1   **Comprehensive Inter-comparison of Generative AI Models for Super-**
2   **Resolution Precipitation Downscaling Across Hydroclimatic Regimes**

3

4   Shivam Singh*[1,2], Simon Michael Papalexiou[3,4], Hebatallah M. Abdelmoaty[5,6],
5   Tom Hartvigsen[7], Antonios Mamalakis[2,7]
6   [1]Environmental Institute, University of Virginia, Charlottesville, VA, USA;
7   [2]Department of Environmental Sciences, University of Virginia, Charlottesville, VA, USA;
8   [3]Institute for Global Water Security, Hamburg University of Technology, Germany.
9   [4]Faculty of Environmental Sciences, Czech University of Life Sciences Prague;
10  [5]Department of Civil Engineering, Schulich School of Engineering, University of Calgary,
11  Canada;
12  [6]Irrigation and Hydraulics Department, Faculty of Engineering, Cairo University, Egypt
13  [7]School of Data Science, University of Virginia, Charlottesville, VA, USA;
14  * Corresponding authors; Shivam Singh (wpa8me@virginia.edu), Antonios Mamalakis
15  (npa4tg@virginia.edu)

16  **Abstract**

17  High-resolution precipitation information is essential for hydrologic modeling, flood
18  forecasting, and climate-risk assessment, yet global weather and climate models operate at
19  spatial resolutions too coarse to resolve storm structure, intermittency, and extremes. Deep-
20  learning-based statistical downscaling provides a computationally efficient alternative to
21  dynamical downscaling, but deterministic convolutional neural networks often yield overly
22  smooth predictions and underestimate fine-scale variability and extreme events. Generative
23  deep-learning models, including generative adversarial networks and diffusion models, offer
24  a promising alternative by enabling stochastic downscaling and explicit representation of
25  uncertainty. This study presents a systematic, hydrologically oriented comparison of three
26  representative deep-learning frameworks for precipitation super-resolution: a
27  convolutional U-NET, a conditional Wasserstein GAN (WGAN), and a conditional denoising
28  diffusion probabilistic model (DDPM). Using a perfect-model experimental design based on
29  ERA5-Land precipitation over distinct hydroclimatic regions of the United States, we
30  evaluate performance under 8-times (8×) and 16-times (16×) downscaling tasks within a
31  unified training and evaluation framework. Models are evaluated using diagnostics that
32  examine precipitation distributions, wet–dry occurrence, extremes, spatial structure, storm
33  morphology, mass consistency, ensemble variability, and computational cost. All three
34  models preserve aggregate rainfall mass despite the absence of explicit physical constraints.
35  Differences arise primarily at fine spatial scales and in the representation of extremes,
36  spatial dependence, and uncertainty. U-NET provides stable and computationally efficient
37  predictions but smooths small-scale variability. WGAN improves fine-scale structure and

38    heavy-tail behavior at the expense of increased noise. The DDPM yields physically coherent
39    ensemble members and an explicit representation of uncertainty, at a substantially higher
40    computational cost.

41    Keywords: Precipitation downscaling; deep learning; generative models; super-resolution;
42    hydrologic extremes; uncertainty quantification.

## 1. Introduction

44    Global climate models are fundamental tools for projecting future hydroclimate, yet their
45    typical horizontal spatial resolution (often on the order of ∼100–200 km) remains too coarse
46    to represent the mesoscale and storm-scale processes that govern precipitation
47    intermittency and extremes (Feser et al., 2011; Palmer, 2014; Schär, 2019). In contrast,
48    hydrologic impact modeling, flood-risk assessment, and climate adaptation planning
49    commonly require precipitation information at finer resolution (∼10 km), where localized
50    gradients, orographic forcing, land–sea contrasts, and convective organization must be
51    adequately represented (Lucas-Picher et al., 2021; Nishant et al., 2023; Piani et al., 2010;
52    Stephens, 2017). As a result, the direct application of coarse-resolution model output is
53    inadequate for hydrologic impact assessment, flood risk estimation, and climate-risk
54    analysis, especially in situations characterized by strong spatial intermittency and extremes
55    (Giorgi & Gutowski, 2015; Tabari et al., 2021; Wood et al., 2004). To address this scale
56    mismatch and provide high-resolution information required for hydrologic and climate-
57    impact applications, downscaling techniques are employed to infer fine-scale fields from
58    coarse-resolution model output. Dynamical downscaling, which relies on physics-based
59    regional climate models (RCMs), has been widely used to improve the representation of
60    mesoscale processes and precipitation extremes (Coppola et al., 2020; Giorgi & Gutowski,
61    2015; Giorgi & Mearns, 1991; Maraun et al., 2010). However, the substantial computational
62    cost of RCMs severely constrains ensemble size, limits the exploration of uncertainty, and
63    restricts their applicability for large multi-model or multi-scenario studies (Deser et al.,
64    2012; Gao et al., 2012; Tomasi et al., 2025). These limitations have motivated growing
65    interest in empirical downscaling approaches, including statistical and machine-learning-
66    based methods, which offer orders-of-magnitude reductions in computational cost (Baño-
67    Medina et al., 2020; Hobeichi et al., 2023; Lange, 2019; Mamalakis et al., 2017; Vrac et al.,
68    2007).

69    Recent advances in deep learning have substantially reshaped empirical precipitation
70    downscaling. Convolutional neural networks (CNNs), particularly end-to-end architectures
71    such as the U-NET, have demonstrated strong skill in reproducing mean precipitation
72    patterns, spatial continuity, and wet–dry occurrence with stable training and fast inference,
73    making them attractive for large-scale and operational applications (Baño-Medina et al.,
74    2020; Höhlein et al., 2020; Vandal et al., 2017). However, deterministic CNNs are typically

75    optimized using pixel-wise loss functions that favor conditional mean solutions, resulting in
76    overly smooth precipitation fields, reduced small-scale variability, and systematic
77    underrepresentation of extremes, especially at high spatial resolutions and large
78    downscaling factors (Abdelmoaty et al., 2025; Ravuri et al., 2021). To address these
79    limitations, stochastic generative models have been increasingly explored to represent the
80    inherent non-uniqueness of fine-scale precipitation conditioned on coarse inputs (Rampal et
81    al., 2024). Generative adversarial networks (GANs), including Wasserstein GANs, have
82    shown improved representation of fine-scale structure and extreme intensities relative to
83    deterministic models, but their training is sensitive to hyperparameter choices and can
84    suffer from instabilities and mode collapse, complicating robustness and calibration
85    (Arjovsky et al., 2017; Gulrajani et al., 2017; Harris et al., 2022). More recently, diffusion-
86    based generative models have emerged as an alternative stochastic framework, offering
87    improved training stability and flexible uncertainty representation through a forward–
88    reverse diffusion process, with latent diffusion variants improving computational efficiency
89    by operating in a compressed feature space (Khader et al., 2023; Lyu et al., 2024; Tomasi et
90    al., 2025). Early applications suggest that diffusion models can outperform both
91    deterministic CNNs and GANs in capturing multiscale variability and spatial organization,
92    though at substantially higher inference cost and with sensitivity to diffusion configuration
93    (Mardani et al., 2025; X. Wang et al., 2025). Despite these advances, existing studies typically
94    develop and evaluate deterministic, adversarial, and diffusion-based approaches in isolation,
95    using case-specific designs and single target resolutions, leaving key questions unresolved
96    regarding their relative performance, uncertainty representation, and computational
97    scalability for hydrologically relevant diagnostics.

98        In this study, we address these gaps through a comprehensive, hydrologically
99    oriented comparison of three representative deep-learning frameworks for daily
100   precipitation super-resolution downscaling: a deterministic U-NET, a conditional
101   Wasserstein GAN (WGAN), and a conditional denoising diffusion probabilistic model
102   (DDPM). Using a common training and evaluation framework, we assess not only mean
103   predictive accuracy but also wet–dry occurrence, storm morphology, spatial dependence,
104   extreme precipitation behavior, and rainfall mass consistency across scales. We further
105   exploit multi-seed realizations of the generative models to quantify ensemble variability and
106   compare it with deterministic behavior, while explicitly evaluating computational
107   requirements for training and inference. By placing deterministic and generative approaches
108   on equal footing and emphasizing physically meaningful diagnostics, this work aims to
109   clarify the relative strengths and limitations of contemporary deep-learning downscaling
110   frameworks and to provide actionable guidance for hydrologic and climate-risk applications
111   that require reliable extremes, coherent spatial structure, and interpretable uncertainty
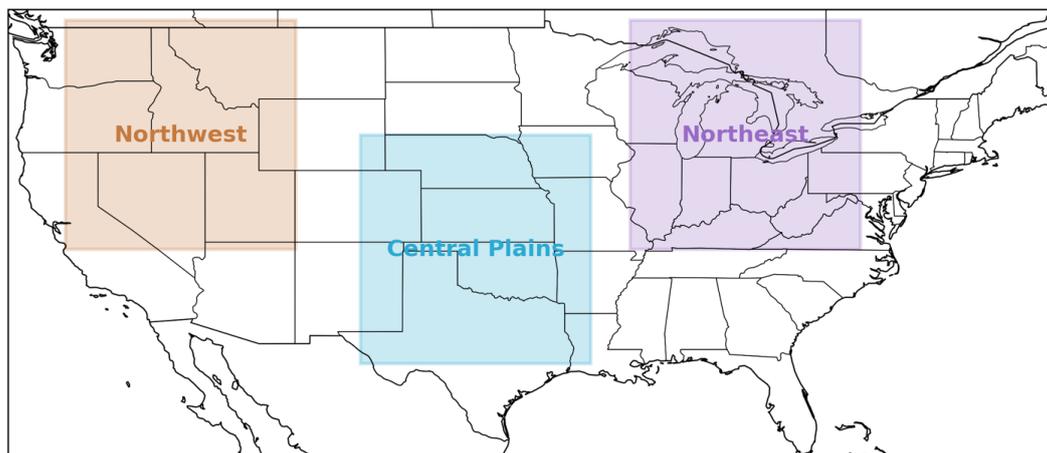112   under increasing resolution demands.

## 2. Data and Methodology

### 2.1 Data

We use precipitation fields from the ERA5-Land reanalysis dataset (Muñoz-Sabater et al., 2021), which provides hourly accumulated precipitation at 0.1° (~9 km) spatial resolution over global land surfaces. Hourly precipitation is aggregated to daily totals and extracted over the contiguous United States (CONUS) for the period 1980–2014. This 35-year record is sufficiently long to sample a wide range of storm types, hydroclimatic variability, and interannual fluctuations. ERA5-Land is selected because it provides physically consistent land-surface precipitation estimates that are widely used in hydrologic modeling and downscaling evaluation.

To assess model behavior across distinct precipitation climatologies, we define three non-overlapping 128 × 128 grid domains (~1150 × 1150 km at mid-latitudes) representing major U.S. hydroclimatic regimes (Figure 1): (i) Central Plains (30.2°N–43.0°N, 105.4°W–92.6°W), dominated by warm-season convective precipitation; (ii) Pacific Northwest (36.6°N–49.4°N, 121.8°W–109.0°W), characterized by wintertime orographic enhancement and strong seasonal contrast; and (iii) Northeast (36.6°N–49.4°N, 90.4°W–77.6°W), influenced by synoptic-scale cyclones and frontal systems.



**Figure 1.** Geographic domains used in this study, showing the three 128 × 128 grid regions over the contiguous United States representing distinct hydroclimatic regimes: the Pacific Northwest, Central Plains, and Northeast

A cross-regional evaluation strategy is adopted to test generalization across hydroclimatic regimes. Models are trained using samples from the Central Plains and Pacific Northwest. Validation samples are drawn from the Central Plains and a subset of the Northeast region to guide model selection and early stopping under distribution shift, while

138   the remaining Northeast samples are reserved for independent testing. This setup ensures
139   that the Northeast region is not used for parameter learning during training, while enabling
140   evaluation in an out-of-training-regime context (Rampal et al., 2024; Vandal et al., 2017).
141   Daily precipitation values below 1 mm/day are treated as dry and set to zero, following
142   commonly used wet-day thresholds in hydroclimatological analyses and reporting
143   conventions (Teutschbein & Seibert, 2012; Trenberth et al., 2015). To ensure that learning
144   is driven by meaningful spatial rainfall structure rather than near-empty scenes, days with
145   fewer than 1% wet pixels (<164 wet pixels in a 128 × 128 domain) are excluded. After
146   filtering, 11,025 daily samples are retained for the Central Plains, 11,747 for the Pacific
147   Northwest, and 12,348 for the Northeast.
148        Low-resolution inputs are generated through block averaging, producing 8× and 16×
149   aggregated precipitation fields while conserving storm-total rainfall volume. Block
150   averaging is preferred over interpolation because it preserves physical mass consistency and
151   avoids introducing artificial spatial correlations (Hsu et al., 2024; Kumar et al., 2023; Stengel
152   et al., 2020).  For the 8× configuration, 128 × 128 fields (~9 km) are aggregated to 16 × 16
153   (~72 km) and models are trained to reconstruct the corresponding 128 × 128 target. For the
154   16× configuration, targets are reconstructed from 8 × 8 (~144 km) inputs. This design
155   defines a perfect-model super-resolution framework in which inputs and targets originate
156   from the same dataset, enabling controlled evaluation of spatial refinement independent of
157   predictor mismatch or bias-correction effects.

