# PreDiff: Precipitation Nowcasting with Latent Diffusion Models

**Zhihan Gao**[*]
Hong Kong University of Science and Technology
zhihan.gao@connect.ust.hk

**Xingjian Shi**[†]
Boson AI
xshiab@connect.ust.hk

**Boran Han**
AWS
boranhan@amazon.com

**Hao Wang**
AWS AI Labs
howngz@amazon.com

**Xiaoyong Jin**
Amazon
jxiaoyon@amazon.com

**Danielle Maddix**
AWS AI Labs
dmmaddix@amazon.com

**Yi Zhu**[†]
Boson AI
yi@boson.ai

**Mu Li**[†]
Boson AI
mu@boson.ai

**Yuyang Wang**
AWS AI Labs
yuyawang@amazon.com

## Abstract

Earth system forecasting has traditionally relied on complex physical models that are computationally expensive and require significant domain expertise. In the past decade, the unprecedented increase in spatiotemporal Earth observation data has enabled data-driven forecasting models using deep learning techniques. These models have shown promise for diverse Earth system forecasting tasks but either struggle with handling uncertainty or neglect domain-specific prior knowledge, resulting in averaging possible futures to blurred forecasts or generating physically implausible predictions. To address these limitations, we propose a two-stage pipeline for probabilistic spatiotemporal forecasting: 1) We develop PreDiff, a conditional latent diffusion model capable of probabilistic forecasts. 2) We incorporate an explicit knowledge control mechanism to align forecasts with domain-specific physical constraints. This is achieved by estimating the deviation from imposed constraints at each denoising step and adjusting the transition distribution accordingly. We conduct empirical studies on two datasets: $N$-body MNIST, a synthetic dataset with chaotic behavior, and SEVIR, a real-world precipitation nowcasting dataset. Specifically, we impose the law of conservation of energy in $N$-body MNIST and anticipated precipitation intensity in SEVIR. Experiments demonstrate the effectiveness of PreDiff in handling uncertainty, incorporating domain-specific prior knowledge, and generating forecasts that exhibit high operational utility.

## 1 Introduction

Earth's intricate climate system significantly influences daily life. Precipitation nowcasting, tasked with delivering accurate rainfall forecasts for the near future (e.g., 0-6 hours), is vital for decision-making across numerous industries and services. Recent advancements in data-driven deep learning (DL) techniques have demonstrated promising potential in this field, rivaling conventional numerical methods [9, 5] with their advantages of being more skillful [5], efficient [37], and scalable [3]. However, accurately predicting the future rainfall remains challenging for data-driven algorithms. The state-of-the-art Earth system forecasting algorithms [47, 58, 41, 37, 9, 65, 2, 29, 3] frequently

---

[*]Work conducted during an internship at Amazon. [†]Work conducted while at Amazon.

Preprint. Under review.

generates blurry predictions. This is caused by the high variability and complexity inherent to Earth's climatic system. Even minor differences in initial conditions can lead to vastly divergent outcomes that are difficult to predict. Most methods adopt a point estimation of the future rainfall and are trained by minimizing pixel-wise loss functions (e.g., mean-squared error). These methods lack the capability of capturing multiple plausible futures and will generate blurry forecasts which lose important operational details.

Therefore, what are needed instead are probabilistic models that can represent the uncertainty inherent in stochastic systems. The probabilistic models can capture multiple plausible futures, generating diverse high-quality predictions that better align with real-world data.

The emergence of diffusion models (DMs) [22] has enabled powerful probabilistic frameworks for generative modeling. DMs have shown remarkable capabilities in generating high-quality images [40, 45, 43] and videos [15, 23]. As a likelihood-based model, DMs do not exhibit mode collapse or training instabilities like GANs [11]. Compared to autoregressive (AR) models [53, 46, 59, 39, 61] that generate images pixel-by-pixel, DMs can produce higher resolution images faster and with higher quality. They also provide better uncertainty modeling without drawbacks like exposure bias in AR models. Latent diffusion models (LDMs) [42, 52] further improve on DMs by separating the model into two phases, only applying the costly diffusion in a compressed latent space. This alleviates the computational costs of DMs without significantly impairing performance.

Despite DMs' success in image and video generation [42, 15, 62, 35, 32, 56], its application to precipitation nowcasting and Earth system forecasting is in early stages [16]. One of the major concerns is that this purely data-centric approach lacks constraints and controls from prior knowledge about the dynamic system. Some spatiotemporal forecasting approaches have incorporated domain knowledge by modifying the model architecture or adding extra training losses [12, 1, 37]. This enables them to be aware of prior knowledge and generate physically plausible forecasts. However, these approaches still face challenges, such as requiring to design new model architectures or retrain the entire model from scratch when constraints change.

Inspired by recent success in controllable generative models [64, 24, 4, 33, 6], we propose a general two-stage pipeline for training data-driven earth system forecasting model. 1) In the first stage, we focus on capturing the intrinsic semantics in the data by training an LDM. To capture Earth's long-term and complex changes, we instantiate the LDM's core neural network as a UNet-style architecture based on Earthformer [9]. 2) In the second stage, we inject prior knowledge of the Earth system by training a knowledge control network that guides the sampling process of the LDM. Specifically the control network parameterizes an energy function that adjusts the transition probabilities during each denoising step. This encourages the generation of physically plausible intermediate latent states while suppressing those likely to violate the given domain knowledge. We summarize our main contributions as follows:

- We introduce a novel LDM based approach *PreDiff* for precipitation nowcasting.
- We propose a general two-stage pipeline for training data-driven earth system forecasting models. Specifically, we develop a knowledge control mechanism to guide the latent denoising diffusion processes of PreDiff. This ensures that the generated predictions align with domain-specific prior knowledge better, enhancing the reliability of the forecasts.
- Our method achieves state-of-the-art performance on the $N$-body MNIST [9] dataset and attains state-of-the-art perceptual quality on the SEVIR [55] dataset.

## 2 Related Work

**Deep learning for precipitation nowcasting**   In recent years, the field of DL has experienced remarkable advancements, revolutionizing various domains of study, including Earth science. One area where DL has particularly made significant strides is in the field of Earth system forecasting, especially precipitation nowcasting. Precipitation nowcasting benefits from the success of DL architectures such as convolutional neural networks (CNNs), recurrent neural networks (RNNs), and Transformers, which have demonstrated their effectiveness in handling spatiotemporal tensors, the typical formulation for Earth system observation data. ConvLSTM [47], a pioneering approach in DL for precipitation nowcasting, combines the strengths of CNNs and LSTMs processing spatial and temporal data. PredRNN [58] builds upon ConvLSTM by incorporating a spatiotemporal memory

flow structure. E3D-LSTM [57] integrates 3D CNN to LSTM to enhance long-term high-level relation modeling. PhyDNet [12] incorporated partial differential equation (PDE) constraints in the latent space. MetNet [49] and its successor, MetNet-2 [5], propose architectures based on ConvLSTM and dilated CNN, enabling skillful precipitation forecasts up to twelve hours ahead. DGMR [41] takes an adversarial training approach to generate sharp and accurate nowcasts, addressing the issue of blurry predictions.

In addition to precipitation nowcasting, there has been a surge in the modeling of global weather and medium-range weather forecasting due to the availability of extensive Earth observation data, such as the European Centre for Medium-Range Weather Forecasts (ECMWF)'s ERA5 [19] dataset. Several DL-based models have emerged in this area. FourCastNet [37] proposes an architecture with Adaptive Fourier Neural Operators (AFNO) [13] as building blocks for autoregressive weather forecasting. FengWu [3] introduces a multi-model Transformer-based global medium-range weather forecast model that achieves skillful forecasts up to ten days ahead. GraphCast [29] combines graph neural networks with convolutional LSTMs to tackle sub-seasonal forecasting tasks, representing weather phenomena as spatiotemporal graphs. Pangu-Weather [2] proposes a 3D Transformer model with Earth-specific priors and a hierarchical temporal aggregation strategy for medium-range global weather forecasting. While recent years have seen remarkable progress in DL for precipitation nowcasting, existing methods still face some limitations. Some methods are deterministic, failing to capture uncertainty and resulting in blurry generation. Others lack the capability of incorporating prior knowledge, which is crucial for controllable generation. In contrast, PreDiff captures the uncertainty in the underlying data distribution with a denoising diffusion process, avoiding simply averaging all possibilities into blurry forecasts. Our knowledge control mechanism facilitates controlled generation under the guidance of prior knowledge.

**Diffusion models** Diffusion models (DMs) [22] are a class of generative models that have become increasingly popular in recent years. DMs learn the data distribution by constructing a forward process that adds noise to the data, and then approximating the reverse process to remove the noise. Latent diffusion models (LDMs) [42] are a variant of DMs that are trained on latent vector outputs from a variational autoencoder. LDMs have been shown to be more efficient to train and capable of generating higher quality images compared to original DMs. Building on the success of DMs in image generation, DMs have also been adopted for video generation. MCVD [56] trains a DM by randomly masking past and/or future frames in blocks and conditioning on the remaining frames. It generates long videos by autoregressively sampling blocks of frames in a sliding window manner. PVDM [62] projects videos into low-dimensional latent space as 2D vectors, and presents a joint training of unconditional and frame conditional video generations. LFDM [35] employs a flow predictor to estimate latent flows between video frames and learns a DM for temporal latent flow generation. VideoFusion [32] decomposes the transition noise in DMs into per-frame noise and the noise along time axis, and trains two networks jointly to match the noise decomposition. While DMs have demonstrated impressive performance in video synthesis, its applications to precipitation nowcasting and other Earth science tasks have not been well explored. Hatanaka et al. [16] uses DMs to super-resolve coarse numerical predictions for solar forecast. Concurrent to our work, LDCast [30] applies LDMs for precipitation nowcasting. However, LDCast has not studied how to integrate prior knowledge to the DM, which is a unique advantage and novelty of PreDiff.

**Conditional controls on diffusion models** Another key advantage of DMs is the ability to condition generation on text, class labels, and other modalities for controllable and diverse output. For instance, ControlNet [64] enables fine-tuning a pretrained DM by freezing the base model and training a copy end-to-end with conditional inputs. Composer [24] decomposes images into representative factors used as conditions to guide the generation. Beyond text and class labels, conditions in other modalities, including physical constraints, can also be leveraged to provide valuable guidance. TopDiff [33] constrains topology optimization using load, boundary conditions, and volume fraction. Physdiff [63] trains a physics-based motion projection module with reinforcement learning to project denoised motions in diffusion steps into physically plausible ones. Nonetheless, while conditional control has proven to be a powerful technique in various domains, its application in DL for precipitation nowcasting remains an unexplored area.
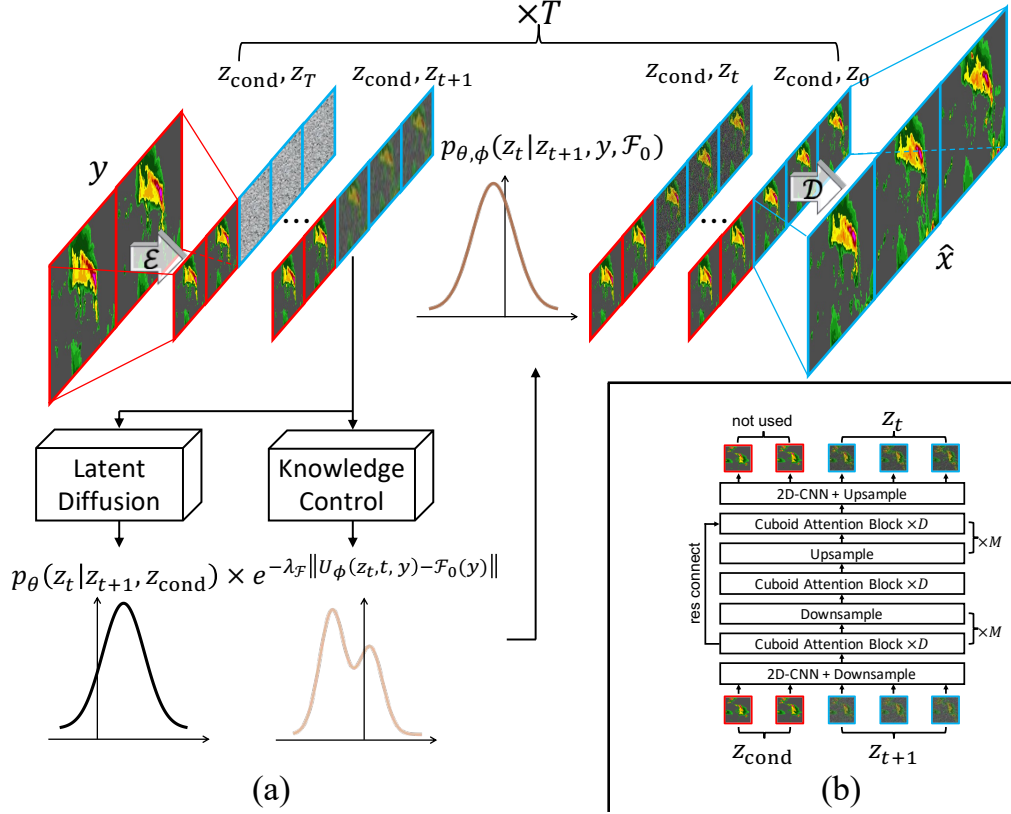
Figure 1: **PreDiff overview**. (a) PreDiff inference with knowledge control. An observation sequence $y$ is encoded into a latent context $z_{cond}$ by the frame-wise encoder $\mathcal{E}$. The latent diffusion model $p_\theta(z_t|z_{t+1}, z_{cond})$ then generates the latent future $z_0$ by autoregressively denoising Gaussian noise $z_T$ conditioned on $z_{cond}$. The transition distribution of each step from $z_{t+1}$ to $z_t$ can be further refined via knowledge control, according to auxiliary prior knowledge. $z_0$ is decoded back to pixel space by the frame-wise decoder $\mathcal{D}$ to produce the final prediction $\widehat{x}$. (b) The Earthformer-UNet architecture used in PreDiff. PreDiff employs an Earthformer-UNet as the backbone for parameterizing the latent diffusion model $p_\theta(z_t|z_{t+1}, z_{cond})$. It takes the concatenation of the latent context $z_{cond}$ (in the red border) and the previous-step noisy latent future $z_{t+1}$ (in the blue border) as input, and outputs $z_t$. (Best viewed in color).

