

Generating time-consistent dynamics with discriminator-guided image diffusion models

Philipp Hess^{1,2}, Maximilian Gelbrecht^{1,2}, Christof Schötz^{1,2}, Michael Aich^{1,2}
Yu Huang^{1,2}, Shangshang Yang^{1,2}, Niklas Boers^{1,2}

¹Technical University of Munich, ²Potsdam Institute for Climate Impact Research,

{philipp.hess, maximilian.gelbrecht, christof.schoetz,
michael.aich, y.huang, shangshang.yang, n.boers}@tum.de

Abstract

Realistic temporal dynamics are crucial for many video generation, processing and modelling applications, e.g. in computational fluid dynamics, weather prediction, or long-term climate simulations. Video diffusion models (VDMs) are the current state-of-the-art method for generating highly realistic dynamics. However, training VDMs from scratch can be challenging and requires large computational resources, limiting their wider application. Here, we propose a time-consistency discriminator that enables pretrained image diffusion models to generate realistic spatiotemporal dynamics. The discriminator guides the sampling inference process and does not require extensions or finetuning of the image diffusion model. We compare our approach against a VDM trained from scratch on an idealized turbulence simulation and a real-world global precipitation dataset. Our approach performs equally well in terms of temporal consistency, shows improved uncertainty calibration and lower biases compared to the VDM, and achieves stable centennial-scale climate simulations at daily time steps.

1 Introduction

Generating time-consistent sequences of images is important to many video generation and synthesis tasks [1, 2, 3, 4, 5], for example in computational fluid dynamics [6, 7, 8], probabilistic weather forecasts [9, 10] or climate simulations [11, 12, 13, 14].

The success of image diffusion models (IDMs) [15, 16, 17] has sparked a large interest in extending their generation to time-consistent videos, achieving remarkable results [1, 2, 18, 19, 20, 3, 21, 22]. However, training video diffusion models (VDMs) from scratch is challenging and requires large amounts of computational resources [23]. Moreover, recent state-of-the-art VDMs are not always released open source [3], limiting their adaptability to a wider scientific community.

Therefore, efforts have been made to leverage pretrained image models for video editing tasks such as style-transfer or inpainting [24, 5, 25, 26, 27, 28]. Video editing relies on full or partial temporal information in the source video that can then be combined with inference-level guidance techniques to preserve temporal consistency during the editing process. Such video processing tasks are also important to many scientific applications, for example, in data reconstruction using inpainting methods or downscaling applications using super-resolution techniques in fluid dynamics [29, 30], meteorology [31] and climate science [32, 30, 33, 34, 35, 36].

Generating videos with IDMs without relying on a source video or a given encoding of the dynamics is much more challenging. Most approaches rely on finetuning an IDM on video data, e.g., by inserting additional temporal layers into the architecture [37, 38, 39, 40], which can still be computationally demanding and requires a deep understanding of the architecture.

We propose a novel guidance approach, inspired by temporal discriminators in generative adversarial networks [41, 42, 43], for the generation of realistic, time-consistent, and stable spatiotemporal dynamics with pretrained IDMs. Our discriminator guidance is lightweight and efficient, adding only about 3%-8% to the generation time, and is trained independently of the IDM, making extensions of different IDMs to new downstream tasks straightforward. We perform a comprehensive evaluation on challenging datasets with high-dimensional chaotic dynamics, including 2D Navier-Stokes turbulence simulations and global precipitation reanalysis, using the extensive catalog of established metrics from fluid dynamics and Earth system science. We find that our method performs similarly well as a VDM trained from scratch in terms of temporal dynamics, while achieving better uncertainty calibration and lower biases. Moreover, our guidance approach enables stable climate simulations for more than 100 years, while the VDM exhibits unstable drifts in global averages.

2 Related work

Video GANs. Generative adversarial networks (GANs) have been widely explored for synthesizing temporally-consistent videos. Earlier work [44, 45] introduced the idea of using an adversarial discriminator to distinguish between real and generated video frames, which was improved in following studies [46, 47, 48]. DVD-GAN [41] proposed two separate discriminators for spatial and time domains, the latter being similarly motivated as our time-consistency discriminator.

Video prediction GANs with temporal discriminators have shown great success in turbulence modelling [49] and probabilistic weather predictions [42, 43]. However, while temporal discriminators provide powerful tools that enable the generation of dynamically consistent videos in GANs, adversarial training is generally prone to instabilities and mode collapse, making GANs challenging to optimize.

Video diffusion models. Generative diffusion models (DMs) [15, 16, 50, 17], have largely superseded GANs owing to their improved training stability, high-fidelity output and iterative sampling process, which enables downstream tasks without retraining [15, 16].

Video diffusion models [1] have achieved state-of-the-art performance [2, 18, 19, 20, 3, 21, 22], e.g., through latent VDMs [51, 52, 53], and improved training strategies [18, 54, 21]. Classifier-free guidance has also been explored to enable variable-length conditioning on past video frames with VDMs [55].

The ability of VDMs to model uncertainties and to produce sharp outputs makes them powerful tools, e.g., for weather prediction [9, 10, 56], super-resolution (downscaling) [11], reconstructing spatiotemporal dynamics from sparse sensor measurements [57], emulating precipitation dynamics directly from remote sensing observations [58], or climate model simulations [12, 13, 14].

However, VDMs require large computational resources for training [23, 59], limiting their applicability.

Video synthesis with image diffusion models. Due to the high computational costs and lack of open source availability of VDMs, recent efforts have focused on utilizing available pretrained image diffusion models (IDMs) for video processing and editing tasks. In video processing tasks, the temporal dynamics are usually given in a source video that needs to be transformed in a time-consistent manner, for applications such as style-transfer, inpainting, or super-resolution [24]. Approaches to preserve temporal consistency include correlated (warped) noise [5, 25], or transitioning from spatial to temporal-attention blocks [26, 27, 28, 60, 61, 4].

IDMs have been adapted to applications in fluid dynamics, weather prediction and climate modelling. Some applications take temporal dynamics explicitly into account, e.g., in data-assimilation [62] and spatio-temporal downscaling [11, 63]. Sampling guidance from a numerical weather prediction model has also been used to improve the weather forecast from a VDM [64]. Further, IDMs have been combined with a deterministic forecast neural network to produce dynamically consistent simulations from weather to climate time scales [56].

Many applications, however, employ IDMs to process dynamical simulations in each time step without taking time consistency into account, e.g., for downscaling (super-resolution) climate simulations [30, 33, 34, 31, 35, 36], data-assimilation [65], or data reconstruction [32], which could potentially be improved with our method.

When generating videos with IDMs, e.g., from a given starting frame, most work relies on finetuning a pretrained model by inserting temporal-attention layers [37, 38, 39, 40]. Notably, [66] use a temporal

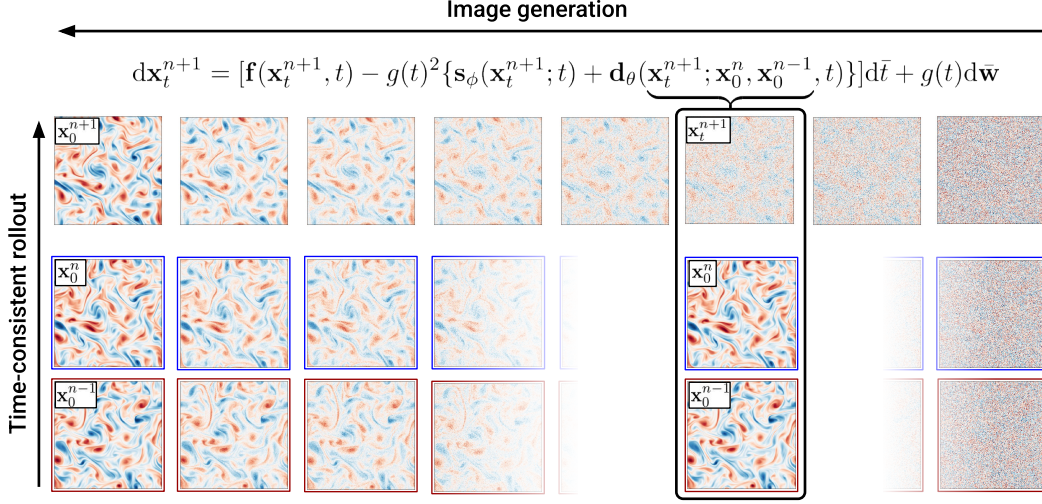


Figure 1: **Overview sketch** of the time-consistency discriminator guidance for generating images in a dynamically realistic sequence. The discriminator guidance $\mathbf{d}_\theta(\cdot)$ uses the current and past time frames, \mathbf{x}^n and \mathbf{x}^{n-1} , to guide the denoising generation of the next \mathbf{x}^{n+1} .

discriminator for finetuning the extended IDM. Our approach, in contrast, is agnostic of the IDM architecture and does not require finetuning.

Discriminator guidance. Inspired by GANs, discriminators have been employed during diffusion model training to improve the performance [67, 68, 69] or enhance sampling speed with adversarial distillation [70, 71]. Discriminators have also been proposed as purely inference-level guidance to improve the image quality of IDMs [72], to pair separately trained video and audio diffusion models [73], or to generate molecular graphs [74, 75]. Similarly, our time-consistency discriminator is only applied during inference as guidance. A notable difference to [72, 73], is that our training does not require samples generated with an IDM, which can be computationally costly.

3 Methods

Diffusion models. Diffusion models [15, 16, 17] learn to generate data from a target distribution $\mathbf{x}_0 \sim p_{\text{data}}(\mathbf{x})$ with a time-reversed denoising process that starts with an initial noise sample, e.g. Gaussian white noise $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \sigma_{\text{max}}^2 \mathbf{I})$, and can be formulated as a reversed stochastic differential equation (SDE) [76],

$$d\mathbf{x}_t = [\mathbf{f}(\mathbf{x}_t, t) - g(t)^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)] d\bar{t} + g(t) d\bar{\mathbf{w}}, \quad (1)$$

where $\mathbf{x}_t \in \mathbb{R}^{B \times C \times H \times W}$ is a noised image of batch size B with C channels and H, W pixels in height and width dimension, respectively. The drift term is given by $\mathbf{f}(\cdot)$, $d\bar{\mathbf{w}}$ adds Wiener noise where the bar denotes a time reversal, and $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ is the score function of the noised target distribution. The noise strength σ_t with lower and upper bounds σ_{min} and σ_{max} , respectively, decreases during the reverse processes, following a prescribed schedule $g(t)$ (e.g. see Eq. 8). The score function in Eq. 1 is typically intractable but can be learned with a neural network $\mathbf{S}_\phi(\cdot)$ [77],

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) \approx \frac{\mathbf{S}_\phi(\mathbf{x}_t; t) - \mathbf{x}_t}{\sigma_t^2} =: \mathbf{s}_\phi(\mathbf{x}_t; t), \quad (2)$$

using the loss function $\mathcal{L}(\phi) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x}), \epsilon_t \sim \mathcal{N}(\mathbf{0}, \sigma_t^2 \mathbf{I})} [w(t) \|\mathbf{S}_\phi(\mathbf{x}_0 + \epsilon_t; t) - \mathbf{x}_0\|_2^2]$, where $w(t)$ is a weighting function (see details in Appendix A).