**2.2 Models**

159   In this study we used a deterministic convolutional U-NET as a baseline and two generative
160   models WGAN and DDPM for downscaling precipitation data. A schematic of these deep-
161   learning architectures is presented in Figure 2 and detailed architectures are included in
162   supplementary information, Figure S1.

**2.2.1 U-NET**

164   A U-NET architecture is widely used in the deep learning based super-resolution
165   downscaling experiments (Abdelmoaty et al., 2025; Papalexiou & Mamalakis, 2025; Wang et
166   al., 2021). The model takes as input a coarse precipitation field $x_{LR}$ of size 16 × 16 for the 8×
167   downscaling task or 8 × 8 for the 16× downscaling task, together with a spatially
168   uncorrelated noise tensor $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ of identical dimensions. The noise input is included
169   solely to maintain architectural compatibility with the stochastic generator used in the
170   WGAN framework, but since the network is trained independently using a mean-squared
171   error (MSE) loss, it converges to the conditional mean of the high-resolution target given the
172   coarse input and therefore exhibits deterministic behavior during inference
173   (Lakshminarayanan et al., 2017; Yan et al., 2019). As demonstrated in prior studies,
174   minimization of squared error causes the network to ignore injected noise and produce

175    nearly identical outputs across different noise realizations (Abdelmoaty et al., 2025;
176    Papalexiou & Mamalakis, 2025). The MSE loss is expressed as:

177
$$\mathcal{L}_{\mathrm{MSE}} = \mathbb{E}[\| \, G_\theta(x_{\mathrm{LR}}, \mathbf{z}) - x_{\mathrm{HR}} \, \|_2^2] \tag{1}$$

178    where $G_\theta$ denotes the U-NET mapping parameterized by $\theta$, and $x_{\mathrm{HR}}$ represents the
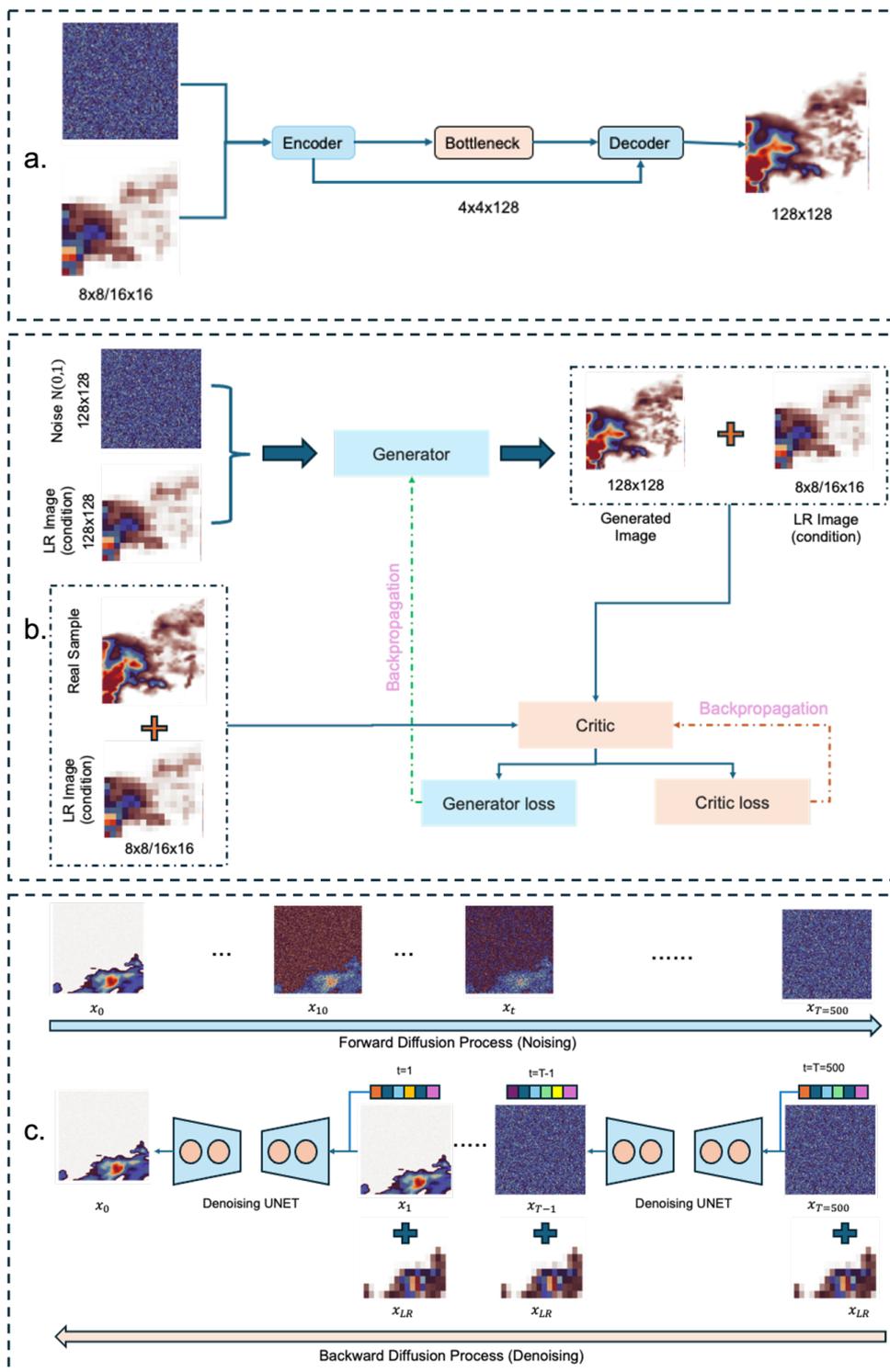179    corresponding high-resolution target field.

180        Architecturally, the encoder consists of two down-sampling stages with 32 and 64
181    filters, each containing pairs of $3 \times 3$ convolutional layers with ReLU activations, followed
182    by $2 \times 2$ max-pooling layers. The bottleneck includes two convolutional layers with 128
183    filters. The decoder employs transposed convolution up-sampling and incorporates skip
184    connections from the encoder to preserve fine-scale spatial organization. To reach the target
185    $128 \times 128$ resolution from the coarse input, three progressive up-sampling stages are
186    applied ($32 \rightarrow 64 \rightarrow 128$), using $2 \times 2$ transposed convolutions with Leaky ReLU
187    activations, followed by refinement convolutions at full resolution. The output layer uses a
188    linear activation to produce precipitation intensity in mm/day. The model is trained end-to-
189    end using the Adam optimizer (learning rate $1 \times 10^{-4}$, batch size 32) with early stopping
190    when validation loss fails to improve for ten consecutive epochs (patience =10), and the best
191    checkpoint based on validation set MSE is retained for evaluation.

192    **2.2.2 Wasserstein GAN (WGAN)**

193    To enable stochastic and spatially realistic precipitation downscaling, we implement a
194    conditional Wasserstein GAN (WGAN), following (Arjovsky et al., 2017; Gulrajani et al., 2017;
195    Papalexiou & Mamalakis, 2025). The generator, which is the same noise-conditional U-NET
196    described above, maps a coarse precipitation field $x_{\mathrm{LR}}$ together with a spatial noise tensor
197    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ to a high-resolution field $\hat{x}_{\mathrm{HR}}$ on the $128 \times 128$ grid. Unlike the MSE-trained U-
198    NET, which converges to the conditional mean, adversarial optimization enables sampling
199    from a distribution of plausible high-resolution rainfall structures conditioned on the same
200    coarse input. The critic $C_\psi$ serves as a conditional discriminator and evaluates the realism of
201    the generated field given the coarse rainfall context. It consists of two pathways: an encoder
202    that reduces the $128 \times 128$ field through strided convolutions, and an embedding of the low-
203    resolution input into a matching spatial feature representation. These feature streams are
204    fused and reduced to a scalar score, enabling the critic to assess global storm morphology
205    while remaining aware of the large-scale meteorological state.

206    Let $x_{\mathrm{HR}} \sim P_r$ denote real high-resolution samples and $\hat{x}_{\mathrm{HR}} = G_\theta(x_{\mathrm{LR}}, \mathbf{z}) \sim P_g$ denote
207    generated samples. The conditional Wasserstein critic loss is:

208    $$\mathcal{L}_{\mathrm{C}} = \mathbb{E}_{\hat{x}_{\mathrm{HR}} \sim P_g}\big[C_\psi(\hat{x}_{\mathrm{HR}}, x_{\mathrm{LR}})\big] - \mathbb{E}_{x_{\mathrm{HR}} \sim P_r}\big[C_\psi(x_{\mathrm{HR}}, x_{\mathrm{LR}})\big] + \lambda\, \mathbb{E}_{\tilde{x} \sim P_{\tilde{x}}}\Big[\big(\| \, \nabla_{\tilde{x}} C_\psi(\tilde{x}, x_{\mathrm{LR}}) \, \|_2 - 1\big)^2\Big] \tag{2}$$

210 **Figure 2:** Schematic of the deep-learning architectures used for precipitation super-
211 resolution downscaling. (a) Deterministic U-NET with an encoder–decoder structure and
212 skip connections producing a single high-resolution output. (b) Conditional Wasserstein
213 GAN (WGAN), where a stochastic generator produces high-resolution fields conditioned on
214 low-resolution input and noise and is trained adversarially against a critic. (c) Conditional
215 denoising diffusion probabilistic model (DDPM), which reconstructs high-resolution
216 precipitation through an iterative, noise-to-signal denoising process conditioned on the
217 coarse input.

218 The generator is trained to maximize the critic score, corresponding to:

$$\mathcal{L}_{\mathrm{G}} = -\mathbb{E}_{\hat{x}_{\mathrm{HR}} \sim P_g}\big[C_\psi(\hat{x}_{\mathrm{HR}}, x_{\mathrm{LR}})\big] \tag{3}$$

219 The model is trained using a gradient penalty weight $\lambda = 10$, updating the critic three times
220 for each generator update. Both networks use the Adam optimizer with learning rate
221 $1 \times 10^{-4}$, $\beta_1 = 0.0$, and $\beta_2 = 0.9$. To characterize stochastic variability, ten WGAN models
222 independently initialized with a random seed are trained for 200 epochs, producing
223 ensembles of generated samples for each coarse precipitation input.