## 3 Method

We follow [47, 48, 55, 1, 9] to formulate precipitation nowcasting as a spatiotemporal forecasting problem. The observation is represented as a spatiotemporal sequence $y = [y^j]_{j=1}^{L_{in}} \in \mathbb{R}^{L_{in} \times H \times W \times C}$, where $H$ and $W$ denote the spatial resolution, and $C$ denotes the number of measurements at each space-time coordinate. Probabilistic forecasting aims to model the conditional probabilistic distribution $p(x|y)$ of the $L_{out}$-step-ahead future $x = [x^j]_{j=1}^{L_{out}} \in \mathbb{R}^{L_{out} \times H \times W \times C}$, given the $L_{in}$-step observation $y$. In what follows, we will present the parameterization of $p(x|y)$ by a conditional LDM and the two-stage pipeline for its training.

### 3.1 Preliminary: Diffusion Models

Diffusion models (DMs) learn the data distribution $p(x)$ by training a model to reverse a pre-defined noising process that progressively corrupts the data. Specifically, the noising process is defined as $q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1-\alpha_t)I), 1 \leq t \leq T$, where $x_0 \sim p(x)$ is the true data, and $x_T \sim \mathcal{N}(0, I)$ is random noise. The coefficients $\alpha_t$ follow a fixed schedule over the timesteps $t$. DMs factorize the joint distribution over the data $x_0$ and noisy latents $x_i$ as

$p(x_{0:T}) = p(x_T) \prod_{t=1}^{T} p(x_{t-1}|x_t)$. Each step of the reverse denoising process is parameterized as a Gaussian distribution $p_\theta(x_{t-1}|x_t) = \mathcal{N}(\mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$, which is trained to recover $x_{t-1}$ from $x_t$.

To apply DMs for spatiotemporal forecasting, $p(x|y)$ is factorized and parameterized as $p(x|y) = \int p(x_{0:T}|y)dx_{1:T} = \int p(x_T) \prod_{t=1}^{T} p_\theta(x_{t-1}|x_t, y)dx_{1:T}$, where $p_\theta(x_{t-1}|x_t, y)$ represents the conditional denoising transition with the condition $y$.

## 3.2 Conditional Diffusion in Latent Space

To improve the computational efficiency of DM training and inference, our *PreDiff* uses a two-phase training process that leverages the benefits of lower-dimensional latent representations, following LDM.

PreDiff training consists of two sequential phases: 1) Training a frame-wise variational autoencoder (VAE) [28] that encodes pixel space into a lower-dimensional latent space, and 2) Training a conditional DM that generates predictions in this acquired latent space.

**Frame-wise autoencoder.** We follow [7], we train a frame autoencoder using a combination of the pixel-wise loss (e.g. L2 loss) and an adversarial loss. Unlike [7], we exclude the perceptual loss since there are no standard pretrained models for perception on Earth observation data. Specifically, the encoder $\mathcal{E}$ is trained to encode a data frame $x^j \in \mathbb{R}^{H \times W \times C}$ to a latent representation $z^j = \mathcal{E}(x^j) \in \mathbb{R}^{H_z \times W_z \times C_z}$. The decoder $\mathcal{D}$ learns to reconstruct the data frame $\widehat{x}^j = \mathcal{D}(z^j)$ from the encoded latent. In the following, we denote $z \sim p_\mathcal{E}(z|x) \in \mathbb{R}^{L \times H_z \times W_z \times C_z}$ as equivalent to $z = [z^j] = [\mathcal{E}(x^j)]$, representing encoding a sequence of frames in pixel space into latent spatiotemporal states. And $x \sim p_\mathcal{D}(x|z)$ denotes decoding a latent spatiotemporal sequence.

**Latent diffusion.** With the context $y$ being encoded by the frame-wise encoder $\mathcal{E}$ into the learned latent space as $z_{\text{cond}} \in \mathbb{R}^{L_{\text{in}} \times H_z \times W_z \times C_z}$ as (1). The conditional distribution $p(z_{0:T}|z_{\text{cond}})$ of the latent future $z_i \in \mathbb{R}^{L_{\text{out}} \times H_z \times W_z \times C_z}$ given $z_{\text{cond}}$ is factorized and parameterized as (2):

$$z_{\text{cond}} \sim p_\mathcal{E}(z_{\text{cond}}|y), \tag{1}$$

$$p(z_{0:T}|z_{\text{cond}}) = p(z_T) \prod_{t=1}^{T} p_\theta(z_{t-1}|z_t, z_{\text{cond}}). \tag{2}$$

where $z_T \sim p(z_T) = \mathcal{N}(0, I)$. As proposed by [22, 45], an equivalent parameterization is to have the DMs learn to match the transition noise $\epsilon_\theta(z_t, t)$ of step $t$ instead of directly predicting $z_{t-1}$. The training objective of PreDiff is simplified as shown in (3):

$$L_{\text{CLDM}} = \mathbb{E}_{(x,y), t, \epsilon \sim \mathcal{N}(0, I)} \|\epsilon - \epsilon_\theta(z_t, t, z_{\text{cond}})\|_2^2. \tag{3}$$

where $(x, y)$ is a sampled context sequence and target sequence data pair, and given that, $z_t \sim q(z_t|z_0)p_\mathcal{E}(z_0|x)$ and $z_{\text{cond}} \sim p_\mathcal{E}(z_{\text{cond}}|y)$.

**Instantiating** $p_\theta(z_{t-1}|z_t, z_{\text{cond}})$. Compared to images, modeling spatiotemporal observation data in precipitation nowcasting poses greater challenges due to their higher dimensionality. We propose replacing the UNet backbone in LDM with *Earthformer-UNet*, derived from Earthformer's encoder [9], which is known for its ability to model intricate and extensive spatiotemporal dependencies in the Earth system.

As illustrated in Fig. 1(b), Earthformer-UNet adopts a hierarchical UNet architecture with self cuboid attention [9] as the building blocks, excluding the bridging cross-attention in the encoder-decoder architecture of Earthformer. We find Earthformer-UNet to be more stable and effective at modeling the transition distribution $p_\theta(z_{t-1}|z_t, z_{\text{cond}})$. It takes the concatenation of the encoded latent context $z_{\text{cond}}$ and the noisy latent future $z_t$ as input, and predicts the one-step-ahead noisy latent future $z_{t-1}$ (in practice, the transition noise $\epsilon$ from $z_t$ to $z_{t-1}$ is predicted as shown in (3)). More implementation details are provided in Appendix A.1.

## 3.3 Incorporating Knowledge Controls

Though DMs hold great promise for diverse and realistic generation, the generated predictions may violate physical behaviors by producing implausible forecasts, or disregard domain-specific expertise, thereby fail to give non-trivial results [14, 44]. One possible reason for this is that DMs are not necessarily trained on data full compliant with domain knowledge. when trained on curated data, there is no guarantee the generations sampled from the learned distribution will remain physically realizable. The causes may also stem from the stochastic nature of chaotic systems, the approximation error in denoising steps, etc.

---

**Algorithm 1** One training step of the knowledge control network $U_\phi(\cdot)$

---

1: $(x, y)$ sampled from data
2: $t \sim \text{Uniform}(0, T)$
3: $z_t \sim q(z_t|z_0)p_\mathcal{E}(z_0|x)$
4: $L_U \leftarrow \|U_\phi(z_t, t, y) - \mathcal{F}(x, y)\|$

---

To address this issue, we propose *knowledge control* to incorporate auxiliary prior knowledge:

$$\mathcal{F}(\widehat{x}, y) = \mathcal{F}_0(y) \in \mathbb{R}^d, \tag{4}$$

into the diffusion generation process. The knowledge control imposes a constraint $\mathcal{F}(\cdot)$ on the forecast $\widehat{x}$, optionally with the observation $y$, based on domain expertise. E.g., for an isolated physical system, a knowledge control $E(\widehat{x}, \cdot) = E_0(y^{L_\text{in}}) \in \mathcal{R}$ could enforce conservation of energy by requiring the generation $\widehat{x}$ to have the same total energy $E_0(y^{L_\text{in}})$ as the last observation. The violation $\|\mathcal{F}(\widehat{x}, y) - \mathcal{F}_0(y)\|$ quantifies the deviation of a prediction $\widehat{x}$ from prior knowledge. The larger violation indicates $\widehat{x}$ diverges further from the constraints. Knowledge control hence aims to suppress the probability of generating predictions with large knowledge control violation. Notice that even the target futures $x$ from training data may violate the knowledge control, i.e. $\mathcal{F}(x, y) \neq \mathcal{F}_0(y)$, due to noise in data collection or simulation.

Inspired by classifier guidance [4], we achieve knowledge control by training a neural network $U_\phi(z_t, t, y)$ to estimate $\mathcal{F}(\widehat{x}, y)$ from the intermediate latent $z_t$ at noising step $t$. The key idea is to adjust the transition probability distribution $p_\theta(z_{t-1}|z_t, z_\text{cond})$ in (2) during each latent denoising step to reduce the likelihood of sampling $z_t$ values expected to violate the constraints:

$$p_{\theta,\phi}(z_t|z_{t+1}, y, \mathcal{F}_0) \propto p_\theta(z_t|z_{t+1}, z_\text{cond}) \cdot e^{-\lambda_\mathcal{F}\|U_\phi(z_t, t, y) - \mathcal{F}_0(y)\|}, \tag{5}$$

where $\lambda_\mathcal{F}$ is a guidance scale factor. The knowledge control network is trained by optimizing the objective $L_U$ in Alg. 1. According to [4], (5) can be approximated by shifting the predicted mean of the denoising transition $\mu_\theta(z_{t+1}, t, z_\text{cond})$ by $-\lambda_\mathcal{F}\Sigma_\theta \nabla_{z_t}\|U_\phi(z_t, t, y) - \mathcal{F}_0(y)\|$, where $\Sigma_\theta$ is the variance of the original transition distribution $p_\theta(z_t|z_{t+1}, z_\text{cond}) = \mathcal{N}(\mu_\theta(z_{t+1}, t, z_\text{cond}), \Sigma_\theta(z_{t+1}, t, z_\text{cond}))$. See Appendix B for a full derivation.

The training procedure of knowledge control is outlined in Alg. 1. To obtain the noisy latent $z_t$, we sample using the frame-wise encoder $\mathcal{E}$ and the forward noising process $q(z_t|z_0)$. It's important to note that this trainging procedure is independent of the LDM training. At inference time, knowledge control is applied as a plug-in, without impacting the trained VAE and the LDM. This modular approach allows training lightweight knowledge control networks $U_\phi$ to impose different constraints, without retraining the full model. The key advantage over incorporating constraints into model architectures or training losses is the flexibility to quickly explore diverse domain knowledge.

## 4 Experiments

We conduct empirical studies and compare PreDiff with other state-of-the-art spatiotemporal forecasting models on a synthetic dataset $N$-body MNIST [9] and a real-world precipitation nowcasting benchmark SEVIR[2] [55] to verify the effectiveness of PreDiff in handling the dynamics and uncertainty in complex spatiotemporal systems and generating high quality, accurate forecasts. We impose data-specific knowledge controls: **energy conservation** on $N$-body MNIST and **anticipated precipitation intensity** on SEVIR. Experiments demonstrate that PreDiff under the guidance of knowledge control (PreDiff-KC) is able to generate predictions that comply with domain expertise without severely sacrificing fidelity.

---

[2]Dataset is available at `https://sevir.mit.edu/`

## 4.1 $N$-body MNIST Digits Motion Forecasting

**Dataset.** The Earth is a chaotic system with complex dynamics. The real-world Earth observation data, such as radar echo maps and satellite imagery, are usually not physically complete. We are unable to directly verify whether certain domain knowledge, like conservation laws of energy and momentum, is satisfied or not. This makes it difficult to verify if a method is really capable of modeling certain dynamics. To address this, we follow [9] to generate a synthetic dataset called $N$-body MNIST[3], which is an extension of MovingMNIST [50]. The dataset contains sequences of digits moving subject to the gravitational force from other digits. The governing equation for the motion is $\frac{d^2 \boldsymbol{x}_i}{dt^2} = -\sum_{j \neq i} \frac{G m_j (\boldsymbol{x}_i - \boldsymbol{x}_j)}{(|\boldsymbol{x}_i - \boldsymbol{x}_j| + d_{\text{soft}})^r}$, where $\boldsymbol{x}_i$ is the spatial coordinates of the $i$-th digit, $G$ is the gravitational constant, $m_j$ is the mass of the $j$-th digit, $r$ is a constant representing the power scale in the gravitational law, and $d_{\text{soft}}$ is a small softening distance that ensures numerical stability. The motion occurs within a $64 \times 64$ frame. When a digit hits the boundaries of the frame, it bounces back by elastic collision. We use $N = 3$ for chaotic 3-body motion [34]. The forecasting task is to predict 10-step ahead future frames $x \in \mathbb{R}^{10 \times 64 \times 64 \times 1}$ given the length-10 context $y \in \mathbb{R}^{10 \times 64 \times 64 \times 1}$. We generate 20,000 sequences for training, 1,000 sequences for validation and 1,000 sequences for testing. Empirical studies on such a synthetic dataset with known dynamics helps provide useful insights for model development and evaluation.

**Evaluation.** In addition to standard metrics MSE, MAE and SSIM, we also report the continuous ranked probability score (CRPS) [10] for assessing probabilistic forecasts, as well as the scores of Fréchet Video Distance (FVD) [51], a metric for evaluating the visual quality of generated videos. CRPS measures the discrepancy between the predicted distribution and the true distribution. When the forecast distribution collapses into a single value, as in deterministic forecasting models, CRPS reduces to Mean Absolute Error (MAE). A lower CRPS value indicates higher forecast accuracy. Similar to Fréchet Inception Distance (FID) [20] for evaluating image generation, FVD estimates the distance between the learned distribution and the true data distribution by comparing the statistics of feature vectors extracted from the generations and the real data. The inception network used in FVD for feature extraction is pre-trained on video classification and is not specifically adapted for processing "unnatural videos" such as spatiotemporal observation data in Earth systems. Consequently, the FVD scores on the $N$-body MNIST and SEVIR datasets cannot be directly compared with those on natural video datasets. Nevertheless, the relative ranking of the FVD scores remains a meaningful indicator of model ability to achieve high visual quality, as FVD has shown consistency with expert evaluations across various domains beyond natural images [38, 26].

