ocean boundary
474 effects and large variability in precipitation. The northeastern (Songhua and Liaohe River Basin) and
475 northwestern (Continental Basin) regions exhibit high consistency between TWSA_{DS} and TWSA_{GRACE},
476 with CC approaching 0.90 and RMSE within 2 cm, indicating robust performance of the deep-learning
477 downscaling framework across inland and continental regions in China.

478 Compared to TWSA_{DS}, the physically informed results (TWSA_{PI}) demonstrate lower performance in
479 all basins, with a median CC of 0.80 (Fig. 5c) and median RMSE of 5.28 cm (Fig. 5d). These results
480 highlight the role of GRACE constraints in enhancing both temporal variability and signal amplitude in
481 the deep learning-based downscaling approach. The largest improvement in temporal variability is
482 observed in the Continental Basin of Northwest China, where CC increases from 0.47 (TWSA_{PI} vs.
483 TWSA_{GRACE}) to 0.88 (TWSA_{DS} vs. TWSA_{GRACE}). In addition, amplitude discrepancies are markedly
484 reduced in the Yangtze and Southwest basins, where RMSE decreases by approximately 90%. In the
485 Yangtze basin, RMSE decreases from 5.28 cm (TWSA_{PI} vs. TWSA_{GRACE}) to 0.89 cm (TWSA_{DS} vs.
486 TWSA_{GRACE}), while in the Southwest basin, it declines from 9.25 cm (TWSA_{PI} vs. TWSA_{GRACE}) to 0.98
487 cm (TWSA_{DS} vs. TWSA_{GRACE}).



488

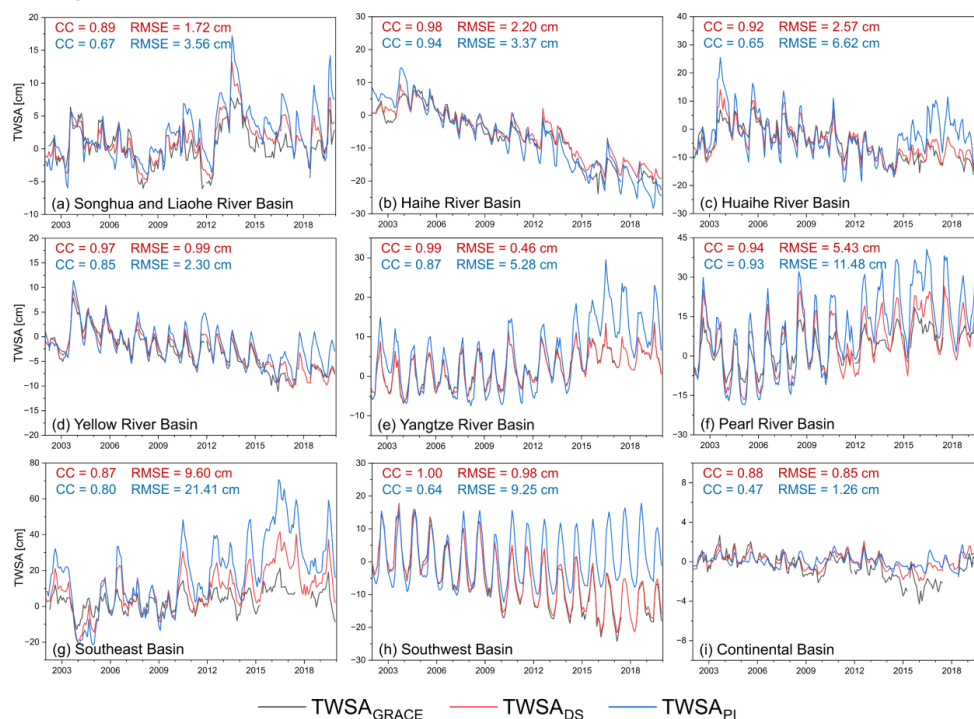
489 **Fig. 5** Basin-scale evaluation of (a)–(b) downscaled TWSA and (c)–(d) physically informed TWSA relative to
490 GRACE TWSA over nine river basins across China during 2002–2019, including correlation coefficient (left
491 column) and root mean square error (right column). Median values of the evaluation index over all basins are
492 labelled in bottom-right of each panel.

493 Basin-scale time series provide more details to understand the performance of downscaled and
494 physically informed TWSA in characterizing seasonal amplitude, interannual variability, and long-term
495 trends in TWSA (Fig. 6). Across nine large river basins in China, the Haihe, Huaihe, Yellow, and
496 Southwest basins show significant TWS declines during 2002–2019. Compared with the physically
497 informed simulations, the downscaled TWSA better reproduces both magnitude and temporal evolution
498 of these declining trends. For example, in the Southwest basin (Fig. 6h), $TWSA_{PI}$ mainly captures
499 seasonal variations but fails to represent the long-term decreasing trend observed by GRACE. In
500 contrast, $TWSA_{DS}$ successfully reproduces this declining trend, demonstrating the benefit of
501 incorporating GRACE constraints into the deep learning framework. As for the Haihe (Fig. 6b), Huaihe
502 (Fig. 6c), and Yellow (Fig. 6d) river basins, $TWSA_{PI}$ shows noticeable deviations from $TWSA_{GRACE}$
503 after 2015, whereas the $TWSA_{DS}$ better reproduces the observed signals. Similar deviations between
504 $TWSA_{GRACE}$ and $TWSA_{PI}$ after 2015 are also evident in the Yangtze basin with an increasing TWS trend
505 (Fig. 6e), whereas $TWSA_{DS}$ closely follows the GRACE observations, successfully reproducing both
506 the increasing trend and seasonal variations.

507 We observe an overestimation of variability amplitude in the physical simulations in the Songhua and
508 Liaohe River basin (Fig. 6a), the Pearl River basin (Fig. 6f), and the Southeast basin (Fig. 6g), where
509 TWS signals are dominated by seasonal and interannual variability. In these basins, the downscaled
510 TWSA reflects the combined influence of GRACE observations and physical simulations, with its
511 temporal variations generally lying between those of the two. Despite improvements, the downscaled
512 product still faces challenges in fully reproducing GRACE variability in these basins, particularly the
513 coastal regions with high precipitation variability, where both GRACE retrievals and physical
514 simulations show large uncertainty. This suggests that improving the quality of training labels and
515 constraints could be important for deep learning-based downscaling. Finally, the physically informed
516 TWSA shows the weakest temporal coherence with GRACE in the Continental Basin of northwestern
517 China (Fig. 6i), where TWS variability is primarily dominated by interannual changes. GRACE also
518 indicates TWS decline during 2015–2017 followed by a recovery in 2018–2019 in this region. We find
519 that the downscaled TWSA improves the temporal coherence with GRACE and captures part of these
520 changes. However, the magnitude of TWS declines during this period is not fully represented in



521 TWSA_{DS}, which may be related to the challenge in linking hydrologic predictors to TWS changes in
 522 arid regions of Northwest China.



523
 524 **Fig. 6** Monthly time series of TWSA during 2002–2019 across nine river basins in China, estimated from GRACE
 525 (TWSA_{GRACE}; black line), downscaled outputs (TWSA_{DS}; red line), and physically informed estimates (TWSA_{PI};
 526 blue line). The correlation coefficient (CC) and root mean square error (RMSE) between GRACE and two other
 527 estimations are labelled in top-left of each panel, with red text indicating TWSA_{GRACE} vs. TWSA_{DS} and blue text
 528 indicating TWSA_{GRACE} vs. TWSA_{PI}.