### 2.2.3 Denoising Diffusion Probabilistic Model (DDPM)

226 To exploit recent advances in likelihood-based generative modeling for high-resolution
227 precipitation reconstruction, we implement a conditional DDPM following the framework
228 proposed by Ho et al. (2020). In their framework, the model learns to reverse a fixed forward
229 diffusion process that progressively perturbs high-resolution (HR) precipitation fields with
230 Gaussian noise over T(500) steps. During training, the network predicts the added noise at
231 a randomly sampled timestep, conditioned on the corresponding low-resolution (LR)
232 precipitation field (either $16 \times 16$ or $8 \times 8$). At inference, samples are obtained by iteratively
233 denoising pure Gaussian noise while conditioning on the LR field. The forward process is
234 defined as:

$$q(x_{1:T} \mid x_0) = \prod_{t=1}^{T} q(x_t \mid x_{t-1}) \tag{4}$$

$$q(x_t \mid x_{t-1}) = \mathcal{N}\big(x_t; \sqrt{1 - \beta_t}\, x_{t-1}, \beta_t I\big) \tag{5}$$

237 where $\{\beta_t\}_{t=1}^{\mathrm{T}}$ is the variance schedule. We adopt the cosine noise schedule (Song & Dhariwal,
238 2023) to improve sample smoothness and reduce denoising artifacts. The cumulative signal
239 retention is given by

$$\bar{\alpha}_t = \prod_{s=1}^{t}(1 - \beta_s) \tag{6}$$

$$\tilde{\alpha}_t = \frac{\cos^2\left(\frac{T/t+s}{1+s}\cdot\frac{\pi}{2}\right)}{\cos^2\left(\frac{s}{1+s}\cdot\frac{\pi}{2}\right)} \tag{7}$$

241

242    with a small offset $s = 0.008$ to avoid extreme noise ratios. At training time, the model
243    predicts the noise $\epsilon$ added to $x_0$ at timestep $t$, conditioned on the LR field $x_{\mathrm{LR}}$:

244

$$\mathcal{L} = \mathbb{E}_{x_0,\epsilon,t}[\|\epsilon - \epsilon_\theta(x_t, t, x_{\mathrm{LR}})\|_2^2] \tag{8}$$

245    where $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon$. The network backbone is a conditional U-NET augmented
246    with sinusoidal time embeddings and Feature-wise Linear Modulation (FiLM; (Perez et al.,
247    2018) layers to inject timestep context and LR conditioning into intermediate feature
248    representations. Time is encoded using 256-dimensional sinusoidal embeddings,

250

$$\gamma(\mathrm{t}) = \left[\sin\left(\frac{t}{10000^{2i/256}}\right), \cos\left(\frac{t}{10000^{2i/256}}\right)\right]_{i=0}^{127} \tag{9}$$

249

251    and refined using a multilayer perceptron to produce a contextualized embedding $t_{emb}$.
252    Within each convolutional block, FiLM modulates internal activations according to the
253    timestep and LR conditioning:

254

$$\mathrm{FiLM}(h, t_{emb}) = \gamma(t_{emb}) \odot h + \beta(t_{emb}) \tag{10}$$

255    This allows the model to adapt its representations to different stages of the denoising
256    trajectory. To condition DDPM, the LR field is up-sampled to $128 \times 128$ via bilinear
257    interpolation and concatenated with the noisy HR input, ensuring that the coarse-scale
258    spatial context informs the fine-scale reconstruction. All precipitation fields undergo a
259    $\log(1+x)$ transformation followed by min–max normalization to stabilize training under
260    heavy-tailed rainfall distributions. Training minimizes the noise-prediction loss using the
261    AdamW optimizer (learning rate $1 \times 10^{-4}$, weight decay $1 \times 10^{-4}$). During inference, HR
262    precipitation samples are generated by iteratively applying the reverse diffusion update,

263

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta(x_t, t, c)\right) + \sigma_t \mathbf{z}, \qquad \mathbf{z} \sim \mathcal{N}(0, \mathrm{I}) \tag{11}$$

264    for all t from T to 1, injecting noise at each step except the final one. This additional noise
265    makes the denoising process more stochastic and for same LR input, keeping all other
266    parameters same, we get different downscaled output.

267     All three models (U-NET, WGAN, and DDPM) were trained using 10 different random
268     seeds, yielding 10 independently initialized and trained realizations of each architecture. For
269     a given low-resolution input, one prediction was generated from each trained instance,
270     forming a 10-member ensemble. This ensemble primarily reflects epistemic uncertainty
271     arising from random weight initialization and stochastic optimization during training. To
272     ensure a fair comparison across downscaling factors, the core network architectures were
273     kept fixed across experiments. For the deterministic U-NET and the WGAN, the same
274     generator and critic architectures were used for both 8× and 16× configurations. The
275     increased difficulty of the 16× case was introduced solely by providing coarser inputs,
276     obtained by bilinearly interpolating $16 \times 16$ fields to $8 \times 8$ before being passed to the
277     networks, while keeping the target resolution at $128 \times 128$.

278     For the DDPM, the same denoising U-NET architecture was used for both scaling
279     configurations. In the 8× setup, conditioning fields were bilinearly interpolated from $16 \times 16$
280     to $128 \times 128$ through three successive 2× interpolations, consistent with the diffusion
281     model's multi-scale refinement process. In the more challenging 16× setup, conditioning
282     fields originated from $8 \times 8$ inputs and were similarly interpolated to $128 \times 128$ through
283     repeated (four-times) 2× bilinear interpolation steps before being concatenated with the
284     noisy target field at each diffusion timestep. This design ensured that differences in
285     performance across scaling factors reflect the increased information gap in the input, rather
286     than changes in model capacity or architecture. By holding network architectures fixed and
287     varying only the effective resolution of the conditioning input, this experimental setup
288     enables a controlled and equitable comparison of deterministic and generative models
289     across downscaling factors. This implementation allows direct comparison with U-NET and
290     WGAN models under identical input configurations and downscaling ratios (8× and 16×)
291     isolating differences attributable to generative formulation rather than architectural
292     capacity or preprocessing design.

## 3. Performance Evaluation

294     To evaluate the statistical fidelity, spatial realism, and hydrologically relevant characteristics
295     of the downscaled precipitation fields produced by the U-NET, WGAN, and DDPM models, we
296     employ a set of complementary evaluation measures. These measures examine
297     distributional consistency, representation of extremes and dry occurrence, rainfall mass
298     conservation across spatial scales, spatial organization and storm morphology, and binary
299     precipitation detection skill. All analyses are conducted on a test set of 12,616 daily
300     precipitation fields from an unseen Northeast hydroclimatic region. For each modeling
301     framework, results are computed across a 10-member ensemble consisting of independently
302     trained models initialized with different random seeds.

### 303 3.1 Distributional Consistency

### 304 3.1.1 Quantile–Quantile (Q–Q) Analysis

305 Quantile behavior is assessed by comparing the p-quantile of model predictions ($Q_p^{\text{mod}}$) with
306 the corresponding observational quantile ($Q_p^{\text{obs}}$).

$$Q_p^{\text{mod}} = F_{\text{mod}}^{-1}(p) \tag{12}$$

$$Q_p^{\text{obs}} = F_{\text{obs}}^{-1}(p) \tag{13}$$

308 where $p \in [0,1]$ and $F^{-1}$ denotes the empirical inverse CDF. When plotted together,
309 alignment with the 1:1 line indicates agreement in distributional shape, while deviations at
310 high $p$ quantify errors in extreme precipitation intensity representation.

### 312 3.1.2 Exceedance Probability

313 Extreme rainfall representation is examined using the complementary CDF and evaluated
314 across intensity thresholds $x$. This emphasizes tail performance (e.g., > 10 mm/day), which
315 is critical for hydrologic risk estimation. The probability that exceeds the threshold ($x$) is
316 expressed as:

$$\mathbb{P}(X^{\text{mod}} > x) = 1 - F_{\text{mod}}(x) \tag{14}$$

317 where $F(x)$ is cumulative distribution function of rainfall intensity.

### 319 3.1.3 Probability of Zero Precipitation ($P_0$)

320 Dry–wet occurrence skill is evaluated using the probability of zero (or near-zero)
321 precipitation, defined as:

$$P_0^{\text{mod}} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{1}(X_i^{\text{mod}} \leq x_{\text{th}}) \tag{15}$$

323 where $N$ is the total number of evaluated grid cells (pixels), $X_i$ is pixel intensity, $x_{\text{th}} =$
324 1 mm/day. This metric quantifies how often a model predicts dry conditions and is essential
325 for diagnosing dry bias.

### 326 3.1.4 Higher-Order Statistical Moments

327 To evaluate the distributional structure of rainfall intensities predicted from the models, we
328 compute the first four statistical moments for all wet pixels ($X_i > 1$ mm/day): the mean,
329 variance, skewness, and kurtosis. These moments respectively describe the central
330 tendency, spread, asymmetry, and tail-heaviness of the precipitation distribution properties.
331 Model fidelity in capturing these distributional characteristics is assessed through mean bias

332 and RMSE**,** which summarize systematic and random departures from the observed
333 moments.

### 3.2 Mass Conservation

### 3.2.1 Cumulative Mean Rainfall Depth

336 We compute mean precipitation at multiple aggregation scales to assess whether rainfall
337 volume is preserved during spatial refinement. For block size d, the observed and predicted
338 mean depths are:

$$P_d^{\text{obs}} = \frac{1}{M} \sum_{j=1}^{M} X_{d,j}^{\text{obs}} \ \text{and} \ P_d^{\text{mod}} = \frac{1}{M} \sum_{j=1}^{M} X_{d,j}^{\text{mod}} \tag{16}$$

339 where $M$ is the number of non-overlapping blocks and denotes the mean precipitation within
340 block $j$ at scale $d$. Agreement across spatial scales indicates physically consistent
341 redistribution of rainfall mass in downscaling precipitation.

### 3.3 Spatial Structure and Storm Morphology

### 3.3.1 Lagged Spatial Autocorrelation

345 Spatial dependence in precipitation fields is evaluated using lagged spatial autocorrelation,
346 which quantifies the similarity of precipitation values separated by a given spatial lag
347 (Papalexiou et al., 2021). For a specified lag distance $d$, the lagged autocorrelation is
348 computed as the Pearson correlation coefficient between paired precipitation values
349 sampled at locations separated by $d$. Specifically, the statistic is defined as

$$r(d) = \frac{\sum_{i=1}^{N} (X(s_i) - \bar{X}) (X(s_i + d) - \bar{X})}{\sqrt{\sum_{i=1}^{N} (X(s_i) - \bar{X})^2 \ \sum_{i=1}^{N} (X(s_i + d) - \bar{X})^2}}, \tag{17}$$

352 where $X$ denotes either $X^{\text{obs}}$ or $X^{\text{mod}}$, $\bar{X}$ is the sample mean of the precipitation field, $s_i$
353 denotes a spatial location, $d$ is the lag vector, and $N$ is the number of valid grid-point pairs.
354 The statistic is computed for a range of lag distances and averaged over all valid pairs at each
355 lag.

### 3.3.2 Fraction Skill Score (FSS)

357 Spatial coherence and storm organization are assessed using the Fraction Skill Score (FSS),
358 a window-based metric that compares the fractional rainfall coverage in predictions and
359 observations (Gilleland et al., 2009; Roberts & Lean, 2008). For a window of size w, the score
360 is defined as

365
$$\text{FSS}(w) = 1 - \frac{\mathbb{E}[(P_{\text{mod}} - P_{\text{obs}})^2]}{\mathbb{E}[P_{\text{mod}}^2 + P_{\text{obs}}^2] + \varepsilon} \qquad (18)$$

361 where $P_{\text{mod}}$ and $P_{\text{obs}}$ denote the fractional rainfall coverage within windows of size
362 $w$ computed from the model prediction and the reference observation, respectively. FSS
363 ranges from 0 (no skill) to 1 (perfect spatial agreement), making it particularly effective for
364 diagnosing displacement errors, spatial smoothing, and the realism of storm geometry.

366 **3.3.3 Radial Power Spectrum**

367 Scale-dependent spatial variability in downscaled predictions is evaluated using the radially
368 averaged Fourier power spectrum, which characterizes how energy is distributed across
369 spatial wavenumbers (Bednarz & Cherukuri, 2023; Harrison et al., 2025; Skamarock, 2004).
370 The 2-D discrete Fourier transform of the rainfall field,

372
$$F(k_x, k_y) = \mathcal{F}\{X(x, y)\} \qquad (19)$$

371 where $\mathcal{F}\{\cdot\}$ denotes the 2-D discrete Fourier transform, yields the spectral power,

378
$$P(k_x, k_y) = \mid F(k_x, k_y) \mid^2 \qquad (20)$$

373 which is then azimuthally averaged to obtain the one-dimensional spectrum $P(k)$.
374 Agreement between predicted and observed spectra indicates that the model accurately
375 captures multiscale storm structure from large synoptic gradients to mesoscale organization
376 and fine-scale convective patterns. Whereas deviations reveal scale-specific biases such as
377 excessive smoothing, noise amplification, or loss of small-scale variability.

379 **3.3.4 Structural Similarity Index (SSIM)**

380 Structural realism is evaluated using the Structural Similarity Index (SSIM), which jointly
381 measures agreement in luminance, contrast, and local spatial structure between the
382 predicted and observed rainfall fields (Meghani et al., 2023; Singh & Goyal, 2023; Zhou Wang
383 et al., 2004). Unlike pixel-wise metrics, SSIM emphasizes coherent patterns such as storm
384 cores, gradients, and spatial organization. SSIM distributions are visualized using violin/box
385 plots, enabling comparison of structural fidelity across models and highlighting variability
386 in performance across the ensemble.

