1    **A Robust Generative Adversarial Network Approach for Climate Downscaling**
2                                   **and Weather Generation**

3    Neelesh Rampal[1,3], Peter B. Gibson[2], Steven Sherwood[3], Gab Abramowitz[3] and Sanaa
4                                        Hoibechi[3]

5

6         [1]National Institute of Water and Atmospheric Research, Auckland, New Zealand.

7         [2]National Institute of Water and Atmospheric Research, Wellington, New Zealand

8    [3] Climate Change Research Centre & ARC Centre of Excellence for Climate Extremes,
9                      University of New South Wales, Sydney, Australia.

10

11   Corresponding author: Neelesh Rampal (neelesh.rampal@niwa.co.nz)

12

13   **Key Points:**

14   •   Regression-based downscaling methods struggle to capture extreme events and poorly
15       preserve important climate statistical measures.

16   •   Generative Adversarial Networks outperform regression-based methods, but their
17       skill is very sensitive to the selection of hyperparameters.

18   •   Using constraints in the Generative Adversarial Network's loss function improves its
19       robustness and skill of across a wide range of metrics.

20
21

22 **Abstract**

23    Anticipating climate impacts and risks in present or future climates requires
24 predicting the statistics of high-impact weather events at fine-scales. Direct numerical
25 simulations of fine-scale weather are computationally too expensive for many uses. While
26 regression-based (deep-learning or statistical) downscaling of low-resolution climate
27 simulations is several orders of magnitude faster than direct numerical simulations, it suffers
28 from several limitations. These limitations include the tendency to regress to the mean, which
29 produces excessively smooth predictions and underestimates the magnitude of extreme
30 events. Additionally, they also fail to preserve statistical measures that are key for climate
31 research. We use a conditional GAN (c-GAN) architecture to downscale daily precipitation
32 as a Regional Climate Model (RCM) emulator. The c-GAN generates plausible residuals on
33 top of the predictable expectation state produced by a regression-based DL algorithm. The
34 skill of c-GANs is highly sensitive to a hyperparameter known as the weight of the
35 adversarial loss ($\lambda_{adv}$), and the value of $\lambda_{adv}$ required for accurate results varies with season
36 and performance metric, casting doubt on the robustness of c-GANs as usually implemented.
37 But, by applying a simple intensity constraint to the loss function, it is possible to obtain
38 robust performance results across $\lambda_{adv}$ spanning two orders of magnitude. C-GANs are
39 considerably more skillful in capturing climatological statistics including the distribution and
40 spatial characteristics of extreme events. We expect c-GANs with this modification to be
41 readily transferable to other problems and time periods, making them a useful weather
42 generator for representing extreme event statistics in present and future climates.

43 **Plain Language Summary**

44    Climate projections produced by Global Climate Models (GCMs) have a typical
45 resolution of 100-200km, which is too coarse for studying climate impacts at regional scales.
46 Dynamical downscaling involves running a Regional Climate Model (RCM) to simulate
47 physical processes that are not resolved at the resolution of GCM, enabling high-resolution
48 climate projections for studying localized climate change impacts. However, RCMs are
49 computationally expensive, limiting both the number of GCMs that can be downscaled and
50 estimates of uncertainty. Deep learning (DL) methods offer a promising, cost-effective
51 alternative to RCMs, and recent studies have emulated certain aspects of RCMs at a fraction
52 of the computational cost. Generative DL algorithms such as Generative Adversarial
53 Networks (GANs) appear to show promise in accurately emulating RCMs, but their training
54 instability and inconsistent performance across climate contexts raises concerns about their
55 robustness for downscaling climate projections. Here we develop and introduce a simple
56 technique to improve the stability in GAN performance across a wide range of training
57 configurations. This improves robustness and utility in broader climate applications.

58

59

60

61    1. **Introduction**

62    The coarse spatial resolution of Global Climate Models (GCMs) limits their ability to
63    simulate climate changes at regional and local scales, where the impacts of climate change
64    are most directly experienced (Benestad, 2004, 2010; Fowler et al., 2007; Maraun, 2016;
65    Maraun et al., 2010). Dynamical downscaling aims to address this resolution issue by
66    capturing finer-scale aspects of mesoscale circulation and regional climate across different
67    landscapes such as mountain ranges, valleys, and coastal boundaries (Feser et al., 2011;
68    Gensini et al., 2023; Giorgi et al., 1994; Hoogewind et al., 2017; Jones et al., 1995; Liu et al.,
69    2017; Prein et al., 2015; Xu et al., 2019). Dynamical downscaling typically involves running
70    a Regional Climate Model (RCM) from the lateral boundary conditions of a GCM. A major
71    drawback of RCMs is their high computational cost, which today limits their spatial
72    resolution to scales of 12-50km when run operationally in Coordinated Regional climate
73    Downscaling Experiment (CORDEX) type experiments (Giorgi et al., 2009). Additionally,
74    the computational cost of RCMs limits the number of GCMs that can be downscaled.
75    Consequently, this small number of downscaled GCMs means that model structural and
76    internal variability uncertainty are under sampled in regional climate projections, despite its
77    known importance on regional scales (Deser et al., 2012; Deser & Phillips, 2023; Gibson et
78    al., 2024; Hawkins & Sutton, 2009, 2011).

79    Recently, computationally efficient statistical/empirical algorithms have been explored
80    for RCM emulation, including simple multiple linear regression (Holden et al., 2015),
81    multilayer perceptron (Chadwick et al., 2011; Hobeichi et al., 2023; Nishant et al., 2023),
82    statistical analogues (Boé et al., 2023), and normalizing flows (Groenke et al., 2020). In both
83    RCM emulation and other downscaling applications, there has been a shift towards
84    regression-based deep learning computer vision algorithms such as CNNs (Babaousmail et
85    al., 2021; Bano-Medina et al., 2023; Doury et al., 2022; van der Meer et al., 2023). These are
86    better suited to the complex non-linear relationships between large-scale predictors and local-
87    scale climate variables (Rampal et al., 2022) and have generally outperformed traditional
88    statistical and machine learning (ML) techniques (Baño-Medina et al., 2020; Rampal, 2024;
89    Rampal et al., 2022).

90    While regression-based approaches (including deep learning) are skillful in capturing the
91    "mean-state" in instantaneous predictions (i.e. they regress to the mean), they tend to
92    underestimate extreme events and struggle to resolve fine scale details (Harris et al., 2022;
93    Mardani et al., 2023; Rampal, 2024; Reddy et al., 2023; Vosper et al., 2023; J. Wang et al.,
94    2021). Unlike weather forecasting, accurate instantaneous predictions are less useful than
95    climatological metrics (i.e. how often a given weather event occurs) in a climate projection
96    context, as atmospheric variability is chaotic and effectively random beyond a short horizon.
97    This may create a trade-off between accuracy of instantaneous predictions, and the skill in
98    capturing climatological metrics and extreme events (Rampal et al., 2024). This is
99    particularly problematic for extreme events (e.g. convective high intensity short duration
100   rainfall events) which can have the highest societal impact. While there have been a wide

101  variety of algorithm developments to overcome such issues in regression-based approaches,
102  these issues persist (for recent reviews see Rampal et al., 2024; Sun et al., 2024) .

103      Generative Adversarial Networks (GANs) are a recent development in ML that may offer
104  a solution to some of these shortcomings. GANs have been used in many research areas
105  (Goodfellow et al., 2014; Isola et al., 2018; Mirza & Osindero, 2014; X. Wang et al., 2018),
106  and have recently been adapted from the computer vision sub-field of super-resolution (which
107  focuses on enhancing image resolution) to climate downscaling. GANs, also often described
108  as conditional GANs (c-GANs) in this context, have significantly improved regression-based
109  computer-vision algorithms in predicting local-scale extreme events and resolving high-
110  resolution spatial structure in the downscaled predictions (Annau et al., 2023; Brochet et al.,
111  2023; Izumi et al., 2022; Leinonen et al., 2021; Miralles et al., 2022; Oyama et al., 2023;
112  Price & Rasp, 2022; Ravuri et al., 2021a; Saha & Ravela, 2022; Vosper et al., 2023; J. Wang
113  et al., 2021).

114      Unlike traditional regression-based ML algorithms, which optimize for loss functions
115  such as mean squared error (MSE), GANs are generative algorithms that incorporate an
116  adversarial loss (Goodfellow et al., 2014; Mirza & Osindero, 2014) and use stochastic noise
117  to generate an ensemble of predictions for a given set of large-scale predictor variables (i.e.
118  coarse-resolution variables from GCMs).The adversarial loss function drives a competitive
119  process between two CNNs, a generator and discriminator. The generator attempts to
120  generate realistic samples (i.e. pseudo-RCM simulations), while the discriminator tries to
121  distinguish between real data (i.e. RCM simulations) and the generator's output (Goodfellow
122  et al., 2014; Mirza & Osindero, 2014). This competition leads to the generator implicitly
123  learning through a powerful loss function that goes beyond traditional pixel-wise
124  comparisons, encouraging the generation of outputs to be distributionally and structurally
125  similar to the real data (Gulrajani et al., 2017).

126      The effectiveness of GANs for climate downscaling in present-day or future climates has
127  not been well-assessed (Rampal, et al., 2024). Existing research mainly focuses on using
128  traditional error metrics such as root-mean-squared error (Rampal et al., 2024; Sun et al.,
129  2024) instead of climatological metrics. Additionally, GANs are notoriously unstable and
130  challenging to train, where stability is often determined by selecting the correct
131  hyperparameters (Arjovsky et al., 2017; Goodfellow et al., 2014; Gulrajani et al., 2017; Mirza
132  & Osindero, 2014).

133      One particularly important hyperparameter is the weighting of the adversarial loss
134  function ($\lambda_{adv}$) during training (refer to section 2.1 for more details), which determines the
135  strength of the adversarial loss during training. While studies have analyzed the impact of
136  model architecture and loss function choices on generated output quality, this research has
137  been limited to computer vision applications (Abu-Srhan et al., 2022; Isola et al., 2018; Ledig
138  et al., 2017; X. Wang et al., 2018). For example, Isola et al., (2018) highlighted that values
139  too large would often hallucinate and generate artifacts (i.e. $\lambda_{adv} = 1$), and found optimal
140  performance when $\lambda_{adv} = 0.01$ for image-to-image translation. Existing studies in

141 downscaling applications have only conducted their research using a specific value of $\lambda_{adv}$
142 (i.e. Annau et al., 2023; Harris et al., 2022; Leinonen et al., 2021; Vosper et al., 2023), with
143 limited exploration of how the strength of the adversarial loss affects climate downscaling.

144     Our study therefore aims to focus on two aspects of evaluating GANs. Firstly, we assess
145 whether GANs add value over regression-based RCM emulators. Secondly, we explore the
146 robustness of GAN performance for RCM emulation, by varying the hyperparameter $\lambda_{adv}$.
147 Our study uses a comprehensive set of evaluation metrics to ensure that GANs are useful in a
148 variety of climate downscaling contexts. These metrics assess the emulator's ability to learn
149 various climate statistics, such as the climatology of precipitation, extreme events, and the
150 persistence of dry spells. We also evaluate the skill of GANs to generate ensembles, which
151 have significant implications for uncertainty quantification in climate science and weather
152 forecasting.

153 **2. Materials and Methods**

154 2.1 Training and Evaluation Data
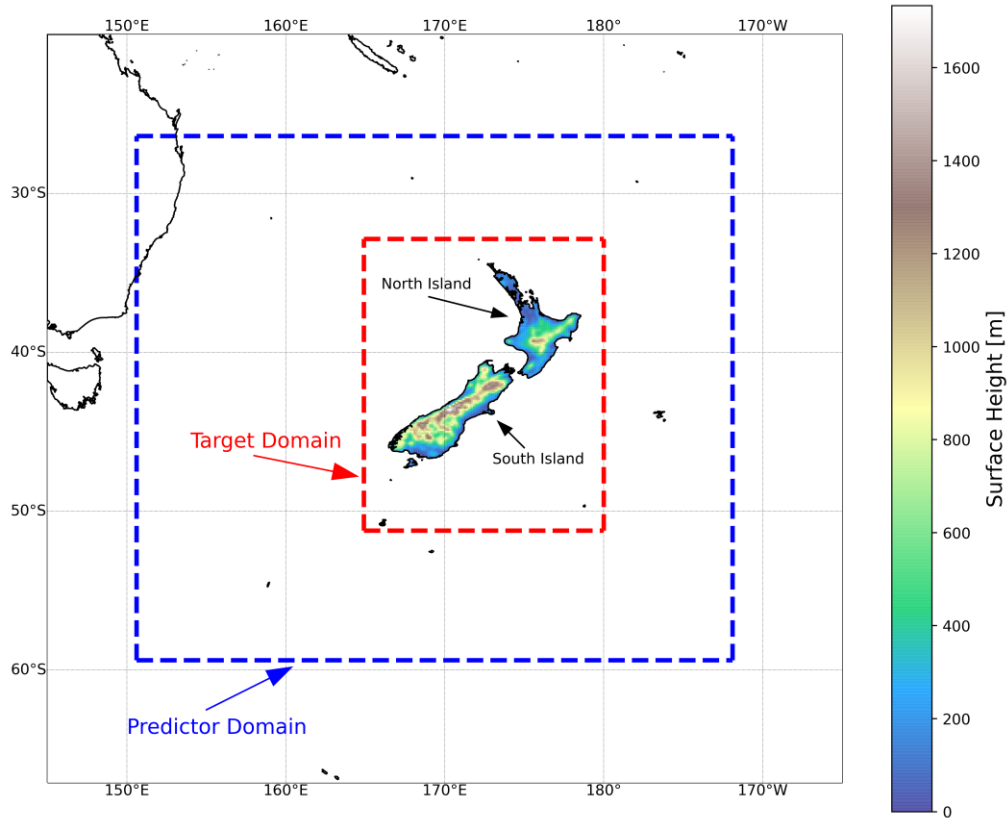
155 2.1.1 Regional Climate Model Configuration

156     Our RCM emulator was trained using predictor and target variables from the Conformal
157 Cubic Atmospheric Model (CCAM), a global non-hydrostatic atmospheric model with a
158 variable-resolution cubic grid (Chapman et al., 2023; Gibson et al., 2023; McGregor & Dix,
159 2008; Thatcher & McGregor, 2009). In contrast to commonly used RCMs like the Weather
160 Research and Forecasting Model (WRF), which rely on lateral boundary conditions from
161 reanalysis or CMIP6 GCMs, CCAM is run as a global variable-resolution model (McGregor
162 & Dix, 2008). CCAM is run globally with spectral nudging to input fields from GCM
163 atmospheric variables. A detailed evaluation of CCAM is presented in Gibson et al. (2023)
164 for this region, which used a very similar version of CCAM (i.e. model grid and physics
165 configuration).

166     Although CCAM is a global model, our emulation efforts concentrate on the New
167 Zealand region (165°E-184°W, 33°S-51°S) as shown in Figure 1 (target domain), where the
168 highest resolution face of CCAM is near-uniformly 12km. Due to its diverse array of
169 microclimates, the New Zealand region provides an ideal case study for RCM emulation.
170 These microclimates arise due to New Zealand's complex geography, including coastlines,
171 mountains, and its position in the mid-latitudes. New Zealand is also exposed to weather
172 phenomena such as tropical cyclones, atmospheric rivers, and large-scale climate drivers such
173 as the El Niño-Southern Oscillation (ENSO), and the Southern Annular Mode (SAM) (Refs).
174 While physical processes governing New Zealand's regional climate are generally well
175 captured by physics-based RCMs (Ackerley et al., 2012; Gibson et al., 2023), a key challenge
176 is ensuring that RCM emulators can also learn these processes (Rampal et al., 2024).