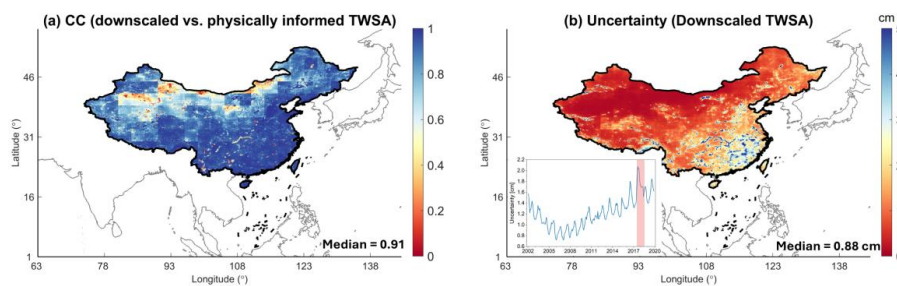
529 4.2 Reliability and uncertainty estimation of downscaled TWSA at fine scales

530 Because accurate in-situ TWS measurements at sub-regional scales in China are not available, we
 531 cannot directly evaluate the absolute magnitude of downscaled TWSA at a high-resolution. Therefore,
 532 we focus on the reliability of downscaled TWSA indicated by the spatial correlation with physically
 533 informed TWSA, uncertainty estimation at the pixel level, and ability to close water balance across
 534 small to medium catchments beyond the native GRACE resolution. We computed the Pearson
 535 correlation between TWSA_{DS} and TWSA_{PI} at each 0.1° pixel over the entire time span of 2002–2019.
 536 The Pearson correlation is invariant to magnitude scaling, which reduces the impact of inaccurate
 537 magnitudes of physical informed estimates. Results show that the overall correlation between
 538 downscaled and physically informed TWSA is high with a median value of 0.91 (Fig. 7a). The relatively
 539 low correlations are mainly found in arid regions in north and northwest China, likely due to weak
 540 hydrologic signals, to which both GRACE observations and physical models have limited sensitivity.

541 The monthly averaged uncertainty during 2002–2019 indicates stable model performance, with a
 542 median uncertainty of 0.88 cm (Fig. 7b). Approximately 82% of pixels have uncertainty lower than the
 543 average uncertainty of GRACE JPL-M over China (1.96 cm). The pixel-scale uncertainty generally
 544 decreases from southeast to northwest across China, which is likely associated with the larger amplitude
 545 of TWS variability in the humid southeastern regions. Relatively high uncertainties (> 4 cm) are mainly
 546 found in the middle and lower Yangtze and southern China, high mountain regions of northwestern and
 547 southwestern China, areas with large lakes, and the highly irrigated regions. These high-uncertainty
 548 pixels appear as scattered and localized patches rather than forming extensive contiguous regions. In

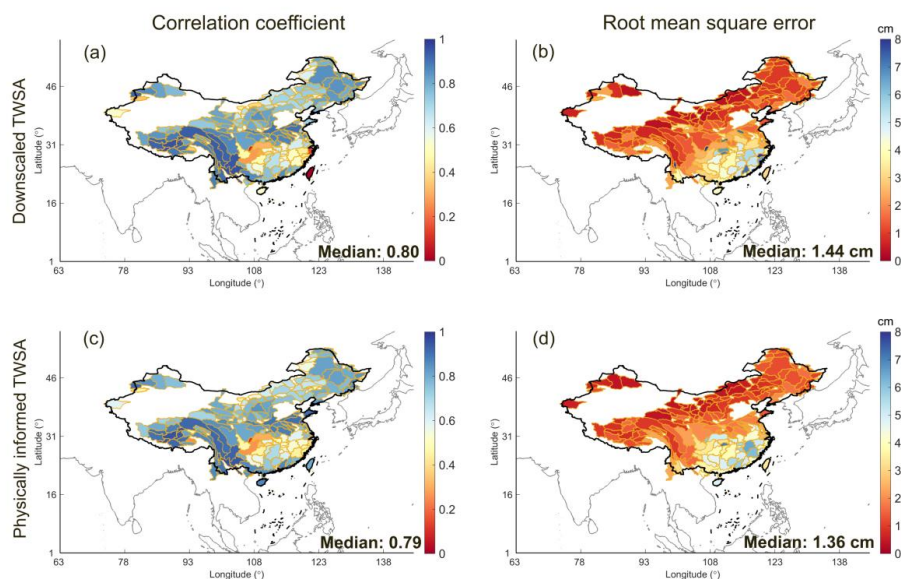


549 addition to spatial patterns, changes in monthly uncertainty shows larger values during the gap period
550 between GRACE and GRACE-FO missions (Jun 2017–May 2018; inserted subplot in Fig. 7b),
551 highlighting the importance of satellite observations in constraining the deep learning framework.



552
553 **Fig. 7** (a) Correlation coefficient between downscaled and physically informed TWSA and (b) monthly averaged
554 uncertainty of downscaled TWSA for each 0.1° pixel during 2002–2019. The inserted subplot in (b) shows the
555 spatially averaged uncertainty across China for each monthly map, with the red rectangle indicating the gap period
556 (Jul 2017–May 2018) between the GRACE and GRACE-FO missions. Missing data during the gap period have
557 been reconstructed by the U-Net deep learning.

558 To evaluate the ability of downscaled TWSA to close the water balance across small to medium
559 catchments beyond GRACE resolution, we compared TWS changes from storage and flux terms across
560 163 small to median catchments ($10,000\text{--}100,000\text{ km}^2$) in China (Section 3.3). The downscaled TWSA
561 shows good agreement with the water budget-derived TWS changes across most catchments. The
562 median CC between TWS changes derived from the water budget and the downscaled TWSA reaches
563 0.80 (Fig. 8a), with a low median $RMSE$ of 1.44 cm (Fig. 8b). Relatively low correlations are observed
564 in the middle reaches of the Yangtze River (Fig. 8a). This may be related to the operation of the Three
565 Gorges Reservoir, where large-scale human water regulation causes storage changes that are not fully
566 represented by the natural water balance. In addition, relatively large $RMSE$ values are found in humid
567 Southeast China (Fig. 8b), where monsoon precipitation leads to high variability in hydrologic fluxes
568 and could introduce additional uncertainty in the water balance estimation. The performance of
569 downscaled TWSA in closing the water balance is comparable to that of the physically informed TWSA,
570 which exhibits a median CC of 0.79 (Fig. 8c) and a median $RMSE$ of 1.36 cm (Fig. 8d). These results
571 suggest that the deep learning-based downscaling framework generally preserves the key hydrologic
572 relationships involved in the water balance.



573

574 **Fig. 8** Catchment-scale evaluation of monthly TWS change derived from (a)–(b) downscaled TWSA and (c)–(d)
575 physically informed TWSA, compared with TWS change from water budget across 163 small to medium
576 catchments in China during 2002–2019. The left column shows correlation coefficients, and the right column
577 shows root mean square errors. Median values of the evaluation index over all basins are shown in bottom-right
578 of each panel.

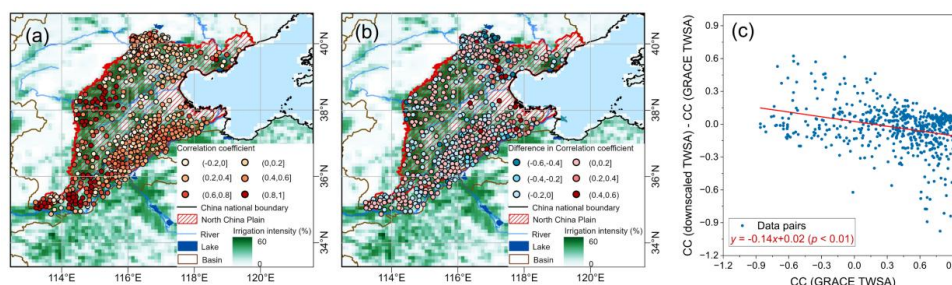
579 4.3 Evaluation across highly irrigated and glacierized regions

580 To further evaluate the physical consistency of the downscaled TWSA, we focus on regions where TWS
581 variability is dominated by specific components, such as groundwater storage in highly irrigated areas
582 and glacier mass balance changes in high mountain regions. In these regions, independent observations
583 of dominant TWS components provide a valuable benchmark for assessing the performance of the
584 downscaled TWSA. Because the comparison reference is not the observed TWS, absolute error metrics
585 such as *RMSE* are less informative, whereas correlation is more appropriate for assessing the
586 consistency between changes in downscaled TWSA and those in independent observations. The NCP
587 is intensively irrigated (Fig. 1) with TWS loss dominated by groundwater consumption. The correlation
588 between the downscaled TWSA and observed groundwater levels indicates generally good performance
589 across 559 groundwater wells spanning 2005–2018. Approximately 80% of the wells show positive
590 correlations, with a median *CC* of 0.41 (Fig. 9a). Notably, the agreement improves in regions with
591 higher irrigation intensity, such as the eastern and southwestern margins of the NCP. The median *CC*
592 increases to 0.48 and 0.65 for wells located in areas with irrigation intensity exceeding 40% and 50%,
593 respectively. This reflects the increasing dominance of groundwater in TWS variability in highly
594 irrigated regions, highlighting the capability of the downscaling framework to capture TWS variations
595 under strong human intervention.