387
$$\text{SSIM}(X^{\text{mod}}, X^{\text{obs}}) \in [-1, 1] \qquad (21)$$

388 **3.3.5 ROC Curve and AUC**

389 Binary precipitation-occurrence skill is assessed using Receiver Operating Characteristic
390 (ROC) analysis based on a random threshold (here, wet/dry threshold of $x = 1\,\text{mm/day}$)
391 (Harris et al., 2022). The Area Under the Curve (AUC) satisfies Eq. (25).

392
$$0.5 \leq \text{AUC} \leq 1 \qquad (22)$$

393 where AUC = 1 denotes perfect discrimination between wet and dry events, and AUC =
394 0.5 indicates no discriminative ability (random chance). This metric evaluates ability of
395 models to correctly identify rainfall occurrence independently of intensity, making it
396 particularly useful for diagnosing dry–wet classification bias.

## 397 4. Computational Requirements

398 All model trainings were conducted on a single NVIDIA A100 GPU (80 GB memory). Each
399 architecture was trained for 200 epochs using 19,200 training samples and evaluated on
400 12,616 test samples. Substantial differences in computational demand were observed across
401 the three model classes. The U-NET baseline completed training in approximately
402 16 minutes, reflecting the efficiency of direct supervised optimization with a single forward–
403 backward pass per batch. In contrast, the WGAN required 1 hour and 12 minutes to train,
404 driven primarily by the need to update the critic multiple times per generator update and to
405 compute the gradient-penalty term that enforces the 1-Lipschitz constraint. The DDPM
406 exhibited the highest training cost, requiring 2 hours and 29 minutes, since each
407 optimization step involves predicting injected noise over a sequence of 500 diffusion
408 timesteps.

409 **Table1:** Comparison of training and inference times for U-NET, WGAN, and DDPM models
410 on a single NVIDIA A100 (80 GB) GPU

| Resource | Model | Training | Inference | Per-Sample Time | Remarks |
|---|---|---|---|---|---|
| | U-NET | ~16 min | ~1.33 sec | ~0.0001sec | Fastest, single forward pass |
| | WGAN (Generator) | ~1h 12min | ~1.32 sec | ~0.0001sec | Multiple critic steps |
| **A100-80GB GPU (1)** | DDPM (T=500) | ~2h 29min | ~1h 42 min | ~0.5 sec | 500 denoising steps per image |
| | DDPM (T=100) | ~2h 23min | ~20 min 53 sec | ~0.1 sec | 100 denoising steps per image |
| | DDPM (T=50) | ~2h 21min | ~10 min 30 sec | ~0.05 sec | 50 denoising steps per image |

411

412       Inference performance showed an even stronger divergence. Because the U-NET and
413      WGAN generators share identical architectures and an equal number of trainable
414      parameters (differing only in training objectives), their inference times were nearly
415      identical, requiring 1.33 seconds and 1.32 seconds, respectively, to process the full test set.
416      In contrast, the DDPM required 1 hour and 42 minutes for inference, since high-resolution
417      precipitation fields are generated through an iterative reverse-diffusion sampling procedure
418      that sequentially refines noise over 500 denoising steps. These differences underscore the
419      computational trade-offs between deterministic convolutional downscaling and likelihood-
420      based generative modeling. While WGAN offers stochastic outputs with moderate added
421      computational burden, DDPM provides calibrated ensemble diversity at substantially higher
422      computational cost, a factor that may strongly influence operational deployment, ensemble
423      forecasting, and climate-model downscaling at scale.

424      **5. Results and Discussion**

425      **5.1 Training behavior**

426      The three downscaling architectures exhibit distinct optimization characteristics that
427      directly reflect their learning objectives and model structures. The deterministic U-NET
428      converges most rapidly under both the 8× and 16× super-resolution configurations, with
429      training and validation losses decreasing sharply and stabilizing within the first few epochs
430      (Figure S2a). The near-perfect overlap between training and validation curves indicates
431      negligible overfitting and strong generalization, consistent with optimization under a mean-
432      squared-error objective that drives the network toward a conditional mean solution. Final
433      loss values are systematically higher for the 16× configuration, reflecting the larger
434      information gap between the coarse 8×8 inputs and the high-resolution targets.

435      The WGAN displays the characteristic oscillatory dynamics associated with
436      adversarial training. For the 8× case, the critic loss stabilizes within a narrow negative range
437      while the generator loss fluctuates around a stationary mean, indicating a sustained
438      adversarial equilibrium and effective enforcement of the 1-Lipschitz constraint through
439      gradient penalty regularization (Figure S2b). Under the more challenging 16× configuration,
440      critic and generator losses exhibit increased variability, reflecting the greater difficulty of
441      discriminating realistic structure when the conditioning input contains extremely limited
442      spatial information. Importantly, the loss trajectories remain bounded across all seeds, with
443      no evidence of divergence or mode collapse, highlighting the stabilizing influence of
444      conditional adversarial training even under severe super-resolution.

445      The DDPM shows the most monotonic and stable optimization behavior among the
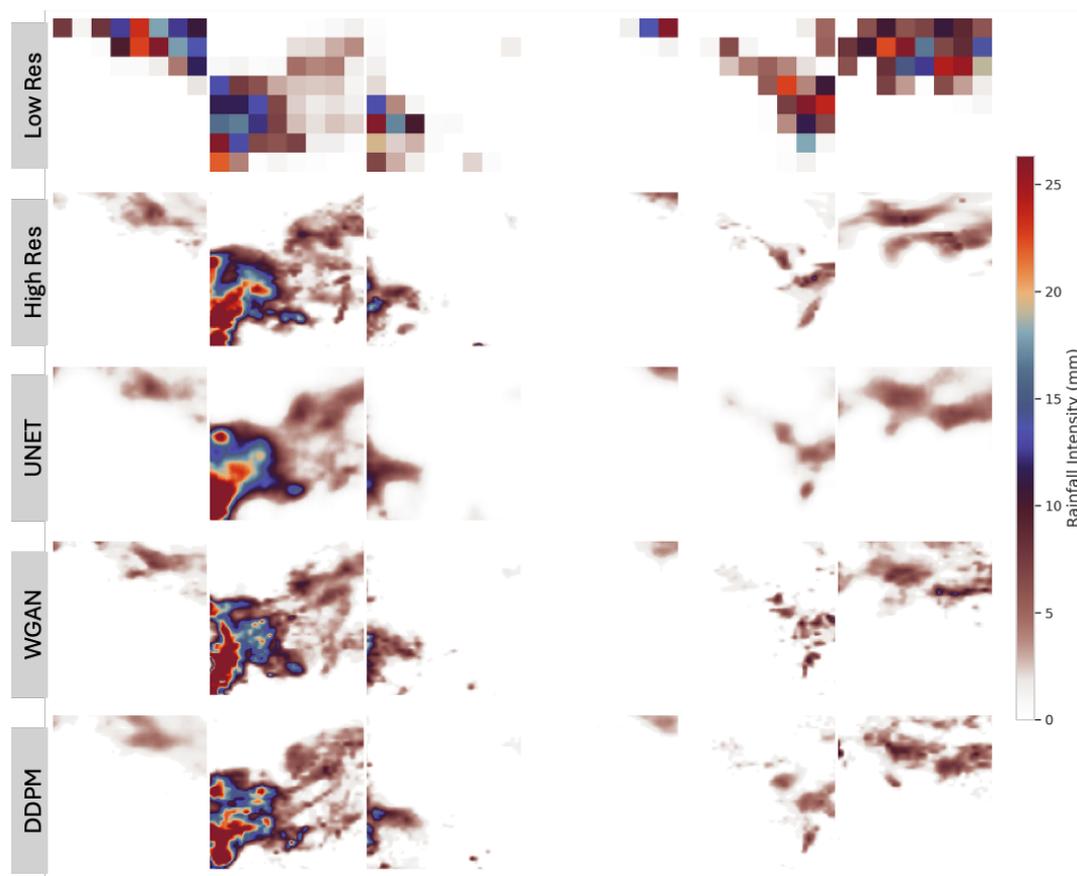446      three models. The noise-prediction loss decreases smoothly for both super-resolution

447 factors, and training and validation curves remain nearly indistinguishable throughout
448 training (Figure S2c), indicating strong generalization. Convergence is slower than for U-NET
449 and WGAN, reflecting the iterative denoising process and timestep conditioning intrinsic to
450 diffusion models. In the 16× configuration, a modest widening between training and
451 validation losses emerges after prolonged training, suggesting that the model approaches
452 representational limits imposed by the extremely coarse input. Nevertheless, training
453 remains stable across all seeds, underscoring the robustness of likelihood-based diffusion
454 models for learning structured precipitation fields under severe information constraints.

455 **5.2 Model Performance Evaluation**

456 **5.2.1 Visual Reconstruction of Precipitation Fields**

457 Visual inspection of reconstructed precipitation fields reveals clear performance differences
458 between the 8× and 16× downscaling tasks (Figure 3; Figure S3). Under the 8× configuration,
459 where coarse inputs retain recognizable storm-scale organization, all three models
460 successfully recover the dominant spatial structure of precipitation events. The U-NET
461 produces smooth, spatially coherent fields but systematically attenuates sharp gradients and
462 localized convective maxima. In contrast, WGAN outputs exhibit sharper boundaries and
463 more textured rainfall patterns, consistent with adversarial training encouraging the
464 reconstruction of high-frequency spatial variability. DDPM reconstructions closely resemble
465 those of WGAN in terms of storm morphology and spatial extent, but with slightly smoother
466 textures attributable to the progressive denoising mechanism. Performance degradation
467 becomes evident for all models under the 16× configuration, reflecting the extremely limited
468 information content of the 8×8 conditioning inputs. In this case, U-NET predictions
469 increasingly collapse toward overly smoothed and diffuse rainfall fields, frequently failing to
470 recover localized storm features. Both generative models retain substantially better spatial
471 organization than U-NET, although fine-scale detail is reduced relative to the 8× case. WGAN
472 continues to generate sharper, more intermittent structures, whereas DDPM maintains
473 coherent storm morphology with comparatively smoother gradients. Overall, these visual
474 results highlight the growing limitations of deterministic regression under extreme
475 downscaling and demonstrate the advantage of generative models in reconstructing
476 physically plausible precipitation patterns when conditioning information becomes severely
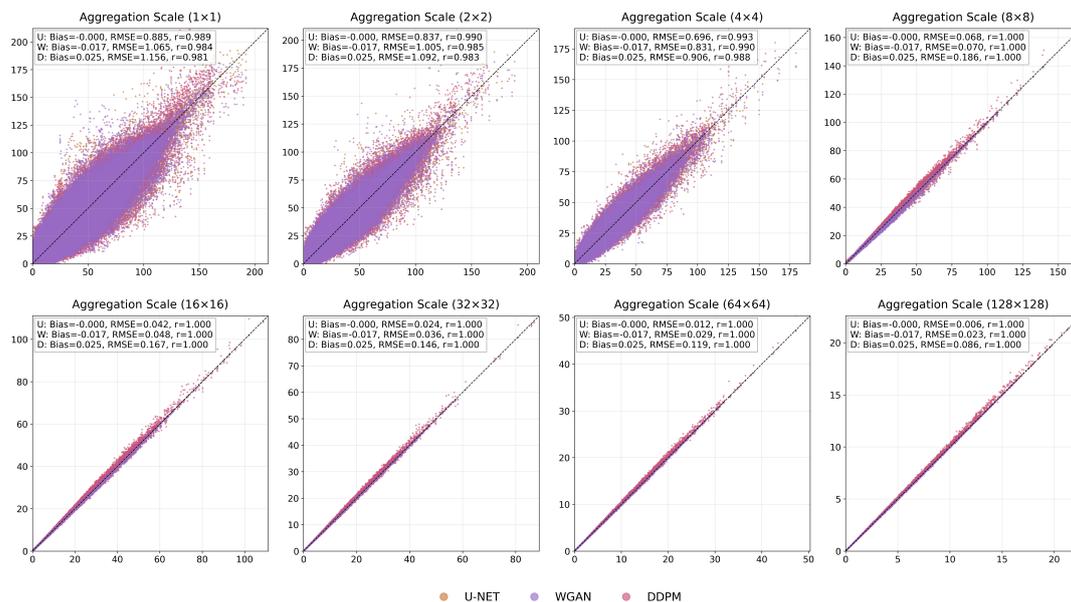477 constrained.