177

178    2.1.2 Training Data

179        The main target variable is daily accumulated high-resolution (~12km) precipitation ($pr$)
180    from CCAM output. Precipitation is logarithmically normalized ($z = \log_e(pr + 0.001)$) to
181    reduce its distributional skewness, as implemented in various weather forecasting and
182    downscaling studies (i.e. Rasp et al., 2020; Renwick et al., 2009). The coarse-resolution
183    predictor variables are daily-averaged large-scale prognostic variables from CCAM, which
184    include zonal wind ($U$), meridional wind ($V$), temperature ($T$) and specific humidity ($Q$) at
185    two pressure levels, 500hPa and 850hPa in the atmosphere. The domain extent of the
186    predictor variables is slightly larger than the target variable (151°E-188°W, 26°S-59°S) as
187    illustrated in Figure 1, and was chosen to prevent information scarcity at the boundaries of
188    the target domain (Bailie et al., 2024; Rampal et al., 2022). These predictor variables are re-
189    gridded from 12km to a resolution of 1.5° (~150km) using conservative remapping. The
190    predictor variables are normalized relative to the mean and standard deviation computed over
191    the entire training dataset as implemented in Rampal et al., (2022) and Rasp et al., (2020).
192    The rationale for using daily-aggregated predictor and target variables instead of sub-daily is
193    to both speed up model training and inference time but also reduce CPU/GPU memory usage.
194    Using daily input fields also ensures that the emulator can be applied to a much larger
195    number of GCMs, since the availability of daily data is much greater than sub-daily data
196    across the CMIP6 archive. It is important to note, that daily-aggregation also incurs a loss of
197    temporal information, making the problem somewhat more challenging than using
198    instantaneous fields (i.e. hourly).

199        Our study focuses solely on evaluating and training DL algorithms on historical RCM
200    simulations. Our evaluation framework does not focus on out-of-distribution performance
201    temporally (i.e. to future climates), but rather tests whether the emulator can be applied more
202    broadly to other un-seen GCM/RCM simulations from training. All emulators were trained
203    on 55 years of simulation (~21,000 days) from the CMIP6 ACCESS-CM2 (1960-2014). We
204    assess the performance of all emulators using ground-truth downscaled simulations from
205    CCAM, configured identically from two additional CMIP6 GCMs (EC-Earth3 and
206    NorESM2-MM). This out-of-sample evaluation covers a 20-year historical period from 1986
207    to 2005 (~7300 days). Here, our emulator is applied to the CCAM-coarsened predictor fields
208    (perfect framework) from these simulations. Doing so provides a true out-of-sample test of
209    the emulator, testing the performance (and ability to generalize) on additional driving fields
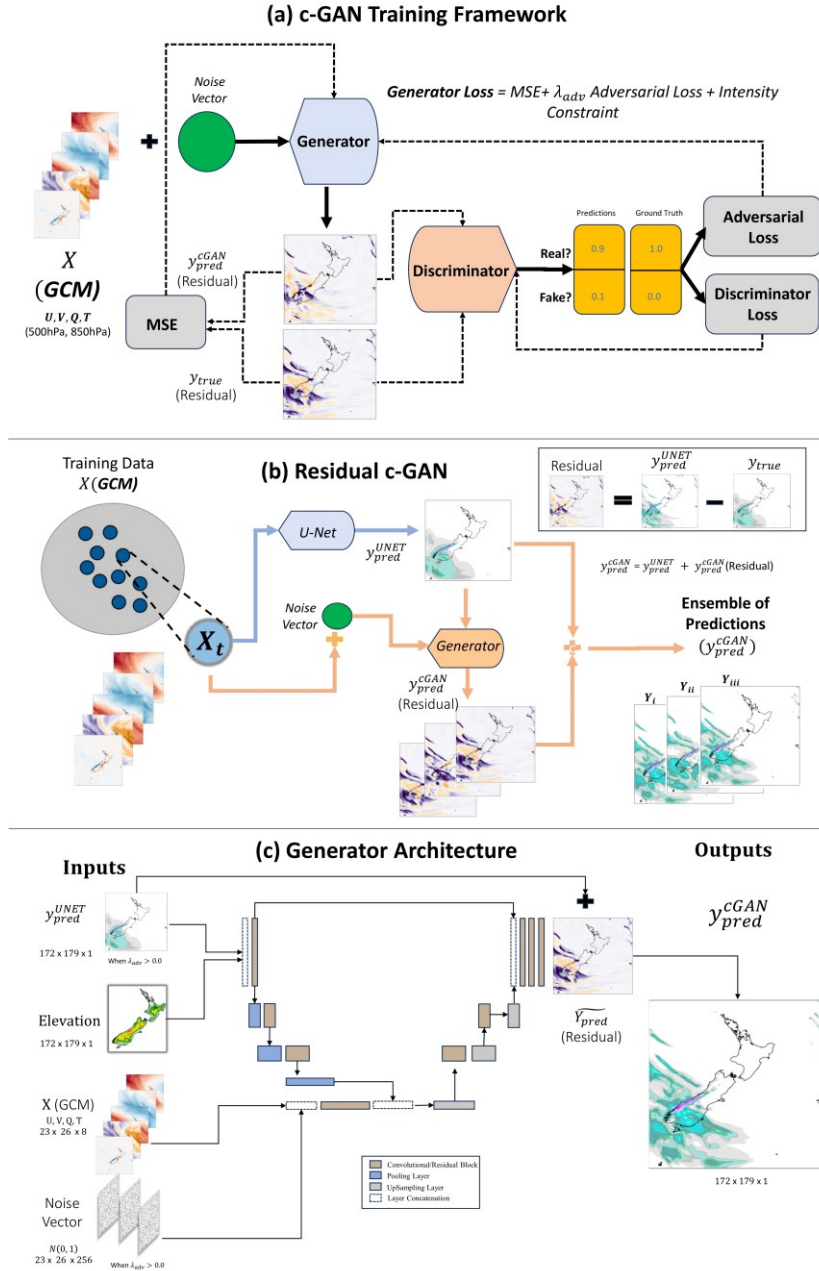210    from GCMs which were unseen in training.

211

212

**Figure 1**: A depiction of the domain extent of the predictor variables (blue) and target variables (red) across the New Zealand region, with the color scale representing the region's surface elevation

216   2.1.3 Training framework

217   We have used CCAM-coarsened predictor variables as opposed to variables from the
218   GCM directly. This training strategy is known as the perfect framework. It differs from the
219   imperfect training framework, which uses GCM fields as predictor variables directly
220   (Rampal, et al., 2024). Training an emulator through the imperfect framework is more
221   challenging as the RCM's mean state can significantly deviate from the GCM (Bartók et al.,
222   2017; Boé et al., 2020; Sørland et al., 2018). An emulator trained in the imperfect framework
223   needs to learn both the deviations between the RCM and GCMs mean state, and the finer
224   scales of RCM (Rampal et al., 2024), whereas the perfect framework emulator is only
225   required to learn the latter.

226   Additionally, emulators trained in the imperfect framework have been shown to learn a
227   relationship that is unique to a specific GCM/RCM pair (Bano-Medina et al., 2023; Boé et
228   al., 2023) and thus is less portable across the wider GCM/RCM matrix. Conversely, training
229   in the perfect framework has very little dependence on the RCM simulation used in training
230   (as it does not have to account for differences in circulation between RCM and GCM). While
231   there is an ongoing discussion about which framework is optimal in an out-of-sample
232   operational setting, training and evaluation in the perfect framework is simpler and involves

8

233  fewer considerations (Rampal, et al., 2024). This is advantageous for the purposes of testing
234  different configurations of GANs in this study, and more broadly applying the emulator to
235  downscaling multiple GCMs to generate high-resolution "pseudo" simulation.

236



237

238  **Figure 2**: (a) An illustration of the training feedback loop of a GAN (c-GAN). The
239  generator creates high-resolution predictions from low-resolution inputs and a noise vector.
240  The discriminator then measures how realistic the generator's outputs are through a
241  discriminator loss. (b) A residual GAN for downscaling. The residual GAN consists of two
242  steps; first, a regression-based U-Net is trained to produce a deterministic prediction ($y_{pred}^{UNET}$)
243  for a given X (predictor fields). This prediction and a noise vector are then input into the

244  generator, which is trained to predict the residual between the U-Net and the ground truth in
245  logarithmic space. (c) Generator Architecture: This flowchart depicts the architecture of the
246  generator within the GAN. It shows how multiple inputs, including a U-Net prediction
247  ($y_{pred}^{UNET}$), elevation data, low-resolution GCM data (X), and a noise vector pass through a
248  series of layers and processes to ultimately produce the high-resolution climate field
249  prediction ($y_{pred}^{cGAN}$) used by the GAN.

250

251  2.2 Generative Adversarial Networks

252      The GAN architecture for downscaling consists of two main components: a generator and
253  a discriminator (also known as a critic). The generator aims to create a high-resolution
254  climate field from a "low-resolution" climate field as an input (the condition), and a
255  discriminator evaluates whether the generated image is likely real (ground truth high-
256  resolution simulations) or fake (synthetic high-resolution fields generated by the generator
257  that may have characteristic artefacts). There are two main loss functions in training a GAN:
258  the generator loss ($G_{loss}$) and the discriminator (critic) loss.

259  2.2.1 Generator Loss

260      In this study we train c-GANs with two different loss function configurations. In a
261  downscaling or image super-resolution context, the generator loss usually consists of
262  traditional loss functions such as the MSE and an adversarial loss function, $G_{adv}$, which is
263  weighted by some constant factor $\lambda_{adv}$, as shown in Equation 1.

264  $$(1) \quad G_{loss}(y_{true}, y_{pred}) = MSE(y_{true}, y_{pred}) + \lambda_{adv} * G_{adv}(D(y_{pred})),$$

265  $$G_{adv}(y_{pred}) = \overline{- D(y_{pred})}$$

266      Here, $y_{true}$ and $y_{pred}$ refer to the ground truth RCM simulations and generated samples
267  from the emulator, respectively. The adversarial loss function ($G_{adv}$) is calculated by taking
268  the negative average of the discriminator's ($D$) output on generated samples $D(y_{pred})$. In
269  simpler terms, the adversarial loss increases when the discriminator is not fooled by the
270  generated images, penalizing the current weight set in the generator. The generator loss
271  shown in Equation 1 is one of the two main loss function configurations explored in this
272  study. It is widely used in many super-resolution and downscaling studies (i.e. Harris et al.,
273  2022; Leinonen et al., 2021; Vosper et al., 2023). Note we use the MSE loss function as
274  opposed to the MAE loss as it is more sensitive to errors in extreme events (not shown).  It is
275  important to note that training with an $\lambda_{adv}$ too large is often unstable (Isola et al., 2018;
276  Vosper et al., 2023), and the majority of existing studies generally use values of $\lambda_{adv}$ less
277  than 0.005 (Harris et al., 2022; Izumi et al., 2022; Leinonen et al., 2021; Vosper et al., 2023;
278  X. Wang et al., 2018).

279     We also explore a second loss function configuration that incorporates an intensity
280 constraint (IC), analogous to Ravuri et al. (2021) and Price & Rasp., (2022). The intensity
281 constraint penalizes both the model's maximum precipitation intensity over the regional
282 domain ($Y^{max}$) at each timestep, and its batch-averaged precipitation rate ($Y^{mean}$) for each
283 location, as shown in Equation 2. The maximum precipitation intensity constraint prevents
284 precipitation intensities from growing too large, and the batch-averaged precipitation (where
285 the batch size is 32) is a proxy for conserving monthly precipitation averages. Note that
286 during training, the batches are randomly shuffled at each epoch.

287     (2): $G_{loss} = MSE(y_{true}, y_{pred}) + \lambda_{adv} * G_{Adv}(y_{pred}) + IC(y_{true}, y_{pred})$

288     where $IC(y_{true}, y_{pred}) = MSE(Y^{max}_{true}, Y^{max}_{pred}) + MSE(Y^{mean}_{true}, Y^{mean}_{pred})$

289 2.2.2 Discriminator Loss

290     Similar to previous studies (Gulrajani et al., 2017; Harris et al., 2022; Leinonen et al.,
291 2021; Vosper et al., 2023), we use the 1-Wasserstein distance ($D_{adv}$) as an discriminator or
292 critic loss function (yielding what are often known as Wasserstein-GANs), where
293     $D_{adv}(y_{true}, y_{pred}) = \overline{D(y_{true})} - \overline{D(y_{pred})}$.

294     We also use a gradient penalty value of 10 (Gulrajani et al., 2017; Harris et al., 2022;
295 Leinonen et al., 2021; Vosper et al., 2023). As implemented in these studies, we also train the
296 discriminator three times as frequently as the generator. Overall, these refinements to the
297 discriminator have been shown to improve training stability and a reduction of sensitivity to
298 the choice of architecture and hyperparameters in c-GANs (Arjovsky et al., 2017).

299 2.2.3 Adversarial Parameter Selection

300     Our study examines how the solutions produced by GANs to the contribution of the
301 adversarial loss weight ($\lambda_{adv}$). Increasing $\lambda_{adv}$ allows the solutions from the GAN to diverge
302 from the regression baseline as the adversarial loss becomes increasingly important. We
303 explore seven different values of $\lambda_{adv}$: 0.0, 0.0001, 0.00125, 0.0025, 0.005, 0.01 and 0.1.
304 Here, $\lambda_{adv} = 0$ refers to the regression baseline. The range of $\lambda_{adv}$ was chosen to encompass
305 the wide variety of values used in climate downscaling / weather forecasting literature.

306 2.3 Algorithm Architectures

307     In this study, we train two types of emulators: a regression baseline in which $\lambda_{adv} =$
308 0.0 and a residual GAN (Figure 2b). For the residual GAN, we test two different loss
309 function configurations: with (Equation 2) and without an additional intensity constraint
310 (Equation 1).

311 2.3.1 Regression Baseline

312       The regression baseline is based on the widely used U-Net deep learning model
313    (Ronneberger et al., 2015), as illustrated in Figure 2c. The U-Net architecture consists of a
314    contracting path, extracting information from the input predictor variables into a lower
315    dimensional latent space. The expansive path involves reconstructing the high-resolution
316    output (precipitation) from the latent space. The U-Net regression model consists of two input
317    data streams: normalized high-resolution elevation data (12km) from CCAM and the large-
318    scale prognostic predictor variables (1.5°). Our model uses residual convolutional layers (or
319    residual blocks) with batch normalization, which have shown better performance than
320    traditional convolutional layers and help address instability issues in deep-learning models
321    (Rampal, et al., 2024; Sun et al., 2024). Following several residual convolutional blocks and
322    pooling layers, the two input streams are concatenated and mixed to form the latent space of
323    the model. Then, there are a series of upsampling (increasing the spatial resolution) and
324    residual convolutional blocks until the output reaches the desired shape. Additionally, we
325    repeated our experiments with and without batch normalization within our residual blocks,
326    which had a minimal impact on our results.