Time-consistency guidance. We propose a discriminator that guides the reverse diffusion process in Eq. 1 to generate time-consistent sequences of images (see Fig. 1). A discrete, temporally ordered

time series of N images is denoted as $\{\mathbf{x}_t^n | n = 1, 2, \dots, N\}$, where the superscript represents the physical time step and the subscript the diffusion time. We train a discriminator to distinguish between images that sampled in a time-consistent manner, i.e. that are conditioned on a noise-free ($t = 0$) sequence of the current and m previous time steps, $\{\mathbf{x}_0^{n-m}, \mathbf{x}_0^{n-m+1}, \dots, \mathbf{x}_0^n\} = \mathbf{x}_0^{(n-m):n}$, following $\mathbf{x}_t^{n+1} \sim p_t(\mathbf{x}_t^{n+1} | \mathbf{x}_0^{(n-m):n})$, and random samples without temporal ordering $\mathbf{x}_t^{n+1} \sim p_t(\mathbf{x}_t^{n+1})$. A optimal discriminator has then the form [78, 72],

$$D_\theta(\mathbf{x}_t^{n+1}; \mathbf{x}_0^{(n-m):n}, t) = \frac{p(\mathbf{x}_t^{n+1} | \mathbf{x}_0^{(n-m):n})}{p(\mathbf{x}_t^{n+1} | \mathbf{x}_0^{(n-m):n}) + p(\mathbf{x}_t^{n+1})}. \quad (3)$$

Computing the scores, by applying the logarithm and gradient, on both sides of Eq. 3,

$$\nabla_{\mathbf{x}_t^{n+1}} \log \left(\frac{D_\theta(\mathbf{x}_t^{n+1}; \mathbf{x}_0^{(n-m):n}, t)}{1 - D_\theta(\mathbf{x}_t^{n+1}; \mathbf{x}_0^{(n-m):n}, t)} \right) = \nabla_{\mathbf{x}_t^{n+1}} \log \left(\frac{p(\mathbf{x}_t^{n+1} | \mathbf{x}_0^{(n-m):n})}{p(\mathbf{x}_t^{n+1})} \right),$$

allows us to define the time-consistency guidance as

$$\mathbf{d}_\theta(\mathbf{x}_t^{n+1}; \mathbf{x}_0^{(n-m):n}, t) := \nabla_{\mathbf{x}_t^{n+1}} \log \left(\frac{D_\theta(\mathbf{x}_t^{n+1}; \mathbf{x}_0^{(n-m):n}, t)}{1 - D_\theta(\mathbf{x}_t^{n+1}; \mathbf{x}_0^{(n-m):n}, t)} \right). \quad (4)$$

The guidance term in Eq. 4 can then be added to the unconditional score to enable time-consistent sampling of images in an autoregressive manner using the reverse SDE in Eq. 1,

$$d\mathbf{x}_t^{n+1} = [\mathbf{f}(\mathbf{x}_t^{n+1}, t) - g(t)^2 \{\mathbf{s}_\phi(\mathbf{x}_t^{n+1}; t) + \lambda_t \mathbf{d}_\theta(\mathbf{x}_t^{n+1}; \mathbf{x}_0^{(n-m):n}, t)\}]d\bar{t} + g(t)d\bar{\mathbf{w}} \quad (5)$$

where the strength of the guidance is controlled through the parameter λ_t (see Appendix B.1 for a more detailed discussion). In our experiments, we find that using the current and previous time step ($m = 1$) works best for conditioning the discriminator, which is in line with typical ODE solvers and ML weather models [59].

Discriminator training. The discriminator is trained as a binary classifier $D_\theta : (\mathbf{x}_t^k; \mathbf{x}_0^{(n-m):n}, t) \mapsto q$, conditioned on previous, denoised time frames $\mathbf{x}_0^{(n-m):n}$, and the diffusion noise time t , to predict the probability q of a noised image $\mathbf{x}_t^k = \mathbf{x}_0^k + \epsilon_t$, $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \sigma_t^2 \mathbf{1})$, being temporally consistent with the current and m previous time frames, i.e., whether $\mathbf{x}_t^k = \mathbf{x}_t^{n+1}$. The same noise schedule as for training the diffusion model is used for σ_t (see Appendix A). We use the standard cross entropy loss as a training objective [72],

$$\mathcal{L}_{CE}(\theta) = -\mathbb{E}_{n,l} \left[\log D_\theta(\mathbf{x}_t^{n+1}; \mathbf{x}_0^{(n-m):n}, t) + \log(1 - D_\theta(\mathbf{x}_t^{n+l}; \mathbf{x}_0^{(n-m):n}, t)) \right], \quad (6)$$

with $\mathbb{E}_{n,l} := \mathbb{E}_{n \sim \mathcal{U}(1,N), l \sim \mathcal{N}_{\mathbb{Z}}(\mu, \sigma_{\text{step}}^2) \setminus \{1\}}$, where we uniformly sample a time step n from the dataset of N samples, and introduce an importance sampling for non-time consistent samples ($l \neq 1$) to prioritize time steps from the vicinity of the next time step, $k = n + 1$, using a normal distribution of integers $l \sim \mathcal{N}_{\mathbb{Z}}(\mu, \sigma_{\text{step}}^2) \setminus \{1\}$, and we set $\mu = 1$, $\sigma_{\text{step}} = 2$. The motivation is that fields close to the next time step ahead are hardest to distinguish for the network, due to their high correlation. We find that random cropping of the images further improves the results as it forces the discriminator to focus on different spatial scales. Fig. 2 shows the time consistency prediction of the trained discriminator network during inference. See Appendix B.2 for training and architecture details.

4 Experiments

Data. We evaluate our method on two challenging datasets: an idealized fluid dynamical Navier-Stokes simulation and real-world observational precipitation data from the ERA5 reanalysis [79].

A two-dimensional vorticity simulation is performed by numerically solving the Navier-Stokes equation in vorticity formulation with a 4th-order Runge-Kutta solver on a 256×256 grid with periodic boundary conditions and stochastic forcing. We use 47k samples for training and 13k for validation and test set, respectively (see Appendix C.1 for details). The time-consistency evaluation in the following is performed over 4k samples.

For the second experiment, we use global precipitation fields from the ERA5 reanalysis, which

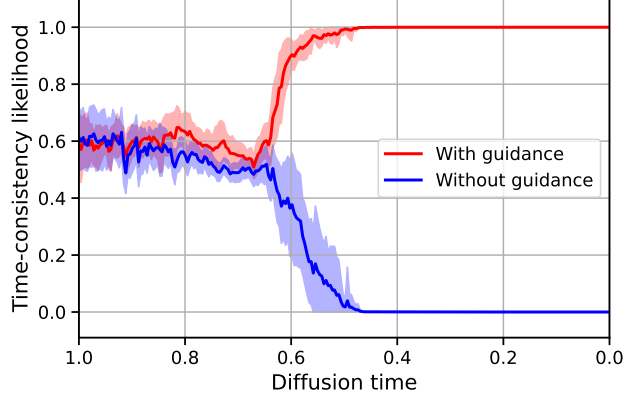


Figure 2: **Time-consistency prediction** of the discriminator network during sampling of vorticity fields with guidance switched on (red) or off (blue). The mean over 50 samples is given by the solid line, and the shaded area shows the standard deviation. With decreasing noise scales in the reserve diffusion process ($t_{\max} = 1 \rightarrow t_{\min} = 0$), the discriminator network reliably predicts whether samples are time-consistent or not.

combines high-resolution numerical weather model simulations with different observational sources using data assimilation (see details in Appendix C.2), as a complementary and real-world dataset that is challenging due to its spatiotemporal intermittency, inherent stochasticity, and skewed distributions. The horizontal spatial resolution is 1° degree, which corresponds to 180×360 grid cells in latitude (height) and longitude (width) direction, respectively. We split the daily data into periods of 1979-2000, 2001-2010 and 2011-2020 for training, validation and testing.

Baselines. We compare our time-consistency discriminator guidance method to sampling from the unconditional DM without guidance, and a video DM baseline trained from scratch. The video DM is set up in an autoregressive manner [1, 59], with a conditional score network $s_\psi(\mathbf{x}_t^{n+1}; \mathbf{x}_0^n, \mathbf{x}_0^{n-1}, t) \approx \nabla_{\mathbf{x}_t^{n+1}} \log p_t(\mathbf{x}_t^{n+1} | \mathbf{x}_0^n, \mathbf{x}_0^{n-1})$. We find that the same hyperparameters work well for both the video DM and unconditional DMs (see Appendix A for details).

Sampling. We use the stochastic EDM sampler [77] for the unconditional and video DMs, with the same parameters for both models (see Tab. 1). We apply the discriminator guidance Eq. 5 in both the first and second-order solver steps (see algorithm 1), which we find to be important to achieve good performance.

5 Results

We evaluate our time-consistency guidance approach against the baselines on Navier-Stokes fluid dynamical simulations of turbulence, and observational daily precipitation from state-of-the-art reanalysis data (ERA5) [79], using established metrics in dynamical systems theory and Earth system science (see Appendix D for definitions).

2D Navier-Stokes turbulence. A qualitative comparison of samples is shown in Fig. 9 for the first five and last time frames of the ground truth and generated simulations. A video of the generated dynamics is also provided¹. While the pairing between generated and ground truth samples is quickly lost due to the chaotic non-linear dynamics, both video DM and discriminator-guided DM produce realistic dynamics in contrast to the unconditional DM. All generative DMs remain sharp over the entire 4000-step rollout.

To better visualize the dynamics, we compute Hovmöller diagrams [80], showing the average over a vertical band of 10 grid columns in the center of the fields (Fig. 3). Our guidance approach is able to reproduce the elongated wave-like structures over multiple time steps that can be seen in the ground

¹https://youtu.be/JMwsZi_b-uk

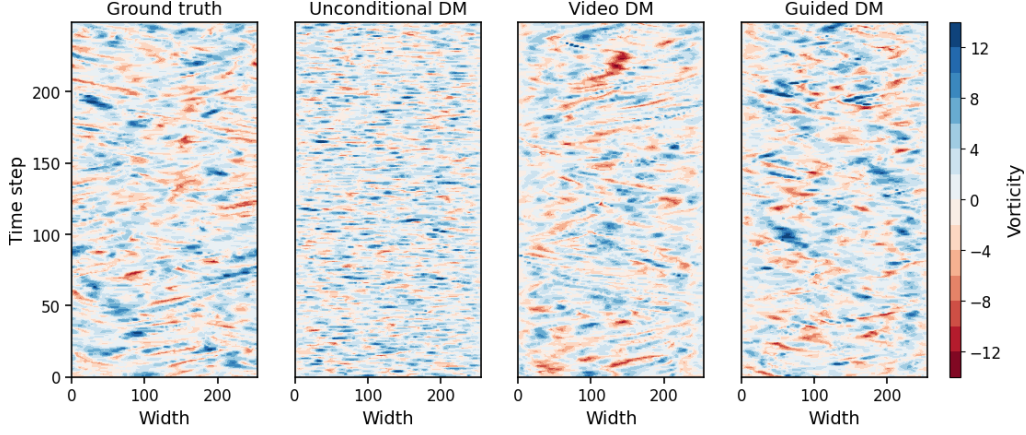


Figure 3: **Hovmöller diagrams**, often used to visualize spatiotemporal dynamics and, in particular, the propagation of waves in fluid dynamics and meteorology, are computed for the 2D vorticity simulation as the mean over a vertical band of grid columns for (from left to right) the ground truth numerical simulation, the unconditional DM, the video DM, and our guidance approach. The guidance method and video DM generate dynamics indistinguishable from the ground truth.

truth, video DM and guided DM simulations, but are absent in the unconditional DM output. We quantify the similarity in the dynamics seen in Fig. 3 by computing the Wasserstein-1 distance between consecutive rows in the Hovmöller diagram and compare their distributions in Fig. 4a and errors in Fig. 14a. A close match between the distributions of the ground truth, video and guided DM output can be seen.