### 4.1.1 Comparison with the State of the Art

We evaluate seven deterministic spatiotemporal forecasting models: **UNet** [55], **ConvLSTM** [47], **PredRNN** [58], **PhyDNet** [12], **E3D-LSTM** [57], **Rainformer** [1] and **Earthformer** [9], as well as two probabilistic spatiotemporal forecasting models: **VideoGPT** [61] and **LDM** [42]. All baselines are trained using default configurations in their officially released code. Results in Table 1 show that PreDiff outperforms these baselines by a large margin in both conventional video prediction metrics (i.e., MSE, MAE, SSIM), and a perceptual quality metric, FVD. The example predictions in Fig. 2 demonstrate that PreDiff generate predictions with sharp and clear digits in accurate positions. In contrast, deterministic baselines resort to generating blurry predictions to accommodate uncertainty. Probabilistic baselines, though producing sharp strokes, either predict *incorrect* positions or *fail to reconstruct* the digits. The performance gap between LDM [42] and PreDiff serves as an ablation study that highlights the importance of the latent backbone's spatiotemporal modeling capacity. Specifically, the Earthformer-UNet utilized in PreDiff demonstrates superior performance compared to the UNet in LDM [42].

### 4.1.2 Knowledge Control: Energy Conservation

In the $N$-body MNIST simulation, digits move in a closed system under Newton's law of gravity and bounce with elastic collisions. This system obeys the law of conservation of energy. The total energy of the whole system $E(x^j)$ at any future time step $j$ during evolution should equal the total energy at the last observation time step $E(y^{L_{\text{in}}})$.

---

[3]Code available at `https://github.com/amazon-science/earth-forecasting-transformer/tree/main/src/earthformer/datasets/nbody`
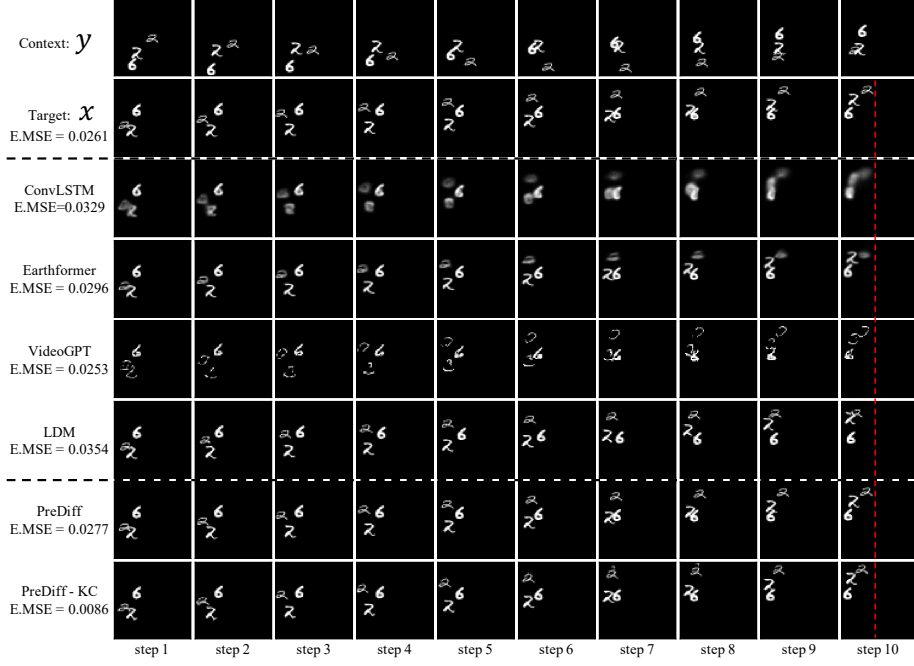
Figure 2: A set of example predictions on the $N$-body MNIST test set. From top to bottom: context sequence $y$, target sequence $x$, predictions by ConvLSTM [47], Earthformer [9], VideoGPT [61], LDM [42], PreDiff, and PreDiff with knowledge control. E.MSE denotes the average error between the total energy (kinetic + potential) of the predictions $E(\widehat{x}^j)$ and the total energy of the last context frame $E(y^{L_{in}})$. The red dashed line is to help the reader to judge the position of the digit "2" in the last frame.

We impose the law of conservation of energy as the knowledge control for PreDiff on $N$-body MNIST in the form of (4) :

$$\mathcal{F}(\widehat{x}, y) \equiv [E(\widehat{x}^1), \ldots, E(\widehat{x}^{L_{out}})]^T, \tag{6}$$

$$\mathcal{F}_0(y) \equiv [E(y^{L_{in}}), \ldots, E(y^{L_{in}})]^T. \tag{7}$$

The ground-truth values of the total energy $E(y^{L_{in}})$ and $E(x^j)$ are directly accessible since $N$-body MNIST is a synthetic dataset from simulation. The total energy can be derived from the velocities (kinetic energy) and positions (potential energy) of the moving digits. A knowledge control network $U_\phi(\cdot)$ is trained following Alg. 1 to guide the PreDiff to generate forecasts $\widehat{x}$ that conserve the same energy as the initial step $E(y^{L_{in}})$.

To verify the effectiveness of the knowledge control on guiding the generations to comply with the law of conservation of energy, we train an energy detector $E_{det}(\widehat{x})$[4] that detects the total energy of the forecasts $\widehat{x}$. We evaluate the energy error between the forecasts and the initial energy with E.MSE$(\widehat{x}, y) \equiv$ MSE$(E_{det}(\widehat{x}), E(y^{L_{in}}))$ and E.MAE$(\widehat{x}, y) \equiv$ MAE$(E_{det}(\widehat{x}), E(y^{L_{in}}))$. In this evaluation, we exclude the methods that generate blurred predictions with ambiguous digit positions. We only focus on the methods that are capable of producing clear digits in precise positions.

As illustrated in Table 1, PreDiff-KC substantially outperforms all baseline methods and PreDiff without knowledge control in E.MSE and E.MAE. This demonstrates that the forecasts of PreDiff-KC comply much better with the law of conservation of energy, while still maintaining high visual quality with an FVD score of $4.063$.

---

[4]The test MSE of the energy detector is $5.56 \times 10^{-5}$, which is much smaller than the scores of E.MSE shown in Table 1. This indicates that the energy detector has high precision and reliability for verifying energy conservation in the model forecasts.

Table 1: Performance comparison on $N$-body MNIST. We report conventional frame quality metrics (MSE, MAE, SSIM), along with Fréchet Video Distance (FVD) [51] for assessing visual quality. Energy conservation is evaluated via E.MSE and E.MAE between the energy of predictions $E_{\text{det}}(\widehat{x})$ and the initial energy $E(y^{L_{\text{in}}})$. Lower values on the energy metrics indicate better compliance with conservation of energy.

| Model | #Param. (M) | Frame Metrics | | | | Energy Metrics | |
|---|---|---|---|---|---|---|---|
| | | MSE ↓ | MAE ↓ | SSIM ↑ | FVD ↓ | E.MSE ↓ | E.MAE ↓ |
| Target | - | 0.000 | 0.000 | 1.0000 | 0.000 | 0.0132 | 0.0697 |
| Persistence | - | 104.9 | 139.0 | 0.7270 | 168.3 | - | - |
| UNet [55] | 16.6 | 38.90 | 94.29 | 0.8260 | 142.3 | - | - |
| ConvLSTM [47] | 14.0 | 32.15 | 72.64 | 0.8886 | 86.31 | - | - |
| PredRNN [58] | 23.8 | 21.76 | 54.32 | 0.9288 | 20.65 | - | - |
| PhyDNet [12] | 3.1 | 28.97 | 78.66 | 0.8206 | 178.0 | - | - |
| E3D-LSTM [57] | 12.9 | 22.98 | 62.52 | 0.9131 | 22.28 | - | - |
| Rainformer [1] | 19.2 | 38.89 | 96.47 | 0.8036 | 163.5 | - | - |
| Earthformer [9] | 7.6 | <u>14.82</u> | <u>39.93</u> | <u>0.9538</u> | 6.798 | - | - |
| VideoGPT [61] | 92.2 | 53.68 | 77.42 | 0.8468 | 39.28 | 0.0228 | 0.1092 |
| LDM [42] | 410.3 | 46.29 | 72.19 | 0.8773 | <u>3.432</u> | 0.0243 | 0.1172 |
| PreDiff | 120.7 | **9.492** | **25.01** | **0.9716** | **0.987** | <u>0.0226</u> | <u>0.1083</u> |
| PreDiff-KC | 129.4 | 21.90 | 43.57 | 0.9303 | 4.063 | **0.0039** | **0.0443** |

Furthermore, we detect energy errors in the target data sequences. The first row of Table 1 indicates that even the target from the training data may not strictly adhere to the prior knowledge. This could be due to discretization errors in the simulation. Table 1 shows that all baseline methods and PreDiff have larger energy errors than the target, meaning purely data-oriented approaches cannot eliminate the impact of noise in the training data. In contrast, guided by the law of conservation of energy, PreDiff-KC overcomes the intrinsic defects in the training data, resulting in lower energy errors than the target.

A typical example shown in Fig. 2 demonstrates that while PreDiff precisely reproduces the ground-truth position of digit "2" in the last frame (aligned to the red dashed line), resulting in nearly the same energy error (E.MSE = 0.0277) as the ground-truth's (E.MSE = 0.0261), PreDiff-KC successfully corrects the motion of digit "2", providing it with physically plausible velocity and position (slightly off the red dashed line). The knowledge control ensures that the generation complies better with the law of conservation of energy, resulting in a much lower E.MSE = 0.0086. On the contrary, none of the evaluated baselines can overcome the intrinsic noise from the data, resulting in energy errors similar to or larger than that of the ground-truth.

Notice that the pixel-wise scores MSE, MAE and SSIM are less meaningful for evaluating PreDiff-KC, since correcting the noise of the energy results in changing the velocities and positions of the digits. A minor change in the position of a digit can cause a large pixel-wise error, even though the digit is still generated sharply and in high quality as shown in Fig. 2.

### 4.2 SEVIR Precipitation Nowcasting

**Dataset.** The Storm EVent ImageRy (SEVIR) [55] is a spatiotemporal Earth observation dataset which consists of $384\,\mathrm{km} \times 384\,\mathrm{km}$ image sequences spanning over 4 hours. Images in SEVIR are sampled and aligned across five different data types: three channels (C02, C09, C13) from the GOES-16 advanced baseline imager, NEXRAD Vertically Integrated Liquid (VIL) mosaics, and GOES-16 Geostationary Lightning Mapper (GLM) flashes. The SEVIR benchmark supports scientific research on multiple meteorological applications including precipitation nowcasting, synthetic radar generation, front detection, etc. Due to computational resource limitations, we adopt a downsampled version of SEVIR for benchmarking precipitation nowcasting. The task is to predict the future VIL up to 60 minutes (6 frames) given 70 minutes of context VIL (7 frames) at a spatial resolution of $128 \times 128$, i.e. $x \in \mathbb{R}^{6 \times 128 \times 128 \times 1}$, $y \in \mathbb{R}^{7 \times 128 \times 128 \times 1}$.
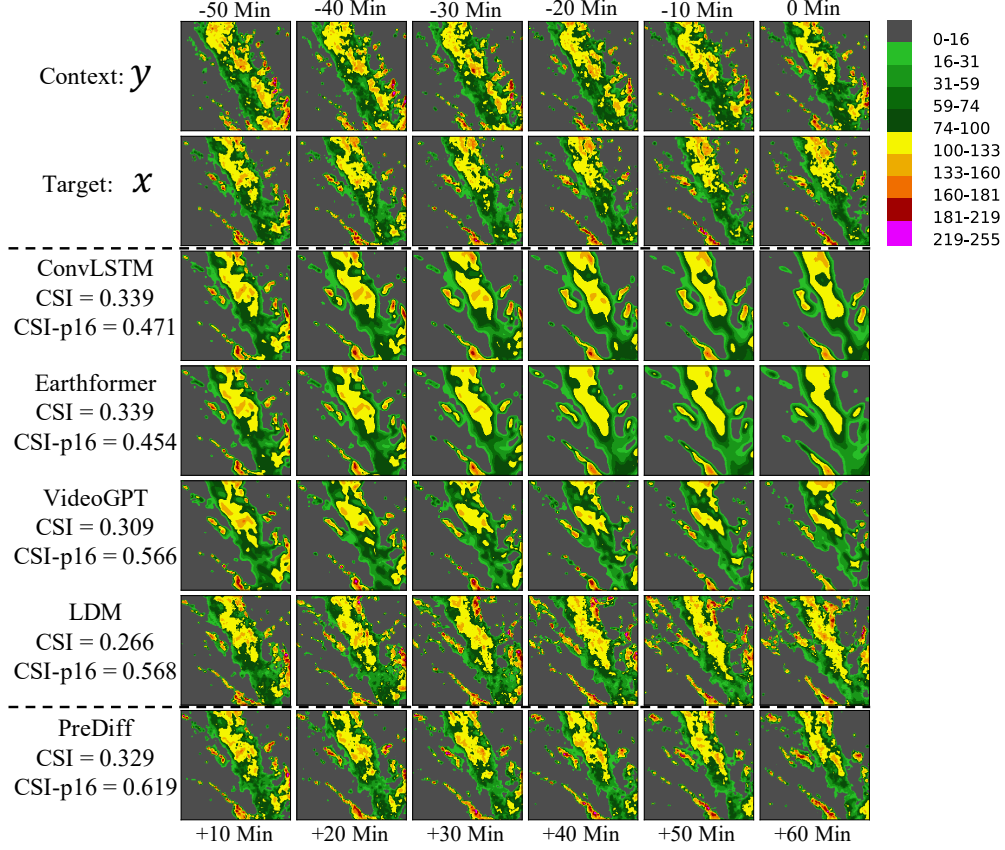
Figure 3: Perceptual quality of forecasts from baselines and PreDiff on the SEVIR test set. From top to bottom: context sequence $y$, target sequence $x$, forecasts from ConvLSTM [47], Earthformer [9], VideoGPT[61], LDM [42], PreDiff.

**Evaluation.** Following [55, 9], we adopt the `Critical Success Index` (CSI) for evaluation, which is commonly used in precipitation nowcasting and is defined as CSI = $\frac{\#\texttt{Hits}}{\#\texttt{Hits}+\#\texttt{Misses}+\#\texttt{F.Alarms}}$. To count the $\#\texttt{Hits}$ (truth=1, pred=1), $\#\texttt{Misses}$ (truth=1, pred=0) and $\#\texttt{F.Alarms}$ (truth=0, pred=1), the prediction and the ground-truth are rescaled back to the range $0-255$ and binarized at thresholds $[16, 74, 133, 160, 181, 219]$. We also follow [41] to report the `CSI` at pooling scale $4 \times 4$ and $16 \times 16$, which evaluate performance on neighborhood aggregations at multiple spatial scales. These pooled `CSI` metrics assess the models' ability to capture local pattern distributions. Additionally, we incorporate FVD [51] and CRPS [10] for assessing the visual quality and uncertainty modeling capabilities of the investigated methods.