327    2.3.2 Residual GAN

328       The residual GAN is trained to predict residuals ($r = \widetilde{y_{\lambda_{adv}=0}} - y_{true}$) between a
329    regression baseline ($\lambda_{adv} = 0$) and the ground truth CCAM, as illustrated in Figure 2b. This
330    residual methodology adapted from Mardani et al. (2023), who employed a similar approach
331    in training a different type of generative model for downscaling called diffusion models. The
332    regression baseline learns the expectation of all possible outcomes (the predictable
333    component) from the RCM simulations, which tend to be smooth in both space and time
334    (large-scale precipitation structures). This allows the residual GAN to focus on generating
335    plausible fluctuations around this expectation, which include high-frequency variations and
336    potentially some larger-scale contributions. The architecture of the generator in the residual
337    GAN is nearly identical to the regression baseline, with two additional predictors: high-
338    resolution prediction of precipitation from the regression baseline ($y_{\lambda_{adv}=0}$), and a stochastic
339    noise vector, as inputs (see Figure 2c). Both the regression baseline and the residual GAN
340    have approximately 3.5 million trainable parameters.

341       The discriminator or critic evaluates the perceptual realism of the residuals (either ground
342    truth or predictors from the residual-GAN) conditioned on the regression baseline
343    precipitation predictions ($y_{\lambda_{adv}=0}$), topography, and large-scale meteorological predictors
344    ($x$), as shown in Figure 1a. The discriminator architecture features two input data streams
345    analogous to the generator architecture: one for low-resolution fields with four convolutional
346    layers, and another for high-resolution fields consisting of five convolutional layers. The two
347    input data streams re subsequently concatenated in lower layers of the network. Both input
348    data streams to the discriminator use strided convolutional layers for dimensionality
349    reduction. To reduce model complexity and computational cost we excluded residual blocks
350    from the discriminator architecture, which had negligible impact on our results (not shown).

351    In the discriminator and generator architectures, we use Leaky Rectified Linear Unit
352    (ReLU) activation functions in all layers, as implemented in Leinonen et al. (2021). ReLU
353    activation functions have been suggested to improve stability in training both the generator
354    and discriminator. For the final output activation of the residual GAN, we use the LeakyRelu
355    function ($r_\beta(x)$).

356
$$r_\beta(x) = \begin{Bmatrix} x & x > 0.0 \\ 0.5\,x & x \le 0.0 \end{Bmatrix}.$$

357    In addition to LeakyRelu, we experimented with other output activation functions, such as
358    a modified hyperbolic tangent (Tanh) function, which had a similar skill across all evaluation
359    metrics used in this study.

360    Similar to previous studies (Gulrajani et al., 2017; Leinonen et al., 2020), both the
361    generator and discriminator are trained with an initial learning rate of 2 x $10^{-4}$, and a batch
362    size of 32. The regression baseline is trained with a learning rate of 7 x $10^{-4}$. We also
363    explored smaller learning rates (i.e. $1\ x\ 10^{-6}$), which overall produced similar results but
364    increased algorithm training times (not shown).  To control overfitting and improve stability
365    during training, we use learning decay for training the generator and the regression baseline
366    with decay rates of 0.9945 and 0.989 per 1000 iterations, respectively. Each model was
367    trained for 240 epochs, which equates to approximately 48 hours of training on a single
368    NVIDIA A100 GPU with 80GB RAM. Additionally, learning rate decay also stabilizes GAN
369    performance across different epochs (i.e. similar results were obtained from using 200 epochs
370    instead of 240), addressing fluctuations in performance reported in prior studies. (Harris et
371    al., 2022).

372    Predictions from the residual GAN are added to the regression baseline and inverse
373    transformed ($pr = \exp(Y_{\lambda_{adv}=0..0} + Y_{\lambda_{adv}})$-0.001) to produce daily precipitation fields. Each
374    experiment was repeated three times with a different random seed to ensure the consistency
375    of results, and a separate regression baseline (U-Net) was trained with and without the
376    intensity constraint. Generating a single simulation (one member) of 20-year daily
377    precipitation (7300-time steps) record takes approximately 20 seconds on an A100 GPU.

378    2.4 Evaluation Metrics

379    Analogous to many RCM historical evaluation studies (i.e. Chapman et al., 2023; Di
380    Virgilio et al., 2019, 2020; Isphording et al., 2023), we use three common climatological
381    evaluation metrics to assess the out-of-sample performance over 20 years from 1986-2005.
382    The first metric assesses the ability to capture seasonal averages in precipitation. This
383    assessment for emulators is particularly important for New Zealand, where significant shifts
384    in large-scale circulations, such as the subtropical and polar jet, occur between summer and
385    winter and affect seasonal precipitation. Here, we use the summer and winter periods for
386    evaluation: December-February (DJF) and June – August (JJA), respectively. We also use
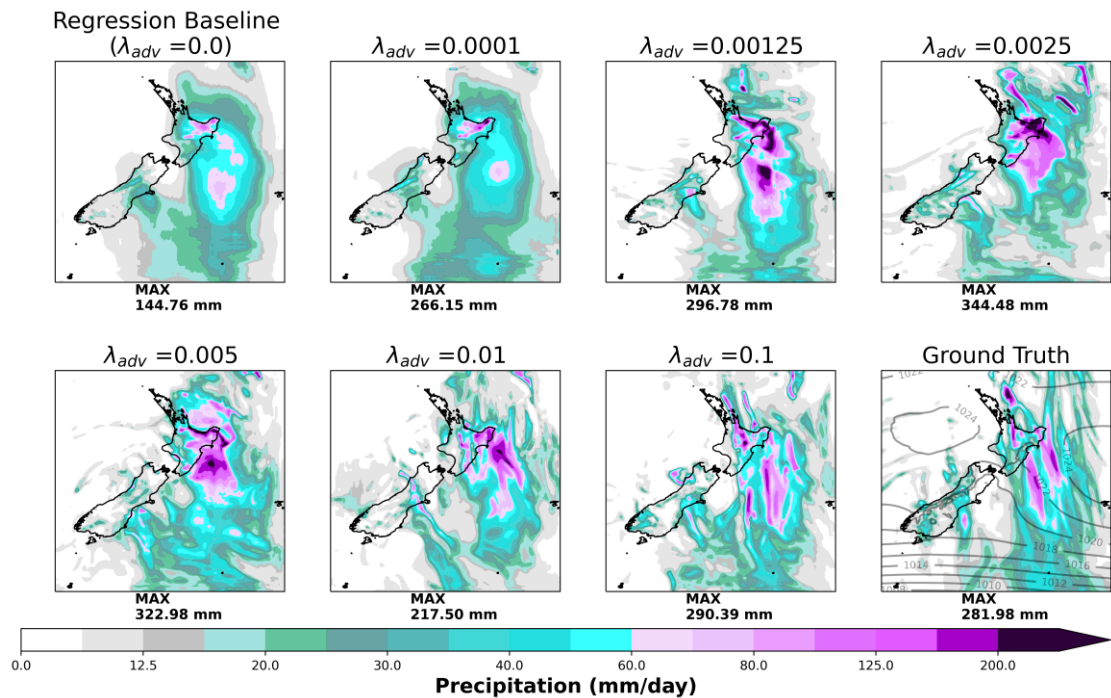387    two other ETCCDI metrics that assess the performance of our emulator on capturing the

388   climatology of extreme events (Isphording et al., 2023; Rampal et al., 2024; Zhang et al.,
389   2011) the wettest day of the year (RX1Day) and the average number of consecutive dry days
390   (CDD) per year.

391   **3. Results**

392   3.1 The Adversarial Effect on Local-scale Extremes

393   3.1.1 Case Studies of Extreme Events

394        To understand how the $\lambda_{adv}$ affects the ability to resolve mesoscale structures and
395   precipitation intensity; we present a case study of an extreme precipitation event simulated in
396   EC-Earth3 over the New Zealand region.  The emulator's predictions of precipitation across
397   all $\lambda_{adv}$ values (including the regression baseline) are spatially aligned with CCAM's
398   precipitation patterns and associated low-pressure centers, as depicted in Figure 3. Overall,
399   this demonstrates the emulator's proficiency in learning the effects of mesoscale circulation
400   on extreme rainfall. This result is also consistent without the intensity constraint (Figure S1)
401   and amongst other case studies (Figure S2-S3).

402



403
404   **Figure 3:** Example of daily precipitation predictions from GAN with the intensity constraint
405   for a simulated extreme event from EC-Earth3 (2002-02-27), relative to the ground truth
406   (CCAM downscaling EC-Earth3). The maximum precipitation intensity across the domain is
407   shown in the text below the plot. The contours show CCAM's Mean Sea Level Pressure
408   (MSLP) patterns for the same event.

409   The regression baseline ($\lambda_{adv} = 0$) significantly underestimates the maximum
410   precipitation intensity and overly smooths mesoscale precipitation structures for the event
411   depicted in Figure 3 relative to ground truth CCAM (also shown in Supplementary Figure S3-
412   S4). However, when $\lambda_{adv} \geq 0.00125$, the emulator can better resolve mesoscale structures
413   and more accurately estimate the maximum precipitation rates across the domain. When the
414   intensity constraint is not used, there are instances where the maximum precipitation intensity
415   is significantly overestimated. Most notably, when $\lambda_{adv} = 0.01$ or $0.1$, the intensity is
416   overestimated by over 200% (see Supplementary Figure S1).

417   3.1.2 Precipitation Distribution

418   To quantify performance more generally, we examine the distribution of precipitation
419   across the entire New Zealand region (including land and ocean) via a one-dimensional
420   histogram of precipitation for all grid points and daily timesteps as shown for both loss
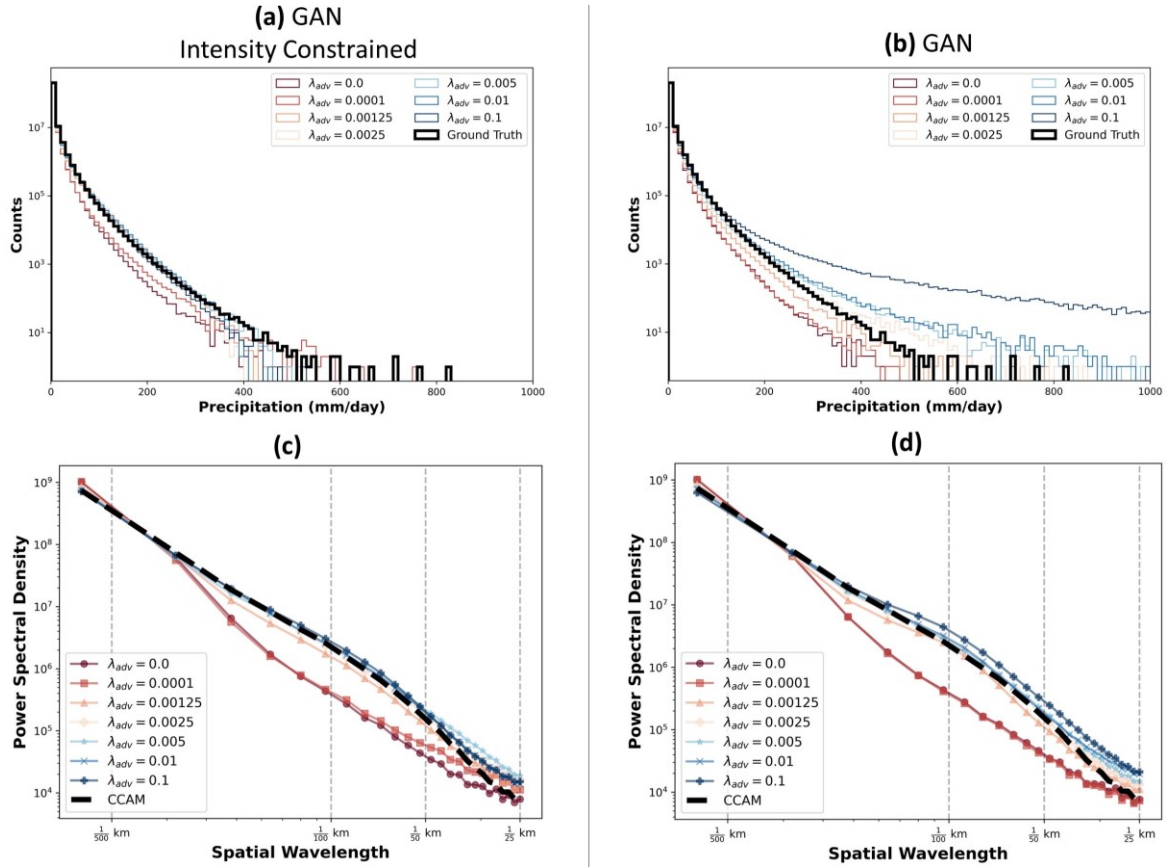421   function configurations (Figure 4).

422   The regression-baseline captures the mean of the precipitation distribution relatively well
423   (the lowest intensity histogram bins) but underestimates the frequency of the most extreme
424   events (i.e. >200mm/day), as shown in Figure 4a-4b. GANs do not always outperform
425   regression models in capturing the precipitation distribution. Rather, their performance
426   depends heavily on the specific loss function configuration (with or without intensity
427   constraints) and the weighting of the adversarial loss ($\lambda_{adv}$).

428   Overall, varying $\lambda_{adv}$ has a minimal effect on the precipitation distribution when the
429   intensity constraint is used (Equation 2; Figure 4a). For nearly all values of $\lambda_{adv}$ the
430   precipitation distribution closely matches CCAM's – albeit slightly underestimating the most
431   extreme precipitation events (>500mm), as illustrated in Figure 4a. In contrast, when no
432   intensity constraint is used (Equation 1), varying $\lambda_{adv}$ has a strong effect on the precipitation
433   distribution (Figure 4b). Here, the regression baseline and $\lambda_{adv} = 0.0001$ case, both
434   underestimate precipitation frequency at all intensities relative to CCAM, whereas when
435   $\lambda_{adv} = 0.1$ there is a significant overestimation of precipitation frequency across all
436   intensities, including a maximum of over 1,000,000 mm/day. Unphysically large precipitation
437   values have also been reported in previous studies (Harris et al., 2022; Vosper et al., 2023).
438   Further evaluation using Quantile-Quantile (Q-Q) plots is shown in Supplementary Figure
439   S4-S6.

440   3.1.3 Mesoscale Variability

441   To evaluate the emulator's skill in resolving finer scale aspects of precipitation, we
442   computed the Power Spectral Density (PSD) on predictions from the 200 rainiest days on
443   average across the domain (although we obtain similar results using all days). The PSD is
444   computed on each day's two-dimensional field of precipitation, and then averaged across all
445   days. Here, the PSD is the integrated Fourier Transform as a function of the spatial
446   wavelength ($K = \sqrt{(k_x^2 + k_y^2)}$, where $k_x$ and $k_y$ are the wavelengths in the $x$ and $y$
447   directions, respectively. We normalized each day's precipitation so that that the PSD receives
448   equal weight from all included days.