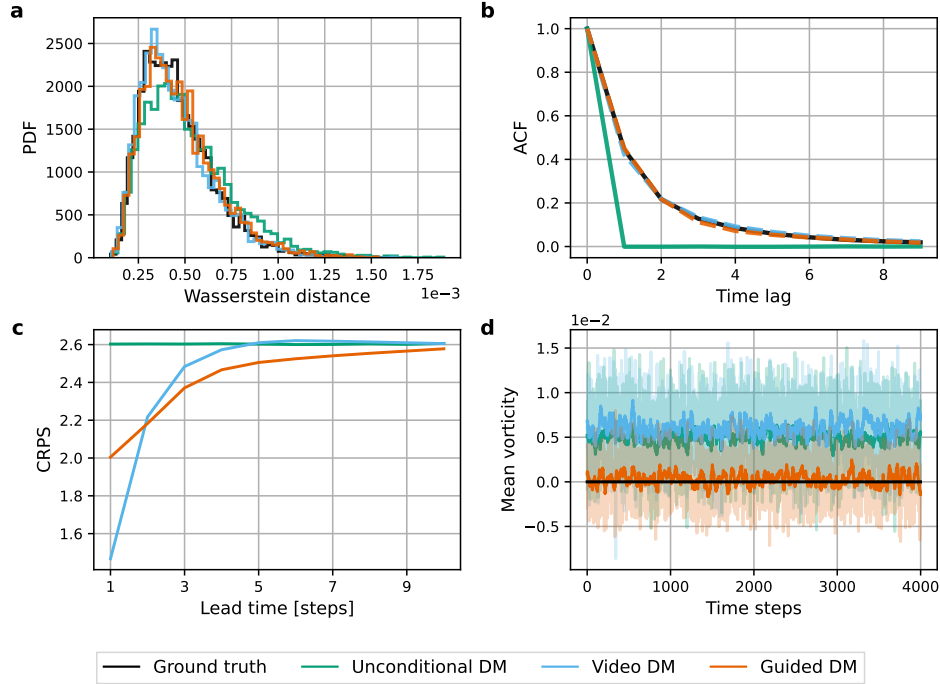


Figure 4: **Quantitative evaluation of 2D Navier-Stokes turbulent vorticity dynamics** in terms of Wasserstein distances between consecutive rows of the Hovmöller diagram in Fig. 3 (a), autocorrelation function (ACF) (b), continuous ranked probability score (CRPS) (c) and running window spatial mean as solid line with the actual time series shown as shades (d), for the ground truth simulation (black), the unconditional DM (green), video DM (blue) and the guided DM (red). Note that only the guided DM achieves an unbiased representation of the vorticity.

We compute autocorrelation functions (ACFs) with a time lag of up to 10 time steps (Fig. 4b, Fig. 14b). Both the video and guided DM achieve very accurate ACFs that are indistinguishable from the ground truth, whereas the unconditional DM generates uncorrelated samples as expected. Forecast skill is compared in terms of the continuous ranked probability score (CRPS) [81] using a 50-member ensemble, 10-step lead times, and 100 forecasts (Fig. 4c, Fig. 14c). We find that the video DM outperforms the guided DM in terms of forecast skill for the first two lead times, while the guided DM has a better forecast skill for longer lead times. We compute the spread skill ratio and find that the guided DM shows a better calibration over all lead times. (Fig. 11) In terms of global mean vorticity, both the unconditional and video DM show significant biases. The guidance method, in contrast, achieves a substantially lower bias (Fig. 4d, Fig. 14d).

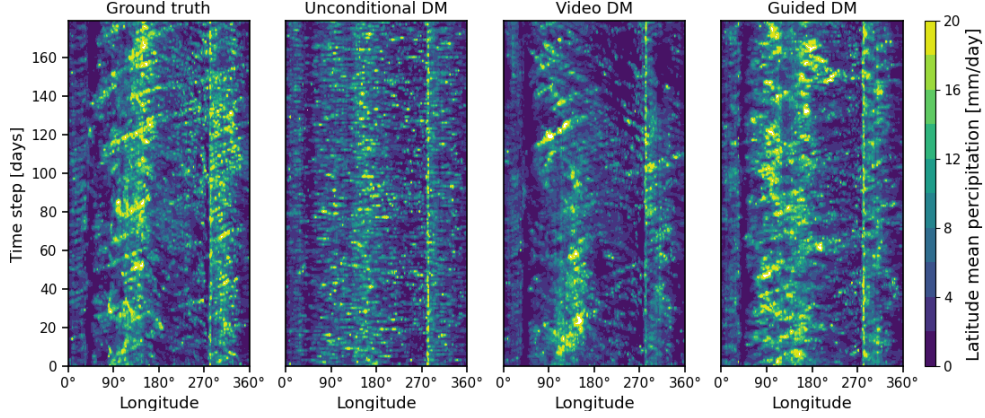


Figure 5: **Hovmöller diagrams** of the global daily precipitation simulation (from left to right) are computed for 180 days as a mean over the latitude band from 10°S to 10°N for the ground truth ERA5, unconditional DM, video DM, and our guidance approach.

ERA5 global precipitation. The first two and last daily precipitation fields from the ground truth test set and generated time series are shown in Fig. 10, videos of the simulations are also provided². All DMs produce realistic spatial patterns and remain sharp for the entire 10-year rollout, while the video DM and guided DM show qualitatively much more realistic dynamics than the unconditional DM.

The generated dynamics are again compared with Hovmöller diagrams over 180 days, using a latitude band from 10°S to 10°N. The unconditional DM produces visible random patterns that are distinctively different from the ERA5 target data (Fig. 5). Our guidance approach enables the unconditional DM to produce realistic dynamics, similar to the video DM and target data, with characteristic west-to-east wave-like patterns, which are challenging to capture even for state-of-the-art climate models [82].

We use the Wasserstein-1 distance again to quantify the similarity in the dynamics seen in the Hovmöller diagram in Fig. 6a (see errors in Fig. 15a). We find that the unconditional DM produces a flat distribution with larger distances shifted to the right of the target data distribution, which is narrower. Both the guided and the video DM capture the target distribution more accurately. Our guidance method produces a very accurate autocorrelation function, slightly outperforming the video DM for longer lags (Fig. 6b, Fig. 15b). We again perform 100 ensemble forecasts with an ensemble size of 50 members and compute the CRPS to evaluate the skill (Fig. 6c, Fig. 15c). We find that the video DM has a slightly better forecast skill than the guided DM for one and two-day ahead predictions. We compute the spread skill ratio of the forecasts and find an improved calibration in the guided DM with respect to the video DM (Fig. 12). To assess a critical characteristic of precipitation dynamics, we compute the waiting times between extreme events above the 95th percentile (Fig. 6d, Fig. 15d). We find that the unconditional DM significantly underestimates the frequency of waiting times larger than 100 days. The video DM captures waiting times less than 100 days accurately, while the guided DM generates slightly more accurate waiting times that are larger than 300 days.

We compute the first three empirical orthogonal functions (EOFs) using principal component analysis (PCA) up to an explained variance greater than 1% and find that the guided DM is able to accurately

²<https://youtu.be/noRxb9trpQ>

capture the first three leading EOFs (Fig. 13), while the video DM shows notable differences in the 2nd and 3rd EOF. We evaluate the stability of the guided and video DM with 10 simulations each over 100 years, as well as a single 170-year guided DM run, and find that the video DM develops instabilities in terms of a shifting mean. Our guidance approach, on the other hand, is stable on centennial time scales and has a much lower global mean difference to the ERA5 ground truth (Fig. 7). We compare the spatial bias in the generated precipitation time series and find that the unconditional DM produces the smallest global mean bias. The guidance method outperforms the video DM, the latter having the largest overall bias.

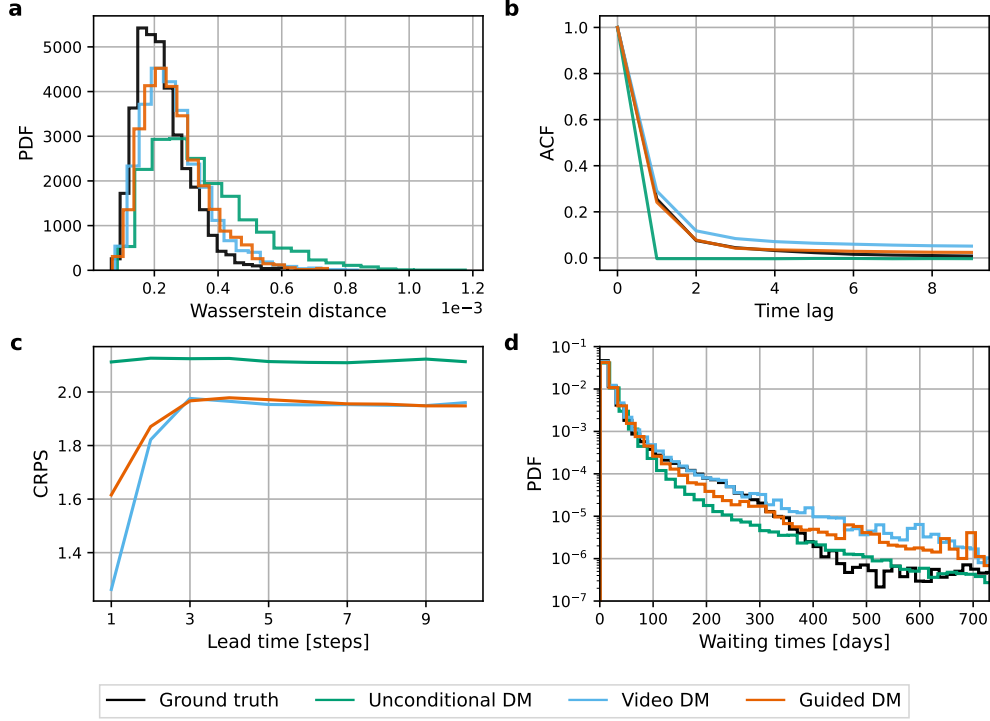


Figure 6: **Quantitative evaluation of daily precipitation dynamics** in terms of Wasserstein distances between consecutive rows of the Hovmöller diagram in Fig. 5 (a), autocorrelation functions (ACFs) (b), CRPS forecast skill (c), extreme event waiting time distributions (d), for the ground truth ERA5 (black), the unconditional DM (green), video DM (blue) and our guidance approach (red).

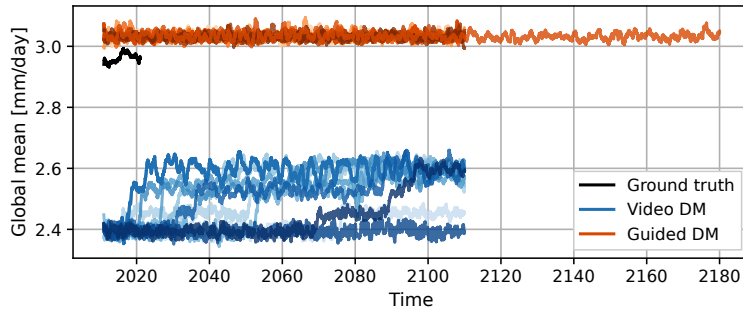


Figure 7: **Long-term precipitation simulations** are shown as an annual rolling global mean for the ERA5 test set (black), the video DM (blue) and our guidance method (red). Shadings of one color denote different ensemble members, showing that the video DM exhibits randomly occurring drifts, whereas the guided DM remains stable.