### 4.2.1 Comparison to the State of the Art

We adjust the configurations of involved baselines accordingly and tune some of the hyperparameters for adaptation on the SEVIR dataset. More implementation details of baselines are provided in Appendix A.2. The experiment results listed in Table 2 show that probabilistic spatiotemporal forecasting methods are not good at achieving high `CSI` scores. However, they are more powerful at capturing the patterns and the true distribution of the data, hence achieving much better FVD scores and `CSI-pool16`. Qualitative results shown in Fig. 3 demonstrate that `CSI` is not aligned with human perceptual judgement. For such a complex system, deterministic methods give up capturing the real patterns and resort to averaging the possible futures, i.e. blurry predictions, to keep the
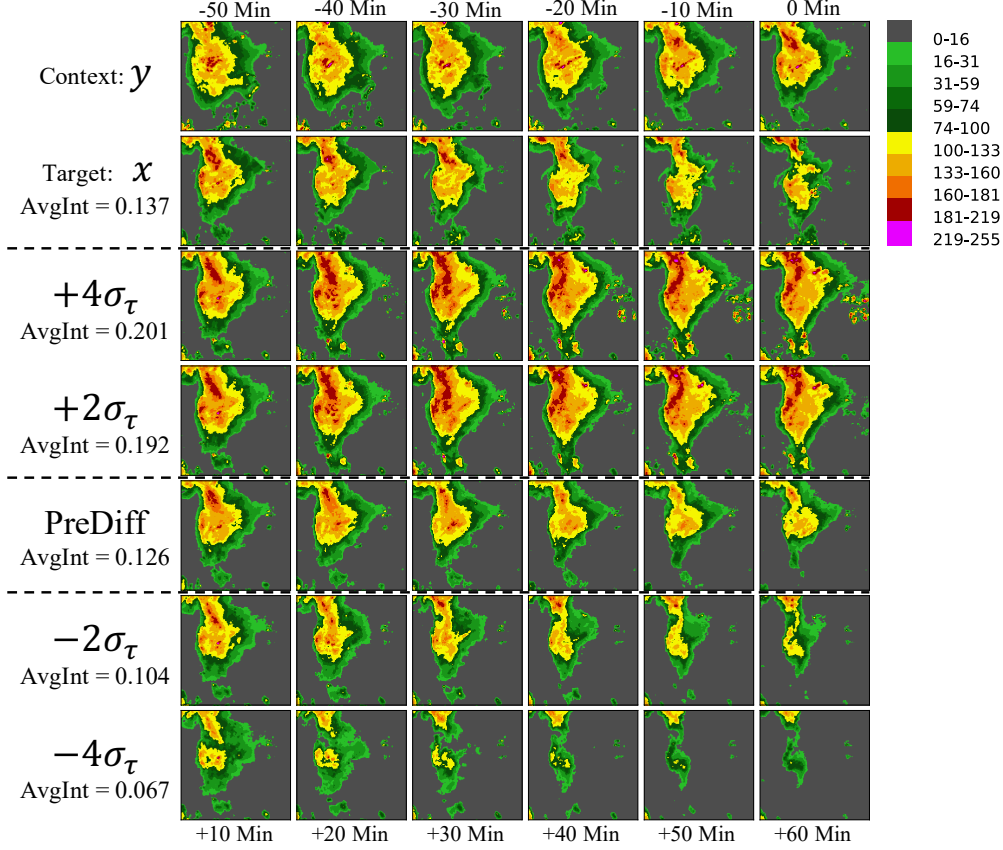
Figure 4: A set of example forecasts from PreDiff-KC, i.e., PreDiff under the guidance of anticipated average intensity. From top to bottom: context sequence $y$, target sequence $x$, forecasts from PreDiff-KC showcasing different levels of anticipated future intensity ($\mu_\tau + n\sigma_\tau$), where $n$ takes the values of $4, 2, 0, -2, -4$.

scores from appearing too inaccurate. Probabilistic approaches, of which PreDiff is the best, though are not favored by per-pixel metrics, perform better at capturing the data distribution within a local area, resulting in higher `CSI-pool16`, lower CRPS, and succeed in keeping the correct local patterns, which can be crucial for recognizing weather events.

### 4.2.2 Knowledge Control: Anticipated Average Intensity

Earth system observation data like radar echo maps are usually not physically complete, posing challenges for directly incorporating physical laws for guidance. However, with highly flexible knowledge control, guided forecasting becomes feasible through auxiliary controls.

Specifically for precipitation nowcasting on SEVIR, we use anticipated precipitation intensity as knowledge control to simulate possible extreme weather events. We denote the average intensity of a data sequence as $I(x) \in \mathbb{R}^+$. In order to estimate the conditional quantiles of future intensity, we train a simple probabilistic time series forecasting model with a parametric (Gaussian) distribution $p_\tau(I(x)|[I(y^j)]) = \mathcal{N}(\mu_\tau([I(y^j)]), \sigma_\tau([I(y^j)]))$ that predict the distribution of the average future intensity $I(x)$ given the average intensity of each context frame $[I(y^j)]_{j=1}^{L_{\text{in}}}$ (abbreviated as $[I(y^j)]$). By defining $\mathcal{F}(\hat{x}, y) \equiv I(\hat{x})$ and $\mathcal{F}_0(y, n\sigma_\tau) \equiv \mu_\tau + n\sigma_\tau$, PreDiff is able to generate samples for extreme cases like anticipated intensity beyond $2\sigma_\tau$, under the guidance of knowledge control.

Table 2: Performance comparison on SEVIR. The `Critical Success Index`, also known as the intersection over union (IoU), is calculated at different precipitation thresholds and denoted as `CSI`-$thresh$. `CSI` reports the mean of `CSI`-$[16, 74, 133, 160, 181, 219]$. `CSI-pool`$s$ with $s = 4$ and $s = 16$ report the `CSI` at pooling scales of $4 \times 4$ and $16 \times 16$. Besides, we include the continuous ranked probability score (CRPS) for probabilistic forecast assessment, and the scores of Fréchet Video Distance (FVD) for evaluating visual quality.

| Model | #Param. (M) | Metrics | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | FVD ↓ | CRPS ↓ | CSI ↑ | CSI-pool4 ↑ | CSI-pool16 ↑ |
| Persistence | - | 525.2 | 0.0526 | 0.2613 | 0.3702 | 0.4690 |
| UNet [55] | 16.6 | 753.6 | 0.0353 | 0.3593 | 0.4098 | 0.4805 |
| ConvLSTM [47] | 14.0 | 659.7 | 0.0332 | 0.4185 | 0.4452 | 0.5135 |
| PredRNN [58] | 46.6 | 663.5 | 0.0306 | 0.4080 | 0.4497 | 0.5005 |
| PhyDNet [12] | 13.7 | 723.2 | 0.0319 | 0.3940 | 0.4379 | 0.4854 |
| E3D-LSTM [57] | 35.6 | 600.1 | 0.0297 | 0.4038 | 0.4492 | 0.4961 |
| Rainformer [1] | 184.0 | 760.5 | 0.0357 | 0.3661 | 0.4232 | 0.4738 |
| Earthformer [9] | 15.1 | 690.7 | 0.0304 | **0.4419** | 0.4567 | 0.5005 |
| DGRM [41] | 71.5 | 485.2 | 0.0435 | 0.2675 | 0.3431 | 0.4832 |
| VideoGPT [61] | 99.6 | 261.6 | 0.0381 | 0.3653 | 0.4349 | 0.5798 |
| LDM [42] | 438.6 | 133.0 | 0.0280 | 0.3580 | 0.4022 | 0.5522 |
| PreDiff | 220.5 | **33.05** | **0.0246** | 0.4072 | **0.4624** | **0.6244** |
| PreDiff-KC ($\in [-2\sigma_\tau, 2\sigma_\tau]$) | 229.4 | 34.18 | - | - | - | - |

Fig. 4 shows a set of generations with knowledge control on anticipated future intensity $\mu_\tau + n\sigma_\tau$, $n \in \{-4, -2, 2, 4\}$. This qualitative example demonstrates that PreDiff is not only capable of capturing the distribution of the future, but also flexible at highlighting possible extreme cases like rainstorms and droughts, which is crucial for decision-making and precaution.

According to Table 2, the FVD score of PreDiff-KC $34.18$ is only slightly higher than the FVD score of PreDiff $33.05$ without knowledge control. This indicates that knowledge control effectively guides the generations while maintaining fidelity and adherence to the true data distribution.

## 5    Conclusions and Broader Impacts

In this paper, we propose PreDiff, a novel latent diffusion model for precipitation nowcasting. We also introduce a general two-stage pipeline for training DL models for Earth system forecasting. Specificaaly, we develop a knowledge control mechanism that is capable of guiding PreDiff to generate forecasts in compliance with domain-specific prior knowledge. Experiments demonstrate that our method achieves state-of-the-art performance on $N$-body MNIST and SEVIR datasets.

Our work has certain limitations: 1) Benchmark datasets and evaluation metrics for precipitation nowcasting and Earth system forecasting are still maturing compared to the computer vision domain. While we utilize conventional precipitation forecast skill scores, visual quality evaluation, and expert judgement, aligning these assessments remains an open challenge. 2) Effective integration of physical principles and domain knowledge into DL models for precipitation nowcasting remains an active research area. Close collaboration between DL researchers and domain experts in meteorology and climatology will be key to developing hybrid models that effectively leverage both data-driven learning and scientific theory. 3) While Earth system observation data have grown substantially in recent years, high-quality data remain scarce in many domains. This scarcity can limit PreDiff's ability to accurately capture the full distribution, occasionally resulting in unrealistic forecast hallucinations under the guidance of prior knowledge as it attempts to circumvent its knowledge control mechanism. Further research on enhancing the sample efficiency of PreDiff and knowledge control is needed.

In conclusion, PreDiff represents a promising advance in knowledge-guided DL for Earth system forecasting, but work remains to improve benchmarking, incorporate scientific knowledge, and boost model robustness through collaborative research between AI and domain experts.

# References

[1] Cong Bai, Feng Sun, Jinglin Zhang, Yi Song, and Shengyong Chen. Rainformer: Features extraction balanced network for radar-based precipitation nowcasting. *IEEE Geoscience and Remote Sensing Letters*, 19:1–5, 2022.

[2] Kaifeng Bi, Lingxi Xie, Hengheng Zhang, Xin Chen, Xiaotao Gu, and Qi Tian. Accurate medium-range global weather forecasting with 3d neural networks. *Nature*, pages 1–6, 2023.

[3] Kang Chen, Tao Han, Junchao Gong, Lei Bai, Fenghua Ling, Jing-Jia Luo, Xi Chen, Leiming Ma, Tianning Zhang, Rui Su, et al. Fengwu: Pushing the skillful global medium-range weather forecast beyond 10 days lead. *arXiv preprint arXiv:2304.02948*, 2023.

[4] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.

[5] Lasse Espeholt, Shreya Agrawal, Casper Sønderby, Manoj Kumar, Jonathan Heek, Carla Bromberg, Cenk Gazen, Jason Hickey, Aaron Bell, and Nal Kalchbrenner. Skillful twelve hour precipitation forecasts using large context neural networks. *arXiv preprint arXiv:2111.07470*, 2021.

[6] Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis. Structure and content-guided video synthesis with diffusion models. *arXiv preprint arXiv:2302.03011*, 2023.

[7] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021.

[8] William Falcon and The PyTorch Lightning team. PyTorch Lightning, March 2019.

[9] Zhihan Gao, Xingjian Shi, Hao Wang, Yi Zhu, Yuyang Wang, Mu Li, and Dit-Yan Yeung. Earthformer: Exploring space-time transformers for earth system forecasting. In *NeurIPS*, 2022.

[10] Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical Association*, 102(477):359–378, 2007.

[11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.

[12] Vincent Le Guen and Nicolas Thome. Disentangling physical dynamics from unknown factors for unsupervised video prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11474–11484, 2020.

[13] John Guibas, Morteza Mardani, Zongyi Li, Andrew Tao, Anima Anandkumar, and Bryan Catanzaro. Adaptive fourier neural operators: Efficient token mixers for transformers. *arXiv preprint arXiv:2111.13587*, 2021.

[14] Derek Hansen, Danielle C. Maddix, Shima Alizadeh, Gaurav Gupta, and Michael W. Mahoney. Learning physical models that can respect conservation laws. In *Proceedings of the $40^{th}$ of International Conference on Machine Learning*, volume 202, 2023.

[15] William Harvey, Saeid Naderiparizi, Vaden Masrani, Christian Weilbach, and Frank Wood. Flexible diffusion modeling of long videos. *arXiv preprint arXiv:2205.11495*, 2022.

[16] Yusuke Hatanaka, Yannik Glaser, Geoff Galgon, Giuseppe Torri, and Peter Sadowski. Diffusion models for high-resolution solar forecasts. *arXiv preprint arXiv:2302.00170*, 2023.

[17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[18] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.

[19] Hans Hersbach, Bill Bell, Paul Berrisford, Shoji Hirahara, András Horányi, Joaquín Muñoz-Sabater, Julien Nicolas, Carole Peubey, Raluca Radu, Dinand Schepers, et al. The era5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730):1999–2049, 2020.

[20] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.

[21] Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhut-dinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.

[22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.

[23] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022.

[24] Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. Composer: Creative and controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778*, 2023.

[25] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr, 2015.

[26] Kevin Kilgour, Mauricio Zuluaga, Dominik Roblek, and Matthew Sharifi. Fr\'echet audio distance: A metric for evaluating music enhancement algorithms. *arXiv preprint arXiv:1812.08466*, 2018.

[27] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[28] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[29] Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirnsberger, Meire Fortunato, Alexander Pritzel, Suman Ravuri, Timo Ewalds, Ferran Alet, Zach Eaton-Rosen, et al. Graphcast: Learning skillful medium-range global weather forecasting. *arXiv preprint arXiv:2212.12794*, 2022.

[30] Jussi Leinonen, Ulrich Hamann, Daniele Nerini, Urs Germann, and Gabriele Franch. Latent diffusion models for generative precipitation nowcasting with accurate uncertainty quantification. *arXiv preprint arXiv:2304.12891*, 2023.

[31] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

[32] Zhengxiong Luo, Dayou Chen, Yingya Zhang, Yan Huang, Liang Wang, Yujun Shen, Deli Zhao, Jingren Zhou, and Tieniu Tan. Videofusion: Decomposed diffusion models for high-quality video generation. *arXiv e-prints*, pages arXiv–2303, 2023.

[33] Franccois Mazé and Faez Ahmed. Topodiff: A performance and constraint-guided diffusion model for topology optimization. *arXiv preprint arXiv:2208.09591*, 2022.

[34] Valtonen MJ, Mauri Valtonen, and Hannu Karttunen. *The three-body problem*. Cambridge University Press, 2006.