596 We do not find a substantial difference in overall TWS-groundwater correlation after and before
597 downscaling: 51% of wells show an increase in *CC*, while 49% exhibit a decrease (Fig. 9b). A consistent
598 case was reported across irrigation zones in India by Wang et al. (2024), who demonstrated that the
599 overall correlation between variations in TWSA and those in groundwater levels did not change much
600 after TWS downscaling. This may be related to the fact that groundwater level variations observed at
601 individual wells may not be representative for the entire aquifer that is subject to intensive groundwater
602 extraction. As a result, the spatial scale at which point-based groundwater observations are aggregated
603 or compared with gridded TWSA may introduce scale-related uncertainties in the evaluation. However,
604 the effect of downscaling is not uniform. We find a significant negative relationship between the change



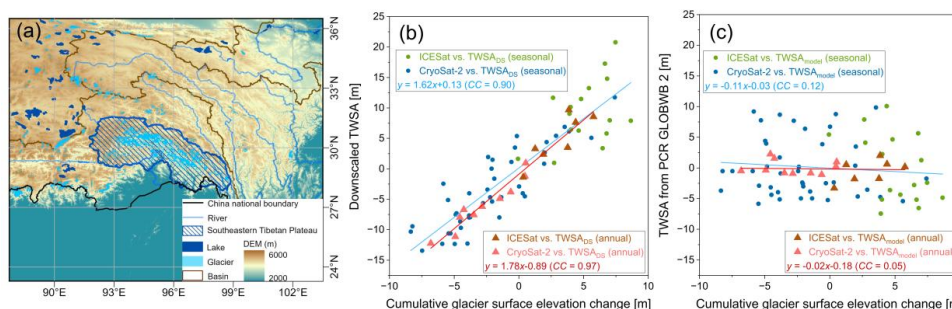
605 in correlation ($CC_{\text{downscaled, obs}} - CC_{\text{GRACE, obs}}$) and the baseline correlation ($CC_{\text{GRACE, obs}}$) (Fig. 9c). This
 606 suggests that in regions with low initial agreement between GRACE TWSA and in situ groundwater
 607 levels, likely reflecting strong spatial heterogeneity in groundwater level changes, the downscaled
 608 TWSA improves the correlation with groundwater level measurements, thereby better capturing
 609 localised TWS variability. This pattern highlights the ability of the deep learning model to recover
 610 localised TWS signals in complex, human-impacted regions.



611

612 **Fig. 9** In-situ validation of downscaled TWSA across the North China Plain. **(a)** Correlation coefficient (CC)
 613 between groundwater level observations and downscaled TWSA ($CC_{\text{downscaled, obs}}$) at 559 wells during 2005–2018.
 614 **(b)** Difference in correlation between TWSA and groundwater observations after and before downscaling
 615 ($CC_{\text{downscaled, obs}} - CC_{\text{GRACE, obs}}$). **(c)** Scatterplots of the difference in CC ($CC_{\text{downscaled, obs}} - CC_{\text{GRACE, obs}}$) versus the
 616 CC before downscaling ($CC_{\text{GRACE, obs}}$). Geographical information and irrigation intensity (%) of the NCP are
 617 shown in panels (a) and (b).

618 In addition to groundwater-dominated TWS variations, we evaluated the reliability of the downscaled
 619 TWSA over a representative glacierized region, the SETP in Southwest China (Fig. 1; Fig. 10a). Based
 620 on satellite altimetry observations from ICESat (2003–2009) and CryoSat-2 (2010–2019), Zhao et al.
 621 (2022) estimated regionally averaged changes in cumulative surface elevation over the SETP. We
 622 compared region-averaged TWSA from the downscaled product and cumulative glacier surface
 623 elevation changes at both seasonal and annual time scales (Fig. 10b). All estimates were converted to
 624 anomalies relative to the 2003–2019 mean to ensure consistency. We find strong agreement between
 625 the downscaled TWSA and satellite altimetry observations at the seasonal scale across both ICESat and
 626 CryoSat-2 periods (CC : 0.90). When aggregated to the annual scale, which reduces short-term
 627 variability and observational noise, the correlation further increases to 0.97 (Fig. 10b). This represents
 628 a substantial improvement compared to conventional high-resolution TWS estimation relied on physical
 629 models, which often poorly simulate glacierized regions due to the lack of glacier modules. Specifically,
 630 the averaged TWSA over the SETP simulated by PCR-GLOBWB 2 shows a low CC value of 0.12 and
 631 0.05 at the seasonal and annual scale, respectively (Fig. 10c). These results further demonstrate that the
 632 downscaled product improves the representation of TWS variability in glacierized areas, thereby
 633 extending its applicability to cryospheric regions.



634

635 **Fig. 10** Validation of downscaled TWSA over the Southeastern Tibetan Plateau (SETP) using independent satellite
 636 altimetry retrievals. **(a)** Geographical information over the SETP. **(b)–(c)** Scatterplots of cumulative glacier
 637 surface elevation change and region-mean TWSA, estimated from **(b)** the downscaled product and **(c)** PCR-