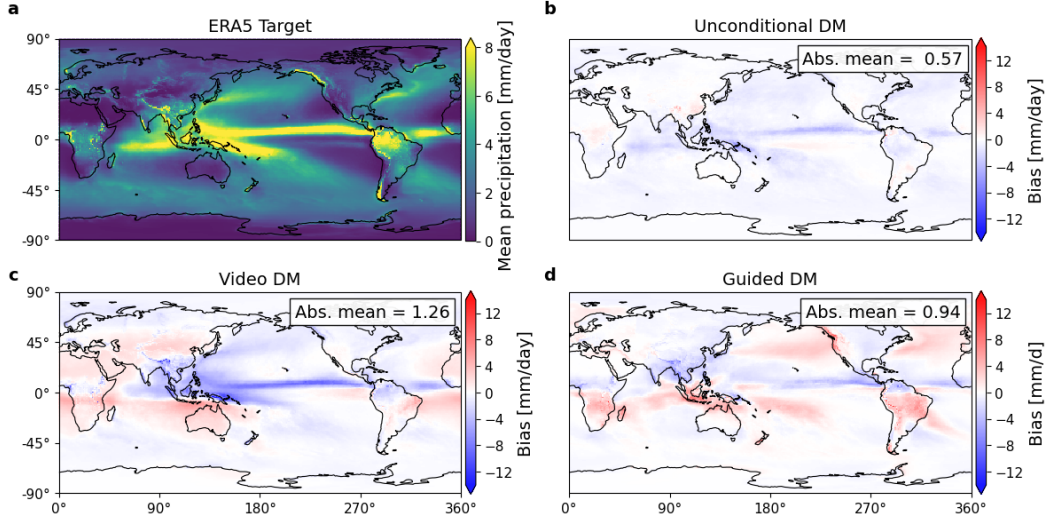


Figure 8: **Global mean bias** (see Appendix D for definition) comparison showing, (a) the test set mean of the ERA5 ground truth, (b) the bias of the unconditional DM, (c) the video DM, (d) our guided DM. Mean absolute bias values are given in the top right.

6 Discussion

We propose a time consistency discriminator that guides the sampling process of unconditionally trained image diffusion models (DMs) to generate time-consistent image sequences, i.e., dynamically realistic simulations and videos. Our discriminator is trained separately from the diffusion model on the target data only and hence independent of the DM architecture.

We evaluate our method on two challenging datasets with complex non-linear dynamics, namely 2D Navier-Stokes turbulence simulations and global precipitation reanalysis. We find that the discriminator guidance enables the unconditionally trained DM to generate realistic dynamics with comparable skill to a video DM trained from scratch. While the video DM produces more accurate short-term forecasts, the guidance method outperforms on longer lead times with lower biases and improved stability. Our method enables stable rollouts over 100 years or longer and is not subject to unstable mean shifts as the video DM, promising immense potential for climate research applications [83, 13, 84]. Our method is computationally efficient (see Appendix B.2) and increases the reusability of pre-trained image diffusion models to a wider range of tasks that require temporal consistency without the need for costly retraining video models from scratch, making generative modeling more sustainable and accessible. Moreover, our guidance method is promising for long rollouts of video DMs, which we leave for future research.

We only consider univariate simulations here, but, we believe that extensions to multiple variables are straightforward. While our method produces more realistic long-term simulations, it has a lower short-term forecast skill than the video DM. Further explorations in terms of the discriminator architecture and training might enable improvements in that respect. We apply our discriminator guidance method only to one type of diffusion model and sampler [77], but it should, in principle, be applicable to others as well [72]. This study focuses on video *generation*, however, the guidance method is also applicable to enforce time-consistency in video *processing* tasks such as super-resolution or inpainting, which is crucial in related applications in weather and climate [32, 31, 36, 33, 34]. We hope that our results will encourage further research on video synthesis with discriminator-guided diffusion models.

Broader impact. This work focuses on the spatiotemporal dynamics of highly non-linear and chaotic systems, with characteristics common in computational fluid dynamics, meteorology, and climate science. However, the ability of our method to enable realistic video generation with image diffusion models might also have potentially negative societal effects, including the amplification of disinformation.

Acknowledgments and Disclosure of Funding

The authors would like to thank Sebastian Bathiany, Alistair White and Maha Badri for their helpful comments and discussions of the work. NB and MG acknowledge funding by the Volkswagen Foundation. MA acknowledges funding under the Excellence Strategy of the Federal Government and the Länder through the TUM Innovation Network EarthCare. This is ClimTip contribution #X; the ClimTip project has received funding from the European Union's Horizon Europe research and innovation programme under grant agreement No. 101137601. YH acknowledges funding by the Alexander von Humboldt Foundation. The authors acknowledge the European Regional Development Fund (ERDF), the German Federal Ministry of Education and Research, and the Land Brandenburg for supporting this project by providing resources on the high-performance computer system at the Potsdam Institute for Climate Impact Research.

References

- [1] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J. Fleet. Video Diffusion Models. *Advances in Neural Information Processing Systems*, 35:8633–8646, December 2022.
- [2] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen Video: High Definition Video Generation with Diffusion Models, October 2022. arXiv:2210.02303 [cs].
- [3] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-Sora: Democratizing Efficient Video Production for All, December 2024. arXiv:2412.20404 [cs].
- [4] Wen Wang, Yan Jiang, Kangyang Xie, Zide Liu, Hao Chen, Yue Cao, Xinlong Wang, and Chunhua Shen. Zero-Shot Video Editing Using Off-The-Shelf Image Diffusion Models, January 2024. arXiv:2303.17599 [cs].
- [5] Giannis Daras, Weili Nie, Karsten Kreis, Alexandros G. Dimakis, Morteza Mardani, Nikola B. Kovachki, and Arash Vahdat. Warped Diffusion: Solving Video Inverse Problems with Image Diffusion Models. *Advances in Neural Information Processing Systems*, 37:101116–101143, December 2024.
- [6] Pan Du, Meet Hemant Parikh, Xiantao Fan, Xin-Yang Liu, and Jian-Xun Wang. Conditional neural field latent diffusion model for generating spatiotemporal turbulence. *Nature Communications*, 15(1):10416, November 2024.
- [7] T. Li, L. Biferale, F. Bonaccorso, M. A. Scarpolini, and M. Buzzicotti. Synthetic Lagrangian turbulence by generative diffusion models. *Nature Machine Intelligence*, 6(4):393–403, April 2024.
- [8] Marten Lienen, David Lüdke, Jan Hansen-Palmus, and Stephan Günemann. From Zero to Turbulence: Generative Modeling for 3D Flow Simulation, March 2024. arXiv:2306.01776 [physics].
- [9] Ilan Price, Alvaro Sanchez-Gonzalez, Ferran Alet, Tom R. Andersson, Andrew El-Kadi, Dominic Masters, Timo Ewalds, Jacklynn Stott, Shakir Mohamed, Peter Battaglia, Remi Lam, and Matthew Willson. Probabilistic weather forecasting with machine learning. *Nature*, 637(8044):84–90, January 2025. Publisher: Nature Publishing Group.
- [10] Lizao Li, Robert Carver, Ignacio Lopez-Gomez, Fei Sha, and John Anderson. Generative emulation of weather forecast ensembles with diffusion models. *Science Advances*, 10(13):eadk4489, March 2024. Publisher: American Association for the Advancement of Science.
- [11] Prakhar Srivastava, Ruihan Yang, Gavin Kerrigan, Gideon Dresdner, Jeremy McGibbon, Christopher S. Bretherton, and Stephan Mandt. Precipitation Downscaling with Spatiotemporal Video Diffusion. *Advances in Neural Information Processing Systems*, 37:56374–56400, December 2024.
- [12] Seth Basetti, Brian Hutchinson, Claudia Tebaldi, and Ben Kravitz. DiffESM: Conditional Emulation of Temperature and Precipitation in Earth System Models With 3D Diffusion Models. *Journal of Advances in Modeling Earth Systems*, 16(10):e2023MS004194, October 2024.
- [13] Salva Rühling Cachay, Brian Henn, Oliver Watt-Meyer, Christopher S. Bretherton, and Rose Yu. Probabilistic Emulation of a Global Climate Model with Spherical Diffusion. *Advances in Neural Information Processing Systems*, 37:127610–127644, December 2024.
- [14] Guillaume Couairon, Renu Singh, Anastase Charantonis, Christian Lessig, and Claire Monteleoni. ArchesWeather & ArchesWeatherGen: a deterministic and generative model for efficient ML weather forecasting, December 2024. arXiv:2412.12971 [cs].

- [15] Yang Song and Stefano Ermon. Generative Modeling by Estimating Gradients of the Data Distribution. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. In *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020.
- [17] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-Based Generative Modeling through Stochastic Differential Equations, February 2021. arXiv:2011.13456 [cs, stat].
- [18] Agrim Gupta, Lijun Yu, Kihyuk Sohn, Xiuye Gu, Meera Hahn, Li Fei-Fei, Irfan Essa, Lu Jiang, and José Lezama. Photorealistic Video Generation with Diffusion Models, December 2023. arXiv:2312.06662 [cs].
- [19] Xiang Wang, Hangjie Yuan, Shiwei Zhang, Dayou Chen, Jiuniu Wang, Yingya Zhang, Yujun Shen, Deli Zhao, and Jingren Zhou. VideoComposer: Compositional Video Synthesis with Motion Controllability, June 2023. arXiv:2306.02018 [cs].
- [20] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, Da Yin, Xiaotao Gu, Yuxuan Zhang, Weihai Wang, Yean Cheng, Ting Liu, Bin Xu, Yuxiao Dong, and Jie Tang. CogVideoX: Text-to-Video Diffusion Models with An Expert Transformer, October 2024. arXiv:2408.06072 [cs].
- [21] Haoge Deng, Ting Pan, Haiwen Diao, Zhengxiong Luo, Yufeng Cui, Huchuan Lu, Shiguang Shan, Yonggang Qi, and Xinlong Wang. Autoregressive Video Generation without Vector Quantization, March 2025. arXiv:2412.14169 [cs].
- [22] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, Kathrina Wu, Qin Lin, Junkun Yuan, Yanxin Long, Aladdin Wang, Andong Wang, Changlin Li, Duoju Huang, Fang Yang, Hao Tan, Hongmei Wang, Jacob Song, Jiawang Bai, Jianbing Wu, Jinbao Xue, Joey Wang, Kai Wang, Mengyang Liu, Pengyu Li, Shuai Li, Weiyan Wang, Wenqing Yu, Xincheng Deng, Yang Li, Yi Chen, Yutao Cui, Yuanbo Peng, Zhentao Yu, Zhiyu He, Zhiyong Xu, Zixiang Zhou, Zunnan Xu, Yangyu Tao, Qinglin Lu, Songtao Liu, Dax Zhou, Hongfa Wang, Yong Yang, Di Wang, Yuhong Liu, Jie Jiang, and Caesar Zhong. HunyuanVideo: A Systematic Framework For Large Video Generative Models, January 2025. arXiv:2412.03603 [cs].
- [23] Zijun Deng, Xiangteng He, and Yuxin Peng. Efficiency-optimized Video Diffusion Models. In *Proceedings of the 31st ACM International Conference on Multimedia*, MM '23, pages 7295–7303, New York, NY, USA, October 2023. Association for Computing Machinery.
- [24] Wenhao Sun, Rong-Cheng Tu, Jingyi Liao, and Dacheng Tao. Diffusion Model-Based Video Editing: A Survey, June 2024. arXiv:2407.07111 [cs].
- [25] Pascal Chang, Jingwei Tang, Markus Gross, and Vinicius C. Azevedo. How I Warped Your Noise: a Temporally-Correlated Noise Prior for Diffusion Models. In *The Twelfth International Conference on Learning Representations*, October 2023.
- [26] Duygu Ceylan, Chun-Hao P. Huang, and Niloy J. Mitra. Pix2Video: Video Editing using Image Diffusion. pages 23206–23217, 2023.
- [27] Yabo Zhang, Yuxiang Wei, Dongsheng Jiang, Xiaopeng Zhang, Wangmeng Zuo, and Qi Tian. ControlVideo: Training-free Controllable Text-to-video Generation. October 2023.
- [28] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei, Xintao Wang, Ying Shan, and Qifeng Chen. FateZero: Fusing Attentions for Zero-shot Text-based Video Editing. pages 15932–15942, 2023.
- [29] Zhong Yi Wan, Ricardo Baptista, Yi-fan Chen, John Anderson, Anudhyan Boral, Fei Sha, and Leonardo Zepeda-Núñez. Debias Coarsely, Sample Conditionally: Statistical Downscaling through Optimal Transport and Probabilistic Diffusion Models, May 2023. arXiv:2305.15618 [physics].