[35] Haomiao Ni, Changhao Shi, Kai Li, Sharon X. Huang, and Martin Renqiang Min. Conditional image-to-video generation with latent flow diffusion models, 2023.

[36] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.

[37] Jaideep Pathak, Shashank Subramanian, Peter Harrington, Sanjeev Raja, Ashesh Chattopadhyay, Morteza Mardani, Thorsten Kurth, David Hall, Zongyi Li, Kamyar Azizzadenesheli, et al. FourCastNet: A global data-driven high-resolution weather model using adaptive fourier neural operators. *arXiv preprint arXiv:2202.11214*, 2022.

[38] Kristina Preuer, Philipp Renz, Thomas Unterthiner, Sepp Hochreiter, and Gunter Klambauer. Fréchet chemnet distance: a metric for generative models for molecules in drug discovery. *Journal of chemical information and modeling*, 58(9):1736–1741, 2018.

[39] Ruslan Rakhimov, Denis Volkhonskiy, Alexey Artemov, Denis Zorin, and Evgeny Burnaev. Latent video transformer. *arXiv preprint arXiv:2006.10704*, 2020.

[40] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.

[41] Suman Ravuri, Karel Lenc, Matthew Willson, Dmitry Kangin, Remi Lam, Piotr Mirowski, Megan Fitzsimons, Maria Athanassiadou, Sheleem Kashem, Sam Madge, et al. Skilful precipitation nowcasting using deep generative models of radar. *Nature*, 597(7878):672–677, 2021.

[42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.

[43] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023.

[44] Nadim Saad, Gaurav Gupta, Shima Alizadeh, and Danielle C. Maddix. Guiding continuous operator learning through physics-based boundary constraints. In *Proceedings of the $11^{th}$ International Conference on Learning Representations*, 2023.

[45] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.

[46] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. *arXiv preprint arXiv:1701.05517*, 2017.

[47] Xingjian Shi, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Convolutional LSTM network: A machine learning approach for precipitation nowcasting. In *NeurIPS*, volume 28, 2015.

[48] Xingjian Shi, Zhihan Gao, Leonard Lausen, Hao Wang, Dit-Yan Yeung, Wai-kin Wong, and Wang-chun Woo. Deep learning for precipitation nowcasting: A benchmark and a new model. In *NeurIPS*, volume 30, 2017.

[49] Casper Kaae Sønderby, Lasse Espeholt, Jonathan Heek, Mostafa Dehghani, Avital Oliver, Tim Salimans, Shreya Agrawal, Jason Hickey, and Nal Kalchbrenner. Metnet: A neural weather model for precipitation forecasting. *arXiv preprint arXiv:2003.12140*, 2020.

[50] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of video representations using LSTMs. In *ICML*, pages 843–852. PMLR, 2015.

[51] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. Fvd: A new metric for video generation. In *DGS@ICLR*, 2019.

[52] Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. In *Neural Information Processing Systems (NeurIPS)*, 2021.

[53] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. Conditional image generation with pixelcnn decoders. *Advances in neural information processing systems*, 29, 2016.

[54] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, volume 30, 2017.

[55] Mark Veillette, Siddharth Samsi, and Chris Mattioli. SEVIR: A storm event imagery dataset for deep learning applications in radar and satellite meteorology. *Advances in Neural Information Processing Systems*, 33:22009–22019, 2020.

[56] Vikram Voleti, Alexia Jolicoeur-Martineau, and Christopher Pal. Masked conditional video diffusion for prediction, generation, and interpolation. *arXiv preprint arXiv:2205.09853*, 2022.

[57] Yunbo Wang, Lu Jiang, Ming-Hsuan Yang, Li-Jia Li, Mingsheng Long, and Li Fei-Fei. Eidetic 3D LSTM: A model for video prediction and beyond. In *International conference on learning representations*, 2018.

[58] Yunbo Wang, Haixu Wu, Jianjin Zhang, Zhifeng Gao, Jianmin Wang, Philip Yu, and Mingsheng Long. PredRNN: A recurrent neural network for spatiotemporal predictive learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.

[59] Dirk Weissenborn, Oscar Täckström, and Jakob Uszkoreit. Scaling autoregressive video models. In *International Conference on Learning Representations*, 2019.

[60] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018.

[61] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. VideoGPT: Video generation using vq-vae and transformers. *arXiv preprint arXiv:2104.10157*, 2021.

[62] Sihyun Yu, Kihyuk Sohn, Subin Kim, and Jinwoo Shin. Video probabilistic diffusion models in projected latent space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.

[63] Ye Yuan, Jiaming Song, Umar Iqbal, Arash Vahdat, and Jan Kautz. Physdiff: Physics-guided human motion diffusion model. *arXiv preprint arXiv:2212.02500*, 2022.

[64] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023.

[65] Lu Zhou and Rong-Hua Zhang. A self-attention–based neural network for three-dimensional multivariate modeling and its skillful enso predictions. *Science Advances*, 9(10):eadf2827, 2023.

# A  Implementation Details

All experiments are conducted on machines with NVIDIA A10G GPUs (24GB memoery). All models, including PreDiff, knowledge control networks and the baselines, can fit in a single GPU without gradient checkpointing or without model parallelization.

## A.1  PreDiff

**Frame-wise autoencoder**  We follow [7, 42] to build frame-wise VAEs (not VQVAEs) and train them adversarially from scratch on $N$-body MNIST and SEVIR frames. As shown in Sec. 3.2, on $N$-body MNIST dataset, the spatial downsampling ratio is $4 \times 4$. A frame $x^j \in \mathbb{R}^{64 \times 64 \times 1}$ is encoded to $z^j \in \mathbb{R}^{16 \times 16 \times 3}$ by parameterizing $p_{\mathcal{E}}(z^j | x^j) = \mathcal{N}(\mu_{\mathcal{E}}(x^j) | \sigma_{\mathcal{E}}(x^j))$. On SEVIR dataset, the spatial downsampling ratio is $8 \times 8$. A frame $x^j \in \mathbb{R}^{128 \times 128 \times 1}$ is encoded to $z^j \in \mathbb{R}^{16 \times 16 \times 4}$ similarly.

The detailed configurations of the encoder and decoder of the VAE on $N$-body MNIST are shown in Table 3 and Table 4. The detailed configurations of the encoder and decoder of the VAE on SEVIR are shown in Table 5 and Table 6. The discriminators for adversarial training on $N$-body MNIST and SEVIR datasets share the same configurations, which are shown in Table 7.

**Latent diffusion model that instantiates** $p_\theta(z_{t-1} | z_t, z_{\text{cond}})$  Stemming from Earthformer [9], we build *Earthformer-UNet*, which is a hierarchical UNet with cuboid attention [9] layers as basic building blocks. On $N$-body MNIST, it takes the concatenation along the sequence length axis of $z_{\text{cond}} \in \mathbb{R}^{10 \times 16 \times 16 \times 3}$ and $z_t \in \mathbb{R}^{10 \times 16 \times 16 \times 3}$ as input, and outputs $z_{t-1} \in \mathbb{R}^{10 \times 16 \times 16 \times 3}$. On SEVIR, it takes the concatenation along the temporal axis of $z_{\text{cond}} \in \mathbb{R}^{7 \times 16 \times 16 \times 4}$ and $z_t \in \mathbb{R}^{6 \times 16 \times 16 \times 4}$ as input, and outputs $z_{t-1} \in \mathbb{R}^{6 \times 16 \times 16 \times 4}$. Besides, we add the embedding of the denoising step $t$ to the state in front of each cuboid attention block via an embeding layer `TEmbed`, following [22]. The detailed configurations of the Earthformer-UNet is described in Table 8

**Knowledge control networks**  A knowledge control network parameterizes $U_\phi(z_t, t, y)$ to predict $\mathcal{F}(\widehat{x}, y)$ using the noisy latent $z_t$. In practice, we build an Earthformer encoder [9] with a final pooling block as the knowledge control network to parameterize $U_\phi(z_t, t, z_{\text{cond}})$, which takes $t$, and the concatenation of $z_{\text{cond}}$ and $z_t$, instead of $t$, $y$ and $z_t$ as the inputs. We find this implementation accurate enough when $t$ is small. The detailed configurations of the knowledge control network is described in Table 9

**Optimization**  We train the frame-wise VAEs using the Adam optimizer [27] following [7]. We train the LDMs and the knowledge control networks using the AdamW optimizer [31] following [9]. Detailed configurations are shown in Table 10,  Table 11 and  Table 12 for the frame-wise VAE, the LDM and the knowledge control network, respectively. We adopt data parallel and gradient accumulation to use a larger total batch size while the GPU can only afford a smaller micro batch size.

**Source code**  Example implementations in PyTorch-Lightning [36, 8] for training the frame-wise VAE, the Earthformer-UNet for latent diffusion, and the knowledge control network from scratch on $N$-body MNIST are available at `https://anonymous.4open.science/r/prediff`. The full implementations and the pretrained weights of PreDiff and PreDiff-KC will be released upon acceptance.

Table 3: The details of the encoder of the frame-wise VAE on $N$-body MNIST frames. It encodes an input frame $x^j \in \mathbb{R}^{64\times64\times1}$ into a latent $z^j \in \mathbb{R}^{16\times16\times3}$. Conv3 × 3 is the 2D convolutional layer with $3 \times 3$ kernel. GroupNorm32 is the Group Normalization (GN) layer [60] with 32 groups. SiLU is the Sigmoid Linear Unit activation layer [18] with function $\texttt{SiLU}(x) = x \cdot \texttt{sigmoid}(x)$. The Attention is the self attention layer [54] that first maps the input to queries $Q$, keys $K$ and values $V$ by three Linear layers, and then does self attention operation: $\texttt{Attention}(x) = \texttt{Softmax}(QK^T/\sqrt{C})V$.

| Block | Layer | Resolution | Channels |
|---|---|---|---|
| Input $x^j$ | - | $64 \times 64$ | 1 |
| 2D CNN | Conv3 × 3 | $64 \times 64$ | $1 \to 128$ |
| ResNet Block ×2 | GroupNorm32 | $64 \times 64$ | 128 |
|  | Conv3 × 3 | $64 \times 64$ | 128 |
|  | GroupNorm32 | $64 \times 64$ | 128 |
|  | Conv3 × 3 | $64 \times 64$ | 128 |
|  | SiLU | $64 \times 64$ | 128 |
| Downsampler | Conv3 × 3 | $64 \times 64 \to 32 \times 32$ | 128 |
| ResNet Block ×2 | GroupNorm32 | $32 \times 32$ | 128 |
|  | Conv3 × 3 | $32 \times 32$ | $128 \to 256, 256$ |
|  | GroupNorm32 | $32 \times 32$ | 256 |
|  | Conv3 × 3 | $32 \times 32$ | 256 |
|  | SiLU | $32 \times 32$ | 256 |
| Downsampler | Conv3 × 3 | $32 \times 32 \to 16 \times 16$ | 256 |
| ResNet Block ×2 | GroupNorm32 | $16 \times 16$ | 256 |
|  | Conv3 × 3 | $16 \times 16$ | $256 \to 512, 512$ |
|  | GroupNorm32 | $16 \times 16$ | 512 |
|  | Conv3 × 3 | $16 \times 16$ | 512 |
|  | SiLU | $16 \times 16$ | 512 |
| Self Attention Block | GroupNorm32 | $16 \times 16$ | 512 |
|  | Attention | $16 \times 16$ | 512 |
|  | Linear | $16 \times 16$ | 512 |
| ResNet Block ×2 | GroupNorm32 | $16 \times 16$ | 512 |
|  | Conv3 × 3 | $16 \times 16$ | 512 |
|  | GroupNorm32 | $16 \times 16$ | 512 |
|  | Conv3 × 3 | $16 \times 16$ | 512 |
|  | SiLU | $16 \times 16$ | 512 |
| Output Block | GroupNorm32 | $16 \times 16$ | 512 |
|  | SiLU | $16 \times 16$ | 512 |
|  | Conv3 × 3 | $16 \times 16$ | $512 \to 6$ |
|  | Conv3 × 3 | $16 \times 16$ | 6 |

Table 4: The details of the decoder of the frame-wise VAE on $N$-body MNIST frames. It decodes a latent $z^j \in \mathbb{R}^{16 \times 16 \times 3}$ back to a frame in pixel space $x^j \in \mathbb{R}^{64 \times 64 \times 1}$. `Conv3 × 3` is the 2D convolutional layer with $3 \times 3$ kernel. `GroupNorm32` is the Group Normalization (GN) layer [60] with 32 groups. `SiLU` is the Sigmoid Linear Unit activation layer [18] with function $\mathtt{SiLU}(x) = x \cdot \mathtt{sigmoid}(x)$. The `Attention` is the self attention layer [54] that first maps the input to queries $Q$, keys $K$ and values $V$ by three `Linear` layers, and then does self attention operation: $\mathtt{Attention}(x) = \mathtt{Softmax}(QK^T/\sqrt{C})V$.

| Block | Layer | Resolution | Channels |
|---|---|---|---|
| Input $z^j$ | - | $16 \times 16$ | 3 |
| 2D CNN | Conv3 × 3 | $16 \times 16$ | 3 |
| | Conv3 × 3 | $16 \times 16$ | $3 \rightarrow 512$ |
| Self Attention Block | GroupNorm32 | $16 \times 16$ | 512 |
| | Attention | $16 \times 16$ | 512 |
| | Linear | $16 \times 16$ | 512 |
| ResNet Block ×3 | GroupNorm32 | $16 \times 16$ | 512 |
| | Conv3 × 3 | $16 \times 16$ | 512 |
| | GroupNorm32 | $16 \times 16$ | 512 |
| | Conv3 × 3 | $16 \times 16$ | 512 |
| | SiLU | $16 \times 16$ | 512 |
| Upsampler | Conv3 × 3 | $16 \times 16 \rightarrow 32 \times 32$ | 512 |
| ResNet Block ×3 | GroupNorm32 | $32 \times 32$ | 512 |
| | Conv3 × 3 | $32 \times 32$ | $512 \rightarrow 256, 256, 256$ |
| | GroupNorm32 | $32 \times 32$ | 256 |
| | Conv3 × 3 | $32 \times 32$ | 256 |
| | SiLU | $32 \times 32$ | 256 |
| Upsampler | Conv3 × 3 | $32 \times 32 \rightarrow 64 \times 64$ | 256 |
| ResNet Block ×3 | GroupNorm32 | $64 \times 64$ | 256 |
| | Conv3 × 3 | $64 \times 64$ | $256 \rightarrow 128, 128, 128$ |
| | GroupNorm32 | $64 \times 64$ | 128 |
| | Conv3 × 3 | $64 \times 64$ | 128 |
| | SiLU | $64 \times 64$ | 128 |
| Output Block | GroupNorm32 | $64 \times 64$ | 128 |
| | SiLU | $64 \times 64$ | 128 |
| | Conv3 × 3 | $64 \times 64$ | $128 \rightarrow 1$ |