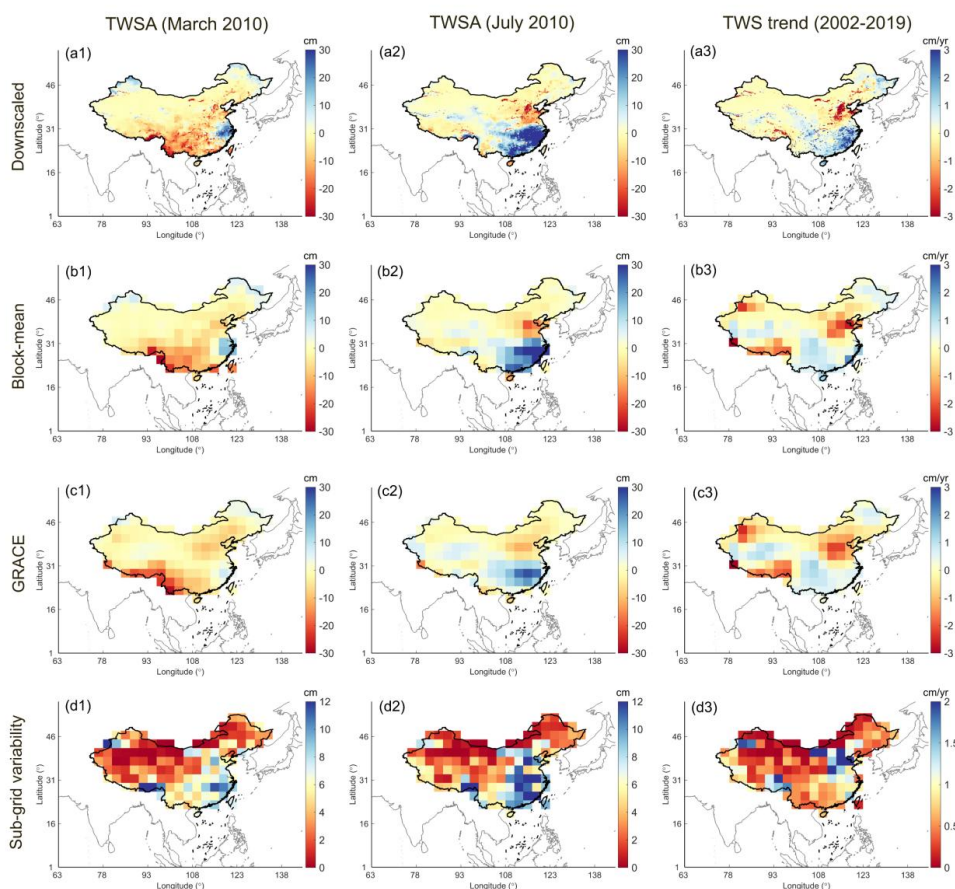


638 GLOBWB 2. In (b) and (c), seasonal mean and annual mean scatterplots are represented by dots and triangles,
639 respectively. Data pairs during the ICESat period are shown in green (seasonal) and brown (annual), whereas
640 those during the CrySat-2 period are shown in blue (seasonal) and red (annual). The expression of linear regression
641 and correlation coefficient are shown in light blue (seasonal) and dark red (annual).

642 4.4 High-resolution spatial details

643 We present spatial patterns of downscaled TWSA for a representative dry month (March 2010), a wet
644 month (July 2010), and the long-term trend during 2002–2019 to illustrate high-resolution spatial details
645 (Fig. 11a1–a3). The downscaled outputs provide enhanced representation of fine-scale spatial
646 variability and TWS trends, offering more intricate details and effectively pinpointing locations of
647 extreme TWSA values. In March 2010 (Fig. 11a1), the downscaled TWSA shows pronounced negative
648 anomalies over Southwest China, consistent with the peak of the severe 2009–2010 drought in this
649 region (Long et al., 2014; Lu et al., 2011). In contrast, in the wet month of July 2010 (Fig. 11a2), large
650 positive TWSA is reproduced across the middle and lower reaches of the Yangtze River Basin,
651 associated with intense Meiyu rainfall and widespread flooding (Gao et al., 2016). The spatial pattern
652 of TWS trends during 2002–2019 (Fig. 11a3) reveals hotspots at sub-regional scales of hundreds to
653 thousands of square kilometres, consistent with previous studies based on in situ observations and/or
654 model simulations. These hotspots include significant TWS declines in the Tianshan Mountains in
655 Northwest China, the SETP and the Himalayas in Southwest China, the Hetao Irrigation District in the
656 Yellow River Basin, the NCP in North China, and the Songhua-Nenjiang Irrigation Plain in Northeast
657 China. In contrast, substantial TWS increases are observed in the Karakoram and western Kunlun
658 Mountains in western China, the endorheic basins of the Tibetan Plateau, and humid southeastern China.
659 The identified hotspots from downscaled TWSA are consistent with previously reported hydrologic
660 processes, including glacier mass changes in high mountain regions (Farinotti et al., 2020; Farinotti et
661 al., 2015; Hugonnet et al., 2021; Zhao et al., 2022), irrigation-induced groundwater depletion in
662 intensively cultivated areas (Feng et al., 2013; Feng et al., 2025), lake expansion in the endorheic
663 Tibetan Plateau (Li et al., 2019c), and flood-driven water storage increases in humid Southeast China
664 (Slater et al., 2021). The agreement with independent studies strongly supports the physical plausibility
665 of the downscaled TWSA patterns.

666 In addition, downscaled TWSA was aggregated to a coarse resolution of 3° for the selected dry and wet
667 months as well as the long-term trend (Fig. 11b1–b3) to compare with GRACE observations. Except
668 for a few coastal grid cells, the block-mean TWSA from the downscaled product closely reproduces
669 GRACE estimates (Fig. 11c1–c3), capturing both spatial patterns and extreme values and demonstrating
670 spatial reliability. To assess the extent to which fine-scale spatial variability is represented, the standard
671 deviation of all 0.1° pixels within each 3° grid cell (900 pixels in each grid cell) was calculated as an
672 indicator of sub-grid variability (Fig. 11d1–d3). While we acknowledge that standard deviation is not a
673 direct measure of spatial details, a higher pixel-wise standard deviation, given a reasonable large-scale
674 spatial pattern, generally indicates richer fine-scale spatial variability. Overall, higher sub-grid
675 variability in both dry (Fig. 11d1) and wet (Fig. 11d2) months is observed in humid and topographically
676 complex regions, while relatively low variability is found in arid and semi-arid regions of Northwest
677 China. The median standard deviation of the dry month (March 2010; Fig. 11d1) is 9.31 cm, whereas
678 that of the wet month (July 2010; Fig. 11d2) reaches 12.62 cm. The spatial pattern of variability in the
679 long-term trend (Fig. 11d3) highlights regions with heterogeneous TWS changes, particularly in the
680 NCP, Hetao Irrigation District in the Yellow River Basin, Tianshan Mountains and SETP. These patterns
681 coincide with complex physical processes under sub-grid heterogeneity, indicating the potential of
682 downscaling TWSA in identifying localized extremes and sub-regional hotspots of water storage change.



683

684 **Fig. 11** Spatial pattern of TWSA in a representative dry month (March 2010; the first column), a wet month (July
685 2010; the second column), long-term trend during 2002–2019 (the third column), and sub-grid spatial variability
686 in each 3° grid cell (the fourth column).

687 5 Discussion

688 5.1 Comparison with representative downscaling studies

689 Two representative recent studies on TWS downscaling were selected to compare with our results. Gou
690 and Soja (2024) developed a CNN-based framework to downscale JPL-M from 3° to 0.5° at the global
691 scale, incorporating physical simulations from the WaterGAP Global Hydrologic Model (WGHM). This
692 study was selected because its model framework is conceptually comparable to our study, as both
693 frameworks utilize convolutional architectures. The other representative study was conducted by Xiong
694 et al. (2025), which employed a joint inversion-based downscaling approach to generate 0.5° TWSA
695 over China. The downscaling results based on JPL-M were compared to be consistent with the GRACE
696 data source of our study. This work represents a recent effort specifically designed for regional TWSA
697 reconstruction over China. Together, these two studies provide representative benchmarks for
698 evaluating the performance of our downscaled results.

699 The proposed downscaling framework in our study generally shows better performance compared to
700 these two studies, not only improving large-scale consistency with GRACE observations, but also
701 enhancing physical consistency and spatial detail across diverse hydrologic regimes (Table 2). At the 3°
702 grid scale, the downscaled TWSA of our study achieves a higher correlation coefficient (CC : 0.95) and



703 a substantially lower *RMSE* (1.38 cm) compared with Gou and Soja (2024) (*CC*: 0.80, *RMSE*: 3.85 cm)
 704 and Xiong et al. (2025) (*CC*: 0.57, *RMSE*: 5.00 cm). Across large basins, our results maintain
 705 comparable correlation with Gou and Soja (2024) (*CC*: 0.94 for both) while further reducing *RMSE*
 706 (1.72 cm in this study vs. 2.63 cm in Gou and Soja (2024)). These results outperform those from Xiong
 707 et al. (2025), which has a *CC* of 0.93 and *RMSE* of 3.07 cm for large river basins in China. In addition,
 708 the uncertainty level of the downscaled TWSA (0.88 cm) remains comparable to that of Gou and Soja
 709 (2024) (0.81 cm), demonstrating stable model performance at the pixel scale. The advantage of the
 710 proposed framework is more evident in process-based evaluations at finer scales. In particular, for water
 711 balance closure across small to medium catchments, the downscaled TWSA in this study shows higher
 712 correlation (*CC*: 0.80) and lower *RMSE* (1.44 cm) than both reference studies (*CC*: 0.70 and *RMSE*:
 713 1.69 cm in Gou and Soja (2024); *CC*: 0.35 and *RMSE*: 2.56 cm in Xiong et al. (2025)). This
 714 improvement could partly be attributed to the higher spatial resolution of generated TWSA of this study
 715 (0.1° compared to 0.5° in the reference studies), which allows better representation of sub-basin
 716 variability.