- [30] Tobias Bischoff and Katherine Deck. Unpaired Downscaling of Fluid Flows with Diffusion Bridges. *Artificial Intelligence for the Earth Systems*, 3(2), May 2024. Publisher: American Meteorological Society Section: Artificial Intelligence for the Earth Systems.
- [31] Morteza Mardani, Noah Brenowitz, Yair Cohen, Jaideep Pathak, Chieh-Yu Chen, Cheng-Chin Liu, Arash Vahdat, Karthik Kashinath, Jan Kautz, and Mike Pritchard. Residual Diffusion Modeling for Km-scale Atmospheric Downscaling, January 2024. ISSN: 2693-5015.
- [32] Étienne Pléziat, Robert J. H. Dunn, Markus G. Donat, and Christopher Kadow. Artificial intelligence reveals past climate extremes by reconstructing historical records. *Nature Communications*, 15(1):9191, October 2024. Publisher: Nature Publishing Group.
- [33] Philipp Hess, Michael Aich, Baoxiang Pan, and Niklas Boers. Fast, scale-adaptive and uncertainty-aware downscaling of Earth system model fields with generative machine learning. *Nature Machine Intelligence*, 7(3):363–373, March 2025. Publisher: Nature Publishing Group.
- [34] Michael Aich, Philipp Hess, Baoxiang Pan, Sebastian Bathiany, Yu Huang, and Niklas Boers. Conditional diffusion models for downscaling & bias correction of Earth system model precipitation, April 2024. arXiv:2404.14416 [physics].
- [35] Fenghua Ling, Zeyu Lu, Jing-Jia Luo, Lei Bai, Swadhin K. Behera, Dachao Jin, Baoxiang Pan, Huidong Jiang, and Toshio Yamagata. Diffusion model-based probabilistic downscaling for 180-year East Asian climate reconstruction. *npj Climate and Atmospheric Science*, 7(1):1–11, June 2024. Publisher: Nature Publishing Group.
- [36] Henry Addison, Elizabeth Kendon, Suman Ravuri, Laurence Aitchison, and Peter AG Watson. Machine learning emulation of precipitation from km-scale regional climate simulations using a diffusion model, July 2024. arXiv:2407.14158 [physics].
- [37] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-A-Video: Text-to-Video Generation without Text-Video Data. September 2022.
- [38] Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis. Structure and Content-Guided Video Synthesis with Diffusion Models. pages 7346–7356, 2023.
- [39] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-A-Video: One-Shot Tuning of Image Diffusion Models for Text-to-Video Generation. pages 7623–7633, 2023.
- [40] Omer Bar-Tal, Hila Chefer, Omer Tov, Charles Herrmann, Roni Paiss, Shiran Zada, Ariel Ephrat, Junhwa Hur, Guanghui Liu, Amit Raj, Yuanzhen Li, Michael Rubinstein, Tomer Michaeli, Oliver Wang, Deqing Sun, Tali Dekel, and Inbar Mosseri. Lumiere: A Space-Time Diffusion Model for Video Generation. In *SIGGRAPH Asia 2024 Conference Papers*, SA ’24, pages 1–11, New York, NY, USA, December 2024. Association for Computing Machinery.
- [41] Aidan Clark, Jeff Donahue, and Karen Simonyan. Adversarial Video Generation on Complex Datasets, September 2019. arXiv:1907.06571 [cs].
- [42] Suman Ravuri, Karel Lenc, Matthew Willson, Dmitry Kangin, Remi Lam, Piotr Mirowski, Megan Fitzsimons, Maria Athanassiadou, Sheleem Kashem, Sam Madge, Rachel Prudden, Amol Mandhane, Aidan Clark, Andrew Brock, Karen Simonyan, Raia Hadsell, Niall Robinson, Ellen Clancy, Alberto Arribas, and Shakir Mohamed. Skilful precipitation nowcasting using deep generative models of radar. *Nature*, 597(7878):672–677, 2021. arXiv: 2104.00954 Publisher: Springer US.
- [43] Puja Das, August Posch, Nathan Barber, Michael Hicks, Kate Duffy, Thomas Vandal, Debjani Singh, Katie van Werkhoven, and Auroop R. Ganguly. Hybrid physics-AI outperforms numerical weather prediction for extreme precipitation nowcasting. *npj Climate and Atmospheric Science*, 7(1):1–15, November 2024. Publisher: Nature Publishing Group.

- [44] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating Videos with Scene Dynamics. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [45] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error, February 2016. arXiv:1511.05440 [cs].
- [46] Masaki Saito, Eiichi Matsumoto, and Shunta Saito. Temporal Generative Adversarial Nets With Singular Value Clipping. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2830–2839, 2017.
- [47] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. MoCoGAN: Decomposing Motion and Content for Video Generation. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1526–1535, Salt Lake City, UT, June 2018. IEEE.
- [48] Masaki Saito, Shunta Saito, Masanori Koyama, and Sosuke Kobayashi. Train Sparsely, Generate Densely: Memory-Efficient Unsupervised Training of High-Resolution Temporal GAN. *International Journal of Computer Vision*, 128(10):2586–2606, November 2020.
- [49] You Xie, Erik Franz, Mengyu Chu, and Nils Thuerey. tempoGAN: a temporally coherent, volumetric GAN for super-resolution fluid flow. *ACM Trans. Graph.*, 37(4):95:1–95:15, July 2018.
- [50] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in neural information processing systems*, volume 34, pages 8780–8794. Curran Associates, Inc., 2021.
- [51] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent Video Diffusion Models for High-Fidelity Long Video Generation, March 2023. arXiv:2211.13221 [cs].
- [52] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin Rombach. Stable Video Diffusion: Scaling Latent Video Diffusion Models to Large Datasets, November 2023. arXiv:2311.15127 [cs].
- [53] Xin Ma, Yaohui Wang, Gengyun Jia, Xinyuan Chen, Ziwei Liu, Yuan-Fang Li, Cunjian Chen, and Yu Qiao. Latte: Latent Diffusion Transformer for Video Generation, January 2024. arXiv:2401.03048 [cs].
- [54] Yang Jin, Zhicheng Sun, Ningyuan Li, Kun Xu, Kun Xu, Hao Jiang, Nan Zhuang, Quzhe Huang, Yang Song, Yadong Mu, and Zhouchen Lin. Pyramidal Flow Matching for Efficient Video Generative Modeling, October 2024. arXiv:2410.05954 [cs].
- [55] Kiwhan Song, Boyuan Chen, Max Simchowitz, Yilun Du, Russ Tedrake, and Vincent Sitzmann. History-Guided Video Diffusion, February 2025. arXiv:2502.06764 [cs].
- [56] Shangshang Yang, Congyi Nai, Xinyan Liu, Weidong Li, Jie Chao, Jingnan Wang, Leyi Wang, Xichen Li, Xi Chen, Bo Lu, Ziniu Xiao, Niklas Boers, Huiling Yuan, and Baoxiang Pan. Generative assimilation and prediction for weather and climate, March 2025. arXiv:2503.03038 [cs].
- [57] Zeyu Li, Wang Han, Yue Zhang, Qingfei Fu, Jingxuan Li, Lizi Qin, Ruoyu Dong, Hao Sun, Yue Deng, and Lijun Yang. Learning spatiotemporal dynamics with a pretrained generative model. *Nature Machine Intelligence*, 6(12):1566–1579, December 2024. Publisher: Nature Publishing Group.
- [58] Jason Stock, Jaideep Pathak, Yair Cohen, Mike Pritchard, Piyush Garg, Dale Durran, Morteza Mardani, and Noah Brenowitz. DiffObs: Generative Diffusion for Global Forecasting of Satellite Observations, April 2024. arXiv:2404.06517.
- [59] Ilan Price, Alvaro Sanchez-Gonzalez, Ferran Alet, Tom R. Andersson, Andrew El-Kadi, Dominic Masters, Timo Ewalds, Jacklynn Stott, Shakir Mohamed, Peter Battaglia, Remi Lam, and Matthew Willson. GenCast: Diffusion-based ensemble forecasting for medium-range weather, May 2024. arXiv:2312.15796 [physics].

- [60] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. TokenFlow: Consistent Diffusion Features for Consistent Video Editing. In *The Twelfth International Conference on Learning Representations*, October 2023.
- [61] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2Video-Zero: Text-to-Image Diffusion Models are Zero-Shot Video Generators. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15954–15964, 2023.
- [62] François Rozet and Gilles Louppe. Score-based Data Assimilation for a Two-Layer Quasi-Geostrophic Model, November 2023. arXiv:2310.01853 [stat].
- [63] Jonathan Schmidt, Luca Schmidt, Felix Strnad, Nicole Ludwig, and Philipp Hennig. Spatiotemporally Coherent Probabilistic Generation of Weather from Climate, January 2025. arXiv:2412.15361 [cs].
- [64] Zhanxiang Hua, Yutong He, Chengqian Ma, and Alexandra Anderson-Frey. Weather Prediction with Diffusion Guided by Realistic Forecast Processes, February 2024. arXiv:2402.06666 [physics].
- [65] Langwen Huang, Lukas Gianinazzi, Yuejiang Yu, Peter D. Dueben, and Torsten Hoefler. DiffDA: a diffusion model for weather-scale data assimilation. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *ICML’24*, pages 19798–19815, Vienna, Austria, July 2024. JMLR.org.
- [66] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align Your Latents: High-Resolution Video Synthesis With Latent Diffusion Models. pages 22563–22575, 2023.
- [67] Zhendong Wang, Huangjie Zheng, Pengcheng He, Weizhu Chen, and Mingyuan Zhou. Diffusion-GAN: Training GANs with Diffusion. September 2022.
- [68] Runhui Huang, Jianhua Han, Guansong Lu, Xiaodan Liang, Yihan Zeng, Wei Zhang, and Hang Xu. DiffDis: Empowering Generative Diffusion Model with Cross-Modal Discrimination Capability. pages 15713–15723, 2023.
- [69] Myeongjin Ko, Euiyeon Kim, and Yong-Hoon Choi. Adversarial Training of Denoising Diffusion Model Using Dual Discriminators for High-Fidelity Multi-Speaker TTS. *IEEE Open Journal of Signal Processing*, 5:577–587, 2024. Conference Name: IEEE Open Journal of Signal Processing.
- [70] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial Diffusion Distillation. In Aleš Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol, editors, *Computer Vision – ECCV 2024*, pages 87–103, Cham, 2024. Springer Nature Switzerland.
- [71] Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, Eli Shechtman, Fredo Durand, and Bill Freeman. Improved Distribution Matching Distillation for Fast Image Synthesis. *Advances in Neural Information Processing Systems*, 37:47455–47487, December 2024.
- [72] Dongjun Kim, Yeongmin Kim, Se Jung Kwon, Wanmo Kang, and Il-Chul Moon. Refining Generative Process with Discriminator Guidance in Score-based Diffusion Models, June 2023. arXiv:2211.17091 [cs].
- [73] Akio Hayakawa, Masato Ishii, Takashi Shibuya, and Yuki Mitsufuji. Discriminator-Guided Cooperative Diffusion for Joint Audio and Video Generation, May 2024. arXiv:2405.17842.
- [74] Filip Ekström Kelvinius and Fredrik Lindsten. Discriminator Guidance for Autoregressive Diffusion Models. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, pages 3403–3411. PMLR, April 2024. ISSN: 2640-3498.
- [75] Thomas J. Kerby and Kevin R. Moon. Training-Free Guidance for Discrete Diffusion Models for Molecular Generation, September 2024. arXiv:2409.07359 [stat].