449    Our results indicate that for very small values of $\lambda_{adv}$ ($\leq$0.00125) including the
450    regression baseline, variability at small spatial is underestimated regardless of the intensity
451    constraint. Here, GANs do not fully resolve mesoscale structures in a similar capacity to
452    CCAM as shown in Figure 4(c-d). For all other values of $\lambda_{adv}$ The emulator's PSD closely
453    follows CCAM for both loss functions. However, there is one major exception when $\lambda_{adv}$ =
454    0.1, and variability is overestimated across all spatial wavelengths, leading to an exaggerated
455    representation of large-scale and mesoscale variability. It is important to note that at very
456    small spatial scales (~1/25km), there is generally good agreement across all $\lambda_{adv}$, including
457    the regression baseline. This agreement is primarily attributed to incorporating topography as
458    a predictor variable (not shown), enabling the algorithm to account for the influence of
459    orographic precipitation (Bailie et al., 2024). We thus conclude that GANs can relatively
460    robustly capture the range of spatial scales, and that this is not much affected by the intensity
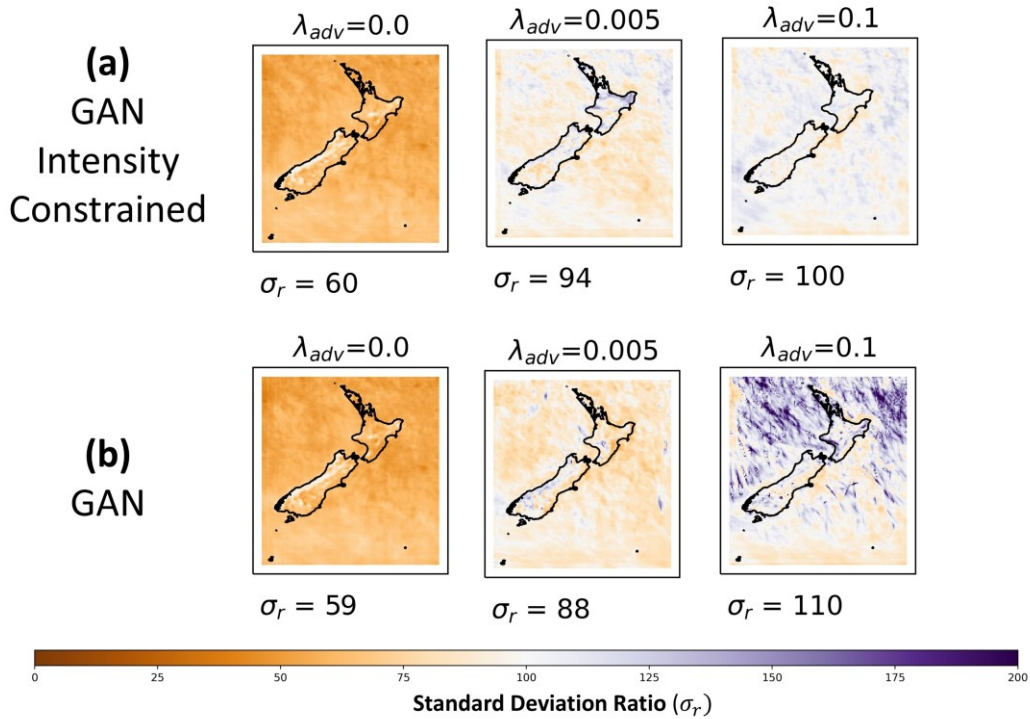461    constraint.



462

463    **Figure 4**: The precipitation distribution as a function of $\lambda_{adv}$ when the RCM emulator is
464    applied out-of-sample to EC-Earth3. (a) the histogram with the intensity constraint, and (b)
465    without. Here, precipitation counts in the histogram are aggregated across all locations over
466    the domain. The black curve highlights ground truth CCAM.

467    3.2 Temporal Variability

468       To evaluate performance on temporal variability (un-normalized precipitation), we
469    compute the ratio ($\sigma_r$) in temporal standard deviation between each emulator's precipitation
470    field ($\sigma_e$) and ground truth CCAM ($\sigma_{CCAM}$), where $\sigma_r = \frac{100 * \sigma_e}{\sigma_{CCAM}}$. Ratios less than 100%
471    underestimate temporal variability, while values greater than 100% overestimate variability.
472    This ratio is computed for each grid cell. Precipitation values exceeding 2000mm/day were
473    excluded when computing standard deviation, which removes the contribution of
474    unphysically large precipitation values, but this only affects the case when $\lambda_{adv} = 0.1$
475    without the intensity constraint.

476       The regression baseline and GAN (Figure 5a-b, left panel) substantially underestimates
477    temporal variability by over 40%, regardless of generative loss configuration (Equation 1 or
478    2). However, as $\lambda_{adv}$ increases further the ratio increases, as illustrated in the center panel in
479    Figure 5a-b and Supplementary Figure S7. With the intensity constraint and when $\lambda_{adv}$=0.1,
480    the emulator performs exceptionally well at capturing CCAM's temporal variability with an
481    average ratio of 100%. When no intensity constraint is used, best performance is achieved
482    when $\lambda_{adv}$ =0.01, with an average 97% ratio (Supplementary Figure S7). However, when
483    $\lambda_{adv} = 0.1$ without the intensity constraint the average ratio exceeds 110%, and in several
484    individual grid points it exceeds 800%. Note if values exceed 2000mm, the ratio exceeds
485    1000%. Thus, for capturing temporal variability robustly, training with an intensity constraint
486    appears important. However, without the intensity constraint, when $\lambda_{adv}$ is large ($\geq$ 0.1), the
487    behavior appears unstable, likely due to the overestimation in extreme precipitation (Figure
488    2b) which inflates the temporal standard deviation.



489

490    **Figure 5**: The percentage ratio of RCM emulated to ground truth temporal standard
491    deviation in CCAM for the EC-Earth3 simulation. (a) shows the percentage ratio for the
492    LeakyReLU activation with an intensity constraint applied and (b) without the constraint for
493    three values of $\lambda_{adv}$. The variance ratio is calculated per grid pixel relative to the CCAM
494    ground truth. The text below each Figure shows the average ratio ($\sigma_r$) across the entire
495    domain.

496    3.3 Climate Statistics

497    This section evaluates the skill in emulating climate statistics/metrics and conventional
498    error metrics such as MAE.

499    3.3.1 Seasonal Precipitation

500    Figure 6 shows the 10-member ensemble average out-of-sample emulator performance in
501    representing four key climate metrics: (a) DJF and (b) JJA climatological precipitation, (c)
502    RX1Day, and (d) CDD - each averaged the domain. The regression-baseline and the $\lambda_{adv} =$
503    0.0001 case have the highest MAE across all out-of-sample evaluation metrics. Increasing
504    $\lambda_{adv}$ improves the skill in reproducing the spatial patterns of seasonal precipitation (DJF and
505    JJA) where the lowest MAEs are observed for $\lambda_{adv} \geq 0.01$. The regression baseline has an
506    overall dry bias and increasing $\lambda_{adv}$ better captures seasonal precipitation rates over the New
507    Zealand region, as illustrated in Figure 7a(i-ii). A similar result is also shown for JJA
508    climatological precipitation (Figure 7b(i-ii)), where the improvement is even more notable.

509    3.3.2 Rx1day Climatology

510    For the Rx1Day climatology, the regression baseline and $\lambda_{adv} = 0.0001$ cases again have
511    the highest MAE (Figure 6c) and are generally dry-biased (Figure 8a(i-ii)) for both loss
512    function configurations. The MAE decreases for higher lambda, except at $\lambda_{adv} = 0.1$ with no
513    intensity constraint where there is a sharp increase in MAE (250%) and the RX1Day
514    climatology is significantly overestimated (Figure 8a(ii) (rightmost panel)). On the other
515    hand, the lowest MAE is achieved with this same $\lambda_{adv} = 0.1$ with the intensity constraint.
516    The spatial patterns in the RX1Day climatology for $\lambda_{adv} = 0.1$ also match the ground truth.
517    Overall, the RX1Day climatology performance seems to most robust across $\lambda_{adv}$ values when
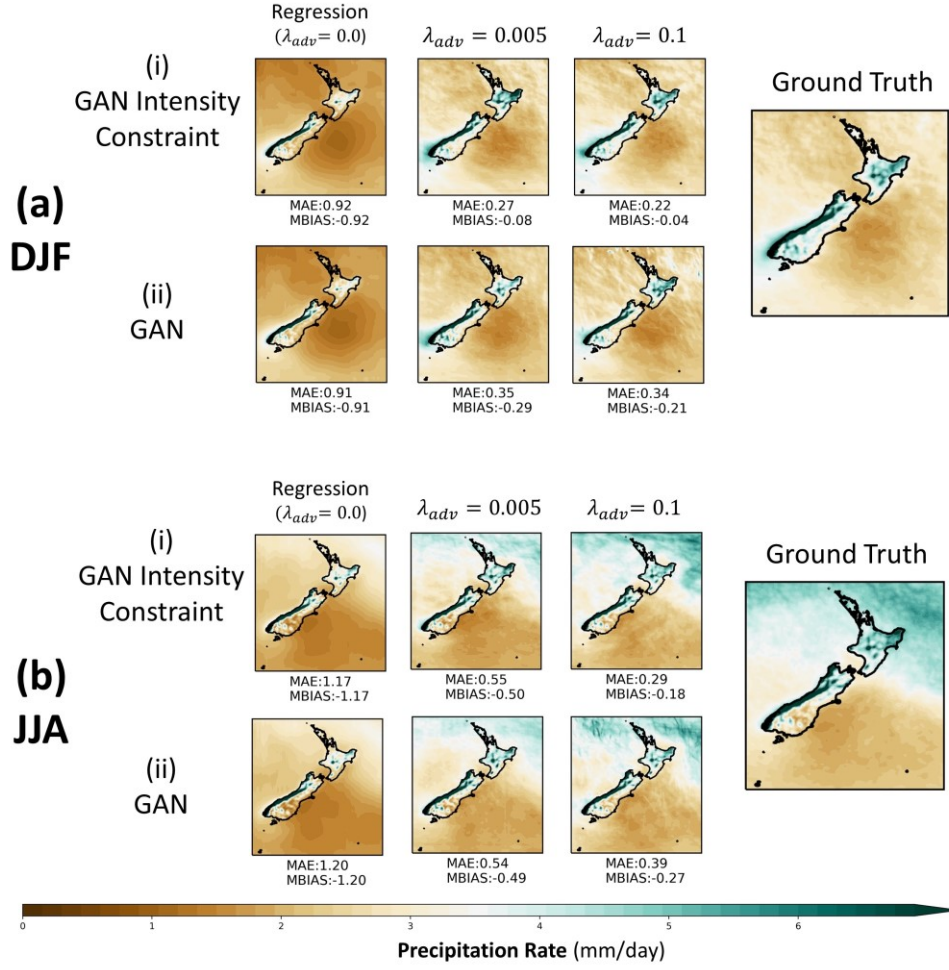518    the intensity constraint is applied.

519

**Figure 6**: The MAE as a function of $\lambda_{adv}$ for the GAN trained with (green) and without (red) the intensity constraint across four key statistics — mean DJF (a) and JJA (b) precipitation, RX1Day (c), and CDD (d) — relative to ground truth CCAM RCM simulation from EC-Earth3. The performance of the regression baseline is shown as the dashed line, both with (green) and without (red) the intensity constraint.

### 3.3.3 Consecutive Dry Days

The results for CDD show the same trends as those for Rx1Day, except at $\lambda_{adv}$=0.1 for where the MAE abruptly increases for both loss function configurations. Upon visual inspection in Figure 8b, the MAE increase appears to be due to an overestimation in CDD over the ocean, particularly on the eastern coast of the South Island and the northern region of the North Island of New Zealand. Interestingly, the configuration with the intensity constraint

532    appears to have a larger MAE across all values of $\lambda_{adv}$ compared to the configuration
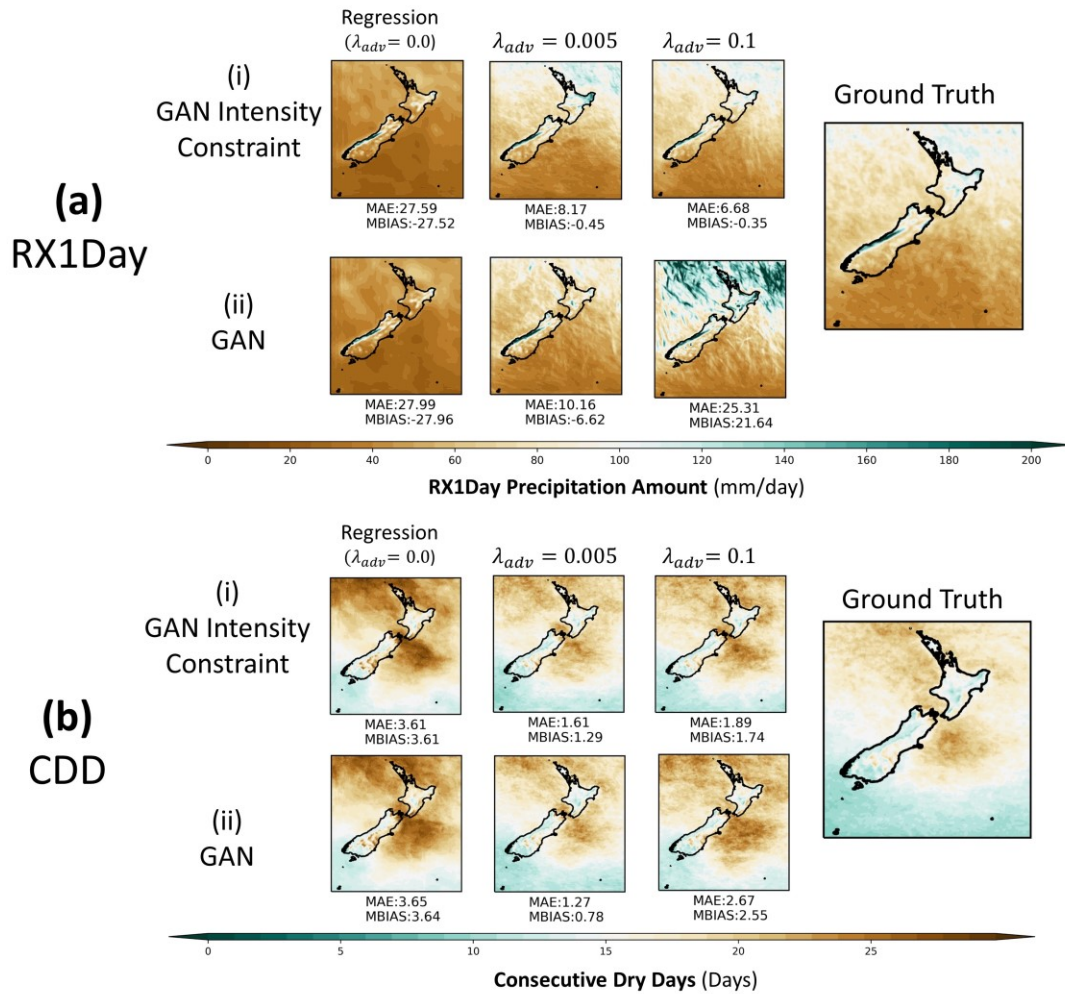533    without it.



**Figure 7**: The performance of the two GAN loss function configurations as a function of $\lambda_{adv}$; with (i) and without the intensity constraint (ii) in generating DJF and JJA climatological precipitation relative to ground truth CCAM RCM simulations (EC-Earth3) for a single ensemble member. The regression baseline is indicated by $\lambda_{adv}= 0.0$. The text for each subplot shows the MAE and the mean bias (MBIAS) relative to ground truth.