717 Across the NCP, where TWS variations are dominated by groundwater variations, the results of this
 718 study are comparable to those of Gou and Soja (2024). The comparable performance is reflected in the
 719 similar correlation with in-situ groundwater levels. Gou and Soja (2024) reports a *CC* of 0.42–0.63
 720 under increasing irrigation intensity, while our study achieves *CC* values of 0.41–0.65. Results of Xiong
 721 et al. (2025) show weaker agreement with a *CC* of 0.36–0.49. This could be related to the fact that both
 722 Gou and Soja (2024) and this study integrate outputs from multiple hydrologic simulations (WGHM
 723 and PCR-GLOBWB 2, respectively), which provide essential constraints on human-driven water
 724 storage changes, particularly relevant in the NCP. However, the absence of glacier processes in Gou and
 725 Soja (2024) results in poor performance over glacierized regions, with low correlations with glacier
 726 surface elevation changes (*CC*: 0.53 at seasonal scale and 0.63 at annual scale). Xiong et al. (2025)
 727 incorporated glacier processes in their inversion framework, and thus got a better performance (*CC*:
 728 0.61 at seasonal scale and 0.81 at annual scale). In comparison, downscaled TWSA in our study achieves
 729 substantially higher correlations with glacier surface elevation changes (*CC*: 0.90 at seasonal scale and
 730 0.97 at annual scale), showing clear improvements in representing glacier-dominant TWS changes.
 731 Furthermore, our results outperform these two studies in recovering fine-scale spatial variability. The
 732 standard deviation within each 3° grid cell of our downscaled TWSA is comparable to that of Xiong et
 733 al. (2025) and much higher than that of Gou and Soja (2024), indicating a stronger ability to resolve
 734 sub-grid heterogeneity than Gou and Soja (2024).

735 **Table 2** Performance comparison of downscaled TWSA in this study with two representative studies across
 736 multiple evaluation metrics. All metrics represent the median values calculated from the evaluated samples.

Evaluation index	Gou and Soja (2024)	Xiong et al. (2025)	This study
<i>Comparison with GRACE TWSA of 3° grid cells across China</i>			
<i>CC</i>	0.80	0.57	0.95
<i>RMSE</i> (cm)	3.85	5.00	1.38
<i>Comparison with GRACE TWSA of nine large river basins</i>			
<i>CC</i>	0.94	0.93	0.94
<i>RMSE</i> (cm)	2.63	3.07	1.72
<i>Pixel-wise uncertainty estimation</i>			
<i>Uncertainty</i> (cm)	0.81	–	0.88
<i>Comparison with TWS change derived from water budget in small to medium catchments</i>			
<i>CC</i>	0.70	0.35	0.80
<i>RMSE</i> (cm)	1.69	2.56	1.44
<i>Correlation with groundwater levels across the NCP</i>			
<i>Positive CC</i> (%)	80	79	80
<i>CC</i>	0.42	0.36	0.41
<i>CC</i> over areas with AEI > 40%	0.46	0.39	0.48
<i>CC</i> over areas with AEI > 50%	0.63	0.49	0.65

Correlation with regionally averaged glacier surface elevation change across the SETP



CC (seasonal)	0.53	0.61	0.90
CC (annual)	0.63	0.81	0.97
<i>Recovery of spatial details within 3° grid cell in representative months</i>			
SD (cm; dry month March 2010)	2.51	11.51	9.31
SD (cm; wet month July 2010)	2.79	12.31	12.62

737 *Note:* Xiong et al. (2025) does not provide uncertainty estimation associated with their downscaling product.
738 Regarding the underlying GRACE data, this study and Gou and Soja (2024) are based on the JPL mascon solution,
739 whereas Xiong et al. (2025) is based on the JPL spherical harmonic solution. AEI: area equipped irrigation; SD:
740 standard deviation.

741 5.2 Applications, limitations and future work

742 Important applications of downscaled TWSA can be summarized in three main aspects. First, the high-
743 resolution TWSA enables detailed assessment of spatial heterogeneity in TWS changes, facilitating the
744 identification of sub-regional hotspots associated with climate variability and human activities. For
745 example, high mountain regions are clearly resolved in the downscaled TWSA patterns, and human-
746 induced TWS changes, such as large TWS depletion in the Hetao Irrigation District within the Yellow
747 River Basin, can be explicitly captured compared to the smoothed patterns in the original GRACE
748 observations (Fig. 11b3). This improved spatial representation allows for more accurate characterization
749 of water risk and resource dynamics at sub-regional scales, tracking impacts from both climate
750 variability and human activities on China's water availability.

751 Second, the enhanced spatial resolution improves the identification of hydrologic extremes by capturing
752 TWSA variability within individual grid cells. For example, during the severe drought in Southwest
753 China in 2009–2010, the downscaled TWSA shows a region-averaged negative anomaly of 16 cm in
754 Yunnan Province (97°–106°E, 21°–29°N) at the drought peak in March 2010. This is comparable to the
755 GRACE estimate (-16.6 cm) at the region-averaged scale. However, the downscaled TWSA reveals a
756 clear northeast-southwest gradient in drought severity across Southwest China (Fig. 11a1). The most
757 pronounced negative anomaly reaches -57 cm in southwestern Yunnan, nearly twice the corresponding
758 GRACE estimate (-29 cm). This highlights the potential of the downscaled TWSA to resolve sub-
759 regional heterogeneity and capture localized extremes beyond the native resolution of GRACE.

760 Third, the dataset is provided at 0.1° resolution, consistent with the spatial resolution of commonly used
761 climate forcing (e.g., ERA5-Land) and hydrologic flux estimation (e.g., 0.1° precipitation derived from
762 Integrated Multisatellite Retrievals for Global Precipitation Mission). This consistency allows direct
763 multi-variable analysis without additional resampling. High-resolution TWSA facilitates integrated
764 analyses of interactions between water storage and hydrologic fluxes at fine scales, whose relationships
765 remain unclear (Tiwari et al., 2025). Further, the downscaled TWSA can serve as a storage reference
766 for evaluating performance of hydrologic model, reducing the risk of equifinality associated with flux-
767 only calibration (Arsenault and Brissette, 2014; Beven, 2006).

768 One caveat of the deep learning framework is that it primarily learns spatial relationships, with limited
769 explicit constraint on temporal dynamics. This is inherent to the convolutional architecture, while the
770 spatial construction of training and validation samples at each monthly time step further reinforces this
771 tendency. In this study, the temporal continuity of input physical variables mitigates this limitation, and
772 the downscaled TWSA still shows good performance in temporal evaluations, such as the reliability of
773 TWSA time series within fixed spatial scale. However, temporal dependencies are not explicitly
774 constrained in the current framework. Future work could couple the U-Net framework with temporal
775 modelling architectures, such as ConvLSTM (Shi et al., 2015), to jointly learn spatial patterns and
776 temporal dependencies. This extension would improve the robustness of downscaled outputs to capture
777 the spatiotemporal dynamics of TWSA and enhance its capability for predictions and projections.