- [76] Brian D. O. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, May 1982.
- [77] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the Design Space of Diffusion-Based Generative Models, October 2022. arXiv:2206.00364 [cs, stat].
- [78] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- [79] Hans Hersbach, Bill Bell, Paul Berrisford, Shoji Hirahara, András Horányi, Joaquín Muñoz-Sabater, Julien Nicolas, Carole Peubey, Raluca Radu, Dinand Schepers, Adrian Simmons, Cornel Soci, Saleh Abdalla, Xavier Abellan, Gianpaolo Balsamo, Peter Bechtold, Gionata Biavati, Jean Bidlot, Massimo Bonavita, Giovanna De Chiara, Per Dahlgren, Dick Dee, Michail Diamantakis, Rossana Dragani, Johannes Flemming, Richard Forbes, Manuel Fuentes, Alan Geer, Leo Haimberger, Sean Healy, Robin J. Hogan, Elías Hólm, Marta Janisková, Sarah Keeley, Patrick Laloyaux, Philippe Lopez, Cristina Lupu, Gabor Radnoti, Patricia de Rosnay, Iryna Rozum, Freja Vamborg, Sebastien Villaume, and Jean-Noël Thépaut. The ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730):1999–2049, 2020. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/qj.3803>.
- [80] Ernest Hovmöller. The Trough-and-Ridge diagram. *Tellus*, 1(2):62–66, 1949. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.2153-3490.1949.tb01260.x>.
- [81] Tilmann Gneiting, Fadoua Balabdaoui, and Adrian E. Raftery. Probabilistic forecasts, calibration and sharpness. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(2):243–268, 2007. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-9868.2007.00587.x>.
- [82] Min-Seop Ahn, Daehyun Kim, Daehyun Kang, Jiwoo Lee, Kenneth R. Sperber, Peter J. Gleckler, Xianan Jiang, Yoo-Geun Ham, and Hyemi Kim. MJO Propagation Across the Maritime Continent: Are CMIP6 Models Better Than CMIP5 Models? *Geophysical Research Letters*, 47(11):e2020GL087250, 2020. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1029/2020GL087250>.
- [83] Ashesh Chattopadhyay, Y. Qiang Sun, and Pedram Hassanzadeh. Challenges of learning multi-scale dynamics with AI weather models: Implications for stability and one solution, December 2024. arXiv:2304.07029 [physics].
- [84] Bo Wang, Li-Lian Wang, and Ziqing Xie. Accurate calculation of spherical and vector spherical harmonic expansions via spectral element grids. *Advances in Computational Mathematics*, 44(3):951–985, June 2018.
- [85] Navid C. Constantinou, Gregory LeClaire Wagner, Lia Siegelman, Brodie C. Pearson, and André Palóczy. GeophysicalFlows.jl: Solvers for geophysical fluid dynamics problems in periodic domains on CPUs & GPUs. *Journal of Open Source Software*, 6(60):3053, April 2021.
- [86] V. Fortin, M. Abaza, F. Anctil, and R. Turcotte. Why Should Ensemble Spread Match the RMSE of the Ensemble Mean? *Journal of Hydrometeorology*, 15(4):1708–1713, August 2014. Publisher: American Meteorological Society Section: Journal of Hydrometeorology.

A Diffusion models

For all diffusion models trained in this work, we use the DDPM++ UNet [17] with hyperparameters adapted as shown in Tab. 1. The UNets have around 33.8M parameters for the vorticity configuration and 25.7M parameters in the precipitation configuration. For inference, we use weights from an exponential moving average (EMA). We extend the architecture by applying periodic padding in both spatial dimensions for the vorticity simulation and in the longitude direction for the ERA5 precipitation data. We use the EDM preconditioning from Karras et al. 2022 [77], which is given for the unconditional DM by

$$\mathcal{S}_\phi(\mathbf{x}_t; t) := c_{\text{skip}}(t)\mathbf{x}_t + c_{\text{out}}(t)\mathbf{f}_\phi(c_{\text{in}}(t)\mathbf{x}_t; c_{\text{noise}}(t)), \quad (7)$$

where $\mathbf{f}_\phi(\cdot)$ is the UNet denoiser network and the coefficients are defined in [77]. For training, we use a log-normal distribution to sample the noise levels $\ln(\sigma) \sim \mathcal{N}(P_{\text{mean}}, P_{\text{std}}^2)$ [77]. For sampling, we also use the EDM noise schedule [77], defined as

$$\sigma_i := \left(\sigma_{\text{max}}^{1/\rho} + \frac{i}{N-1} \left(\sigma_{\text{min}}^{1/\rho} - \sigma_{\text{max}}^{1/\rho} \right) \right)^\rho, \quad (8)$$

where $i \in \{0, \dots, N-1\}$, N is the number of sampling steps and σ_{min} , σ_{max} , ρ are defined in Tab. 1. The training takes about 21 min per epoch on a H100 Nvidia GPU for the Navier-Stokes dataset and 8 min per epoch for the ERA5 precipitation dataset. We use the stochastic EDM sampler [77] with the parameters given in Tab. 1 for the video, unconditional and guided diffusion model. A single generation step takes on average 0.19 and 0.18 seconds for the vorticity and precipitation fields, respectively, on a H100 GPU.

Table 1: Configuration details of the diffusion model architectures, training and sampling.

	2D vorticity	ERA5 precipitation
Architecture		
Input dimension (unconditional)	(B,1,256,256)	(B,1,180,360)
Input dimension (video)	(B,3,256,256)	(B,3,180,360)
Output dimension	(B,1,256,256)	(B,1,180,360)
Num. Resnet blocks	3	2
Num. attention blocks	3	2
Attention resolution	8,4	8,4
Channel	(1,2,2)	(1,2,2)
Channel multiplier	128	128
Training		
Batch size (B)	2	2
Learning rate (LR)	10^{-4}	10^{-4}
Optimizer	AdamW	AdamW
Epochs	350	250
EMA rate	0.9999	0.9999
Sampling		
σ_{min}	0.002	0.002
σ_{max}	80	80
σ_{data}	0.5	0.5
ρ	7	7
P_{mean}	-1.2	-1.2
P_{std}	1.2	1.2
S_{tmin}	0	0
S_{tmax}	1000	1000
S_{noise}	1.005	1.0
S_{churn}	55	80
Num. steps	50	100
λ (guidance strength)	14	68

B Discriminator model

B.1 Time-consistency guidance

Given a time series of length N , $\{\mathbf{x}^n | n = 1, 2, \dots, N\}$, containing images $\mathbf{x}^n \in \mathbb{R}^{B \times C \times H \times W}$ that are temporally ordered, we aim to train a discriminator network to classify whether a shown sample \mathbf{x}^{n+1} is in temporal order with respect to the current and a sequence of m previous time frames $\{\mathbf{x}^{n-m}, \mathbf{x}^{n-m+1}, \dots, \mathbf{x}^n\} = \mathbf{x}^{(n-m):n}$ or not. In other words, whether a sample is drawn from the conditional distribution $\mathbf{x}^{n+1} \sim p(\mathbf{x}^{n+1} | \mathbf{x}^{(n-m):n})$ or the unconditional distribution $\mathbf{x}^{n+1} \sim p(\mathbf{x}^{n+1})$. Here and in the following we drop the explicit dependency on the noise time t . As shown in [78], an optimal discriminator with parameters θ can then be written as

$$D_\theta(\mathbf{x}^{n+1}; \mathbf{x}^{(n-m):n}) = \frac{p(\mathbf{x}^{n+1} | \mathbf{x}^{(n-m):n})}{p(\mathbf{x}^{n+1} | \mathbf{x}^{(n-m):n}) + p(\mathbf{x}^{n+1})}, \quad (9)$$

which we can rewrite as

$$\frac{D_\theta(\mathbf{x}^{n+1}; \mathbf{x}^{(n-m):n})}{1 - D_\theta(\mathbf{x}^{n+1}; \mathbf{x}^{(n-m):n})} = \frac{p(\mathbf{x}^{n+1} | \mathbf{x}^{(n-m):n})}{p(\mathbf{x}^{n+1})}. \quad (10)$$

Taking the log and computing the gradient with respect to \mathbf{x}^{n+1} gives

$$\nabla_{\mathbf{x}^{n+1}} \log \left(\frac{D_\theta(\mathbf{x}^{n+1}; \mathbf{x}^{(n-m):n})}{1 - D_\theta(\mathbf{x}^{n+1}; \mathbf{x}^{(n-m):n})} \right) = \nabla_{\mathbf{x}^{n+1}} \log \left(\frac{p(\mathbf{x}^{n+1} | \mathbf{x}^{(n-m):n})}{p(\mathbf{x}^{n+1})} \right), \quad (11)$$

which we use for our guidance term:

$$\mathbf{d}_\theta(\mathbf{x}^{n+1}; \mathbf{x}^{(n-m):n}) := \nabla_{\mathbf{x}^{n+1}} \log \left(\frac{p(\mathbf{x}^{n+1} | \mathbf{x}^{(n-m):n})}{p(\mathbf{x}^{n+1})} \right). \quad (12)$$

Using the expression

$$p(\mathbf{x}^{n+1} | \mathbf{x}^{(n-m):n}) = p(\mathbf{x}^{n+1}) \frac{p(\mathbf{x}^{n+1} | \mathbf{x}^{(n-m):n})}{p(\mathbf{x}^{n+1})}, \quad (13)$$

and computing the score functions gives

$$\begin{aligned} \nabla_{\mathbf{x}^{n+1}} \log p(\mathbf{x}^{n+1} | \mathbf{x}^{(n-m):n}) &= \nabla_{\mathbf{x}^{n+1}} \log p(\mathbf{x}^{n+1}) \\ &\quad + \nabla_{\mathbf{x}^{n+1}} \log \left(\frac{p(\mathbf{x}^{n+1} | \mathbf{x}^{(n-m):n})}{p(\mathbf{x}^{n+1})} \right). \end{aligned} \quad (14)$$

From Eq. 14 we see that we can approximate the conditional score $\nabla_{\mathbf{x}^{n+1}} \log p(\mathbf{x}^{n+1} | \mathbf{x}^{(n-m):n})$ with an unconditional score model $\nabla_{\mathbf{x}^{n+1}} \log p(\mathbf{x}^{n+1}) \approx \mathbf{s}_\theta(\mathbf{x}^{n+1})$ and the guidance in Eq. 12, as

$$\nabla_{\mathbf{x}^{n+1}} \log p(\mathbf{x}^{n+1} | \mathbf{x}^{(n-m):n}) \approx \mathbf{s}_\theta(\mathbf{x}^{n+1}) + \mathbf{d}_\theta(\mathbf{x}^{n+1}; \mathbf{x}^{(n-m):n}). \quad (15)$$

B.2 Network and training

We adapt the noise time-conditioned encoder part of the UNet from [72] for the discriminator model with a two-layer fully connected network as a decoder (see Tab. 1 for hyperparameter configurations). We add periodic padding in both spatial dimensions for the vorticity simulation and in the longitude direction for the ERA5 precipitation data. The discriminator network has 1.7M parameters in the vorticity configuration and 7.5M in the precipitation configuration, making it around 19.8 and 3.4 times smaller, respectively, than the diffusion networks. Training the discriminator takes about 1.2 minutes on 2 H100 GPUs per epoch for the Navier-Stokes dataset and 0.8 min per epoch for the precipitation dataset. The evaluation of the guidance term in Eq. 4 is computationally much cheaper than a generative sampling step, taking on average 0.007 and 0.016 seconds for the vorticity and precipitation fields, respectively, on the H100 GPU. Hence, the discriminator guidance evaluation corresponds to around 3% and 8% of the respective generation time.

B.3 Guided sampling

We adapt the stochastic discriminator guidance sampler [77, 72], which solves the sampling ODE with stochastic churn with a second-order accurate solver. Stochasticity is controlled through the parameters $S_{\text{noise}}, S_{\text{churn}}, S_{t_{\text{max}}}, S_{t_{\text{min}}}$ in the coefficient γ_i with

$$\gamma_i = \begin{cases} \min\left(\frac{S_{\text{churn}}}{T}, \sqrt{2} - 1\right) & \text{if } t_i \in [S_{t_{\text{min}}}, S_{t_{\text{max}}}], \\ 0 & \text{otherwise.} \end{cases} \quad (16)$$

We provide pseudocode of the implementation in Alg. 1.