540    3.3.4 In-sample Performance

541        It is also important to discuss in-sample performance, that is, evaluating the emulator on
542    the same RCM simulation as it was trained on (ACCESS-CM2) between 1986-2005.
543    Differences between in-sample and out-of-sample performance can shed light on the
544    emulator's ability to generalize further (for example, to other GCMs).  The in-sample
545    performance across the four metrics is generally significantly better (lower error) than on EC-
546    Earth3, particularly for the regression baseline and lower values of $\lambda_{adv}$ ($\leq 0.01$). The higher
547    in-sample performance suggests that the algorithm may have slightly overfitted to the
548    ACCESS-CM2 training distribution despite efforts to prevent it. However, for $\lambda_{adv} =0.1$

20

549 with intensity constraint, in-sample and out-of-sample performances are similar for all
550 metrics except for CDD, as illustrated in Supplementary Figure S9-S10. One potential
551 explanation is that adversarial training mitigates overfitting, allowing the algorithm to learn
552 more generalizable relationships, though further research would be required to test this. We
553 also assessed the out-of-sample emulator performance on the NorESM2-MM GCM (i.e. the
554 model trained on ACCESS-CM2 is applied to NorESM2-MM predictor fields). The results
555 were nearly identical to the EC-Earth3 evaluation, as summarized in Supplementary Figure
556 S8. This result is important as is implies a GAN emulator trained only on one RCM/GCM
557 simulation pair can be broadly applied to historical climates from other GCMs, a finding that
558 differs from the common view of GANs as being unstable.



559

560 **Figure 8**: The performance of the two GAN loss function configurations as a function of
561 $\lambda_{adv}$; with (i) and without the intensity constraint (ii) in generating climatological RX1Day
562 and CDD relative to ground truth CCAM RCM simulations (EC-Earth3) for a single
563 ensemble member. The regression baseline is indicated by $\lambda_{adv}$= 0.0. The text for each
564 subplot shows the MAE and the mean bias (MBIAS) relative to ground truth.

565

566    3.3.5 Summary

567    Overall, when considering all climate statistical metrics, the lowest MAE scores occur
568    when $\lambda_{adv}$ is set between 0.05 and 0.1 with the intensity constraint. Note, that for this range
569    of $\lambda_{adv}$ we also see good performance in accurately capturing precipitation distribution. In
570    comparison, without the intensity constraint, the best performance is generally achieved when
571    $\lambda_{adv}$ is between 0.0025 to 0.01, with the lowest scores notably at $\lambda_{adv} = 0.01$. However, the
572    larger values of $\lambda_{adv}$ within this range (0.005, 0.01) do not accurately capture precipitation
573    distribution (as detailed in Section 3.2), making $\lambda_{adv} = 0.0025$ the only viable option.

574    3.4 Ensemble Statistics

575    Moving beyond standard downscaling metrics, this section assesses the ensemble spread
576    produced by GANs. It aims to determine whether GANs can skillfully generate ensembles
577    that capture the "true" variability of potential outcomes that is essential for uncertainty
578    quantification in a downscaling or weather generation context. This study uses the spread-
579    error relationship (section 3.4.1) and the Continuous Ranked Probability Score (CRPS;
580    section 3.4.2) metrics, which are commonly used for evaluating ensemble weather forecasts
581    (i.e. Doblas-Reyes et al., 2005; Leutbecher & Palmer, 2008; Palmer et al., 2008), and more
582    recently for DL-based weather forecasts (Harris et al., 2022; Kochkov et al., 2024; Price &
583    Rasp, 2022; Ravuri et al., 2021a; Vosper et al., 2023).

584    3.4.1 Spread-error Relationship

585    The spread-error relationship evaluates an ensemble's dispersion – also commonly known
586    as the ensemble's calibration. The spread-error relationship describes a relation between
587    spread of the ensemble about its mean (RMSS) and the error in the ensemble mean (hereon
588    referred to as RMSE) (Doblas-Reyes et al., 2005; Fortin et al., 2014; Leutbecher & Palmer,
589    2008; Palmer et al., 2008). A well-calibrated (statistically perfect) ensemble of infinite size
590    generally has a linear spread-error relationship (black dashed line in Figure 9a-b), meaning
591    that the average distance between the ground truth and the ensemble mean equals the average
592    distance between individual ensemble members and the ensemble mean. The key
593    characteristic of a well-calibrated ensemble is that individual ensemble members are not
594    statistically distinguishable from the ground truth data. This relationship has been widely
595    used in the ensemble weather forecasting (Leutbecher & Palmer, 2008), and has recently
596    examined in a downscaling context (Vosper et al., 2023).
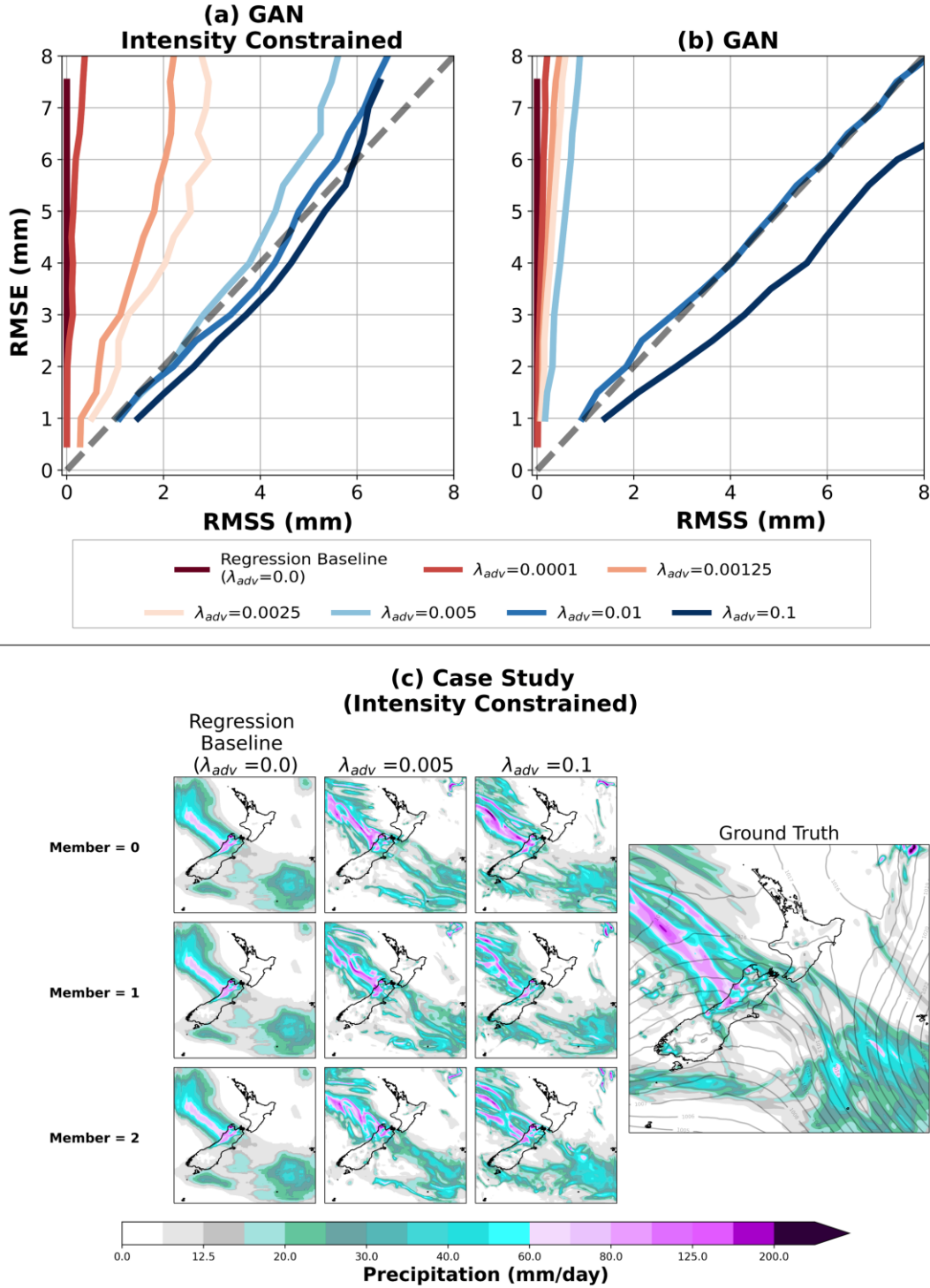
597    The spread-error relationship is computed for a 10-member ensemble spanning the 20-
598    year evaluation period for each $\lambda_{adv}$. Each ensemble member is distinguished by a unique
599    noise vector for the same large-scale predictor variables, as depicted in Figures 2b and c. We
600    also generated a 100-member ensemble spanning one year (i.e. larger ensemble but shorter
601    duration) to understand the impact of ensemble size on the spread-error relationship, which
602    did not alter our findings (not shown). Similar to previous studies (Kochkov et al., 2024;
603    Vosper et al., 2023), to compute spread error curves, we first average RMSS and RMSE

604  values over time, then compute the mean RMSS across all RMSE bins. Due to our smaller
605  ensemble size ($n = 10$), the RMSS and RMSE values are adjusted by the factors 1.11 ($\frac{n}{n-1}$)
606  and 0.9 ($\frac{n}{n+1}$), respectively, as outlined in Vosper et al., (2023) & Leutbecher and Palmer,
607  (2008).

608  In the regression baseline ($\lambda_{adv} = 0.0$) and $\lambda_{adv} = 0.0001$ cases, there is no spread
609  amongst their ensemble members (the RMSS is zero) and thus the slope of the spread-error
610  relationship is infinite (Figure 9a-b). When the intensity constraint is used, increasing $\lambda_{adv}$
611  improves dispersion or calibration, and when $\lambda_{adv}$ is between 0.005 and 0.1 the spread-error
612  curve and its slope are more closely aligned with the perfect ensemble ($y = x$). Visual
613  inspection of an individual case (Figure 9c) likewise shows that for small values of $\lambda_{adv}$ no
614  dispersion is evident among ensemble members, while for $\lambda_{adv}$ greater than 0.005 dispersion
615  becomes more pronounced. Most importantly the dispersion appears perceptually realistic,
616  where each member's precipitation patterns are different but all consistent with large-scale
617  circulation patterns.

618  Conversely, when no intensity constraint is used, for nearly all values of $\lambda_{adv}$ the spread-
619  error curves are primarily under-dispersive or poorly calibrated, as illustrated in Figure 9b.
620  The ensemble is well-calibrated when $\lambda_{adv} = 0.01$ as the spread-error curves are close to a
621  perfect ensemble ($y = x$), but rapidly transitions to being over dispersive for larger
622  $\lambda_{adv} (\geq 0.01)$. Overall, in the absence of the intensity constraint, there is some instability or
623  heightened sensitivity to the spread-error relationship as a function of $\lambda_{adv}$.
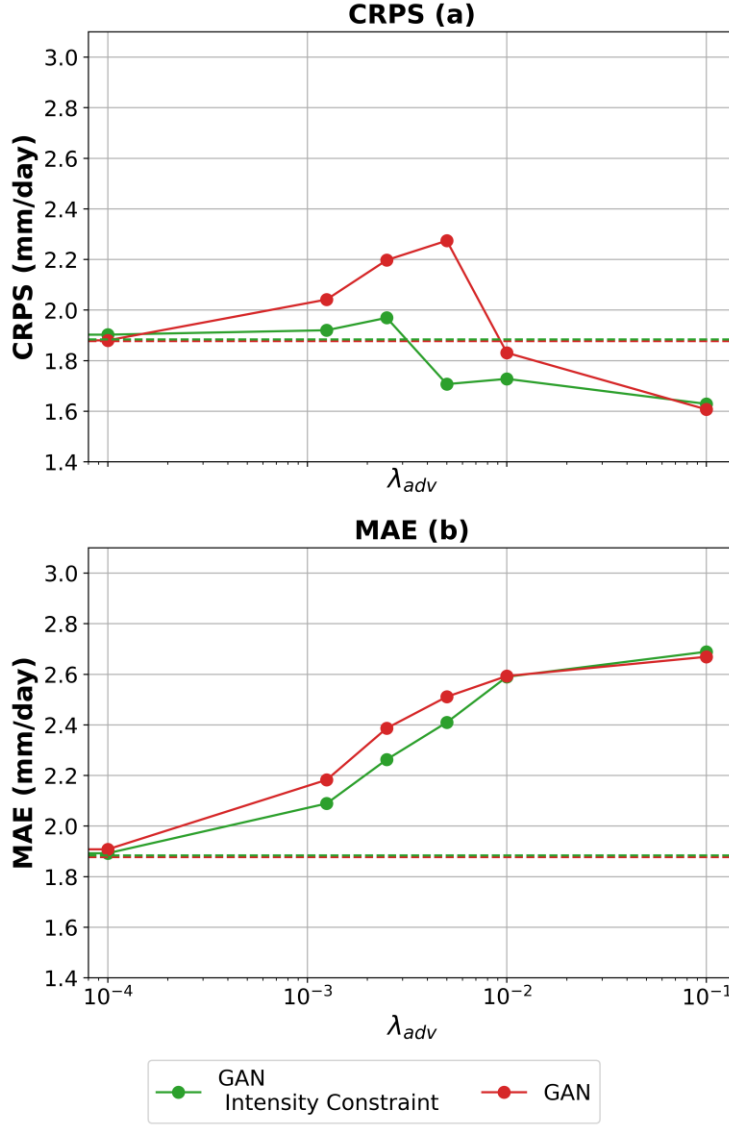
624
625

**Figure 9:** (a) The spread-error relationship as a function of $\lambda_{adv}$ when the intensity constraint is used. (b) The spread-error relationship as a function of $\lambda_{adv}$ when intensity constraint is not used. (c) Examples of solutions from three ensemble members across all values of $\lambda_{adv}$, and its corresponding ground truth CCAM precipitation (b). (c) shows a case study (2004-01-16) from EC-Earth3 to illustrate dispersion across three different ensemble members.

633    3.4.2 Continuous Ranked Probability Score (CRPS)

634    The Continuous Ranked Probability Score (CRPS) is a proper scoring rule that assesses
635    the accuracy of an ensemble in representing the full range of uncertainty within a prediction
636    (Gneiting & Katzfuss, 2014; Gneiting & Raftery, 2007; Hersbach, 2000; Lerch et al., 2017;
637    Matheson & Winkler, 1976). The CRPS measures the distance between the predicted
638    probability distribution and ground truth, but also assesses the ensemble's spread or
639    calibration. The CRPS is often interpreted as a generalization of the MAE (absolute
640    difference between a prediction and ground truth) for probabilistic forecast evaluation. In the
641    case of a deterministic prediction (e.g., regression baselines or a single ensemble member) the
642    CRPS equals the MAE, allowing for a comparison between deterministic and ensemble
643    predictions. One notable advantage of the CRPS is that it is less sensitive to double counting
644    of position (location of where precipitation occurs) and intensity (precipitation amount)
645    errors, which is a widely known limitation of MAE in ensemble forecast evaluation
646    (Hersbach, 2000).

647    Both the MAE and CRPS are computed on an individual prediction basis (per grid cell)
648    and averaged across all timesteps (7300 timesteps), latitudes, and longitudes. Here, the MAE
649    is calculated for each ensemble member and then averaged across all members. To reduce the
650    effect of outlier precipitation values on the computation of the MAE and CRPS, we exclude
651    grid points for a given timestep (and all corresponding members) when at least one ensemble
652    member has a precipitation value exceeding 2000 mm.