778 It is worth noting that the estimated uncertainty in the downscaled product refers to that introduced
779 during the deep learning downscaling. However, it does not include uncertainties in the input labels and
780 predictors, which may contribute to lower performance in coastal regions in Southeast China and arid
781 regions in Northwest China. In coastal regions, GRACE retrievals are affected by signal leakage and
782 reduced accuracy near land-ocean boundaries. Missing values over ocean areas in PCR-GLOBWB 2
783 limit the number of effective training patches, potentially diminishing the robustness of learning



784 processes. In addition, Southeast China is characterized by high variability in precipitation, driven by
785 the East Asian monsoon and ocean-atmosphere coupling in the western Pacific (Ding and Chan, 2005).
786 This may introduce additional uncertainty in climate-related predictors. In arid regions, TWS variability
787 is often weakly coupled with hydrologic predictors. The lower signal-to-noise ratios of GRACE
788 observations and model-based inputs further increase uncertainty in both training and validation,
789 making it more difficult to learn reliable relationships (Gou and Soja, 2024). Future work could integrate
790 uncertainty propagation from input data within the modelling framework, enabling a more
791 comprehensive and systematic characterization of dataset uncertainty.

792 **6 Conclusions**

793 This study downscaled 3° GRACE TWSA to a much higher resolution of 0.1° using a deep learning
794 framework. The downscaled TWSA over China during 2002–2019 demonstrates its high reliability in
795 preserving large-scale GRACE signals (3°) while enhancing fine-scale spatial detail (0.1°). Compared
796 with GRACE TWSA, the downscaled product achieves high *CC* values of 0.95 at the 3° grid scale and
797 0.94 at the basin scale, with corresponding low *RMSE* values of 1.38 cm and 1.72 cm. Independent
798 training runs estimate a median pixel-wise uncertainty of 0.88 cm, lower than the median uncertainty
799 of GRACE JPL-M across China (1.96 cm). Process-based evaluation further confirms the physical
800 consistency of the downscaled TWSA. Water balance closure across 163 small to medium catchments
801 (10,000–100,000 km²) shows good agreement between TWS changes derived from the downscaled
802 TWSA and ERA5L-based hydrologic fluxes (median *CC* of 0.80 and *RMSE* of 1.44 cm). In addition,
803 downscaled TWSA shows median *CC* values of 0.48 and 0.65 with groundwater level observations over
804 the NCP, corresponding to regions with irrigation intensity exceeding 40% and 50%, respectively. Over
805 the glacierized SETP, the downscaled TWSA captures variations in cumulative glacier surface elevation
806 derived from ICESat and CryoSat-2 altimetry, with high *CC* values of 0.90 and 0.97 at the seasonal and
807 annual scales, respectively.

808 Results from our study generally show better performance than those from two recent representative
809 studies on TWS downscaling in China. In addition, the spatial pattern of downscaled TWS trends
810 reveals hotspots at sub-regional scales of hundreds to thousands of square kilometres, such as mountain
811 ranges and irrigation zones, consistent with previous studies based on in situ observations and/or model
812 simulations. The downscaled TWSA provides enhanced spatial detail for identifying fine-scale hotspots
813 and sub-regional heterogeneity in TWS changes associated with hydro-climatic extremes, which are not
814 resolved by the native resolution of GRACE. The resolution of downscaled TWSA (0.1°) aligns well
815 with the commonly used climate forcing and hydrologic flux estimation, enabling direct analyses in
816 storage-flux interactions and evaluation of hydrologic model performance from a storage perspective.

817 Estimating high-resolution TWSA is still challenging. Despite some limitations (e.g., lack of explicit
818 constraints on temporal generalization), this study provides valuable insights for integrating multi-
819 source information to bridge the scale gap between coarse-resolution gravity observations and fine-
820 scale hydrologic processes. The generated datasets not only support investigation of sub-regional TWS
821 variability and associated hydrologic processes in China, but also offer insights for advancing statistical
822 downscaling of TWSA at broader scales.

823 **7 Data and code availability**

824 The new 0.1° TWSA dataset (Li and Sun, 2026a) is archived on *Zenodo* at
825 <https://doi.org/10.5281/zenodo.19502906>. The uploaded data include (1) the standard version covering
826 2002–2019 with the most comprehensive observations for quality control, and (2) the extended version
827 spanning 2020–2023 to support recent hydrologic analyses. This repository also provides the complete
828 training datasets used for model development and reproducibility, consisting of training labels, dynamic
829 predictors, and static predictors.

830 The U-Net model was implemented using standard Python packages. For reproducibility, the core codes
831 used to build the model (Li and Sun, 2026b) are available on *Zenodo* at
832 <https://doi.org/10.5281/zenodo.19502631>.

833 All datasets used in this study are readily accessible. GRACE and GRACE-FO data are provided at
834 <https://grace.jpl.nasa.gov/>. Water storage simulations from PCR-GLOBWB 2 are accessed at



835 <https://globalhydrology.nl/research/models/pcr-globwb-2-0/>. Glacier mass balance produced by
836 Hugonnet et al. (2021) can be accessed at <https://doi.org/10.6096/13>. Input predictors include: ERA5L
837 data available at [https://cds.climate.copernicus.eu/datasets/reanalysis-era5-land-monthly-](https://cds.climate.copernicus.eu/datasets/reanalysis-era5-land-monthly-means?tab=download)
838 [means?tab=download](https://cds.climate.copernicus.eu/datasets/reanalysis-era5-land-monthly-means?tab=download); NDVI derived at <https://modis.gsfc.nasa.gov/data/dataproduct/mod13.php>; DEM
839 data available at <https://search.earthdata.nasa.gov/>; irrigation maps accessed at
840 <https://www.fao.org/aquastat/en/geospatial-information/global-maps-irrigated-areas/latest-version>;
841 and population datasets derived at [https://www.earthdata.nasa.gov/data/catalog/sedac-ciesin-sedac-](https://www.earthdata.nasa.gov/data/catalog/sedac-ciesin-sedac-gpww4-popcount-r11-4.11)
842 [gpww4-popcount-r11-4.11](https://www.earthdata.nasa.gov/data/catalog/sedac-ciesin-sedac-gpww4-popcount-r11-4.11). Independent observations of groundwater levels are accessed at
843 <https://www.geodata.cn/oldindex.html>. Geographic information for the study region includes:
844 shapefiles of China's national boundary, major rivers, and the delineation of river basins accessed from
845 Resource and Environmental Science Data Center at <https://www.resdc.cn/>; glacier distribution
846 accessed from RGI 6.0 at <https://www.glims.org/RGI/andolph60.html>; and location of lakes accessed
847 from HydroLakes at <https://www.hydrosheds.org/products/hydrolakes>.

848 **Author contributions**

849 X.L. conceptualized the study. X.L. and Y.S. developed the methodology and performed the analysis.
850 L.J.S., X.G., N.W., and B.R.S. provided critical input and contributed to the interpretation of results.
851 X.L. drafted the manuscript. All authors discussed the results and revised the manuscript.

852 **Competing interests**

853 The authors declare no competing interests.

854 **Acknowledgements**

855 The authors thank F. Zhao from Tsinghua University for providing cumulative glacier surface elevation
856 over the SETP retrieved from ICESat and CryoSat-2 satellites, as well as J. Gou from ETH Zurich and
857 B. Li from NASA Goddard Space Flight Center for valuable discussions on the U-Net model training.

858 **Financial support**

859 This work is supported by UK Research and Innovation (UKRI) under the Horizon Europe guarantee
860 scheme (EP/Z002729/1), originally funded by the European Commission through the Marie Curie
861 Programme. X.G. acknowledges support from the National Natural Science Foundation of China
862 (Grants 42371041 & U2340230).