Figure 3. Reconstruction of high-resolution precipitation fields for the 16× downscaling task (8×8 → 128×128). Rows show the low-resolution input, high-resolution ground truth, and predictions from U-NET, WGAN, and DDPM respectively.

**5.2.2 Scale-Dependent Rainfall Depth Consistency (Mass Conservation)**

The scatter plots of block-averaged precipitation depth demonstrate strong scale-dependent consistency between predicted and observed rainfall across both super-resolution configurations (Figure 4; Figures S4–S5). Despite the absence of any explicit mass-conservation constraint during training, all three models preserve storm-integrated rainfall depth remarkably well. At the native grid scale (1×1), discrepancies are largest, with RMSE dominated by pixel-scale intensity errors and spatial displacement. In the 8× case, U-NET exhibits the lowest RMSE (0.885 mm/day) and negligible bias, while WGAN and DDPM show comparatively larger errors (1.065 mm/day and 1.156 mm/day, respectively), reflecting increased small-scale variability. Nevertheless, even at this finest scale, correlations remain

492    high, indicating that the total precipitation volume is broadly preserved despite local
493    mismatches.



494
495    Figure 4. Mass (rainfall depth) consistency across spatial aggregation scales for 8×
496    downscaling. Scatter plots compare block-averaged predicted and observed precipitation
497    depths at aggregation scales of 1×1, 2×2, 4×4, 8×8, 16×16, 32×32, 64×64, and 128×128. Each
498    point represents a spatially aggregated precipitation value computed over the
499    corresponding block size. Results are shown for single seed initialized U-NET (orange),
500    WGAN (purple), and DDPM (pink). The dashed black line denotes perfect 1:1 agreement.
501    Panel annotations report model-specific mean bias, root-mean-square error (RMSE), and
502    Pearson correlation coefficient (*r*) at each scale.

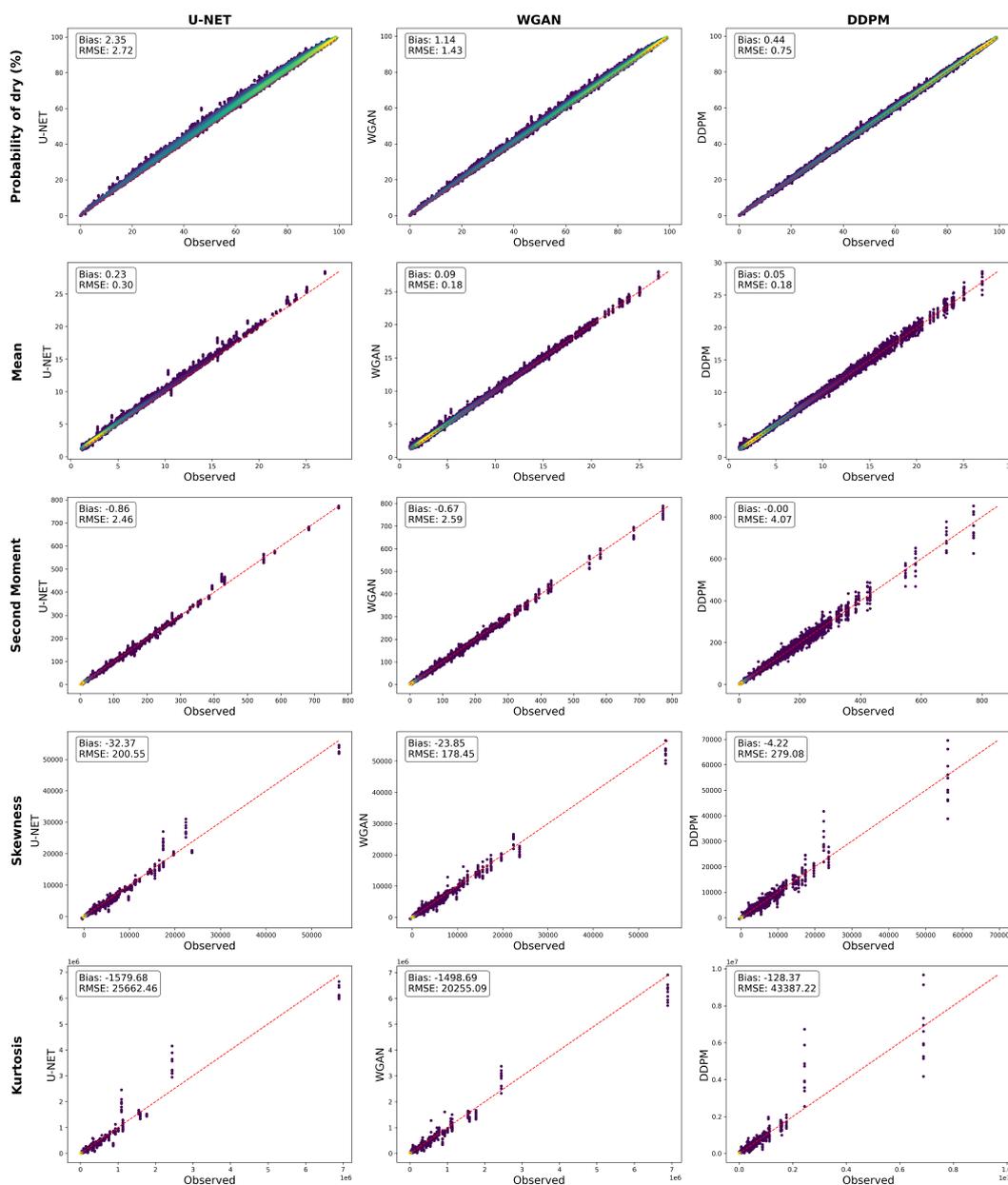503    As aggregation scale increases, scatter collapses rapidly toward the 1:1 line and RMSE
504    decreases sharply for all models. At coarser scales (≥16×16 blocks), biases approach zero
505    and correlations approach unity, demonstrating that errors associated with fine-scale
506    structure largely cancel under spatial averaging. These results indicate that while pixel-level
507    fidelity remains challenging, particularly for generative models, the downscaled fields retain
508    robust depth consistency at hydrologically relevant scales, supporting their applicability for
509    basin-scale water-balance and impact analyses.

### 5.2.3 Statistical Distribution and Storm Morphology

511    To ensure consistent wet–dry classification, predicted precipitation values below 1 mm/day
512    were treated as dry and set to zero. The statistical comparison plots show that U-NET, WGAN

18

513 and DDPM reproduce key distributional properties of the precipitation fields under 8× and
514 16× downscaling (Figure 5 and Supplementary Information, Figure S6 respectively).
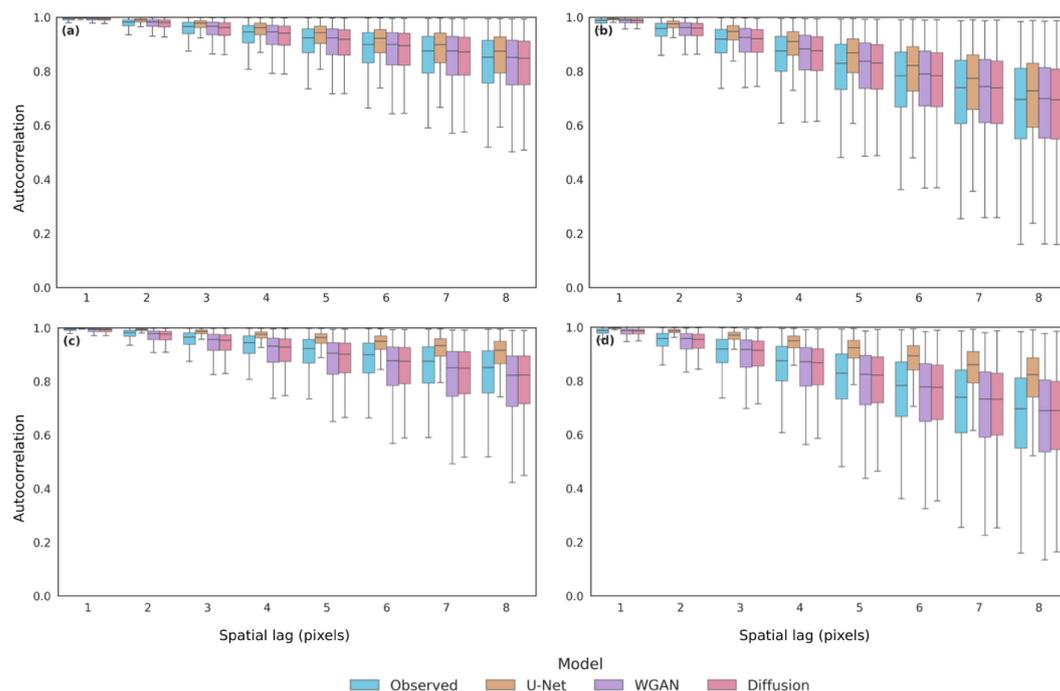


515

516 Figure 5. Comparison of observed versus predicted precipitation statistics for 8×
517 downscaling across U-NET, WGAN, and DDPM. Each panel shows the relationship between

518    observed statistics (x-axis) and model predictions (y-axis) for dry-pixel probability ($P_0$),
519    mean, second moment, skewness, and kurtosis. Bias and RMSE are reported for each metric.

520    All models align closely with the 1:1 reference line for dry-pixel probability and mean
521    rainfall, indicating accurate reconstruction of occurrence frequency and overall storm
522    magnitude. U-NET and WGAN exhibit particularly small deviations for these two metrics,
523    while DDPM achieves the lowest RMSE for $P_0$ in the 16× case. We observe more pronounced
524    differences for higher-order statistics. The second moment is captured reasonably well by
525    all models, although DDPM shows a larger spread at higher values, especially in the 8×
526    configuration. Skewness and kurtosis exhibit the largest errors across models, reflecting the
527    difficulty of recovering tail behavior and extreme-event structure. WGAN generally shows
528    smaller bias in skewness, whereas U-NET and DDPM display higher variance. Despite these
529    challenges at higher-order moments, all models maintain consistent performance trends
530    across both downscaling factors, with U-NET providing stable low-order statistics, WGAN
531    delivering sharper distributional structure, and DDPM capturing overall variability while
532    showing higher spread for extremes.

533        The spatial autocorrelation analysis illustrates how effectively each model preserves
534    fine-scale structure as pixel lag increases. Under the 8× configuration, all three architectures
535    reproduce the observed correlation decay well at short lags ($1 - 3$ pixels), indicating that
536    they can recover local spatial gradients when the coarse input still contains partial storm
537    structure (Figure 6). U-NET yields the highest correlations at the smallest lags, consistent
538    with its tendency to generate smooth and spatially coherent fields. WGAN produces slightly
539    lower short-lag correlations but aligns more closely with the observed decay pattern at
540    intermediate lags ($4 - 6$ pixels), reflecting its capacity to introduce sharper gradients and
541    finer texture. DDPM follows a similar trajectory to WGAN, with correlations slightly below
542    WGAN at short lags but closer to observations across mid-range scales, suggesting a balanced
543    reconstruction of local variability and structural detail. Ensemble spread is wider for WGAN
544    and DDPM than for U-NET, reflecting the stochastic nature of generative sampling and the
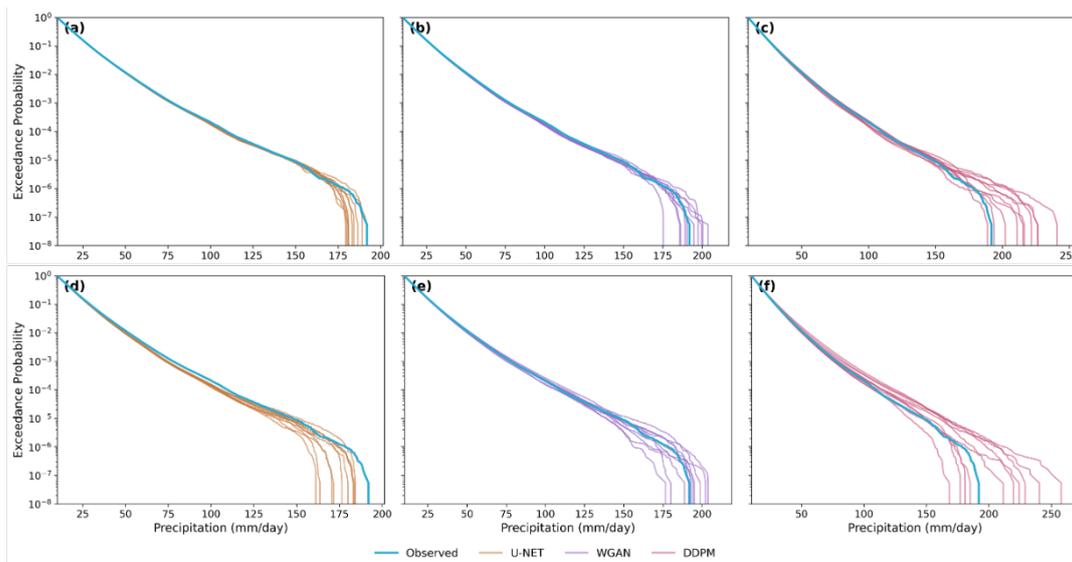545    diversity of high-resolution realizations.