653

654

**Figure 10**: The CRPS (a) and MAE (b) as a function of $\lambda_{adv}$ for both loss function configurations on EC-Earth3 relative to ground truth, with (green) and without the intensity constraint (red).

The MAE (Figure 10b) is lowest for the regression baseline ($\lambda_{adv} = 0$) and increases rapidly as a function of $\lambda_{adv}$, where it is 50% greater when $\lambda_{adv} = 0.1$ for both loss function configurations. The regression baseline's lower MAE is expected, as it directly optimizes for the Mean Squared Error (MSE) during training, which aligns closely with the MAE metric. This relationship between MAE and $\lambda_{adv}$ is also somewhat expected, as by design larger $\lambda_{adv}$ values allow for more deviation from the regression baselines (which are optimized for MSE), leading to an increased MAE. Several studies have also noted that GANs have a higher MAE scores than regression-based DL algorithms (i.e. J. Wang et al., 2021).

As for the CRPS metric, GANs only outperform the regression baseline at certain values of $\lambda_{adv}$. For instance, the GAN's CRPS is larger than the regression baseline when

669    $\lambda_{adv}$ ≤0.0025 with the intensity constraint, and $\lambda_{adv}$ ≤0.005 without it (Figure 10a). On the

670    other hand, the GAN's CRPS is lower than the regression baseline when $\lambda_{adv}$ ≥0.005 with

671    the intensity constraint, and $\lambda_{adv}$ ≥0.01 without it, where in both loss configurations the

672    lowest CRPS is achieved when $\lambda_{adv} = 0.1$. Note, the CRPS scores with the intensity

673    constraint are typically lower than those without across all $\lambda_{adv}$ values, except at $\lambda_{adv} = 0.1$,

674    where the scores are similar.

675    3.4.3 Summary

676    In summary of our results from both CRPS and spread-error metrics, we find that smaller

677    values of $\lambda_{adv}$ (<0.05) tend to generate under-dispersive (poorly calibrated) ensembles with

678    larger or similar CRPS scores to the regression baseline for both loss configurations. When

679    $\lambda_{adv}$ is between 0.005 and 0.1, GANs trained with an intensity constraint generate dispersive

680    (well-calibrated) ensembles with lower CRPS scores than the regression-baseline. However,

681    GANs trained without an intensity constraint produce well-calibrated ensembles and lower

682    CRPS scores than the regression-baseline only when $\lambda_{adv} = 0.01$ and become over-

683    dispersive for larger $\lambda_{adv}$. Additionally, GANs trained with an intensity constraint are more

684    dispersive and have lower CRPS scores across all $\lambda_{adv}$ than those trained without it,

685    suggesting that the intensity constraint improves robustness beyond its intended design.

686    **4 Discussion**

687    4.1 The Importance of Constraints

688    In general, smaller values of $\lambda_{adv}$ (typically below 0.005) are common amongst

689    downscaling (climate and weather) and super-resolution studies when no intensity constraint

690    is used (Harris et al., 2022; Ledig et al., 2017; Leinonen et al., 2021; Vosper et al., 2023; X.

691    Wang et al., 2018). Our optimal range of $\lambda_{adv}$ ($0.00125 \leq \lambda_{adv} < 0.005$) without the

692    intensity constraint aligns with these values. Our findings show that this range of $\lambda_{adv}$

693    balances good performance in capturing rainfall mean variations (Figure 7 & 8) and

694    distribution e.g. for extreme events (Figure 4). While larger values of $\lambda_{adv}$ (>0.005) perform

695    better on the former, they drastically overestimate extreme precipitation events

696    (>200mm/day). As $\lambda_{adv}$ becomes too small, GAN performance converges towards that of

697    regression-based DL algorithm, which generally performs poorly across all metrics.

698    The agreement between our optimal $\lambda_{adv}$ range (without the intensity constraint) and

699    other studies is promising, but one should be cautious about this range of $\lambda_{adv}$ as they have

700    not been properly assessed in literature for their performance in climate settings (i.e. how

701    well they capture climate statistics). Our results demonstrate that GANs within this $\lambda_{adv}$

702    range produce under-dispersive ensembles (Figure 9b), limiting their usefulness for

703    uncertainty quantification (see also section 4.2). Additionally, their errors on climate

704    statistical metrics are much higher than larger $\lambda_{adv}$ values (Figure 6-9). More broadly, there

705    are other challenges with training without the intensity constraint, such as the case for large

706    $\lambda_{adv}$, where precipitation extremes significantly overestimated. This raises concerns about

707    GAN robustness (without the intensity constraint) under climate change, due to potential
708    unreliability in simulating extreme events.

709        Our study also shows that an intensity constraint in the loss function improves the
710    robustness of GANs and allows training with large $\lambda_{adv}$. Larger values of $\lambda_{adv}$ ($\geq$0.005)
711    generate more dispersive ensembles in the results and improve accuracy in climate statistical
712    metrics compared to smaller $\lambda_{adv}$ values (below 0.005), all while accurately representing the
713    precipitation distribution. Several studies have also incorporated intensity constraints into
714    GAN loss functions. These studies have used significantly larger values of $\lambda_{adv}$ (e.g., Ravuri
715    et al., 2021: $\lambda_{adv} = 0.05$; Price & Rasp, 2022: $\lambda_{adv} = 0.1$). They reported substantial
716    improvement over regression-based DL algorithms, focusing primarily on metrics such as
717    CRPS and performance on extreme events in a weather forecasting context. However, they
718    did not directly compare their results to those without an intensity constraint.

719    4.2 Stochastic Weather Generation with GANs

720        The application of GANs as a stochastic weather generator remains both under-utilized
721    and under-evaluated in climate science. Stochastic weather generators can generate large
722    ensembles (or sequences) of climate fields (i.e. Ailliot et al., 2015; Benoit et al., 2018; Furrer
723    & Katz, 2008; Steinschneider et al., 2019; Verdin et al., 2018), which can be used to estimate
724    the likelihood of a certain extreme event occurring (i.e. average recurrence interval), thereby
725    offering valuable insights for disciplines such as catastrophe modeling. Recently, several
726    studies have used generative DL algorithms (including GANs) in a similar capacity to
727    stochastic weather generators (Boulaguiem et al., 2022; Brochet et al., 2023; Peard & Hall,
728    2023; Sha et al., 2024). GANs may have certain benefits over stochastic weather generators,
729    such as their ability to learn complex spatio-temporal relationships (Sha et al., 2024). This
730    may help them better simulate extreme phenomena like cyclones and atmospheric rivers,
731    though further comparison with traditional stochastic weather generators is needed. Although
732    GANs show promise in this context, their success is ultimately hinged on their ability to
733    generate sufficiently dispersive ensembles (that capture the true variability of all possible
734    outcomes).

735        Several studies, which have mainly focused on weather forecasting have assessed the
736    calibration (dispersion) of GAN-generated ensembles (i.e. Harris et al., 2022; Price & Rasp,
737    2022; Ravuri et al., 2021b; Vosper et al., 2023). Collectively, these studies suggest that
738    GANs can produce well-calibrated ensemble predictions across a large range of $\lambda_{adv}$ (0.001-
739    0.1), and thus for this purpose we cannot expect a single value of $\lambda_{adv}$ to work across all
740    problems and regions. However, our study introduces key insights into using GANs for
741    uncertainty quantification not previously detailed in literature. Firstly, our study emphasizes
742    the importance of exploring the $\lambda_{adv}$ parameter, due to its significant impact on ensemble
743    dispersion (calibration). Secondly, incorporating constraints (i.e. intensity constraints) to the
744    loss function can not only improve ensemble dispersion across $\lambda_{adv}$, but also yields more
745    robust performance compared to traditional GAN implementations.

746   4.3 Limitations

747       Our study has only focused on historical training and evaluation. Further research should
748   focus on considering how well GANs extrapolating to future scenarios, especially in warmer
749   climates, may require broader training across both historical and future periods, as well as
750   multiple RCM simulations (Bano-Medina et al., 2023; Chadwick et al., 2011; Doury et al.,
751   2022; Holden et al., 2015). The choice of training simulation may impact the algorithm's
752   ability to extrapolate to future climates across multiple GCMs (Bano-Medina et al., 2023;
753   Rampal et al., 2024). For example, warmer RCM simulations (i.e. with a higher equilibrium
754   climate sensitivity), may offer greater diversity in extreme events and climate variability.
755   When assessing emulator performance in future climates, one should consider evaluating the
756   emulator's ability to reproduce the RCM's climate change signals and non-stationary changes
757   like trends in extreme precipitation. Examples of these evaluation strategies are provided in
758   Bano-Medina et al. (2023), Rampal et al. (2024), and Doury et al. (2022).

759       Further development of statistical constraints incorporated into the loss function should
760   also be considered. In our case, the intensity constraint configuration performs exceptionally
761   well across various evaluation metrics but appears to have a lower skill for CDD. A potential
762   explanation for this lower skill could relate to the concept of metric transitivity (Abramowitz
763   et al., 2019), in which optimizing the algorithm to perform well on specific metrics (i.e.
764   intensity) means it performs slightly worse on other metrics which depend more on the
765   temporal aspects of precipitation (i.e. CDD). In future work, applying additional constraints
766   tailored for CDD could potentially improve the skill for this metric.

767       Our research has only focused on downscaling within the New Zealand domain, and thus
768   it is unclear how generalizable our intensity constraint modification and optimal $\lambda_{adv}$ value is
769   across different domains, especially those larger in size (i.e. CORDEX domains) and
770   involving various variables. Although not detailed here, preliminary evidence, which will be
771   explored in a subsequent study, indicate that this optimal $\lambda_{adv}$ range with the intensity
772   constraint successfully downscales precipitation in different regions and for other variables
773   (i.e. temperature), though further testing is needed to confirm its robustness.

774       Lastly, it is important to highlight common criticisms of GANs, such as "mode collapse"
775   or a lack of diversity in generated samples (Che et al., 2017; Dubiński et al., 2023; Mao et al.,
776   2019; Salimans et al., 2016; Srivastava et al., 2017). While we acknowledge such criticisms,
777   our study suggests that GANs can be very effective in downscaling with careful training
778   strategies (as detailed in Section 2). While diffusion models are an emerging type of
779   downscaling technique (and stable) (Addison et al., 2022, Leinonen et al., 2023), they are
780   significantly slower that GANs in training and inference time.

781

782 **5 Conclusion**

783     This study demonstrates that conditional Generative Adversarial Networks (GANs) can
784 improve upon several of the limitations of regression-based deep learning (DL) algorithms
785 for downscaling in a historical climate setting. We also highlighted the broader potential of
786 GANs for stochastic weather generation, noting their skill in generating ensembles that
787 accurately encompass the full spectrum of possible outcomes.

788     We trained a series of GANs on a single historical RCM simulation (ACCESS-CM2) and
789 tested their performance on two completely unseen GCMs (EC-Earth3 and NorESM2-MM)
790 to assess their generalization potential for downscaling across different GCM/RCM
791 combinations.

792     The best-performing GANs examined here outperformed regression-based DL algorithms
793 across various metrics relative to ground truth RCM simulations. While previous studies have
794 found promising results using GANs for a few problems with limited metrics, we measure
795 skill across a wide range of metrics, extending beyond conventional error metrics (e.g. mean
796 absolute error) used in many DL studies. Crucially we examine climate statistical metrics
797 (climatology of seasonal precipitation, wettest day of the year and length of the longest dry
798 spell), temporal variability, the precipitation intensity distribution, and ensemble calibration
799 (dispersion), which are much more relevant for climate studies.

800     We investigated how the hyperparameter $\lambda_{adv}$ (weighting of the adversarial loss) impacts
801 the skill of the GAN, which has been largely unexplored in literature. GAN performance was
802 strongly dependent on $\lambda_{adv}$ in standard implementations, such that $\lambda_{adv}$ cannot be too big or
803 too small (i.e. there is no convergence to good behavior in either limit), and selecting an
804 optimal value requires trade-offs. Larger values of $\lambda_{adv}$ ($\geq 0.01$) would perform well across
805 most metrics but can drastically overestimate precipitation intensity which diverges
806 monotonically as $\lambda_{adv}$ increases. Smaller values would perform well on precipitation
807 intensity but less well for climate statistics and generating well-calibrated (dispersive)
808 ensembles needed to assess uncertainty. In this situation we cannot be confident that a value
809 of $\lambda_{adv}$ tuned to work well in the historical climate or situation would generalize to an
810 unobserved climate scenario.

811     However, by incorporating a simple intensity constraint into the loss function of the
812 GAN, we significantly improved the robustness of GAN performance, thereby requiring
813 fewer trade-offs when selecting an optimal $\lambda_{adv}$. The intensity constraint allows for the
814 selection of larger $\lambda_{adv}$ ($\geq 0.005$), hence a stronger weighting of the adversarial loss, which
815 performs well across all evaluation metrics, including precipitation intensity, but can also
816 generate well-calibrated (dispersive) ensembles required for stochastic weather generation.

817     While we found an optimal range of $\lambda_{adv}$ between 0.005 and 0.1, we strongly recommend
818 thoroughly exploring and testing this hyperparameter during training in different contexts
819 (i.e. across different regions and variables). We also emphasize the importance of statistical

820 constraints for tailored GAN design and the use of climate-relevant evaluation metrics.
821 Further work will be required to see whether the range of $\lambda_{adv}$ found to succeed here indeed
822 generalizes to both future climates and across multiple variables.

823

830 *Data Availability Statement*

831 The code and datasets supporting this study are accessible to the public. The code can be
832 found on GitHub (https://github.com/nram812/A-Robust-Generative-Adversarial-Network-
833 Approach-for-Climate-Downscaling), and the training and validation datasets are available on
834 Zenodo (https://doi.org/10.5281/zenodo.10889046).

835

836

837 **References**

838 Abu-Srhan, A., Abushariah, M. A. M., & Al-Kadi, O. S. (2022). The effect of loss function

839       on conditional generative adversarial networks. *Journal of King Saud University* -

840       *Computer and Information Sciences*, *34*(9), 6977–6988.

841       https://doi.org/10.1016/j.jksuci.2022.02.018

842 Ackerley, D., Dean, S., Sood, A., & Mullan, A. B. (2012). Regional climate modelling in

843       New Zealand: Comparison to gridded and satellite observations. *Weather and*

844       *Climate*, *32*(1), 3–22. https://doi.org/10.2307/26169722

845 Ailliot, P., Allard, D., Monbet, V., & Naveau, P. (2015). Stochastic weather generators: An

846       overview of weather type models. *Journal de La Société Française de Statistique*,

847       *156*(1), 101–113.

848     Annau, N. J., Cannon, A. J., & Monahan, A. H. (2023). Algorithmic Hallucinations of Near-

849         Surface Winds: Statistical Downscaling with Generative Adversarial Networks to

850         Convection-Permitting Scales. *Artificial Intelligence for the Earth Systems*, *2*(4).

851         https://doi.org/10.1175/AIES-D-23-0015.1

852     Arjovsky, M., Chintala, S., & Bottou, L. (2017). *Wasserstein GAN* (arXiv:1701.07875).

853         arXiv. https://doi.org/10.48550/arXiv.1701.07875

854     Babaousmail, H., Hou, R., Gnitou, G. T., & Ayugi, B. (2021). Novel statistical downscaling

855         emulator for precipitation projections using deep Convolutional Autoencoder over

856         Northern Africa. *Journal of Atmospheric and Solar-Terrestrial Physics*, *218*, 105614.

857         https://doi.org/10.1016/j.jastp.2021.105614

858     Bailie, T., Koh, Y. S., Rampal, N., & Gibson, P. B. (2024). Quantile-Regression-Ensemble:

859         A Deep Learning Algorithm for Downscaling Extreme Precipitation. *Proceedings of*

860         *the AAAI Conference on Artificial Intelligence*, *38*(20), Article 20.