863 **References**

- 864 Arsenault, R. and Brissette, F.P., 2014. Continuous streamflow prediction in ungauged basins: The
865 effects of equifinality and parameter set selection on uncertainty in regionalization approaches.
866 *Water Resources Research*, 50(7): 6135–6153.
- 867 Arshad, A., Mirchi, A., Taghvaeian, S. et al., 2024. Downscaled-GRACE data reveal anthropogenic and
868 climate-induced water storage decline across the Indus Basin. *Water Resources Research*, 60(7):
869 e2023WR035882.
- 870 Bai, P., Liu, X. and Liu, C., 2018. Improving hydrological simulations by incorporating GRACE data
871 for model calibration. *Journal of Hydrology*, 557: 291–304.
- 872 Beven, K., 2006. A manifesto for the equifinality thesis. *Journal of hydrology*, 320(1-2): 18–36.
- 873 Cao, G., Zheng, C., Scanlon, B.R. et al., 2013. Use of flow modeling to assess sustainability of
874 groundwater resources in the North China Plain. *Water Resources Research*, 49(1): 159–175.
- 875 Chen, J., Cazenave, A., Dahle, C. et al., 2022. Applications and challenges of GRACE and GRACE
876 follow-on satellite gravimetry. *Surveys in Geophysics*, 43(1): 305–345.
- 877 Chen, X., Long, D., Hong, Y. et al., 2017. Improved modeling of snow and glacier melting by a
878 progressive two-stage calibration strategy with GRACE and multisource data: How snow and
879 glacier meltwater contributes to the runoff of the U pper Brahmaputra River basin? *Water*
880 *Resources Research*, 53(3): 2431–2466.



- 881 Ding, Y. and Chan, J.C., 2005. The East Asian summer monsoon: an overview. *Meteorology and*
882 *Atmospheric Physics*, 89(1): 117–142.
- 883 Döll, P., Müller Schmied, H., Schuh, C. et al., 2014. Global-scale assessment of groundwater depletion
884 and related groundwater abstractions: Combining hydrological modeling with information from
885 well observations and GRACE satellites. *Water Resources Research*, 50(7): 5698–5720.
- 886 Eicker, A., Schumacher, M., Kusche, J. et al., 2014. Calibration/data assimilation approach for
887 integrating GRACE data into the WaterGAP Global Hydrology Model (WGHM) using an
888 ensemble Kalman filter: First results. *Surveys in Geophysics*, 35(6): 1285–1309.
- 889 Farinotti, D., Immerzeel, W.W., de Kok, R.J. et al., 2020. Manifestations and mechanisms of the
890 Karakoram glacier Anomaly. *Nature geoscience*, 13(1): 8–16.
- 891 Farinotti, D., Longuevergne, L., Moholdt, G. et al., 2015. Substantial glacier mass loss in the Tien Shan
892 over the past 50 years. *Nature Geoscience*, 8(9): 716–722.
- 893 Feng, W., Zhong, M., Lemoine, J.M. et al., 2013. Evaluation of groundwater depletion in North China
894 using the Gravity Recovery and Climate Experiment (GRACE) data and ground-based
895 measurements. *Water Resources Research*, 49(4): 2110–2118.
- 896 Feng, Z., Miao, Q., Shi, H. et al., 2025. Water Saving and Environmental Issues in the Hetao Irrigation
897 District, the Yellow River Basin: Development Perspective Analysis. *Agronomy*, 15(7): 1654.
- 898 Gao, T., Xie, L. and Liu, B., 2016. Association of extreme precipitation over the Yangtze River Basin
899 with global air-sea heat fluxes and moisture transport. *International Journal of Climatology*,
900 36(8): 3020–3038.
- 901 Gerdener, H., Kusche, J., Schulze, K. et al., 2023. The global land water storage data set release 2
902 (GLWS2.0) derived via assimilating GRACE and GRACE-FO data into a global hydrological
903 model. *Journal of Geodesy*, 97(7): 73.
- 904 Gou, J. and Soja, B., 2024. Global high-resolution total water storage anomalies from self-supervised
905 data assimilation using deep learning algorithms. *Nature Water*, 2(2): 139–150.
- 906 Gu, X., Zhang, Q., Li, J. et al., 2019. Attribution of global soil moisture drying to human activities: A
907 quantitative viewpoint. *Geophysical Research Letters*, 46(5): 2573–2582.
- 908 Guan, Y., Gu, X., Dai, A. et al., 2025. Anthropogenic enhancement of subsurface soil moisture droughts.
909 *Nature Climate Change*: 1–8.
- 910 Guan, Y., Gu, X., Slater, L.J. et al., 2023. Increase in ocean-onto-land droughts and their drivers under
911 anthropogenic climate change. *Npj Climate and Atmospheric Science*, 6(1): 195.
- 912 Gupta, H.V., Kling, H., Yilmaz, K.K. et al., 2009. Decomposition of the mean squared error and NSE
913 performance criteria: Implications for improving hydrological modelling. *Journal of hydrology*,
914 377(1-2): 80–91.
- 915 Hugonnet, R., McNabb, R., Berthier, E. et al., 2021. Accelerated global glacier mass loss in the early
916 twenty-first century. *Nature*, 592(7856): 726–731.
- 917 Huss, M., 2013. Density assumptions for converting geodetic glacier volume change to mass change.
918 *The Cryosphere*, 7(3): 877–887.
- 919 Immerzeel, W.W., Lutz, A.F., Andrade, M. et al., 2020. Importance and vulnerability of the world's
920 water towers. *Nature*, 577(7790): 364–369.
- 921 Janzing, J., Wanders, N., van Tiel, M. et al., 2025. Hyper-resolution large-scale hydrological modelling
922 benefits from improved process representation in mountain regions. *Hydrology and Earth*
923 *System Sciences*, 29(23): 7041–7071.
- 924 Kalu, I., Ndehedehe, C.E., Ferreira, V.G. et al., 2024. Statistical downscaling of GRACE terrestrial
925 water storage changes based on the Australian Water Outlook model. *Scientific Reports*, 14(1):
926 10113.



- 927 Kling, H., Fuchs, M. and Paulin, M., 2012. Runoff conditions in the upper Danube basin under an
928 ensemble of climate change scenarios. *Journal of hydrology*, 424: 264–277.
- 929 Landerer, F.W., Dickey, J.O. and Güntner, A., 2010. Terrestrial water budget of the Eurasian pan-Arctic
930 from GRACE satellite measurements during 2003–2009. *Journal of Geophysical Research:*
931 *Atmospheres*, 115(D23).
- 932 Landerer, F.W. and Swenson, S., 2012. Accuracy of scaled GRACE terrestrial water storage estimates.
933 *Water resources research*, 48(4).
- 934 Li, B., Rodell, M., Kumar, S. et al., 2019a. Global GRACE data assimilation for groundwater and
935 drought monitoring: Advances and challenges. *Water Resources Research*, 55(9): 7564–7586.
- 936 Li, F. and Kusche, J., 2026. Reproducing GRACE total water storage change at finer spatial scales.
937 *Geophysical Research Letters*, 53(1): e2025GL119881.
- 938 Li, X., Long, D., Han, Z. et al., 2019b. Evapotranspiration estimation for Tibetan Plateau headwaters
939 using conjoint terrestrial and atmospheric water balances and multisource remote sensing.
940 *Water Resources Research*, 55(11): 8608–8630.
- 941 Li, X., Long, D., Huang, Q. et al., 2019c. High-temporal-resolution water level and storage change data
942 sets for lakes on the Tibetan Plateau during 2000–2017 using multiple altimetric missions and
943 Landsat-derived lake shoreline positions. *Earth System Science Data*, 11(4): 1603–1627.
- 944 Li, X., Long, D., Scanlon, B.R. et al., 2022. Climate change threatens terrestrial water storage over the
945 Tibetan Plateau. *Nature Climate Change*, 12(9): 801–807.
- 946 Li, X., Long, D., Scanlon, B.R. et al., 2025. Retrievals and simulations of terrestrial water storage
947 changes and runoff over the Tibetan Plateau: Challenges and opportunities. *Fundamental*
948 *Research*.
- 949 Li, X. and Sun, Y., 2026a. A 0.1° terrestrial water storage anomaly dataset over China. In: Zenodo.
- 950 Li, X. and Sun, Y., 2026b. Code for: A physically guided deep learning reconstruction of terrestrial
951 water storage anomalies at 0.1° across China. In: Zenodo.
- 952 Liu, J., Yang, H., Gosling, S.N. et al., 2017. Water scarcity assessments in the past, present, and future.
953 *Earth's future*, 5(6): 545–559.
- 954 Long, D., Scanlon, B.R., Longuevergne, L. et al., 2013. GRACE satellite monitoring of large depletion
955 in water storage in response to the 2011 drought in Texas. *Geophysical Research Letters*, 40(13):
956 3395–3401.
- 957 Long, D., Shen, Y., Sun, A. et al., 2014. Drought and flood monitoring for a large karst plateau in
958 Southwest China using extended GRACE data. *Remote Sensing of Environment*, 155: 145–160.
- 959 Long, D., Xu, Y., Cui, Y. et al., 2025. Unprecedented large-scale aquifer recovery through human
960 intervention. *Nature Communications*, 16(1): 7296.
- 961 Long, D., Yang, W., Scanlon, B.R. et al., 2020. South-to-North Water Diversion stabilizing Beijing's
962 groundwater levels. *Nature Communications*, 11(1): 3665.
- 963 Loomis, B., Luthcke, S. and Sabaka, T., 2019. Regularization and error characterization of GRACE
964 mascons. *Journal of geodesy*, 93(9): 1381–1398.
- 965 Lu, E., Luo, Y., Zhang, R. et al., 2011. Regional atmospheric anomalies responsible for the 2009–2010
966 severe drought in China. *Journal of Geophysical Research: Atmospheres*, 116(D21).
- 967 Muñoz-Sabater, J., Dutra, E., Agustí-Panareda, A. et al., 2021. ERA5-Land: A state-of-the-art global
968 reanalysis dataset for land applications. *Earth System Science Data*, 13(9): 4349–4383.
- 969 Pascal, C., Ferrant, S., Selles, A. et al., 2022. Evaluating downscaling methods of GRACE (Gravity
970 Recovery and Climate Experiment) data: a case study over a fractured crystalline aquifer in
971 southern India. *Hydrology and Earth System Sciences*, 26(15): 4169–4186.