Algorithm 1 Time-consistency (TC) guided sampling (adapted from [77, 72])

```

1: input:  $S_\phi, D_\theta, \mathbf{x}^n, \mathbf{x}^{n-1}, \lambda, t_{i \in \{0, \dots, T\}}, \gamma_{i \in \{0, \dots, T-1\}}, S_{\text{noise}}$ .
2: sample  $\mathbf{x}_T^{n+1} \sim \mathcal{N}(\mathbf{0}, t_T^2 \mathbf{I})$ 
3: for  $i = 0$  to  $T$  do
4:   sample  $\epsilon_i \sim \mathcal{N}(\mathbf{0}, S_{\text{noise}}^2 \mathbf{I})$ 
5:    $\hat{t}_i \leftarrow t_i + \gamma_i t_i$ 
6:    $\hat{\mathbf{x}}_i^{n+1} \leftarrow \mathbf{x}_i^{n+1} + \sqrt{\hat{t}_i^2 - t_i^2} \epsilon_i$ 
7:    $\mathbf{s}_i \leftarrow (\hat{\mathbf{x}}_i^{n+1} - S_\phi(\hat{\mathbf{x}}_i^{n+1}; \hat{t}_i)) / \hat{t}_i$ 
8:    $\mathbf{d}_i \leftarrow -t_i \nabla_{\hat{\mathbf{x}}_i^{n+1}} \log \left( \frac{D_\theta(\hat{\mathbf{x}}_i^{n+1}; \mathbf{x}^n, \mathbf{x}^{n-1}, \hat{t}_i)}{1 - D_\theta(\hat{\mathbf{x}}_i^{n+1}; \mathbf{x}^n, \mathbf{x}^{n-1}, \hat{t}_i)} \right)$  ▷ TC guidance
9:    $\mathbf{x}_{i+1} \leftarrow \hat{\mathbf{x}}_i + (t_{i+1} - \hat{t}_i)(\mathbf{s}_i + \lambda \mathbf{d}_i)$ 
10:  if  $t_{i+1} \neq 0$  then
11:     $\mathbf{s}'_i \leftarrow (\mathbf{x}_{i+1}^{n+1} - S_\phi(\mathbf{x}_{i+1}^{n+1}; \hat{t}_{i+1})) / \hat{t}_{i+1}$ 
12:     $\mathbf{d}'_i \leftarrow -t_{i+1} \nabla_{\mathbf{x}_{i+1}^{n+1}} \log \left( \frac{D_\theta(\mathbf{x}_{i+1}^{n+1}; \mathbf{x}^n, \mathbf{x}^{n-1}, t_{i+1})}{1 - D_\theta(\mathbf{x}_{i+1}^{n+1}; \mathbf{x}^n, \mathbf{x}^{n-1}, t_{i+1})} \right)$  ▷ TC guidance
13:     $\mathbf{x}_{i+1}^{n+1} \leftarrow \hat{\mathbf{x}}_{i+1}^{n+1} + (t_{i+1} - \hat{t}_i) \left[ \left(\frac{1}{2} \mathbf{s}_i + \lambda \mathbf{d}_i\right) + \frac{1}{2} (\mathbf{s}'_i + \lambda \mathbf{d}'_i) \right]$ 
14: return  $\mathbf{x}_T^{n+1}$ 

```

Table 2: Discriminator model architecture and training parameters.

	2D vorticity	ERA5 precipitation
Architecture		
Input dimension	(B,3,256,256)	(B,3,180,360)
Output dimension	(B,1)	(B,1)
Num. Resnet blocks	2	2
Num. attention blocks	2	2
Attention resolution	8,4	8,4
Channel	(1,2,2)	(4,2,1)
Channel multiplier	128	64
MLP layer size	2	2
Num. MLP layer	1024	1024
Training		
Batch size (B)	8	8
Learning rate (LR)	10^{-4}	10^{-4}
Optimizer	AdamW	AdamW
Epochs	500	500

C Data

C.1 2D Navier-Stokes experiments

We perform numerical simulations of the two-dimensional incompressible Navier-Stokes equations in vorticity stream function formulation using the GeophysicalFlow.jl Julia package [85]. The simulation

uses periodic boundary conditions, hyperviscosity $\nu = 2e^{-7}$ of second order, a linear drag coefficient of $\mu = 1e^{-1}$ and integration time step $\Delta t = 0.005$. We subsample the simulation saving every fourth time step to disk. We apply stochastic forcing defined with an Ornstein-Uhlenbeck process and a forcing wavenumber $k_f = 6 \cdot 2\pi/L$, where $L = 256$, forcing bandwidth $\delta_f = 1.5 \cdot 2\pi/L$ and an energy input rate of $\epsilon = 0.1$. We wait for 500 steps for the simulation to reach a statistical equilibrium. We then standardize the data by subtracting the mean and dividing by the standard deviation for training the diffusion and discriminator networks.

C.2 Precipitation data

As a challenging, real-world application, we use global daily precipitation fields from the ERA5 reanalysis dataset [79]. The data is openly available for download at the Copernicus Climate Data Store (<https://cds.climate.copernicus.eu/datasets/reanalysis-era5-single-levels?tab=overview>). We regrid the data to 1° horizontal spatial resolution using bilinear interpolation. As additional preprocessing steps we apply a log-transform with $\tilde{x} = \log(x + \epsilon) - \log(\epsilon)$, $\epsilon = 10^{-4}$, and normalization approximately into the range $[-1, 1]$ for the diffusion and discriminator networks. We split the daily data into periods of 1979-2000, 2001-2010 and 2011-2020 for training, validation and testing.

D Evaluation metrics

We denote the ground truth and predicted spatial fields with $y_{k,l}^n$ and $x_{k,l}^n$, respectively, where $n = 1, \dots, N$ is the time index, $k = 1, \dots, K$ is the height (or latitude) index and $l = 1, \dots, L$ is the width (or longitude) index. We weigh the spherical data with a factor

$$w(k) = \frac{\cos(\text{lat}(k))}{\frac{1}{K} \sum_{i=1}^K \cos(\text{lat}(i))}$$

that accounts for the spherical geometry of the precipitation data and set $w(k) = 1$ for the vorticity experiments.

D.1 Deterministic metrics

Root mean square error. We define a spatially weighted root mean square error (RMSE) as

$$\text{RMSE} := \sqrt{\frac{1}{L} \sum_{l=1}^L \frac{1}{K} \sum_{k=1}^K w(k) \sum_{n=1}^N (y_{k,l}^n - x_{k,l}^n)^2}. \quad (17)$$

Bias. The bias at each spatial location is defined as

$$\text{Bias}_{k,l} := \frac{1}{N} \sum_{n=1}^N (y_{k,l}^n - x_{k,l}^n). \quad (18)$$

Autocorrelation function. The global mean over local autocorrelation functions (ACFs) is calculated by first removing the monthly mean for seasonal adjustment and standardizing the time series and then computing the ACF with

$$\text{ACF}(j) := \frac{1}{L} \sum_{l=1}^L \frac{1}{K} \sum_{k=1}^K w(k) \frac{\frac{1}{N} \sum_{n=1}^N [(x_{k,l}^n - \bar{x}_{k,l})(x_{k,l}^{n-j} - \bar{x}_{k,l})]}{\sigma_{k,l}^2}, \quad (19)$$

where the bar denotes the temporal mean and $\sigma_{k,l}^2$ is the variance at a spatial location.

Wasserstein distance. We use the Wasserstein distances to compute changes between consecutive rows in the Hovmöller diagrams, which allow us to quantitatively assess their similarity. We, therefore, treat two rows with consecutive time steps as two tuples of probabilities, (p_1, \dots, p_K) and (q_1, \dots, q_K) ,

by taking their absolute value and normalizing them to sum to 1. We then compute the Wasserstein-1 distance with

$$W_1(P, Q) := \frac{1}{N} \sum_{i=1}^N |F_p(i) - F_q(i)|, \quad (20)$$

where F_p and F_q denote the cumulative distribution functions.

D.2 Probabilistic metrics

Continuous ranked probability score. The continuous ranked probability score (CRPS) is computed for an ensemble size B at a given time step n following [9],

$$\text{CRPS}^n := \frac{1}{L} \sum_{l=1}^L \frac{1}{K} \sum_{k=1}^K w(k) \left(\frac{1}{B} \sum_{b=1}^B |x_{k,l}^{n,b} - y_{k,l}^{n,b}| - \frac{1}{2B^2} \sum_{b=1}^B \sum_{b'=1}^B |x_{k,l}^{n,b} - x_{k,l}^{n,b'}| \right), \quad (21)$$

and a lower CRPS score is better.

Spread skill ratio. The ensemble spread for a single time step n is defined, as in [9],

$$\text{Spread}^n := \sqrt{\frac{1}{L} \sum_{l=1}^L \frac{1}{K} \sum_{k=1}^K w(k) \frac{1}{B-1} \sum_{b=1}^B \left(x_{k,l}^{n,b} - \tilde{x}_{k,l}^n \right)^2}, \quad (22)$$

where $\tilde{x}_{k,l}^n$ is the ensemble mean. The ensemble skill at time step n is then defined as

$$\text{Skill}^n := \sqrt{\frac{1}{L} \sum_{l=1}^L \frac{1}{K} \sum_{k=1}^K w(k) \left(y_{k,l}^n - \tilde{x}_{k,l}^n \right)^2}. \quad (23)$$

Assuming that the ensemble members are all exchangeable, the spread skill ratio is then defined [86, 59],

$$\text{Spread-skill-ratio} := \sqrt{\frac{M+1}{M}} \frac{\text{Spread}}{\text{Skill}}, \quad (24)$$

which should be close to 1 for a perfect forecast.

E Additional analysis

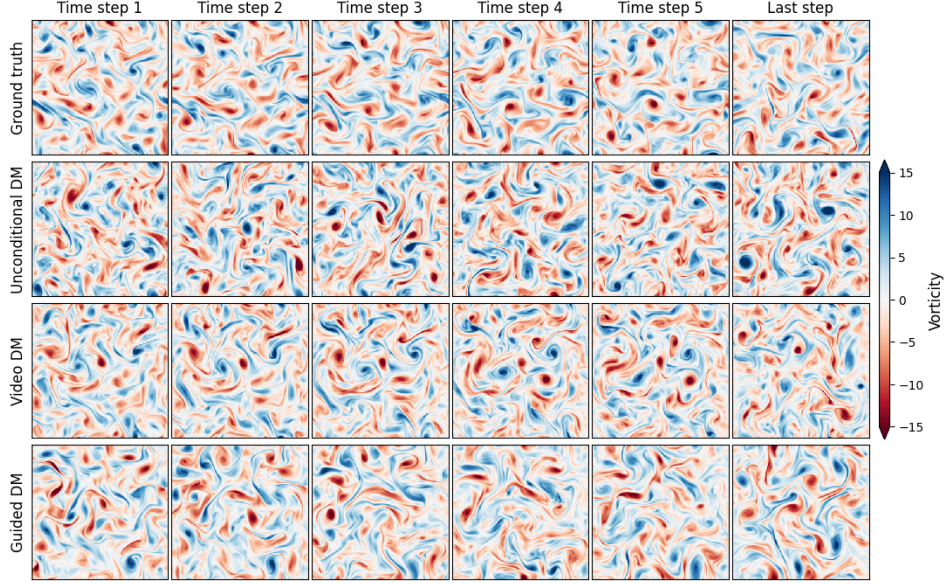


Figure 9: **Qualitative comparison** of the first five and last 2D vorticity fields from the direct numerical Navier-Stokes turbulence simulation (top), unconditional DM (upper middle), video DM (lower middle) and our discriminator guidance DM (bottom). Each row shows a single rollout starting from the same initial condition. Note that the pairing between the generated and ground truth samples decreases due to the chaotic dynamics.