861         https://doi.org/10.1609/aaai.v38i20.30193

862     Bano-Medina, J., Iturbide, M., Fernandez, J., & Gutierrez, J. M. (2023). *Transferability and*

863         *explainability of deep learning emulators for regional climate model projections:*

864         *Perspectives for future applications* (arXiv:2311.03378). arXiv.

865         http://arxiv.org/abs/2311.03378

866     Baño-Medina, J., Manzanas, R., & Gutiérrez, J. M. (2020). Configuration and

867         intercomparison of deep learning neural models for statistical downscaling.

868         *Geoscientific Model Development*, *13*(4), 2109–2124. https://doi.org/10.5194/gmd-

869         13-2109-2020

870     Bartók, B., Wild, M., Folini, D., Lüthi, D., Kotlarski, S., Schär, C., Vautard, R., Jerez, S., &

871         Imecs, Z. (2017). Projected changes in surface solar radiation in CMIP5 global

872      climate models and in EURO-CORDEX regional climate models for Europe. *Climate*

873      *Dynamics*, *49*(7), 2665–2683. https://doi.org/10.1007/s00382-016-3471-2

874    Benestad, R. E. (2004). Empirical-statistical downscaling in climate modeling. *Eos,*

875      *Transactions American Geophysical Union*, *85*(42), 417–422.

876      https://doi.org/10.1029/2004EO420002

877    Benestad, R. E. (2010). Downscaling precipitation extremes: Correction of analog models

878      through PDF predictions. *Theoretical and Applied Climatology*, *100*(1–2), 1–21.

879      https://doi.org/10.1007/s00704-009-0158-1

880    Benoit, L., Allard, D., & Mariethoz, G. (2018). Stochastic Rainfall Modeling at Sub-

881      kilometer Scale. *Water Resources Research*, *54*(6), 4108–4130.

882      https://doi.org/10.1029/2018WR022817

883    Boé, J., Mass, A., & Deman, J. (2023). A simple hybrid statistical–dynamical downscaling

884      method for emulating regional climate models over Western Europe. Evaluation,

885      application, and role of added value? *Climate Dynamics*, *61*(1), 271–294.

886      https://doi.org/10.1007/s00382-022-06552-2

887    Boé, J., Somot, S., Corre, L., & Nabat, P. (2020). Large discrepancies in summer climate

888      change over Europe as projected by global and regional climate models: Causes and

889      consequences. *Climate Dynamics*, *54*(5), 2981–3002. https://doi.org/10.1007/s00382-

890      020-05153-1

891    Boulaguiem, Y., Zscheischler, J., Vignotto, E., Van Der Wiel, K., & Engelke, S. (2022).

892      Modeling and simulating spatial extremes by combining extreme value theory with

893      generative adversarial networks. *Environmental Data Science*, *1*, e5.

894      https://doi.org/10.1017/eds.2022.4

895     Brochet, C., Raynaud, L., Thome, N., Plu, M., & Rambour, C. (2023). Multivariate

896          Emulation of Kilometer-Scale Numerical Weather Predictions with Generative

897          Adversarial Networks: A Proof of Concept. *Artificial Intelligence for the Earth*

898          *Systems*, *2*(4). https://doi.org/10.1175/AIES-D-23-0006.1

899     Chadwick, R., Coppola, E., & Giorgi, F. (2011). An artificial neural network technique for

900          downscaling GCM outputs to RCM spatial scale. *Nonlinear Processes in Geophysics*,

901          *18*(6), 1013–1028. https://doi.org/10.5194/npg-18-1013-2011

902     Chapman, S., Syktus, J., Trancoso, R., Thatcher, M., Toombs, N., Wong, K. K.-H., &

903          Takbash, A. (2023). Evaluation of Dynamically Downscaled CMIP6-CCAM Models

904          Over Australia. *Earth's Future*, *11*(11), e2023EF003548.

905          https://doi.org/10.1029/2023EF003548

906     Che, T., Li, Y., Jacob, A. P., Bengio, Y., & Li, W. (2017). *Mode Regularized Generative*

907          *Adversarial Networks* (arXiv:1612.02136). arXiv.

908          https://doi.org/10.48550/arXiv.1612.02136

909     Deser, C., Phillips, A., Bourdette, V., & Teng, H. (2012). Uncertainty in climate change

910          projections: The role of internal variability. *Climate Dynamics*, *38*(3), 527–546.

911          https://doi.org/10.1007/s00382-010-0977-x

912     Deser, C., & Phillips, A. S. (2023). A range of outcomes: The combined effects of internal

913          variability and anthropogenic forcing on regional climate trends over Europe.

914          *Nonlinear Processes in Geophysics*, *30*(1), 63–84. https://doi.org/10.5194/npg-30-63-

915          2023

916     Di Virgilio, G., Evans, J. P., Di Luca, A., Grose, M. R., Round, V., & Thatcher, M. (2020).

917          Realised added value in dynamical downscaling of Australian climate change.

918    *Climate Dynamics*, *54*(11–12), 4675–4692. https://doi.org/10.1007/s00382-020-

919    05250-1

920  Di Virgilio, G., Evans, J. P., Di Luca, A., Olson, R., Argüeso, D., Kala, J., Andrys, J.,

921    Hoffmann, P., Katzfey, J. J., & Rockel, B. (2019). Evaluating reanalysis-driven

922    CORDEX regional climate models over Australia: Model performance and errors.

923    *Climate Dynamics*, *53*(5–6), 2985–3005. https://doi.org/10.1007/s00382-019-04672-

924    w

925  Doblas-Reyes, F. J., Hagedorn, R., & Palmer, T. N. (2005). The rationale behind the success

926    of multi-model ensembles in seasonal forecasting – II. Calibration and combination.

927    *Tellus A: Dynamic Meteorology and Oceanography*, *57*(3), 234.

928    https://doi.org/10.3402/tellusa.v57i3.14658

929  Doury, A., Somot, S., Gadat, S., Ribes, A., & Corre, L. (2022). Regional climate model

930    emulator based on deep learning: Concept and first evaluation of a novel hybrid

931    downscaling approach. *Climate Dynamics*, *60*(5), 1751–1779.

932    https://doi.org/10.1007/s00382-022-06343-9

933  Dubiński, J., Deja, K., Wenzel, S., Rokita, P., & Trzciński, T. (2023). *Selectively increasing*

934    *the diversity of GAN-generated samples* (arXiv:2207.01561). arXiv.

935    http://arxiv.org/abs/2207.01561

936  Feser, F., Rockel, B., von Storch, H., Winterfeldt, J., & Zahn, M. (2011). Regional Climate

937    Models Add Value to Global Model Data: A Review and Selected Examples. *Bulletin*

938    *of the American Meteorological Society*, *92*(9), 1181–1192.

939    https://doi.org/10.1175/2011BAMS3061.1

940      Fortin, V., Abaza, M., Anctil, F., & Turcotte, R. (2014). Why Should Ensemble Spread

941           Match the RMSE of the Ensemble Mean? *Journal of Hydrometeorology*, *15*(4), 1708–

942           1713. https://doi.org/10.1175/JHM-D-14-0008.1

943      Fowler, H. J., Blenkinsop, S., & Tebaldi, C. (2007). Linking climate change modelling to

944           impacts studies: Recent advances in downscaling techniques for hydrological

945           modelling. *International Journal of Climatology*, *27*(12), 1547–1578.

946           https://doi.org/10.1002/joc.1556

947      Furrer, E. M., & Katz, R. W. (2008). Improving the simulation of extreme precipitation

948           events by stochastic weather generators. *Water Resources Research*, *44*(12).

949           https://doi.org/10.1029/2008WR007316

950      Gensini, V. A., Haberlie, A. M., & Ashley, W. S. (2023). Convection-permitting simulations

951           of historical and possible future climate over the contiguous United States. *Climate*

952           *Dynamics*, *60*(1), 109–126. https://doi.org/10.1007/s00382-022-06306-0

953      Gibson, P. B., Rampal, N., Dean, S. M., & Morgenstern, O. (2024). Storylines for Future

954           Projections of Precipitation Over New Zealand in CMIP6 Models. *Journal of*

955           *Geophysical Research: Atmospheres*, *129*(5), e2023JD039664.

956           https://doi.org/10.1029/2023JD039664

957      Gibson, P. B., Stone, D., Thatcher, M., Broadbent, A., Dean, S., Rosier, S. M., Stuart, S., &

958           Sood, A. (2023). High-Resolution CCAM Simulations Over New Zealand and the

959           South Pacific for the Detection and Attribution of Weather Extremes. *Journal of*

960           *Geophysical Research: Atmospheres*, *128*(14), e2023JD038530.

961           https://doi.org/10.1029/2023JD038530

962      Giorgi, F., Brodeur, C. S., & Bates, G. T. (1994). Regional Climate Change Scenarios over

963           the United States Produced with a Nested Regional Climate Model. *Journal of*

964          *Climate*, *7*(3), 375–399. https://doi.org/10.1175/1520-

965          0442(1994)007<0375:RCCSOT>2.0.CO;2

966     Giorgi, F., Jones, C., & Asrar, G. R. (2009). *Addressing climate information needs at the*

967          *regional level: The CORDEX framework.*

968     Gneiting, T., & Katzfuss, M. (2014). Probabilistic Forecasting. *Annual Review of Statistics*

969          *and Its Application*, *1*(Volume 1, 2014), 125–151. https://doi.org/10.1146/annurev-

970          statistics-062713-085831

971     Gneiting, T., & Raftery, A. E. (2007). Strictly Proper Scoring Rules, Prediction, and

972          Estimation. *Journal of the American Statistical Association*, *102*(477), 359–378.

973          https://doi.org/10.1198/016214506000001437

974     Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville,

975          A., & Bengio, Y. (2014). Generative Adversarial Nets. *Advances in Neural*

976          *Information Processing Systems*, *27*.

977          https://proceedings.neurips.cc/paper_files/paper/2014/hash/5ca3e9b122f61f8f06494c9

978          7b1afccf3-Abstract.html

979     Groenke, B., Madaus, L., & Monteleoni, C. (2020). ClimAlign: Unsupervised statistical

980          downscaling of climate variables via normalizing flows. *Proceedings of the 10th*

981          *International Conference on Climate Informatics*, 60–66.

982          https://doi.org/10.1145/3429309.3429318

983     Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., & Courville, A. C. (2017). Improved

984          Training of Wasserstein GANs. *Advances in Neural Information Processing Systems*,

985          *30*.

986          https://proceedings.neurips.cc/paper_files/paper/2017/hash/892c3b1c6dccd52936e27c

987          bd0ff683d6-Abstract.html

988    Harris, L., McRae, A. T. T., Chantry, M., Dueben, P. D., & Palmer, T. N. (2022). A

989        Generative Deep Learning Approach to Stochastic Downscaling of Precipitation

990        Forecasts. *Journal of Advances in Modeling Earth Systems*, *14*(10), e2022MS003120.

991        https://doi.org/10.1029/2022MS003120

992    Hawkins, E., & Sutton, R. (2009). The Potential to Narrow Uncertainty in Regional Climate

993        Predictions. *Bulletin of the American Meteorological Society*, *90*(8), 1095–1108.

994        https://doi.org/10.1175/2009BAMS2607.1

995    Hawkins, E., & Sutton, R. (2011). The potential to narrow uncertainty in projections of

996        regional precipitation change. *Climate Dynamics*, *37*(1–2), 407–418.

997        https://doi.org/10.1007/s00382-010-0810-6

998    Hersbach, H. (2000). Decomposition of the Continuous Ranked Probability Score for

999        Ensemble Prediction Systems. *Weather and Forecasting*, *15*(5), 559–570.

1000        https://doi.org/10.1175/1520-0434(2000)015<0559:DOTCRP>2.0.CO;2

1001    Hobeichi, S., Nishant, N., Shao, Y., Abramowitz, G., Pitman, A., Sherwood, S., Bishop, C.,

1002        & Green, S. (2023). Using Machine Learning to Cut the Cost of Dynamical

1003        Downscaling. *Earth's Future*, *11*(3), e2022EF003291.

1004        https://doi.org/10.1029/2022EF003291

1005    Holden, P. B., Edwards, N. R., Garthwaite, P. H., & Wilkinson, R. D. (2015). Emulation and

1006        interpretation of high-dimensional climate model outputs. *Journal of Applied*

1007        *Statistics*, *42*(9), 2038–2055. https://doi.org/10.1080/02664763.2015.1016412

1008    Hoogewind, K. A., Baldwin, M. E., & Trapp, R. J. (2017). The Impact of Climate Change on

1009        Hazardous Convective Weather in the United States: Insight from High-Resolution

1010        Dynamical Downscaling. *Journal of Climate*, *30*(24), 10081–10100.

1011        https://doi.org/10.1175/JCLI-D-16-0885.1

1012    Isola, P., Zhu, J.-Y., Zhou, T., & Efros, A. A. (2018). *Image-to-Image Translation with*

1013        *Conditional Adversarial Networks* (arXiv:1611.07004). arXiv.

1014        https://doi.org/10.48550/arXiv.1611.07004

1015    Isphording, R. N., Alexander, L. V., Bador, M., Green, D., Evans, J. P., & Wales, S. (2023).

1016        A Standardized Benchmarking Framework to Assess Downscaled Precipitation

1017        Simulations. *Journal of Climate*, *1*(aop). https://doi.org/10.1175/JCLI-D-23-0317.1

1018    Izumi, T., Amagasaki, M., Ishida, K., & Kiyama, M. (2022). Super-resolution of sea surface

1019        temperature with convolutional neural network- and generative adversarial network-

1020        based methods. *Journal of Water and Climate Change*, *13*(4), 1673–1683.

1021        https://doi.org/10.2166/wcc.2022.291

1022    Jones, R. G., Murphy, J. M., & Noguer, M. (1995). Simulation of climate change over europe

1023        using a nested regional-climate model. I: Assessment of control climate, including

1024        sensitivity to location of lateral boundaries. *Quarterly Journal of the Royal*

1025        *Meteorological Society*, *121*(526), 1413–1449.

1026        https://doi.org/10.1002/qj.49712152610

1027    Kochkov, D., Yuval, J., Langmore, I., Norgaard, P., Smith, J., Mooers, G., Klöwer, M.,

1028        Lottes, J., Rasp, S., Düben, P., Hatfield, S., Battaglia, P., Sanchez-Gonzalez, A.,

1029        Willson, M., Brenner, M. P., & Hoyer, S. (2024). *Neural General Circulation Models*

1030        *for Weather and Climate* (arXiv:2311.07222). arXiv. http://arxiv.org/abs/2311.07222

1031    Ledig, C., Theis, L., Huszar, F., Caballero, J., Cunningham, A., Acosta, A., Aitken, A.,

1032        Tejani, A., Totz, J., Wang, Z., & Shi, W. (2017). *Photo-Realistic Single Image Super-*

1033        *Resolution Using a Generative Adversarial Network* (arXiv:1609.04802). arXiv.