Table 5: The details of the encoder of the frame-wise VAE on SEVIR frames. It encodes an input frame $x^j \in \mathbb{R}^{128 \times 128 \times 1}$ into a latent $z^j \in \mathbb{R}^{16 \times 16 \times 4}$. Conv3 $\times$ 3 is the 2D convolutional layer with $3 \times 3$ kernel. GroupNorm32 is the Group Normalization (GN) layer [60] with 32 groups. SiLU is the Sigmoid Linear Unit activation layer [18] with function $\texttt{SiLU}(x) = x \cdot \texttt{sigmoid}(x)$. The Attention is the self attention layer [54] that first maps the input to queries $Q$, keys $K$ and values $V$ by three Linear layers, and then does self attention operation: $\texttt{Attention}(x) = \texttt{Softmax}(QK^T/\sqrt{C})V$.

| Block | Layer | Resolution | Channels |
|---|---|---|---|
| Input $x^j$ | - | $128 \times 128$ | 1 |
| 2D CNN | Conv3 $\times$ 3 | $128 \times 128$ | $1 \rightarrow 128$ |
| ResNet Block $\times 2$ | GroupNorm32 | $128 \times 128$ | 128 |
| | Conv3 $\times$ 3 | $128 \times 128$ | 128 |
| | GroupNorm32 | $128 \times 128$ | 128 |
| | Conv3 $\times$ 3 | $128 \times 128$ | 128 |
| | SiLU | $128 \times 128$ | 128 |
| Downsampler | Conv3 $\times$ 3 | $128 \times 128 \rightarrow 64 \times 64$ | 128 |
| ResNet Block $\times 2$ | GroupNorm32 | $64 \times 64$ | 128 |
| | Conv3 $\times$ 3 | $64 \times 64$ | $128 \rightarrow 256, 256$ |
| | GroupNorm32 | $64 \times 64$ | 256 |
| | Conv3 $\times$ 3 | $64 \times 64$ | 256 |
| | SiLU | $64 \times 64$ | 256 |
| Downsampler | Conv3 $\times$ 3 | $64 \times 64 \rightarrow 32 \times 32$ | 256 |
| ResNet Block $\times 2$ | GroupNorm32 | $32 \times 32$ | 256 |
| | Conv3 $\times$ 3 | $32 \times 32$ | $256 \rightarrow 512, 512$ |
| | GroupNorm32 | $32 \times 32$ | 512 |
| | Conv3 $\times$ 3 | $32 \times 32$ | 512 |
| | SiLU | $32 \times 32$ | 512 |
| Downsampler | Conv3 $\times$ 3 | $32 \times 32 \rightarrow 16 \times 16$ | 512 |
| ResNet Block $\times 2$ | GroupNorm32 | $16 \times 16$ | 512 |
| | Conv3 $\times$ 3 | $16 \times 16$ | 512 |
| | GroupNorm32 | $16 \times 16$ | 512 |
| | Conv3 $\times$ 3 | $16 \times 16$ | 512 |
| | SiLU | $16 \times 16$ | 512 |
| Self Attention Block | GroupNorm32 | $16 \times 16$ | 512 |
| | Attention | $16 \times 16$ | 512 |
| | Linear | $16 \times 16$ | 512 |
| ResNet Block $\times 2$ | GroupNorm32 | $16 \times 16$ | 512 |
| | Conv3 $\times$ 3 | $16 \times 16$ | 512 |
| | GroupNorm32 | $16 \times 16$ | 512 |
| | Conv3 $\times$ 3 | $16 \times 16$ | 512 |
| | SiLU | $16 \times 16$ | 512 |
| Output Block | GroupNorm32 | $16 \times 16$ | 512 |
| | SiLU | $16 \times 16$ | 512 |
| | Conv3 $\times$ 3 | $16 \times 16$ | $512 \rightarrow 8$ |
| | Conv3 $\times$ 3 | $16 \times 16$ | 8 |

Table 6: The details of the decoder of the frame-wise VAE on SEVIR frames. It decodes a latent $z^j \in \mathbb{R}^{16 \times 16 \times 4}$ back to a frame in pixel space $x^j \in \mathbb{R}^{128 \times 128 \times 1}$. Conv3 $\times$ 3 is the 2D convolutional layer with $3 \times 3$ kernel. GroupNorm32 is the Group Normalization (GN) layer [60] with 32 groups. SiLU is the Sigmoid Linear Unit activation layer [18] with function $\texttt{SiLU}(x) = x \cdot \texttt{sigmoid}(x)$. The Attention is the self attention layer [54] that first maps the input to queries $Q$, keys $K$ and values $V$ by three Linear layers, and then does self attention operation: $\texttt{Attention}(x) = \texttt{Softmax}(QK^T/\sqrt{C})V$.

| Block | Layer | Resolution | Channels |
|---|---|---|---|
| Input $z^j$ | - | $16 \times 16$ | 4 |
| 2D CNN | Conv3 $\times$ 3 | $16 \times 16$ | 4 |
| | Conv3 $\times$ 3 | $16 \times 16$ | $4 \rightarrow 512$ |
| Self Attention Block | GroupNorm32 | $16 \times 16$ | 512 |
| | Attention | $16 \times 16$ | 512 |
| | Linear | $16 \times 16$ | 512 |
| ResNet Block $\times 3$ | GroupNorm32 | $16 \times 16$ | 512 |
| | Conv3 $\times$ 3 | $16 \times 16$ | 512 |
| | GroupNorm32 | $16 \times 16$ | 512 |
| | Conv3 $\times$ 3 | $16 \times 16$ | 512 |
| | SiLU | $16 \times 16$ | 512 |
| Upsampler | Conv3 $\times$ 3 | $16 \times 16 \rightarrow 32 \times 32$ | 512 |
| ResNet Block $\times 3$ | GroupNorm32 | $32 \times 32$ | 512 |
| | Conv3 $\times$ 3 | $32 \times 32$ | 512 |
| | GroupNorm32 | $32 \times 32$ | 512 |
| | Conv3 $\times$ 3 | $32 \times 32$ | 512 |
| | SiLU | $32 \times 32$ | 512 |
| Upsampler | Conv3 $\times$ 3 | $32 \times 32 \rightarrow 64 \times 64$ | 512 |
| ResNet Block $\times 3$ | GroupNorm32 | $64 \times 64$ | 512 |
| | Conv3 $\times$ 3 | $64 \times 64$ | $512 \rightarrow 256, 256, 256$ |
| | GroupNorm32 | $64 \times 64$ | 256 |
| | Conv3 $\times$ 3 | $64 \times 64$ | 256 |
| | SiLU | $64 \times 64$ | 256 |
| Upsampler | Conv3 $\times$ 3 | $64 \times 64 \rightarrow 128 \times 128$ | 256 |
| ResNet Block $\times 3$ | GroupNorm32 | $128 \times 128$ | 256 |
| | Conv3 $\times$ 3 | $128 \times 128$ | $256 \rightarrow 128, 128, 128$ |
| | GroupNorm32 | $128 \times 128$ | 128 |
| | Conv3 $\times$ 3 | $128 \times 128$ | 128 |
| | SiLU | $128 \times 128$ | 128 |
| Output Block | GroupNorm32 | $128 \times 128$ | 128 |
| | SiLU | $128 \times 128$ | 128 |
| | Conv3 $\times$ 3 | $128 \times 128$ | $128 \rightarrow 1$ |

Table 7: The details of the discriminator for the adversarial loss of on $N$-body MNIST and SEVIR frames. `Conv4 × 4` is the 2D convolutional layer with $4 \times 4$ kernel, $2 \times 2$ or $1 \times 1$ stride, and $1 \times 1$ padding. `BatchNorm` is the Batch Normalization (BN) layer [25] . The negative slope in `LeakyReLU` is 0.2.

| Block | Layer | Resolution | | Channels |
| | | $N$-body MNIST | SEVIR | |
|---|---|---|---|---|
| Input $x^j$ | - | $64 \times 64$ | $128 \times 128$ | 1 |
| 2D CNN | Conv4 × 4 | $64 \times 64 \to 32 \times 32$ | $128 \times 128 \to 64 \times 64$ | $1 \to 64$ |
| Downsampler | LeakyReLU | $32 \times 32$ | $64 \times 64$ | 64 |
| | Conv4 × 4 | $32 \times 32 \to 16 \times 16$ | $64 \times 64 \to 32 \times 32$ | $64 \to 128$ |
| | BatchNorm | $16 \times 16$ | $32 \times 32$ | 128 |
| Downsampler | LeakyReLU | $16 \times 16$ | $32 \times 32$ | 128 |
| | Conv4 × 4 | $16 \times 16 \to 8 \times 8$ | $32 \times 32 \to 16 \times 16$ | $128 \to 256$ |
| | BatchNorm | $8 \times 8$ | $16 \times 16$ | 256 |
| Downsampler | LeakyReLU | $8 \times 8$ | $16 \times 16$ | 256 |
| | Conv4 × 4 | $8 \times 8 \to 7 \times 7$ | $16 \times 16 \to 15 \times 15$ | $256 \to 512$ |
| | BatchNorm | $7 \times 7$ | $15 \times 15$ | 512 |
| Output Block | LeakyReLU | $7 \times 7$ | $15 \times 15$ | 512 |
| | Conv4 × 4 | $7 \times 7 \to 6 \times 6$ | $15 \times 15 \to 14 \times 14$ | 1 |
| | AvgPool | $6 \times 6 \to 1$ | $15 \times 15 \to 1$ | 1 |

Table 8: The details of the Earthformer-UNet as the latent diffusion model on $N$-body MNIST and SEVIR datasets. The `ConcatMask` layer for ObservationMask block concatenates one more channel to the input to indicates whether the input is the encoded observation $z_{\mathrm{cond}}$ or the noisy latent $z_t$. 1 for $z_{\mathrm{cond}}$ and 0 for $z_t$. `Conv3 × 3` is the 2D convolutional layer with $3 \times 3$ kernel. `GroupNorm32` is the Group Normalization (GN) layer [60] with 32 groups. If the number of the input data channels is smaller than 32, then the number of groups is set to the number of channels. `SiLU` is the Sigmoid Linear Unit activation layer [18] with function $\mathtt{SiLU}(x) = x \cdot \mathtt{sigmoid}(x)$. The negative slope in `LeakyReLU` is 0.1. `Dropout` is the dropout layer [21] with the probability 0.1 to drop an element to be zeroed. The `FFN` consists of two `Linear` layers separated by a `GeLU` activation layer [18]. `PosEmbed` is the positional embedding layer [54] that adds learned positional embeddings to the input. `TEmbed` is the embedding layer [22] that embeds the denoising step $t$. `PatchMerge` splits a 2D input tensor with $C$ channels into $N$ non-overlapping $p \times p$ patches and merges the spatial dimensions into channels, gets $N$ $1 \times 1$ patches with $p^2 \cdot C$ channels and concatenates them back along spatial dimensions. Residual connections [17] are added from blocks in the downsampling phase to corresponding blocks in the upsampling phase.

| Block | Layer | Spatial Resolution | Channels $N$-body MNIST | SEVIR |
|---|---|---|---|---|
| Input $[z_{\mathrm{cond}}, z_t]$ | - | $16 \times 16$ | 3 | 4 |
| Observation ask | ConcatMask | $16 \times 16$ | $3 \to 4$ | $4 \to 5$ |
| Projector | GroupNorm32 | $16 \times 16$ | 4 | 5 |
| | SiLU | $16 \times 16$ | 4 | 5 |
| | Conv3 × 3 | $16 \times 16$ | $4 \to 256$ | $5 \to 256$ |
| | GroupNorm32 | $16 \times 16$ | 256 | |
| | SiLU | $16 \times 16$ | 256 | |
| | Dropout | $16 \times 16$ | 256 | |
| | Conv3 × 3 | $16 \times 16$ | 256 | |
| Positional Embedding | PosEmbed | $16 \times 16$ | 256 | |
| Cuboid Attention Block ×4 | TEmbed | $16 \times 16$ | 256 | |
| | LayerNorm | $16 \times 16$ | 256 | |
| | Cuboid(T, 1, 1) | $16 \times 16$ | 256 | |
| | FFN | $16 \times 16$ | 256 | |
| | LayerNorm | $16 \times 16$ | 256 | |
| | Cuboid(1, H, 1) | $16 \times 16$ | 256 | |
| | FFN | $16 \times 16$ | 256 | |
| | LayerNorm | $16 \times 16$ | 256 | |
| | Cuboid(1, 1, W) | $16 \times 16$ | 256 | |
| | FFN | $16 \times 16$ | 256 | |
| Downsampler | PatchMerge | $16 \times 16 \to 8 \times 8$ | $256 \to 1024$ | |
| | LayerNorm | $8 \times 8$ | 1024 | |
| | Linear | $8 \times 8$ | 1024 | |
| Cuboid Attention Block ×8 | TEmbed | $8 \times 8$ | 1024 | |
| | LayerNorm | $8 \times 8$ | 1024 | |
| | Cuboid(T, 1, 1) | $8 \times 8$ | 1024 | |
| | FFN | $8 \times 8$ | 1024 | |
| | LayerNorm | $8 \times 8$ | 1024 | |
| | Cuboid(1, H, 1) | $8 \times 8$ | 1024 | |
| | FFN | $8 \times 8$ | 1024 | |
| | LayerNorm | $8 \times 8$ | 1024 | |
| | Cuboid(1, 1, W) | $8 \times 8$ | 1024 | |
| | FFN | $8 \times 8$ | 1024 | |
| Upsampler | NearestNeighborInterp | $8 \times 8 \to 16 \times 16$ | 1024 | |
| | Conv3 × 3 | $16 \times 16$ | $1024 \to 256$ | |
| Cuboid Attention Block ×4 | TEmbed | $16 \times 16$ | 256 | |
| | LayerNorm | $16 \times 16$ | 256 | |
| | Cuboid(T, 1, 1) | $16 \times 16$ | 256 | |
| | FFN | $16 \times 16$ | 256 | |
| | LayerNorm | $16 \times 16$ | 256 | |
| | Cuboid(1, H, 1) | $16 \times 16$ | 256 | |
| | FFN | $16 \times 16$ | 256 | |
| | LayerNorm | $16 \times 16$ | 256 | |
| | Cuboid(1, 1, W) | $16 \times 16$ | 256 | |
| | FFN | $16 \times 16$ | 256 | |
| Output Block | Linear | $16 \times 16$ | $256 \to 3$ | $256 \to 4$ |