546
547 Figure 6. Spatial autocorrelation of observed and predicted precipitation fields for horizontal
548 and vertical pixel lags (1–8 pixels) under 8× (top row) and 16× (bottom row) downscaling.
549 Boxplots summarize ensemble variability across 10 seeds for each model (U-NET, WGAN,
550 DDPM). Observed correlations are shown in green.

551        Under the more demanding 16× configuration, all models show weakened spatial
552 structure, with reduced correlations and greater spread across lags (Figure 6). U-NET
553 continues to produce the highest lagged autocorrelations and fails to reproduce the observed
554 rate of decorrelation (decay in autocorrelation), indicating over smoothing when limited
555 information is available in the $8 \times 8$ inputs. WGAN and DDPM better capture the observed
556 decline in autocorrelation with increasing lag, although WGAN exhibit slightly broader
557 ensemble variability due to uncertainty introduced at this extreme downscaling ratio in
558 vertical direction. WGAN maintains the closest agreement with observed correlations at
559 intermediate lags, preserving sharper spatial transitions, while DDPM yields slightly lower
560 correlations but relatively consistent behavior across lags.  Overall, for 16× downscaling,
561 deterministic models over-smooth the fields, while generative models capture the observed
562 spatial dependence more accurately but with higher ensemble dispersion due to increased
563 uncertainty.

564     **5.2.4 Extreme Precipitation and Tail Behavior**

565     The exceedance probability analysis evaluates the ability of each model to reproduce the
566     upper tail of the precipitation distribution (Figure 7). Under the 8× configuration, all models
567     broadly follow the observed exceedance behavior at moderate intensities, with clear
568     divergence emerging at the highest values. U-NET systematically underestimates extreme
569     precipitation, consistent with regression-induced smoothing. WGAN exhibits the closest
570     agreement with the observed upper tail, maintaining stronger alignment at high intensities,
571     while DDPM performs comparably at low to intermediate values but displays greater
572     ensemble spread at the most extreme quantiles, reflecting increased sampling variability for
573     rare events.



574
575     Figure 7. Exceedance probability (1–CDF) curves of daily precipitation for observations and
576     model predictions at 8× (top row) and 16× (bottom row) spatial downscaling. Panels show
577     results for (a, d) U-NET, (b, e) WGAN, and (c, f) DDPM. The observed precipitation is shown
578     by the thick blue curve, while thin colored curves represent 10 ensemble members for each
579     generative model. All curves are computed using the same test samples and plotted on a
580     logarithmic exceedance scale, highlighting differences in the representation of extreme
581     precipitation tails across downscaling approaches and spatial scales.

582     Under the more challenging 16× configuration, deviations from observations become more
583     pronounced across all models due to the severely limited information content of the coarse
584     inputs. U-NET further suppresses high-intensity events, whereas WGAN retains the closest
585     correspondence to the observed tail despite increased ensemble dispersion. DDPM captures
586     the overall tail shape but exhibits the largest spread among realizations, indicating increased
587     uncertainty in reconstructing extremes at this downscaling ratio. Overall, generative models,

588 particularly WGAN, better preserve the heavy-tailed nature of precipitation, while
589 deterministic regression systematically attenuates extremes.

590 The corresponding Q–Q diagnostics (Figures S7–S8) support these findings, showing
591 consistent underestimation of upper quantiles by U-NET, closer alignment by WGAN at high
592 intensities, and broader dispersion for DDPM in the upper tail. Together, these results
593 indicate that differences among models are dominated by extreme-value behavior:
594 generative approaches better preserve tail behavior but exhibit increased ensemble
595 variability, particularly under severe super-resolution, whereas deterministic methods favor
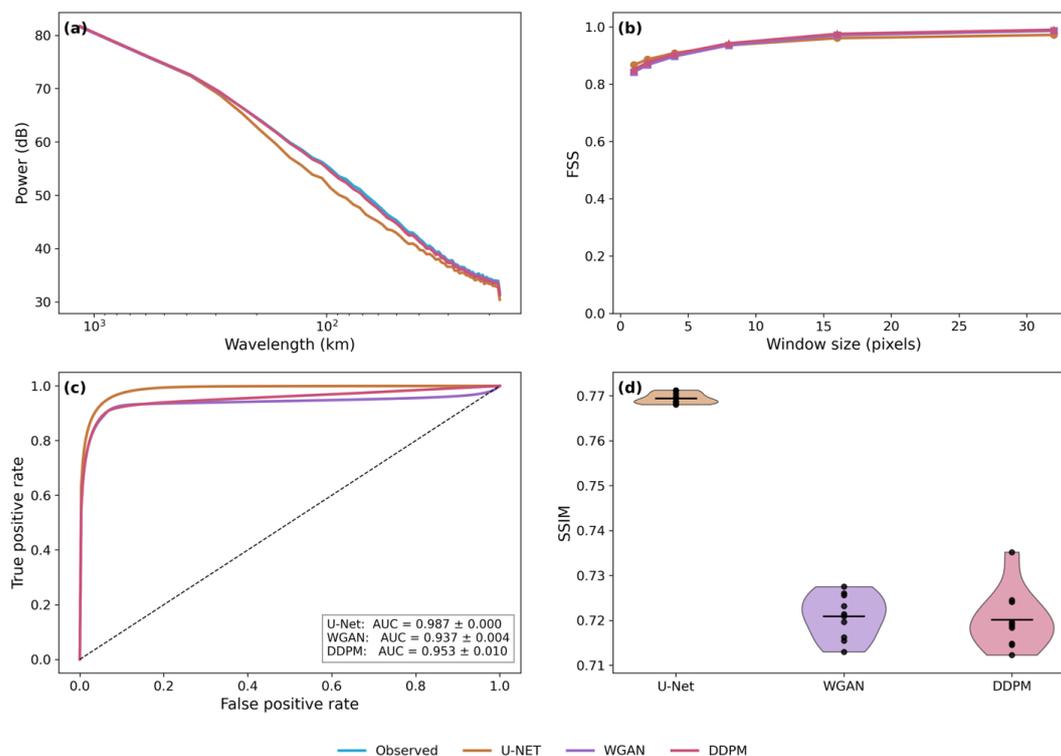596 stability at the expense of extreme-value fidelity.

### 5.2.5 Composite Diagnostics

598 The composite diagnostics provide a complementary evaluation of scale-dependent
599 structure, spatial agreement, event detection, and perceptual similarity using ensemble
600 statistics derived from ten independently trained realizations for each model (Figure 8). The
601 radial power spectra indicate that all three models accurately reproduce the observed large-
602 scale energy content, with close agreement at wavelengths exceeding several hundred
603 kilometers, demonstrating robust preservation of synoptic-scale precipitation organization.
604 At smaller spatial scales, however, U-NET exhibits a pronounced loss of spectral power,
605 reflecting excessive smoothing and suppression of fine-scale variability, whereas WGAN and
606 DDPM retain substantially more energy and remain closer to the observed spectrum,
607 indicating improved representation of small-scale spatial intermittency and storm texture.

608 These scale-dependent differences are consistent with the Fractions Skill Score
609 (Figure 8b), which increases monotonically with window size for all models, reflecting
610 reduced sensitivity to displacement errors at larger spatial scales. U-NET shows marginally
611 higher skill at the smallest windows due to its smoother fields, while WGAN and DDPM
612 converge rapidly and achieve comparable skill at moderate and large window sizes,
613 indicating that enhanced small-scale variability does not compromise spatial agreement at
614 physically meaningful scales. The ROC curves demonstrate strong precipitation occurrence
615 discrimination for all models, with high AUC values indicating reliable separation between
616 wet and dry pixels; U-NET attains the highest AUC, followed by DDPM and WGAN, consistent
617 with the more conservative nature of deterministic predictions versus the sharper, more
618 variable fields produced by generative models (Figure 8c). Finally, the SSIM distributions
619 summarize structural similarity across inference ensembles, showing that U-NET yields the
620 highest median SSIM with minimal spread, while WGAN and DDPM exhibit lower medians
621 and broader distributions, reflecting increased structural diversity and stochasticity. Overall,
622 the diagnostics indicate that while all models perform comparably in capturing large-scale
623 precipitation characteristics, the generative approaches better reproduce fine-scale spatial
624 variability. At the 16× downscaling factor, DDPM achieves a more physically coherent

625 balance between spatial organization and ensemble diversity, whereas WGAN emphasizes
626 sharper extremes with comparatively higher structural variability.



627
628 Figure 8. Multi-metric evaluation of 16× precipitation downscaling performance across 10
629 independently trained models for each architecture. Panels show (a) ensemble-mean radial
630 power spectra, (b) mean Fractions Skill Score (FSS) as a function of spatial window size, (c)
631 mean ROC curves for precipitation occurrence with AUC reported as mean ± standard
632 deviation across models, and (d) distributions of mean SSIM values summarizing structural
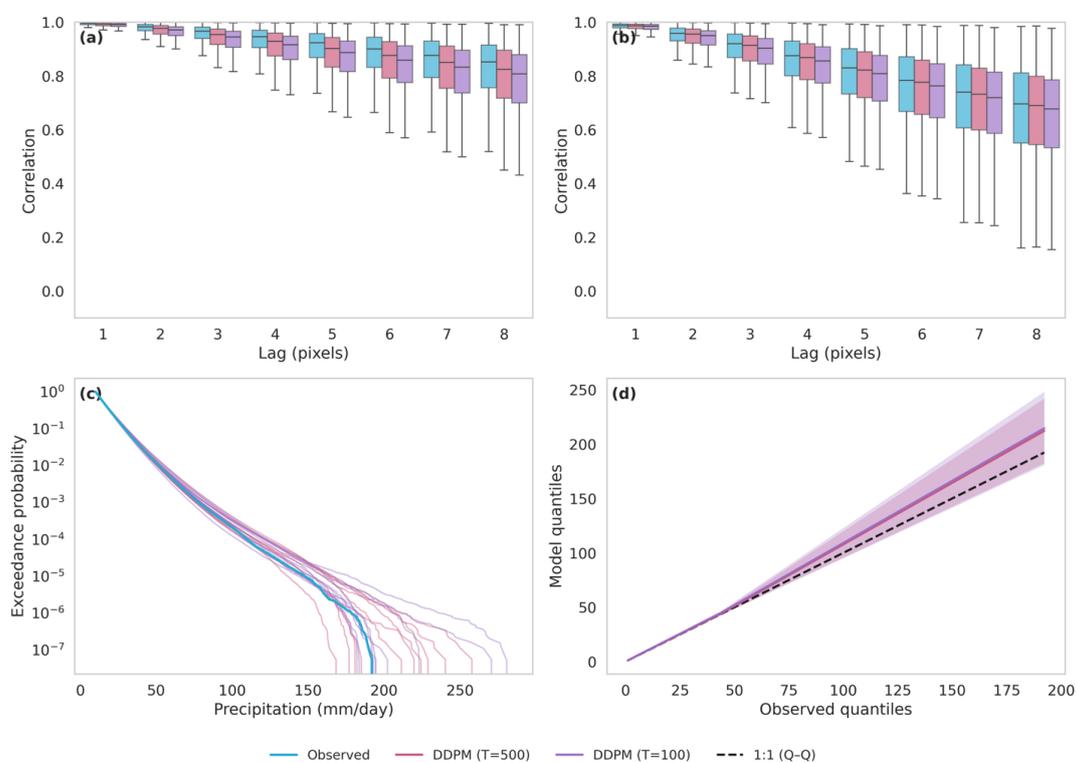633 similarity to observations.