1034        http://arxiv.org/abs/1609.04802

1035    Leinonen, J., Nerini, D., & Berne, A. (2021). Stochastic Super-Resolution for Downscaling

1036        Time-Evolving Atmospheric Fields With a Generative Adversarial Network. *IEEE*

1037        *Transactions on Geoscience and Remote Sensing*, *59*(9), 7211–7223.

1038        https://doi.org/10.1109/TGRS.2020.3032790

1039    Lerch, S., Thorarinsdottir, T. L., Ravazzolo, F., & Gneiting, T. (2017). Forecaster's

1040        Dilemma: Extreme Events and Forecast Evaluation. *Statistical Science*, *32*(1), 106–

1041        127. https://doi.org/10.1214/16-STS588

1042    Leutbecher, M., & Palmer, T. N. (2008). Ensemble forecasting. *Journal of Computational*

1043        *Physics*, *227*(7), 3515–3539. https://doi.org/10.1016/j.jcp.2007.02.014

1044    Liu, C., Ikeda, K., Rasmussen, R., Barlage, M., Newman, A. J., Prein, A. F., Chen, F., Chen,

1045        L., Clark, M., Dai, A., Dudhia, J., Eidhammer, T., Gochis, D., Gutmann, E., Kurkute,

1046        S., Li, Y., Thompson, G., & Yates, D. (2017). Continental-scale convection-

1047        permitting modeling of the current and future climate of North America. *Climate*

1048        *Dynamics*, *49*(1), 71–95. https://doi.org/10.1007/s00382-016-3327-9

1049    Mao, Q., Lee, H.-Y., Tseng, H.-Y., Ma, S., & Yang, M.-H. (2019). *Mode Seeking Generative*

1050        *Adversarial Networks for Diverse Image Synthesis* (arXiv:1903.05628). arXiv.

1051        http://arxiv.org/abs/1903.05628

1052    Maraun, D. (2016). Bias Correcting Climate Change Simulations—A Critical Review.

1053        *Current Climate Change Reports*, *2*(4), 211–220. https://doi.org/10.1007/s40641-016-

1054        0050-x

1055    Maraun, D., Wetterhall, F., Ireson, A. M., Chandler, R. E., Kendon, E. J., Widmann, M.,

1056        Brienen, S., Rust, H. W., Sauter, T., Themeßl, M., Venema, V. K. C., Chun, K. P.,

1057        Goodess, C. M., Jones, R. G., Onof, C., Vrac, M., & Thiele-Eich, I. (2010).

1058        Precipitation downscaling under climate change: Recent developments to bridge the

1059      gap between dynamical models and the end user. *Reviews of Geophysics*, *48*(3).

1060      https://doi.org/10.1029/2009RG000314

1061 Mardani, M., Brenowitz, N., Cohen, Y., Pathak, J., Chen, C.-Y., Liu, C.-C., Vahdat, A.,

1062      Kashinath, K., Kautz, J., & Pritchard, M. (2023). *Generative Residual Diffusion*

1063      *Modeling for Km-scale Atmospheric Downscaling* (arXiv:2309.15214). arXiv.

1064      http://arxiv.org/abs/2309.15214

1065 Matheson, J. E., & Winkler, R. L. (1976). Scoring Rules for Continuous Probability

1066      Distributions. *Management Science*, *22*(10), 1087–1096.

1067      https://doi.org/10.1287/mnsc.22.10.1087

1068 McGregor, J. L., & Dix, M. R. (2008). An Updated Description of the Conformal-Cubic

1069      Atmospheric Model. In K. Hamilton & W. Ohfuchi (Eds.), *High Resolution*

1070      *Numerical Modelling of the Atmosphere and Ocean* (pp. 51–75). Springer.

1071      https://doi.org/10.1007/978-0-387-49791-4_4

1072 Miralles, O., Steinfeld, D., Martius, O., & Davison, A. C. (2022). Downscaling of Historical

1073      Wind Fields over Switzerland Using Generative Adversarial Networks. *Artificial*

1074      *Intelligence for the Earth Systems*, *1*(4). https://doi.org/10.1175/AIES-D-22-0018.1

1075 Mirza, M., & Osindero, S. (2014). *Conditional Generative Adversarial Nets*

1076      (arXiv:1411.1784). arXiv. https://doi.org/10.48550/arXiv.1411.1784

1077 Nishant, N., Hobeichi, S., Sherwood, S. C., Abramowitz, G., Shao, Y., Bishop, C., & Pitman,

1078      A. J. (2023). Comparison of a novel machine learning approach with dynamical

1079      downscaling for Australian precipitation. *Environmental Research Letters*.

1080      https://doi.org/10.1088/1748-9326/ace463

41

1081    Oyama, N., Ishizaki, N. N., Koide, S., & Yoshida, H. (2023). Deep generative model super-

1082       resolves spatially correlated multiregional climate data. *Scientific Reports*, *13*(1),

1083       5992. https://doi.org/10.1038/s41598-023-32947-0

1084    Palmer, T. N., Doblas-Reyes, F. J., Weisheimer, A., & Rodwell, M. J. (2008). Toward

1085       Seamless Prediction: Calibration of Climate Change Projections Using Seasonal

1086       Forecasts. *Bulletin of the American Meteorological Society*, *89*(4), 459–470.

1087       https://doi.org/10.1175/BAMS-89-4-459

1088    Peard, A., & Hall, J. (2023). *Combining deep generative models with extreme value theory*

1089       *for synthetic hazard simulation: A multivariate and spatially coherent approach*

1090       (arXiv:2311.18521). arXiv. https://doi.org/10.48550/arXiv.2311.18521

1091    Prein, A. F., Langhans, W., Fosser, G., Ferrone, A., Ban, N., Goergen, K., Keller, M., Tölle,

1092       M., Gutjahr, O., Feser, F., Brisson, E., Kollet, S., Schmidli, J., van Lipzig, N. P. M.,

1093       & Leung, R. (2015). A review on regional convection-permitting climate modeling:

1094       Demonstrations, prospects, and challenges. *Reviews of Geophysics*, *53*(2), 323–361.

1095       https://doi.org/10.1002/2014RG000475

1096    Price, I., & Rasp, S. (2022). Increasing the accuracy and resolution of precipitation forecasts

1097       using deep generative models. *Proceedings of The 25th International Conference on*

1098       *Artificial Intelligence and Statistics*, 10555–10571.

1099       https://proceedings.mlr.press/v151/price22a.html

1100    Rampal, N. (2024). Enhancing Regional Climate Downscaling Through Advances in

1101       Machine Learning in: Artificial Intelligence for the Earth Systems. *Artificial*

1102       *Intelligence for the Earth Systems*. https://doi.org/10.1175/AIES-D-23-0066.1

1103    Rampal, N., Gibson, P. B., Sood, A., Stuart, S., Fauchereau, N. C., Brandolino, C., Noll, B.,

1104       & Meyers, T. (2022). High-resolution downscaling with interpretable deep learning:

1105    Rainfall extremes over New Zealand. *Weather and Climate Extremes*, *38*, 100525.

1106    https://doi.org/10.1016/j.wace.2022.100525

1107 Rampal, N., Hobeichi, S., Gibson, P. B., Baño-Medina, J., Abramowitz, G., Beucler, T.,

1108    González-Abad, J., Chapman, W., Harder, P., & Gutiérrez, J. M. (2024). Enhancing

1109    Regional Climate Downscaling through Advances in Machine Learning. *Artificial*

1110    *Intelligence for the Earth Systems*, *3*(2). https://doi.org/10.1175/AIES-D-23-0066.1

1111 Rasp, S., Dueben, P. D., Scher, S., Weyn, J. A., Mouatadid, S., & Thuerey, N. (2020).

1112    WeatherBench: A Benchmark Data Set for Data-Driven Weather Forecasting. *Journal*

1113    *of Advances in Modeling Earth Systems*, *12*(11).

1114    https://doi.org/10.1029/2020MS002203

1115 Ravuri, S., Lenc, K., Willson, M., Kangin, D., Lam, R., Mirowski, P., Fitzsimons, M.,

1116    Athanassiadou, M., Kashem, S., Madge, S., Prudden, R., Mandhane, A., Clark, A.,

1117    Brock, A., Simonyan, K., Hadsell, R., Robinson, N., Clancy, E., Arribas, A., &

1118    Mohamed, S. (2021a). Skilful precipitation nowcasting using deep generative models

1119    of radar. *Nature*, *597*(7878), 672–677. https://doi.org/10.1038/s41586-021-03854-z

1120 Ravuri, S., Lenc, K., Willson, M., Kangin, D., Lam, R., Mirowski, P., Fitzsimons, M.,

1121    Athanassiadou, M., Kashem, S., Madge, S., Prudden, R., Mandhane, A., Clark, A.,

1122    Brock, A., Simonyan, K., Hadsell, R., Robinson, N., Clancy, E., Arribas, A., &

1123    Mohamed, S. (2021b). Skilful precipitation nowcasting using deep generative models

1124    of radar. *Nature*, *597*(7878), Article 7878. https://doi.org/10.1038/s41586-021-03854-

1125    z

1126 Reddy, P. J., Matear, R., Taylor, J., Thatcher, M., & Grose, M. (2023). A precipitation

1127    downscaling method using a super-resolution deconvolution neural network with step

1128    orography. *Environmental Data Science*, *2*, e17. https://doi.org/10.1017/eds.2023.18

1129  Renwick, J. A., Mullan, A. B., & Porteous, A. (2009). Statistical Downscaling of New

1130      Zealand Climate. *Weather and Climate*, *29*, 24–44. https://doi.org/10.2307/26169704

1131  Ronneberger, O., Fischer, P., & Brox, T. (2015). U-Net: Convolutional Networks for

1132      Biomedical Image Segmentation. In N. Navab, J. Hornegger, W. M. Wells, & A. F.

1133      Frangi (Eds.), *Medical Image Computing and Computer-Assisted Intervention –*

1134      *MICCAI 2015* (pp. 234–241). Springer International Publishing.

1135      https://doi.org/10.1007/978-3-319-24574-4_28

1136  Saha, A., & Ravela, S. (2022). *Downscaling Extreme Rainfall Using Physical-Statistical*

1137      *Generative Adversarial Learning* (arXiv:2212.01446). arXiv.

1138      http://arxiv.org/abs/2212.01446

1139  Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., & Chen, X. (2016).

1140      *Improved Techniques for Training GANs* (arXiv:1606.03498). arXiv.

1141      https://doi.org/10.48550/arXiv.1606.03498

1142  Sha, J., Chen, X., Chang, Y., Zhang, M., & Li, X. (2024). A spatial weather generator based

1143      on conditional deep convolution generative adversarial nets (cDCGAN). *Climate*

1144      *Dynamics*, *62*(2), 1275–1290. https://doi.org/10.1007/s00382-023-06971-9

1145  Sørland, S. L., Schär, C., Lüthi, D., & Kjellström, E. (2018). Bias patterns and climate change

1146      signals in GCM-RCM model chains. *Environmental Research Letters*, *13*(7), 074017.

1147      https://doi.org/10.1088/1748-9326/aacc77

1148  Srivastava, A., Valkov, L., Russell, C., Gutmann, M. U., & Sutton, C. (2017). VEEGAN:

1149      Reducing Mode Collapse in GANs using Implicit Variational Learning. *Advances in*

1150      *Neural Information Processing Systems*, *30*.

1151      https://proceedings.neurips.cc/paper/2017/hash/44a2e0804995faf8d2e3b084a1e2db1d

1152      -Abstract.html

1153 Steinschneider, S., Ray, P., Rahat, S. H., & Kucharski, J. (2019). A Weather-Regime-Based

1154      Stochastic Weather Generator for Climate Vulnerability Assessments of Water

1155      Systems in the Western United States. *Water Resources Research*, *55*(8), 6923–6945.

1156      https://doi.org/10.1029/2018WR024446

1157 Sun, Y., Deng, K., Ren, K., Liu, J., Deng, C., & Jin, Y. (2024). Deep learning in statistical

1158      downscaling for deriving high spatial resolution gridded meteorological data: A

1159      systematic review. *ISPRS Journal of Photogrammetry and Remote Sensing*, *208*, 14–

1160      38. https://doi.org/10.1016/j.isprsjprs.2023.12.011

1161 Thatcher, M., & McGregor, J. L. (2009). Using a Scale-Selective Filter for Dynamical

1162      Downscaling with the Conformal Cubic Atmospheric Model. *Monthly Weather*

1163      *Review*, *137*(6), 1742–1752. https://doi.org/10.1175/2008MWR2599.1

1164 van der Meer, M., de Roda Husman, S., & Lhermitte, S. (2023). Deep Learning Regional

1165      Climate Model Emulators: A Comparison of Two Downscaling Training

1166      Frameworks. *Journal of Advances in Modeling Earth Systems*, *15*(6),

1167      e2022MS003593. https://doi.org/10.1029/2022MS003593

1168 Verdin, A., Rajagopalan, B., Kleiber, W., Podestá, G., & Bert, F. (2018). A conditional

1169      stochastic weather generator for seasonal to multi-decadal simulations. *Journal of*

1170      *Hydrology*, *556*, 835–846. https://doi.org/10.1016/j.jhydrol.2015.12.036

1171 Vosper, E., Watson, P., Harris, L., McRae, A., Santos-Rodriguez, R., Aitchison, L., &

1172      Mitchell, D. (2023). Deep Learning for Downscaling Tropical Cyclone Rainfall to

1173      Hazard-Relevant Spatial Scales. *Journal of Geophysical Research: Atmospheres*,

1174      *128*(10), e2022JD038163. https://doi.org/10.1029/2022JD038163

1175 Wang, J., Liu, Z., Foster, I., Chang, W., Kettimuthu, R., & Kotamarthi, V. R. (2021). Fast

1176      and accurate learned multiresolution dynamical downscaling for precipitation.

1177    *Geoscientific Model Development*, *14*(10), 6355–6372. https://doi.org/10.5194/gmd-

1178    14-6355-2021

1179    Wang, X., Yu, K., Wu, S., Gu, J., Liu, Y., Dong, C., Loy, C. C., Qiao, Y., & Tang, X. (2018).

1180    *ESRGAN: Enhanced Super-Resolution Generative Adversarial Networks*

1181    (arXiv:1809.00219). arXiv. http://arxiv.org/abs/1809.00219

1182    *WeatherBench: A Benchmark Data Set for Data-Driven Weather Forecasting—Rasp—*

1183    *2020—Journal of Advances in Modeling Earth Systems—Wiley Online Library*. (n.d.).

1184    Retrieved October 22, 2021, from

1185    https://agupubs.onlinelibrary.wiley.com/doi/full/10.1029/2020MS002203

1186    Xu, Z., Han, Y., & Yang, Z. (2019). Dynamical downscaling of regional climate: A review of

1187    methods and limitations. *Science China Earth Sciences*, *62*(2), 365–375.

1188    https://doi.org/10.1007/s11430-018-9261-5

1189    Zhang, X., Alexander, L., Hegerl, G. C., Jones, P., Tank, A. K., Peterson, T. C., Trewin, B.,

1190    & Zwiers, F. W. (2011). Indices for monitoring changes in extremes based on daily

1191    temperature and precipitation data. *WIREs Climate Change*, *2*(6), 851–870.

1192    https://doi.org/10.1002/wcc.147

1193
1194
1195
1196