Table 9: The details of the Earthformer encoders for the parameterization of the knowledge control networks $U_\phi(z_t, t, z_{\text{cond}})$ on $N$-body MNIST and SEVIR datasets. The `ConcatMask` layer for ObservationMask block concatenates one more channel to the input to indicates whether the input is the encoded observation $z_{\text{cond}}$ or the noisy latent $z_t$. 1 for $z_{\text{cond}}$ and 0 for $z_t$. `Conv3 × 3` is the 2D convolutional layer with $3 \times 3$ kernel. `GroupNorm32` is the Group Normalization (GN) layer [60] with 32 groups. If the number of the input data channels is smaller than 32, then the number of groups is set to the number of channels. `SiLU` is the Sigmoid Linear Unit activation layer [18] with function $\text{SiLU}(x) = x \cdot \text{sigmoid}(x)$. The negative slope in `LeakyReLU` is 0.1. `Dropout` is the dropout layer [21] with the probability 0.1 to drop an element to be zeroed. The `FFN` consists of two `Linear` layers separated by a `GeLU` activation layer [18]. `PosEmbed` is the positional embedding layer [54] that adds learned positional embeddings to the input. `TEmbed` is the embedding layer [22] that embeds the denoising step $t$. `PatchMerge` splits a 2D input tensor with $C$ channels into $N$ non-overlapping $p \times p$ patches and merges the spatial dimensions into channels, gets $N$ $1 \times 1$ patches with $p^2 \cdot C$ channels and concatenates them back along spatial dimensions. Residual connections [17] are added from blocks in the downsampling phase to corresponding blocks in the upsampling phase. The `Attention` is the self attention layer [54] with an extra "cls" token for information aggregation. It first flattens the input and concatenates it with the "cls" token. Then it maps the concatenated input to queries $Q$, keys $K$ and values $V$ by three `Linear` layers, and then does self attention operation: $\text{Attention}(x) = \text{Softmax}(QK^T/\sqrt{C})V$. Finally, the value of the "cls" token after self attention operation serves as the layer's output.

| Block | Layer | Spatial Resolution | Channels $N$-body MNIST | SEVIR |
|---|---|---|---|---|
| Input $[z_{\text{cond}}, z_t]$ | - | $16 \times 16$ | 3 | 4 |
| Observation ask | ConcatMask | $16 \times 16$ | $3 \to 4$ | $4 \to 5$ |
| Projector | GroupNorm32 | $16 \times 16$ | 4 | 5 |
| | SiLU | $16 \times 16$ | 4 | 5 |
| | Conv3 × 3 | $16 \times 16$ | $4 \to 64$ | $5 \to 64$ |
| | GroupNorm32 | $16 \times 16$ | 64 | |
| | SiLU | $16 \times 16$ | 64 | |
| | Dropout | $16 \times 16$ | 64 | |
| | Conv3 × 3 | $16 \times 16$ | 64 | |
| Positional Embedding | PosEmbed | $16 \times 16$ | 64 | |
| Cuboid Attention Block | TEmbed | $16 \times 16$ | 64 | |
| | LayerNorm | $16 \times 16$ | 64 | |
| | Cuboid(T, 1, 1) | $16 \times 16$ | 64 | |
| | FFN | $16 \times 16$ | 64 | |
| | LayerNorm | $16 \times 16$ | 64 | |
| | Cuboid(1, H, 1) | $16 \times 16$ | 64 | |
| | FFN | $16 \times 16$ | 64 | |
| | LayerNorm | $16 \times 16$ | 64 | |
| | Cuboid(1, 1, W) | $16 \times 16$ | 64 | |
| | FFN | $16 \times 16$ | 64 | |
| Downsampler | PatchMerge | $16 \times 16 \to 8 \times 8$ | $64 \to 256$ | |
| | LayerNorm | $8 \times 8$ | 256 | |
| | Linear | $8 \times 8$ | 256 | |
| Cuboid Attention Block | TEmbed | $8 \times 8$ | 256 | |
| | LayerNorm | $8 \times 8$ | 256 | |
| | Cuboid(T, 1, 1) | $8 \times 8$ | 256 | |
| | FFN | $8 \times 8$ | 256 | |
| | LayerNorm | $8 \times 8$ | 256 | |
| | Cuboid(1, H, 1) | $8 \times 8$ | 256 | |
| | FFN | $8 \times 8$ | 256 | |
| | LayerNorm | $8 \times 8$ | 256 | |
| | Cuboid(1, 1, W) | $8 \times 8$ | 256 | |
| | FFN | $8 \times 8$ | 256 | |
| Output Pooling Block | GroupNorm32 | $8 \times 8$ | 256 | |
| | Attention | $8 \times 8 \to 1$ | 256 | |
| | Linear | 1 | $256 \to 1$ | |

Table 10: Hyperparameters of the Adam optimizer for training frame-wise VAEs and discriminators on $N$-body MNIST and SEVIR datasets.

| Hyper-parameter of VAE | Value |
| --- | --- |
| Learning rate | $4.5 \times 10^{-6}$ |
| $\beta_1$ | 0.5 |
| $\beta_2$ | 0.9 |
| Weight decay | $10^{-2}$ |
| Batch size | 512 |
| Training epochs | 200 |
| Hyper-parameter of discriminator | Value |
| Learning rate | $4.5 \times 10^{-6}$ |
| $\beta_1$ | 0.5 |
| $\beta_2$ | 0.9 |
| Weight decay | $10^{-2}$ |
| Batch size | 512 |
| Training epochs | 200 |
| Training start step | 50000 |

Table 11: Hyperparameters of the AdamW optimizer for training LDMs on $N$-body MNIST and SEVIR datasets.

| Hyper-parameter of VAE | Value |
| --- | --- |
| Learning rate | $1.0 \times 10^{-3}$ |
| $\beta_1$ | 0.9 |
| $\beta_2$ | 0.999 |
| Weight decay | $10^{-5}$ |
| Batch size | 64 |
| Training epochs | 1000 |
| Warm up percentage | 10% |
| Learning rate decay | Cosine |

Table 12: Hyperparameters of the AdamW optimizer for training knowledge control networks on $N$-body MNIST and SEVIR datasets.

| Hyper-parameter of VAE | Value |
| --- | --- |
| Learning rate | $1.0 \times 10^{-3}$ |
| $\beta_1$ | 0.9 |
| $\beta_2$ | 0.999 |
| Weight decay | $10^{-5}$ |
| Batch size | 64 |
| Training epochs | 200 |
| Warm up percentage | 10% |
| Learning rate decay | Cosine |

## A.2 Baselines

We train baseline algorithms following their officially released configurations and tune the learning rate, learning rate scheduler, working resolution, etc., to optimize their performance on each dataset. We list the modifications we applied to the baselines for each dataset in Table 13.

Table 13: Implementation details of baseline algorithms. Modifications based on the officially released implementations are listed according to different datasets. "-" means no modification is applied. "reverse enc-dec" means adopting the reversed encoder-decoder architecture proposed in [48]. Other terms listed are the hyperparameters in their officially released implementations.

| Model | $N$-body MNIST | SEVIR |
|---|---|---|
| UNet [55] | - | - |
| ConvLSTM [47] | reverse enc-dec [48]<br>conv_kernels = [(7,7),(5,5),(3,3)]<br>deconv_kernels = [(6,6),(4,4),(4,4)]<br>channels=[96, 128, 256] | reverse enc-dec [48]<br>conv_kernels = [(7,7),(5,5),(3,3)]<br>deconv_kernels = [(6,6),(4,4),(4,4)]<br>channels=[96, 128, 256] |
| PredRNN [58] | - | - |
| PhyDNet [12] | - | convcell_hidden = [256, 256, 256, 64] |
| E3D-LSTM [57] | - | - |
| Rainformer [1] | downscaling_factors=[2, 2, 2, 2]<br>hidden_dim=32<br>heads=[4, 4, 8, 16]<br>head_dim=8<br>- | downscaling_factors=[4, 2, 2, 2]<br>-<br>-<br>-<br>- |
| Earthformer [9] | - | - |
| VideoGPT [61] | vqvae_n_codes = 512<br>vqvae_downsample = [1, 4, 4] | vqvae_n_codes = 512<br>vqvae_downsample = [1, 8, 8] |
| LDM [42] | vae: $64 \times 64 \times 1 \rightarrow 16 \times 16 \times 3$<br>conv_dim = 3<br>model_channels = 256 | vae: $128 \times 128 \times 1 \rightarrow 16 \times 16 \times 4$<br>conv_dim = 3<br>model_channels = 256 |
| DGMR [41] | - | - |

# B  Derivation of the Approximation to Knowledge Control Guidance

We derive the approximation to the knowledge control guided denoising transition (5) following [4]. We rewrite (5) to (8) using a normalization constant $Z$ that normalizes $Z \int e^{-\lambda_{\mathcal{F}} \|U_\phi(z_t,t,y) - \mathcal{F}_0(y)\|} dz_t = 1$:

$$p_{\theta,\phi}(z_t|z_{t+1}, y, \mathcal{F}_0) = p_\theta(z_t|z_{t+1}, z_{\text{cond}}) \cdot Z e^{-\lambda_{\mathcal{F}} \|U_\phi(z_t,t,y) - \mathcal{F}_0(y)\|}. \tag{8}$$

In the folowing, we abbreviate $\mu_\theta(z_{t+1}, t, z_{\text{cond}})$ as $\mu_\theta$, and $\Sigma_\theta(z_{t+1}, t, z_{\text{cond}})$ as $\Sigma_\theta$ for brevity. We use $C_i, i = \{1, \ldots, 7\}$ to denote constants.

$$p_\theta(z_t|z_{t+1}, z_{\text{cond}}) = \mathcal{N}(\mu_\theta, \Sigma_\theta),$$
$$\log p_\theta(z_t|z_{t+1}, z_{\text{cond}}) = -\frac{1}{2}(z_t - \mu_\theta)^T \Sigma_\theta^{-1}(z_t - \mu_\theta) + C_1, \tag{9}$$
$$\log Z e^{-\lambda_{\mathcal{F}} \|U_\phi(z_t,t,y) - \mathcal{F}_0(y)\|} = -\lambda_{\mathcal{F}} \|U_\phi(z_t,t,y) - \mathcal{F}_0(y)\| + C_2,$$

By assuming that $\log Z e^{-\lambda_{\mathcal{F}} \|U_\phi(z_t,t,y) - \mathcal{F}_0(y)\|}$ has low curvature compared to $\Sigma_\theta^{-1}$, which is reasonable in the limit of infinite diffusion steps ($\|\Sigma_\theta\| \to 0$), we can approximate it by a Taylor expansion at $z_t = \mu_\theta$

$$\begin{aligned} \log Z e^{-\lambda_{\mathcal{F}} \|U_\phi(z_t,t,y) - \mathcal{F}_0(y)\|} &\approx -\lambda_{\mathcal{F}} \|U_\phi(z_t,t,y) - \mathcal{F}_0(y)\| |_{z_t = \mu_\theta} \\ &\quad - (z_t - \mu_\theta) \lambda_{\mathcal{F}} \nabla_{z_t} \|U_\phi(z_t,t,y) - \mathcal{F}_0(y)\| |_{z_t = \mu_\theta} \\ &= (z_t - \mu_\theta) g + C_3, \end{aligned} \tag{10}$$

where $g = -\lambda_{\mathcal{F}} \nabla_{z_t} \|U_\phi(z_t,t,y) - \mathcal{F}_0(y)\| |_{z_t = \mu_\theta}$. By taking the log of (8) and applying the results from (9) and (10), we get

$$\begin{aligned} \log p_{\theta,\phi}(z_t|z_{t+1}, y, \mathcal{F}_0) &= \log p_\theta(z_t|z_{t+1}, z_{\text{cond}}) + \log Z e^{-\lambda_{\mathcal{F}} \|U_\phi(z_t,t,y) - \mathcal{F}_0(y)\|} \\ &\approx -\frac{1}{2}(z_t - \mu_\theta)^T \Sigma_\theta^{-1}(z_t - \mu_\theta) + (z_t - \mu_\theta) g + C_4 \\ &= -\frac{1}{2}(z_t - \mu_\theta - \Sigma_\theta g)^T \Sigma_\theta^{-1}(z_t - \mu_\theta - \Sigma_\theta g) + \frac{1}{2} g^T \Sigma_\theta g + C_5 \\ &= -\frac{1}{2}(z_t - \mu_\theta - \Sigma_\theta g)^T \Sigma_\theta^{-1}(z_t - \mu_\theta - \Sigma_\theta g) + C_6 \\ &= \log p(z) + C_7, z \sim \mathcal{N}(\mu_\theta + \Sigma_\theta g, \Sigma_\theta). \end{aligned} \tag{11}$$

Therefore, the transition distribution under the guidance of knowledge control shown in (5) can be approximated by a Gaussian similar to the transition without knowledge guidance, but with its mean shifted by $\Sigma_\theta$.

# C More Quantitative Analysis: BIAS on SEVIR

Similar to Critical Success Index (CSI) introduced in Sec. 4.2, $\text{BIAS} = \frac{\#\text{Hits}+\#\text{F.Alarms}}{\#\text{Hits}+\#\text{Misses}}$ is calculated by counting the #Hits (truth=1, pred=1), #Misses (truth=1, pred=0) and #F.Alarms (truth=0, pred=1) of the predictions binarized at thresholds $[16, 74, 133, 160, 181, 219]$. This measurement assesses the model's inclination towards either F.Alarms or Misses.

The results from Table 14 demonstrate that deterministic spatiotemporal forecasting models, such as UNet [55], ConvLSTM [47], PredRNN [58], PhyDNet [12], E3D-LSTM [57], and Earthformer [9], tend to produce predictions with lower intensity. These models prioritize avoiding high-intensity predictions that have a higher chance of being incorrect due to their limited ability to handle such uncertainty effectively. On the other hand, probabilistic spatiotemporal forecasting baselines, including DGMR [41], VideoGPT [61] and LDM [42], demonstrate a more daring approach by predicting possible high-intensity signals, even if it results in lower CSI scores, as depicted in Table 2. Among these baselines, PreDiff achieves the best performance in BIAS. It consistently achieves BIAS scores closest to 1, irrespective of the chosen threshold. These results demonstrate that PreDiff has effectively learned to unbiasedly capture the distribution of intensity.