## 5.3 Sensitivity of DDPM to Diffusion Length

635 We also investigated the sensitivity of DDPM performance to different diffusion lengths ($T =$
636 $500, 100, 50$) by comparing key spatial and distributional characteristics. A primary
637 motivation for this analysis was the substantial difference in inference cost between these
638 models. While U-NET and WGAN generate high-resolution precipitation fields in a single
639 forward pass, DDPM relies on iterative sampling, with inference time scaling approximately
640 linearly with the number of diffusion steps. To assess whether shorter diffusion schedules
641 could offer a practical compromise between computational efficiency and predictive skill, we

642    compared DDPM ensembles trained with T = 500 and T = 100 across 10 independent seeds
643    (Figure 9).

644        The lagged spatial correlation analysis shows that both diffusion configurations
645    reproduce the physically consistent decay of spatial dependence with increasing lag in both
646    the horizontal and vertical directions. However, the T = 500 configuration maintains
647    systematically higher correlations and reduced inter-seed variability at intermediate and
648    larger lags, indicating stronger spatial coherence. In contrast, the T = 100 configuration
649    exhibits modestly lower correlations and a broader spread, reflecting increased variability
650    when fewer diffusion steps are used. Despite these differences, both configurations preserve
651    the dominant spatial dependence structure of the observed precipitation fields, indicating
652    that the overall spatial realism of DDPM is retained even when inference cost is substantially
653    reduced.



654
655    Figure 9. Comparison of DDPM downscaling performance (16×) across 10 random seeds for
656    two diffusion lengths, T = 500 and T = 100. Panels (a–d) summarize complementary
657    aspects of spatial structure and extremes: (a) and (b) show lagged spatial correlation (lags
658    1–8 pixels) in the horizontal and vertical directions, respectively, with boxplots summarizing
659    inter-seed variability and the observed correlations shown for reference. (c) shows

660   exceedance probability (1–CDF) curves for precipitation intensities above 10 mm/day,
661   where colored lines represent individual seeds for each diffusion setting. (d) presents Q–Q
662   plots of precipitation quantiles above 10 mm/day, with shaded bands indicating inter-seed
663   variability.

664   Differences between diffusion lengths are also impacting the representation of
665   extreme precipitation. The exceedance probability curves show that both configurations
666   capture the general shape of the upper tail, but the $T = 100$ ensemble exhibits a slight
667   overestimation at high intensities, with some realizations producing heavier tails. This
668   behavior is also shown in the joint $Q - Q$ plots (Figure 9). These results indicate that
669   reducing diffusion length slightly amplifies small-scale noisiness and extreme-value
670   sampling, rather than altering the central tendency or overall structure of the precipitation
671   distribution.

672   Additional experiments using an even shorter diffusion length ($T = 50$), presented in
673   Figure S9 of the Supplementary Information, further illustrate the trade-off between
674   inference cost and predictive uncertainty. While the $T = 50$ configuration continues to
675   reproduce the overall spatial organization and distributional shape of precipitation, the bias
676   in the upper tail becomes more pronounced. Across diffusion lengths ($T = 500, 100,$ and $50$),
677   DDPM consistently captures the characteristic decay of spatial dependence and the general
678   structure of the precipitation distribution; however, larger diffusion lengths yield
679   systematically closer agreement with observations and reduced ensemble spread,
680   particularly for the 16× downscaling case. These results indicate that diffusion depth
681   primarily controls the strength and coherence of spatial structure rather than altering the
682   fundamental behavior of the model. This analysis underlines a practical trade-off in
683   diffusion-based downscaling that longer diffusion schedules provide more constrained and
684   spatially coherent realizations at higher computational cost, whereas shorter schedules
685   substantially reduce inference time at the expense of increased ensemble variability, and
686   overestimation of extreme precipitation.

## 6. Conclusions

688   This study presents a systematic comparison of deterministic and generative deep-learning
689   approaches for precipitation super-resolution under challenging 8× and 16× downscaling
690   tasks. Using a consistent experimental framework, we evaluated U-Net, WGAN, and DDPM
691   across complementary diagnostics of statistical fidelity, spatial structure, and extreme-value
692   behavior. All three models preserve aggregate rainfall mass despite the absence of explicit
693   conservation constraints, with spatial aggregation showing rapid convergence toward near-
694   perfect agreement at coarse scales. Model differences emerge primarily at fine spatial scales,
695   particularly in the representation of extremes and spatial dependence. Importantly, the
696   cross-regional design provides insight into model robustness across distinct hydroclimatic

697    regimes, highlighting performance under distribution shift between training and evaluation
698    domains.
699        U-NET demonstrates strong stability, structural consistency, and precipitation-
700    occurrence skill, yielding high SSIM and ROC performance and smooth spatial fields that
701    perform well at short spatial lags. However, this deterministic smoothing suppresses small-
702    scale variance, attenuates spectral power, and systematically underestimates extremes,
703    particularly under 16× downscaling. The generative models provide complementary
704    strengths. WGAN more effectively captures fine-scale variability, spatial dependence, and
705    upper-tail behavior, producing sharper precipitation structures and improved extreme-
706    value statistics, though with increased variability and reduced structural consistency at
707    small scales. DDPM offers a balanced alternative, maintaining coherent multi-scale storm
708    morphology while explicitly representing conditional uncertainty through stochastic
709    sampling. This uncertainty manifests as increased ensemble spread in extreme-value
710    diagnostics, especially at high downscaling factors.
711        Sensitivity experiments further show that diffusion length primarily governs a trade-
712    off between computational efficiency and spatial coherence: reducing the number of
713    diffusion steps increases variance and weakens fine-scale structure while preserving the
714    overall distributional and morphological characteristics of precipitation. These results
715    highlight fundamental trade-offs among determinism, spatial realism, uncertainty
716    representation, and computational cost. No single approach is uniformly optimal; instead,
717    model selection should be guided by application-specific priorities. U-NET is well suited for
718    tasks emphasizing stability, occurrence detection, and computational efficiency, whereas
719    generative models are preferable when accurate representation of spatial variability and
720    extremes is critical. Among the generative approaches, WGAN most closely reproduces
721    observed spatial correlations and extreme-value behavior, while DDPM provides physically
722    coherent ensemble realizations with explicit uncertainty representation at higher
723    computational cost.
724        Several limitations necessitate future investigation. The analysis is conducted under
725    a perfect-model framework and focuses on a single region and resolution pair. Diffusion-
726    based inference remains computationally expensive, and further work is needed to assess
727    scalability for operational applications. Future research should explore conditioning on
728    additional physical predictors, hybrid deterministic–generative architectures, explicit
729    physical constraints, and calibration strategies to better control ensemble spread. Overall,
730    this study demonstrates that generative downscaling methods offer clear advantages for
731    representing fine-scale spatial structure and extremes, while deterministic approaches
732    remain valuable for stable and efficient precipitation reconstruction under increasing
733    resolution demands.

## Data Availability

ERA5-Land precipitation data used in this study are publicly available from the Copernicus Climate Change Service (C3S) Climate Data Store. ERA5-Land provides global land-surface variables at approximately 9-km spatial resolution and hourly temporal resolution. The data can be accessed through the C3S Climate Data Store (https://cds.climate.copernicus.eu/datasets/reanalysis-era5-land?tab=download) (Muñoz Sabater, 2019). All preprocessing steps applied in this study are described in the Methods section.

## Author Contributions

SS led the conceptualization, methodology development, model implementation, experiments, analysis, and manuscript writing. SMP, TH, and AM supervised the research, contributed to experimental design, interpretation of results, and manuscript review and editing, and supported funding acquisition. HMA contributed to experimental design and assisted with manuscript review and editing. All authors approved the final manuscript.

## Competing Interests

The author declares no competing interests.

## Acknowledgement

## References

Abdelmoaty, H. M., Papalexiou, S. M., Mamalakis, A., Singh, S., Coia, V., Hairabedian, M., Szeftel, P., & Grover, P. (2025). Generative Adversarial Networks for Downscaling Hourly Precipitation in the Canadian Prairies. *Journal of Geophysical Research: Machine Learning and Computation*, *2*(4). https://doi.org/10.1029/2025JH000678

Arjovsky, M., Chintala, S., & Bottou, L. (2017). *Wasserstein GAN*. http://arxiv.org/abs/1701.07875

Baño-Medina, J., Manzanas, R., & Gutiérrez, J. M. (2020). Configuration and intercomparison of deep learning neural models for statistical downscaling. *Geoscientific Model Development*, *13*(4), 2109–2124. https://doi.org/10.5194/gmd-13-2109-2020

Bednarz, T., & Cherukuri, R. (2023). Use of Physics-Based AI for Simulations and Modeling in the Era of Digital Twins. *SIGGRAPH Asia 2023 Courses*, 1–45. https://doi.org/10.1145/3610538.3614627

Coppola, E., Sobolowski, S., Pichelli, E., Raffaele, F., Ahrens, B., Anders, I., Ban, N., Bastin, S., Belda, M., Belusic, D., Caldas-Alvarez, A., Cardoso, R. M., Davolio, S., Dobler, A., Fernandez, J., Fita, L., Fumiere, Q., Giorgi, F., Goergen, K., … Warrach-Sagi, K. (2020). A first-of-its-kind multi-model convection permitting ensemble for investigating convective

770      phenomena over Europe and the Mediterranean. *Climate Dynamics*, *55*(1–2), 3–34.
771      https://doi.org/10.1007/s00382-018-4521-8

772 Deser, C., Phillips, A., Bourdette, V., & Teng, H. (2012). Uncertainty in climate change
773      projections: the role of internal variability. *Climate Dynamics*, *38*(3–4), 527–546.
774      https://doi.org/10.1007/s00382-010-0977-x

775 Feser, F., Rockel, B., von Storch, H., Winterfeldt, J., & Zahn, M. (2011). Regional Climate Models
776      Add Value to Global Model Data: A Review and Selected Examples. *Bulletin of the*
777      *American Meteorological Society*, *92*(9), 1181–1192.
778      https://doi.org/10.1175/2011BAMS3061.1

779 Gao, X. J., Shi, Y., Zhang, D., Wu, J., Giorgi, F., Ji, Z., & Wang, Y. (2012). Uncertainties in monsoon
780      precipitation projections over China: Results from two high-resolution RCM
781      simulations. *Clim Res*, *52*. https://doi.org/10.3354/cr0108

782 Gilleland, E., Ahijevych, D., Brown, B. G., Casati, B., & Ebert, E. E. (2009). Intercomparison of
783      Spatial Forecast Verification Methods. *Weather and Forecasting*, *24*(5), 1416–1430.
784      https://doi.org/10.1175/2009WAF2222269.1

785 Giorgi, F., & Gutowski, W. J. (2015). Regional dynamical downscaling and the CORDEX
786      initiative. *Annu. Rev. Environ. Resour.*, *40*. https://doi.org/10.1146/annurev-environ-
787      102014-021217

788 Giorgi, F., & Mearns, L. O. (1991). Approaches to the simulation of regional climate change: A
789      review. *Reviews of Geophysics*, *29*(2), 191–216. https://doi.org/10.1029/90RG02636

790 Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., & Courville, A. (2017). *Improved Training*
791      *of Wasserstein GANs*.