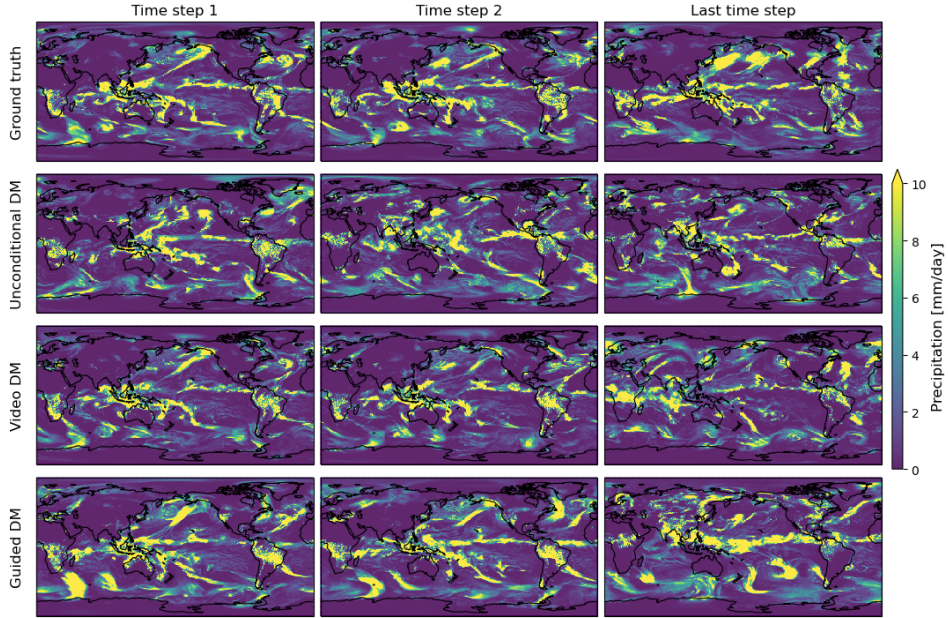


Figure 10: **Qualitative comparison** of the first two and last daily precipitation fields from the ERA5 ground truth (top), unconditional DM (upper middle), video DM (lower middle) and our discriminator guidance DM (bottom).

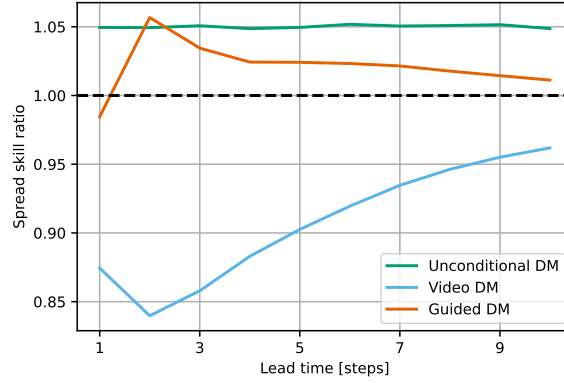


Figure 11: The spread skill ratio (SSR) of vorticity forecasts is shown for 100 ensemble forecasts with 50 members and 10-step lead time for the (blue) video DM and (red) guided DM. A perfect forecast would have a SSR of one.

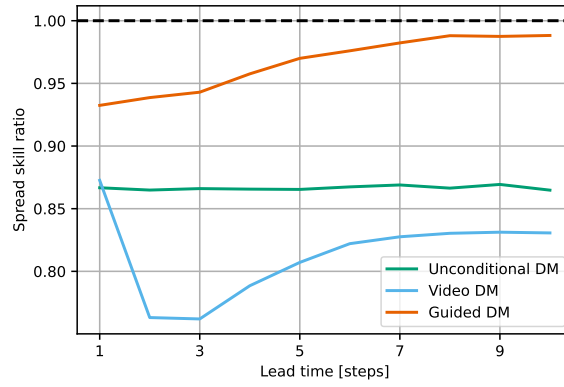


Figure 12: The spread skill ratio (SSR) of precipitation forecast is shown for 100 ensemble forecasts with 50 members and 10-step lead time for the (blue) video DM and (red) guided DM.

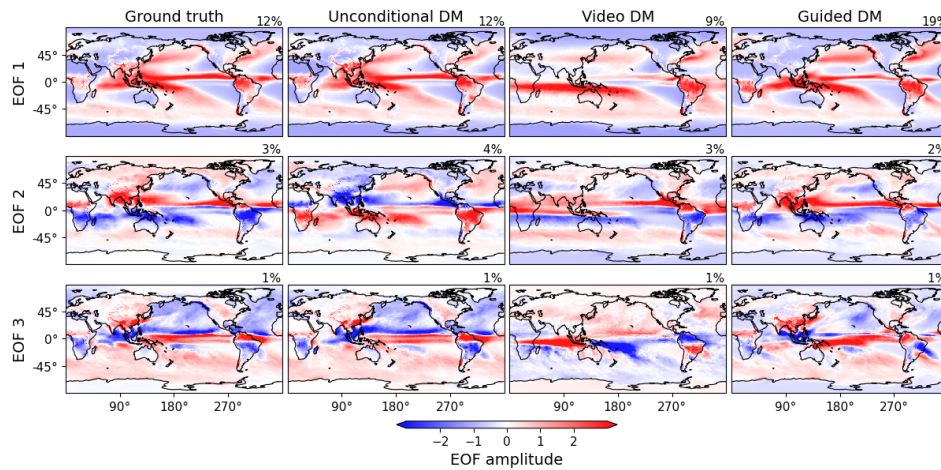


Figure 13: Empirical orthogonal functions (EOFs) are shown for the daily precipitation data for (from left to right) the ERA5 ground truth, the unconditional DM, the video DM and our guidance method. The explained variance is given in the top right of each panel.

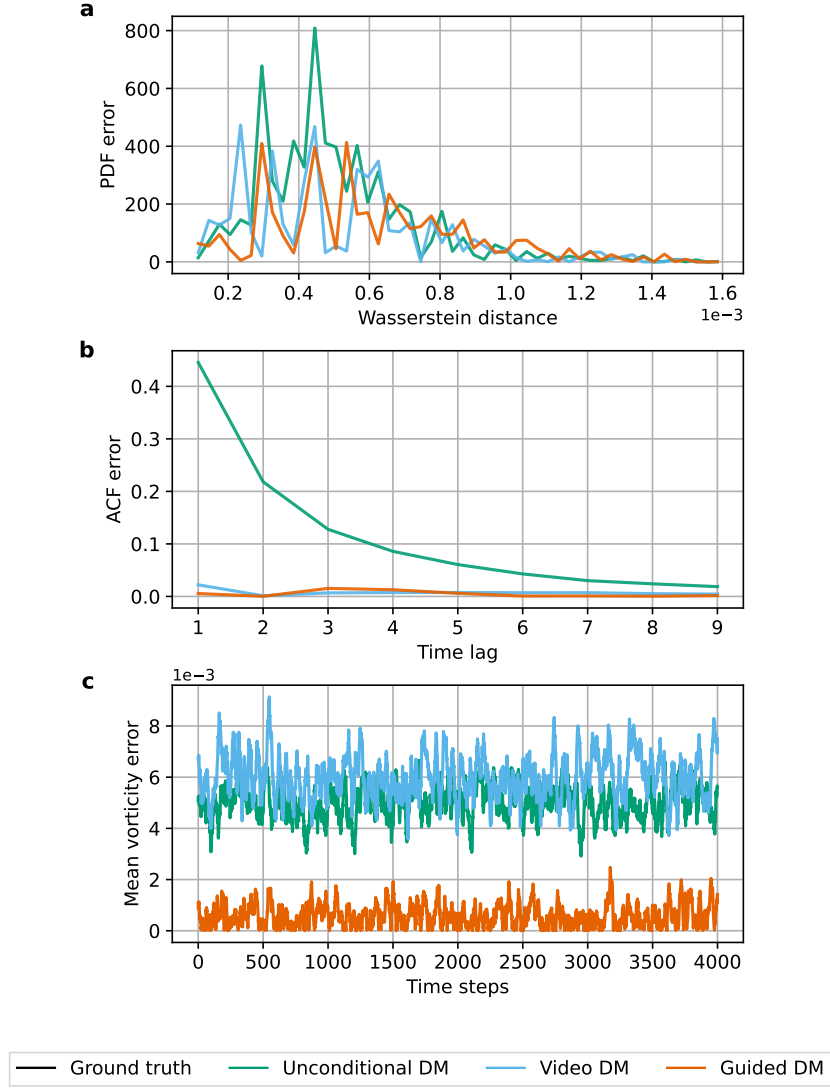


Figure 14: **Absolute errors of the vorticity statistics** shown in Fig. 4, (a) Wasserstein-1 distance, the (b) autocorrelation functions (ACFs), and (c) the global average are shown for the unconditional DM (green), the video DM (blue) and our guidance method (red).

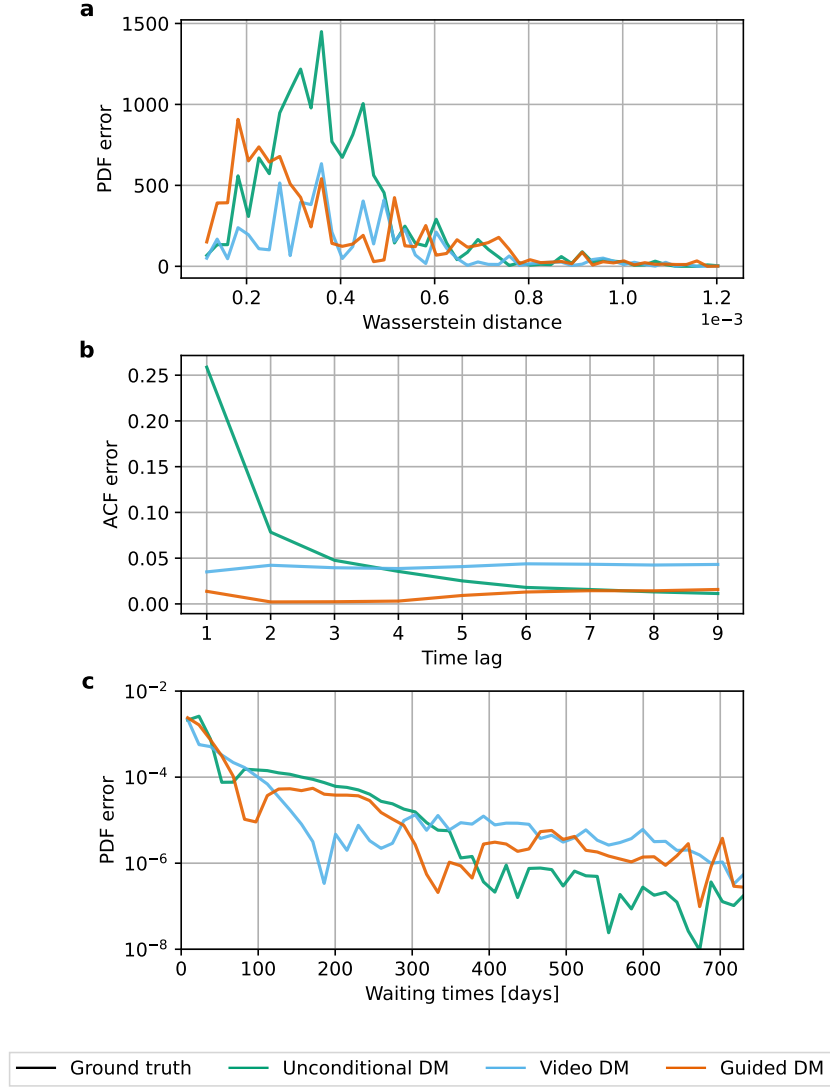


Figure 15: **Absolute errors of the precipitation statistics** shown in Fig. 6, (a) Wasserstein-1 distance, the (b) autocorrelation functions (ACFs), and (c) waiting time distributions are shown for the unconditional DM (green), the video DM (blue) and our guidance method (red).