Table 14: Performance comparison on SEVIR. The BIAS is calculated at different precipitation thresholds and denoted as BIAS-$thresh$. BIAS-m reports the mean of BIAS-$[16, 74, 133, 160, 181, 219]$. A BIAS score closer to 1 indicates that the model is less biased to either F.Alarms or Misses. The best BIAS score is in boldface while the second best is underscored.

| Model | Metrics | | | | | | |
| | BIAS-m | BIAS-219 | BIAS-181 | BIAS-160 | BIAS-133 | BIAS-74 | BIAS-16 |
|---|---|---|---|---|---|---|---|
| Persistence | 1.0177 | 1.0391 | 1.0323 | 1.0258 | 1.0099 | 1.0016 | 0.9983 |
| UNet [55] | 0.6658 | 0.2503 | 0.4013 | 0.5428 | 0.7665 | 0.9551 | 1.0781 |
| ConvLSTM [47] | 0.8341 | 0.5344 | 0.6811 | 0.7626 | <u>0.9643</u> | **0.9957** | 1.0663 |
| PredRNN [58] | 0.6605 | 0.2565 | 0.4377 | 0.4909 | 0.6806 | 0.9419 | 1.1554 |
| PhyDNet [12] | 0.6798 | 0.3970 | 0.6593 | 0.7312 | 1.0543 | 1.0553 | 1.2238 |
| E3D-LSTM [57] | 0.6925 | 0.2696 | 0.4861 | 0.5686 | 0.8352 | 0.9887 | 1.0070 |
| Earthformer [9] | 0.7043 | 0.2423 | 0.4605 | 0.5734 | 0.8623 | 0.9733 | 1.1140 |
| DGRM [41] | 0.7302 | 0.3704 | 0.5254 | 0.6495 | 0.8312 | 0.9594 | 1.0456 |
| VideoGPT [61] | <u>0.8594</u> | 0.6106 | <u>0.7738</u> | <u>0.8629</u> | 0.9606 | 0.9681 | 0.9805 |
| LDM [42] | 1.2951 | <u>1.4534</u> | 1.3525 | 1.3827 | 1.3154 | 1.1817 | 1.0847 |
| PreDiff | **0.9769** | **0.9647** | **0.9268** | **0.9617** | **0.9978** | <u>1.0047</u> | **1.0058** |

# D    More Qualitative Results on $N$-body MNIST

Fig. 5 to Fig. 12 show several sets of example predictions on the $N$-body MNIST test set. In each figure, visualizations from top to bottom are context sequence $y$, target sequence $x$, predictions by ConvLSTM [47], Earthformer [9], VideoGPT [61], LDM [42], PreDiff, PreDiff-KC. E.MSE denotes the average error between the total energy (the sum of kinetic energy and potential energy) of the predictions $E(\widehat{x}^j)$ and the total energy of the last step context $E(y^{L_{\text{in}}})$.
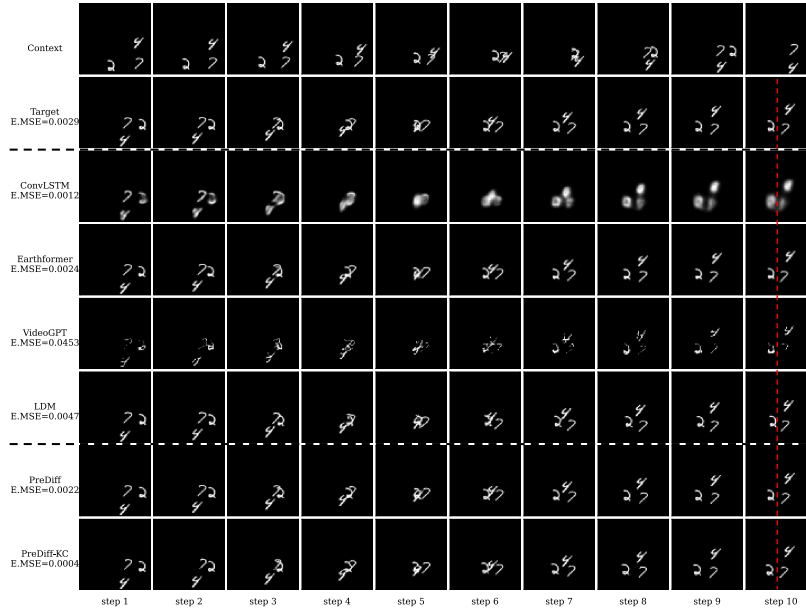


Figure 5: A set of example predictions on the $N$-body MNIST test set. The red dashed line is to help the reader to judge the position of the digit "1" in the last frame.



Figure 6: A set of example predictions on the $N$-body MNIST test set. The red dashed line is to help the reader to judge the position of the digit "0" in the last frame.

Figure 7: A set of example predictions on the $N$-body MNIST test set. The red dashed line is to help the reader to judge the position of the digit "0" in the last frame.



Figure 8: A set of example predictions on the $N$-body MNIST test set. The red dashed line is to help the reader to judge the position of the digit "8" in the last frame.

Figure 9: A set of example predictions on the $N$-body MNIST test set. The red dashed line is to help the reader to judge the position of the digit "4" in the last frame.



Figure 10: A set of example predictions on the $N$-body MNIST test set. The red dashed line is to help the reader to judge the position of the digit "1" in the last frame.

Figure 11: A set of example predictions on the $N$-body MNIST test set. The red dashed line is to help the reader to judge the position of the digit "7" in the last frame.



Figure 12: A set of example predictions on the $N$-body MNIST test set. The red dashed line is to help the reader to judge the position of the digit "7" in the last frame.

# E More Qualitative Results on SEVIR

Fig. 13 to Fig. 18 show several sets of example predictions on the SEVIR test set. In subfigure (a) of each figure, visualizations from top to bottom are context sequence $y$, target sequence $x$, predictions by ConvLSTM [47], Earthformer [9], VideoGPT [61], LDM [42], PreDiff, PreDiff-KC. In subfigure (b) of each figure, visualizations from top to bottom are context sequence $y$, target sequence $x$, predictions by PreDiff-KC with anticipated average future intensity $\mu_\tau + n\sigma_\tau$, $n = 4, 2, 0 - 2, -4$.



Figure 13: A set of example predictions on the SEVIR test set. (a) Comparison of PreDiff with baselines. (b) Predictions by PreDiff-KC under the guidance of anticipated average intensity.
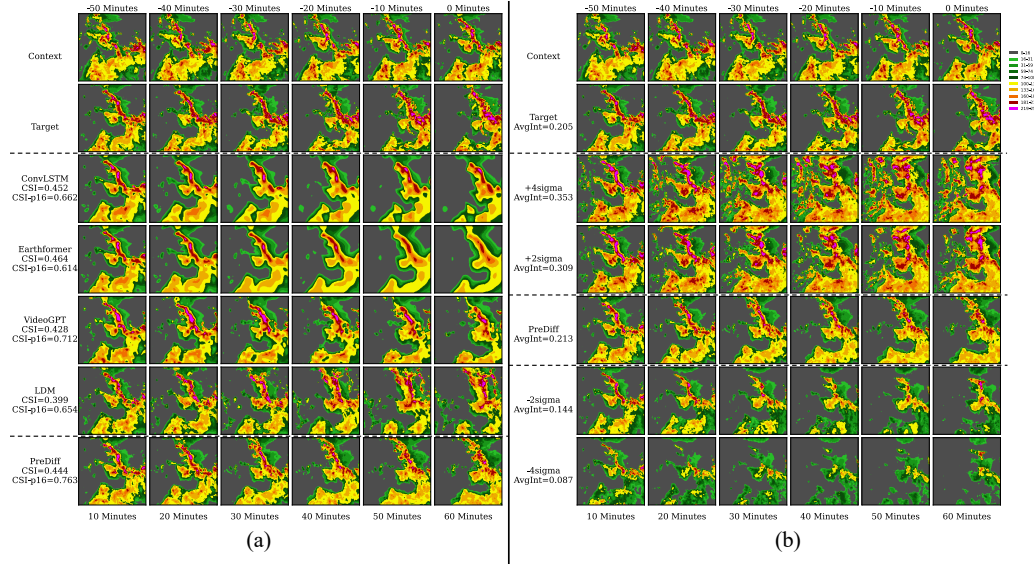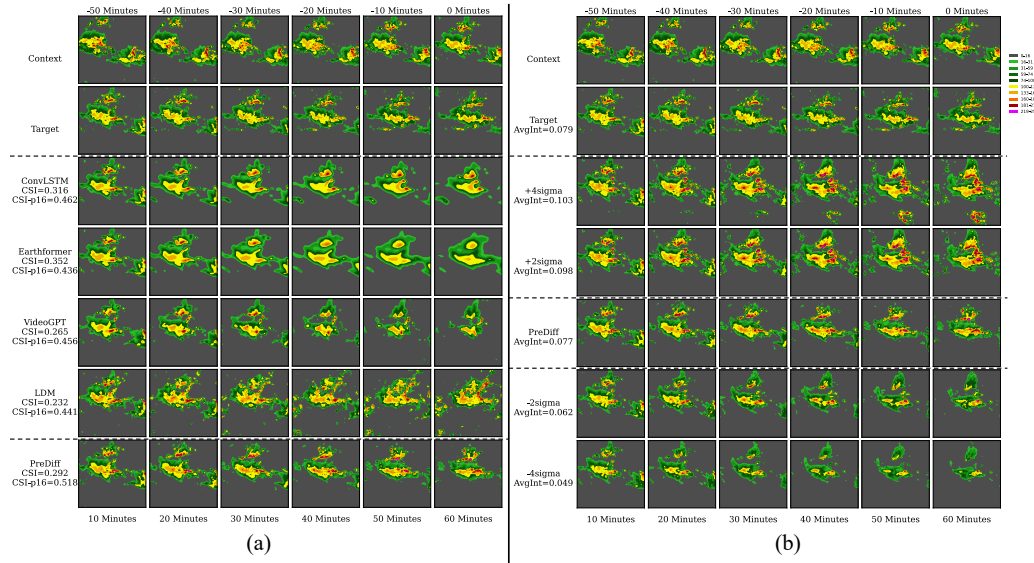


Figure 14: A set of example predictions on the SEVIR test set. (a) Comparison of PreDiff with baselines. (b) Predictions by PreDiff-KC under the guidance of anticipated average intensity.

Figure 15: A set of example predictions on the SEVIR test set. (a) Comparison of PreDiff with baselines. (b) Predictions by PreDiff-KC under the guidance of anticipated average intensity.
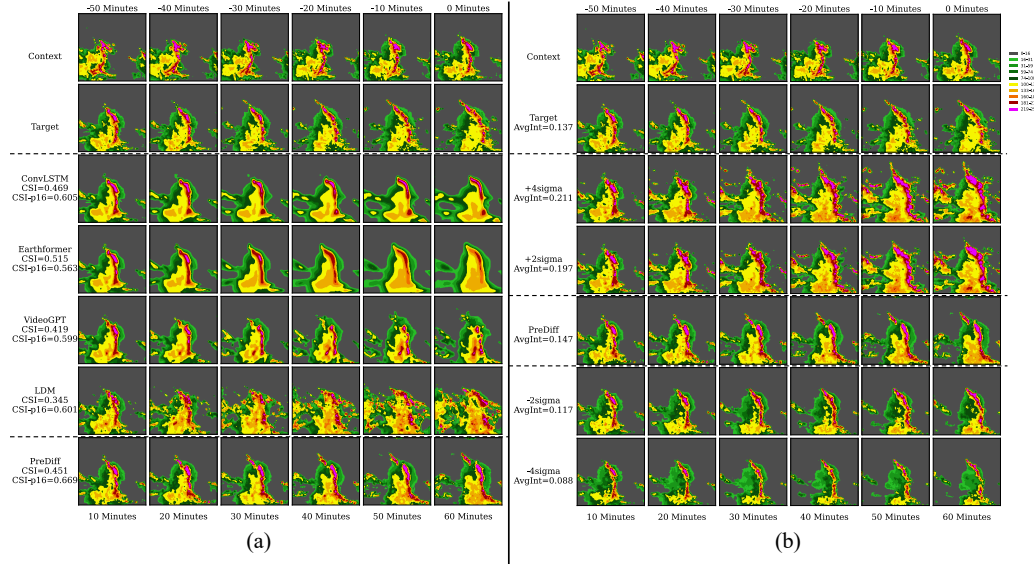


Figure 16: A set of example predictions on the SEVIR test set. (a) Comparison of PreDiff with baselines. (b) Predictions by PreDiff-KC under the guidance of anticipated average intensity.

Figure 17: A set of example predictions on the SEVIR test set. (a) Comparison of PreDiff with baselines. (b) Predictions by PreDiff-KC under the guidance of anticipated average intensity.
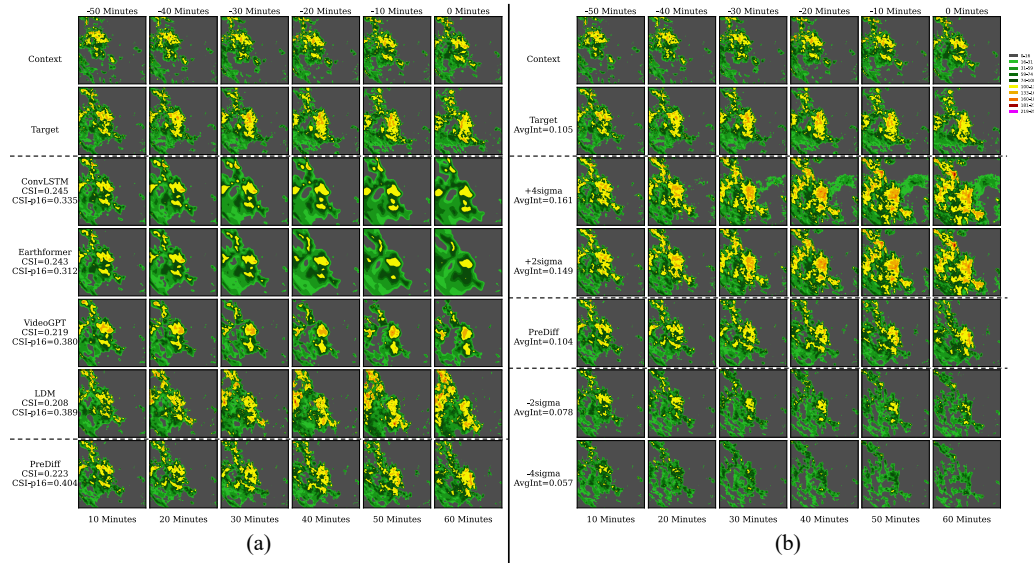


Figure 18: A set of example predictions on the SEVIR test set. (a) Comparison of PreDiff with baselines. (b) Predictions by PreDiff-KC under the guidance of anticipated average intensity.