792 Harris, L., McRae, A. T. T., Chantry, M., Dueben, P. D., & Palmer, T. N. (2022). A Generative
793      Deep Learning Approach to Stochastic Downscaling of Precipitation Forecasts. *Journal*
794      *of Advances in Modeling Earth Systems*, *14*(10).
795      https://doi.org/10.1029/2022MS003120

796 Harrison, D. R., McGovern, A., Karstens, C. D., Bostrom, A., Demuth, J. L., Jirak, I. L., & Marsh,
797      P. T. (2025). An Assessment of How Domain Experts Evaluate Machine Learning in
798      Operational Meteorology. *Weather and Forecasting*, *40*(3), 393–410.
799      https://doi.org/10.1175/WAF-D-24-0144.1

800 Ho, J., Jain, A., & Abbeel, P. (2020). Denoising Diffusion Probabilistic Models. In H. Larochelle,
801      M. Ranzato, R. Hadsell, M. F. Balcan, & H. Lin (Eds.), *Advances in Neural Information*
802      *Processing Systems* (Vol. 33, pp. 6840–6851). Curran Associates, Inc.
803      https://proceedings.neurips.cc/paper_files/paper/2020/file/4c5bcfec8584af0d967f1
804      ab10179ca4b-Paper.pdf

805 Hobeichi, S., Nishant, N., Shao, Y., Abramowitz, G., Pitman, A., Sherwood, S., Bishop, C., &
806      Green, S. (2023). Using Machine Learning to Cut the Cost of Dynamical Downscaling.
807      *Earth's Future*, *11*(3). https://doi.org/10.1029/2022EF003291

808 Höhlein, K., Kern, M., Hewson, T., & Westermann, R. (2020). A comparative study of
809      convolutional neural network models for wind field downscaling. *Meteorological*
810      *Applications*, *27*(6). https://doi.org/10.1002/met.1961

811 Hsu, L.-H., Chiang, C.-C., Lin, K.-L., Lin, H.-H., Chu, J.-L., Yu, Y.-C., & Fahn, C.-S. (2024).
812      Downscaling Taiwan precipitation with a residual deep learning approach. *Geoscience*
813      *Letters*, *11*(1), 23. https://doi.org/10.1186/s40562-024-00340-y

814 Khader, F., Müller-Franzes, G., Tayebi Arasteh, S., Han, T., Haarburger, C., Schulze-Hagen, M.,
815      Schad, P., Engelhardt, S., Baeßler, B., Foersch, S., Stegmaier, J., Kuhl, C., Nebelung, S.,

816  Kather, J. N., & Truhn, D. (2023). Denoising diffusion probabilistic models for 3D medical
817      image generation. *Scientific Reports*, *13*(1), 7303. https://doi.org/10.1038/s41598-
818      023-34341-2
819  Kumar, B., Atey, K., Singh, B. B., Chattopadhyay, R., Acharya, N., Singh, M., Nanjundiah, R. S., &
820      Rao, S. A. (2023). On the modern deep learning approaches for precipitation
821      downscaling. *Earth Science Informatics*, *16*(2), 1459–1472.
822      https://doi.org/10.1007/s12145-023-00970-4
823  Lakshminarayanan, B., Pritzel, A., & Blundell, C. (2017). Simple and Scalable Predictive
824      Uncertainty Estimation using Deep Ensembles. In I. Guyon, U. Von Luxburg, S. Bengio, H.
825      Wallach, R. Fergus, S. Vishwanathan, & R. Garnett (Eds.), *Advances in Neural Information
826      Processing Systems* (Vol. 30). Curran Associates, Inc.
827      https://proceedings.neurips.cc/paper_files/paper/2017/file/9ef2ed4b7fd2c810847ff
828      a5fa85bce38-Paper.pdf
829  Lange, S. (2019). Trend-preserving bias adjustment and statistical downscaling with
830      ISIMIP3BASD (v1.0). *Geoscientific Model Development*, *12*(7), 3055–3070.
831      https://doi.org/10.5194/gmd-12-3055-2019
832  Lucas-Picher, P., Argüeso, D., Brisson, E., Tramblay, Y., Berg, P., Lemonsu, A., Kotlarski, S., &
833      Caillaud, C. (2021). Convection-permitting modeling with regional climate models:
834      Latest developments and next steps. *WIREs Climate Change*, *12*(6).
835      https://doi.org/10.1002/wcc.731
836  Lyu, R., Wang, L., Sun, Y., Bai, H., & Lu, C.-T. (2024). Downscaling Precipitation with Bias-
837      informed Conditional Diffusion Model. *2024 IEEE International Conference on Big Data
838      (BigData)*, 8768–8770. https://doi.org/10.1109/BigData62323.2024.10825056
839  Mamalakis, A., Langousis, A., Deidda, R., & Marrocu, M. (2017). A parametric approach for
840      simultaneous bias correction and high-resolution downscaling of climate model rainfall.
841      *Water Resources Research*, *53*(3), 2149–2170.
842      https://doi.org/10.1002/2016WR019578
843  Maraun, D., Wetterhall, F., Ireson, A. M., Chandler, R. E., Kendon, E. J., Widmann, M., Brienen,
844      S., Rust, H. W., Sauter, T., Themeßl, M., Venema, V. K. C., Chun, K. P., Goodess, C. M., Jones,
845      R. G., Onof, C., Vrac, M., & Thiele-Eich, I. (2010). Precipitation downscaling under climate
846      change: Recent developments to bridge the gap between dynamical models and the end
847      user. *Reviews of Geophysics*, *48*(3), RG3003. https://doi.org/10.1029/2009RG000314
848  Meghani, S., Singh, S., Kumar, N., & Goyal, M. K. (2023). Predicting the spatiotemporal
849      characteristics of atmospheric rivers: A novel data-driven approach. *Global and
850      Planetary Change*, *231*, 104295. https://doi.org/10.1016/j.gloplacha.2023.104295
851  Muñoz Sabater, J. (2019). *ERA5-Land hourly data from 1950 to present*. Copernicus Climate
852      Change Service (C3S) Climate Data Store (CDS).
853      https://doi.org/10.24381/cds.e2161bac
854  Muñoz-Sabater, J., Dutra, E., Agustí-Panareda, A., Albergel, C., Arduini, G., Balsamo, G.,
855      Boussetta, S., Choulga, M., Harrigan, S., Hersbach, H., Martens, B., Miralles, D. G., Piles, M.,
856      Rodríguez-Fernández, N. J., Zsoter, E., Buontempo, C., & Thépaut, J.-N. (2021). ERA5-
857      Land: a state-of-the-art global reanalysis dataset for land applications. *Earth System
858      Science Data*, *13*(9), 4349–4383. https://doi.org/10.5194/essd-13-4349-2021
859  Nishant, N., Hobeichi, S., Sherwood, S., Abramowitz, G., Shao, Y., Bishop, C., & Pitman, A.
860      (2023). Comparison of a novel machine learning approach with dynamical downscaling

861       for Australian precipitation. *Environmental Research Letters*, *18*(9), 094006.
862           https://doi.org/10.1088/1748-9326/ace463
863   Palmer, T. (2014). Build high-resolution global climate models. *Nature*, *515*.
864           https://doi.org/10.1038/515338a
865   Papalexiou, S. M., & Mamalakis, A. (2025). Machine unlearning: bias correction in neural
866           network downscaled storms. *Journal of Hydrology*, 134689.
867           https://doi.org/10.1016/j.jhydrol.2025.134689
868   Papalexiou, S. M., Serinaldi, F., & Porcu, E. (2021). Advancing Space-Time Simulation of
869           Random Fields: From Storms to Cyclones and Beyond. *Water Resources Research*, *57*(8).
870           https://doi.org/10.1029/2020WR029466
871   Perez, E., Strub, F., De Vries, H., Dumoulin, V., & Courville, A. (2018). FiLM: Visual Reasoning
872           with a General Conditioning Layer. *Proceedings of the AAAI Conference on Artificial
873           Intelligence*, *32*(1). https://doi.org/10.1609/aaai.v32i1.11671
874   Piani, C., Haerter, J. O., & Coppola, E. (2010). Statistical bias correction for daily precipitation
875           in regional climate models over Europe. *Theoretical and Applied Climatology*, *99*(1–2),
876           187–192. https://doi.org/10.1007/s00704-009-0134-9
877   Rampal, N., Hobeichi, S., Gibson, P. B., Baño-Medina, J., Abramowitz, G., Beucler, T., González-
878           Abad, J., Chapman, W., Harder, P., & Gutiérrez, J. M. (2024). Enhancing Regional Climate
879           Downscaling through Advances in Machine Learning. *Artificial Intelligence for the Earth
880           Systems*, *3*(2). https://doi.org/10.1175/aies-d-23-0066.1
881   Ravuri, S., Lenc, K., Willson, M., Kangin, D., Lam, R., Mirowski, P., Fitzsimons, M.,
882           Athanassiadou, M., Kashem, S., Madge, S., Prudden, R., Mandhane, A., Clark, A., Brock, A.,
883           Simonyan, K., Hadsell, R., Robinson, N., Clancy, E., Arribas, A., & Mohamed, S. (2021).
884           Skilful precipitation nowcasting using deep generative models of radar. *Nature*,
885           *597*(7878), 672–677. https://doi.org/10.1038/s41586-021-03854-z
886   Roberts, N. M., & Lean, H. W. (2008). Scale-Selective Verification of Rainfall Accumulations
887           from High-Resolution Forecasts of Convective Events. *Monthly Weather Review*, *136*(1),
888           78–97. https://doi.org/10.1175/2007MWR2123.1
889   Schär, C. (2019). Kilometer-scale climate models. Prospects and challenges article. *Am.
890           Meteorol. Soc.*, *101*.
891   Singh, S., & Goyal, M. K. (2023). An innovative approach to predict atmospheric rivers:
892           Exploring convolutional autoencoder. *Atmospheric Research*, *289*, 106754.
893   Skamarock, W. C. (2004). Evaluating Mesoscale NWP Models Using Kinetic Energy Spectra.
894           *Monthly Weather Review*, *132*(12), 3019–3032. https://doi.org/10.1175/MWR2830.1
895   Song, Y., & Dhariwal, P. (2023). *Improved Techniques for Training Consistency Models*.
896           http://arxiv.org/abs/2310.14189
897   Stengel, K., Glaws, A., Hettinger, D., & King, R. N. (2020). Adversarial super-resolution of
898           climatological wind and solar data. *Proceedings of the National Academy of Sciences*,
899           *117*(29), 16805–16815. https://doi.org/10.1073/pnas.1918964117
900   Stephens, G. (2017). Challenges and Advances in Convection-Permitting Climate Modeling.
901           *Bulletin of the American Meteorological Society*, *98*(5), 1027–1030.
902           https://doi.org/10.1175/BAMS-D-16-0263.1
903   Tabari, H., Paz, S. M., Buekenhout, D., & Willems, P. (2021). Comparison of statistical
904           downscaling methods for climate change impact analysis on precipitation-driven
905           drought. *Hydrology and Earth System Sciences*, *25*(6), 3493–3517.
906           https://doi.org/10.5194/hess-25-3493-2021

907    Teutschbein, C., & Seibert, J. (2012). Bias correction of regional climate model simulations
908        for hydrological climate-change impact studies: Review and evaluation of different
909        methods. *Journal of Hydrology*, *456–457*, 12–29.
910        https://doi.org/10.1016/j.jhydrol.2012.05.052
911    Tomasi, E., Franch, G., & Cristoforetti, M. (2025). Can AI be enabled to perform dynamical
912        downscaling? A latent diffusion model to mimic kilometer-scale COSMO5.0_CLM9
913        simulations. *Geoscientific Model Development*, *18*(6), 2051–2078.
914        https://doi.org/10.5194/gmd-18-2051-2025
915    Trenberth, K. E., Fasullo, J. T., & Shepherd, T. G. (2015). Attribution of climate extreme events.
916        *Nature Climate Change*, *5*(8), 725–730. https://doi.org/10.1038/nclimate2657
917    Vandal, T., Kodra, E., Ganguly, S., Michaelis, A., Nemani, R., & Ganguly, A. R. (2017). DeepSD:
918        Generating High Resolution Climate Change Projections through Single Image Super-
919        Resolution. *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge
920        Discovery and Data Mining*, 1663–1672. https://doi.org/10.1145/3097983.3098004
921    Vrac, M., Stein, M. L., Hayhoe, K., & Liang, X. Z. (2007). A general method for validating
922        statistical downscaling methods under future climate change. *Geophysical Research
923        Letters*, *34*(18). https://doi.org/10.1029/2007GL030295
924    Wang, F., Tian, D., Lowe, L., Kalin, L., & Lehrter, J. (2021). Deep Learning for Daily
925        Precipitation and Temperature Downscaling. *Water Resources Research*, *57*(4).
926        https://doi.org/10.1029/2020WR029308
927    Wood, A. W., Leung, L. R., Sridhar, V., & Lettenmaier, D. P. (2004). Hydrologic Implications of
928        Dynamical and Statistical Approaches to Downscaling Climate Model Outputs. *Climatic
929        Change*, *62*(1–3), 189–216. https://doi.org/10.1023/B:CLIM.0000013685.99609.9e
930    Yan, W., Wang, Y., Gu, S., Huang, L., Yan, F., Xia, L., & Tao, Q. (2019). *The Domain Shift Problem
931        of Medical Image Segmentation and Vendor-Adaptation by Unet-GAN* (pp. 623–631).
932        https://doi.org/10.1007/978-3-030-32245-8_69
933    Zhou Wang, Bovik, A. C., Sheikh, H. R., & Simoncelli, E. P. (2004). Image quality assessment:
934        from error visibility to structural similarity. *IEEE Transactions on Image Processing*,
935        *13*(4), 600–612. https://doi.org/10.1109/TIP.2003.819861
936