- 972 Pfeffer, W.T., Arendt, A.A., Bliss, A. et al., 2014. The Randolph Glacier Inventory: a globally complete
973 inventory of glaciers. *Journal of glaciology*, 60(221): 537–552.
- 974 Piao, S., Ciais, P., Huang, Y. et al., 2010. The impacts of climate change on water resources and
975 agriculture in China. *Nature*, 467(7311): 43–51.
- 976 Pokhrel, Y., Felfelani, F., Satoh, Y. et al., 2021. Global terrestrial water storage and drought severity
977 under climate change. *Nature Climate Change*, 11(3): 226–233.
- 978 Reager, J.T., Thomas, B.F. and Famiglietti, J.S., 2014. River basin flood potential inferred using
979 GRACE gravity observations at several months lead time. *Nature Geoscience*, 7(8): 588–592.
- 980 Rodell, M., Famiglietti, J.S., Wiese, D.N. et al., 2018. Emerging trends in global freshwater availability.
981 *Nature*, 557(7707): 651–659.
- 982 Rowlands, D.D., Luthcke, S.B., Klosko, S. et al., 2005. Resolving mass flux at high spatial and temporal
983 resolution using GRACE intersatellite measurements. *Geophysical Research Letters*, 32(4).
- 984 Save, H., Bettadpur, S. and Tapley, B.D., 2016. High-resolution CSR GRACE RL05 mascons. *Journal*
985 *of Geophysical Research: Solid Earth*, 121(10): 7547–7569.
- 986 Scanlon, B.R., Fakhreddine, S., Rateb, A. et al., 2023. Global water resources and the role of
987 groundwater in a resilient water future. *Nature Reviews Earth & Environment*, 4(2): 87–101.
- 988 Scanlon, B.R., Zhang, Z., Save, H. et al., 2018. Global models underestimate large decadal declining
989 and rising water storage trends relative to GRACE satellite data. *Proceedings of the National*
990 *Academy of Sciences*, 115(6): E1080–E1089.
- 991 Shi, X., Chen, Z., Wang, H. et al., 2015. Convolutional LSTM network: A machine learning approach
992 for precipitation nowcasting. *Advances in neural information processing systems*, 28.
- 993 Slater, L., Villarini, G., Archfield, S. et al., 2021. Global changes in 20-year, 50-year, and 100-year river
994 floods. *Geophysical Research Letters*, 48(6): e2020GL091824.
- 995 Sun, A.Y., Save, H., Rateb, A. et al., 2024. Deciphering the role of total water storage anomalies in
996 mediating regional flooding. *Geophysical Research Letters*, 51(16): e2023GL108126.
- 997 Sutanudjaja, E.H., Van Beek, R., Wanders, N. et al., 2018. PCR-GLOBWB 2: a 5 arcmin global
998 hydrological and water resources model. *Geoscientific Model Development*, 11(6): 2429–2453.
- 999 Tapley, B.D., Watkins, M.M., Flechtner, F. et al., 2019. Contributions of GRACE to understanding
1000 climate change. *Nature Climate Change*, 9(5): 358–369.
- 1001 Tiwari, A.D., Pokhrel, Y., Boulange, J. et al., 2025. Similarities and divergent patterns in hydrologic
1002 fluxes and storages simulated by global water models. *Nature Water*: 1–11.
- 1003 Vishwakarma, B.D., Zhang, J. and Sneeuw, N., 2021. Downscaling GRACE total water storage change
1004 using partial least squares regression. *Scientific data*, 8(1): 95.
- 1005 Wahr, J., Swenson, S. and Velicogna, I., 2006. Accuracy of GRACE mass estimates. *Geophysical*
1006 *Research Letters*, 33(6).
- 1007 Wang, Y., Li, C., Cui, Y. et al., 2024. Spatial downscaling of GRACE-derived groundwater storage
1008 changes across diverse climates and human interventions with Random Forests. *Journal of*
1009 *Hydrology*, 640: 131708.
- 1010 Wiese, D.N., Bienstock, B., Blackwood, C. et al., 2022. The mass change designated observable study:
1011 overview and results. *Earth and Space Science*, 9(8): e2022EA002311.
- 1012 Wiese, D.N., Landerer, F.W. and Watkins, M.M., 2016. Quantifying and reducing leakage errors in the
1013 JPL RL05M GRACE mascon solution. *Water Resources Research*, 52(9): 7490–7502.
- 1014 Xiong, Y., Feng, W., Bai, H. et al., 2025. High-resolution terrestrial water storage anomalies and
1015 components in China from GRACE/GFO via joint inversion downscaling. *Water Resources*
1016 *Research*, 61(7): e2024WR038996.



- 1017 Yang, H., Cao, W., Zhi, C. et al., 2021. Evolution of groundwater level in the North China Plain in the
1018 past 40 years and suggestions on its overexploitation treatment. *Geology in China*, 48(4): 1142–
1019 1155.
- 1020 Yang, W., Long, D., Scanlon, B.R. et al., 2022. Human intervention will stabilize groundwater storage
1021 across the North China Plain. *Water Resources Research*, 58(2): e2021WR030884.
- 1022 Yao, T., Bolch, T., Chen, D. et al., 2022. The imbalance of the Asian water tower. *Nature Reviews Earth
1023 & Environment*, 3(10): 618–632.
- 1024 Yazdian, H., Salmani-Dehaghi, N. and Alijanian, M., 2023. A spatially promoted SVM model for
1025 GRACE downscaling: Using ground and satellite-based datasets. *Journal of Hydrology*, 626:
1026 130214.
- 1027 Zhang, Q., Liu, C., Xu, C.-y. et al., 2006. Observed trends of annual maximum water level and
1028 streamflow during past 130 years in the Yangtze River basin, China. *Journal of hydrology*,
1029 324(1-4): 255–265.
- 1030 Zhao, F., Long, D., Li, X. et al., 2022. Rapid glacier mass loss in the Southeastern Tibetan Plateau since
1031 the year 2000 from satellite observations. *Remote Sensing of Environment*, 270: 112853.