SeasonCast: A Masked Latent Diffusion Model for Skillful Subseasonal-to-Seasonal Prediction

Anonymous Author(s)

Affiliation Address email

Abstract

Accurate weather prediction on the subseasonal-to-seasonal (S2S) scale is critical for anticipating and mitigating the impacts of climate change. However, existing data-driven methods struggle beyond the medium-range timescale due to error accumulation in their autoregressive approach. In this work, we propose SeasonCast, a scalable and skillful probabilistic model for S2S prediction. SeasonCast consists of two components, a VAE model that encodes raw weather data into a continuous, lower-dimensional latent space, and a diffusion-based transformer model that generates a sequence of future latent tokens given the initial conditioning tokens. During training, we mask random future tokens and train the transformer to estimate their distribution given conditioning and visible tokens using a per-token diffusion head. During inference, the transformer generates the full sequence of future tokens by iteratively unmasking random subsets of tokens. This joint sampling across space and time mitigates compounding errors from autoregressive approaches. The low-dimensional latent space enables modeling long sequences of future latent states, allowing the transformer to learn weather dynamics beyond initial conditions. SeasonCast performs competitively with leading probabilistic methods at the medium-range timescale while being $10\times$ to $20\times$ faster, and achieves state-of-the-art performance at the subseasonal-to-seasonal scale across accuracy, physics-based, and probabilistic metrics.

1 Introduction

2

4

5

6

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

29

30

31

32

33

Subseasonal-to-seasonal (S2S) weather prediction, which predicts atmospheric conditions on timescales ranging from two to six weeks, is critical for disaster preparedness, resource management, and long-term planning. This timescale bridges the gap between short-term weather forecasts and longer-term climate projections, enabling more informed decision-making for extreme weather events such as droughts, floods, and heatwaves [49, 34, 50, 7]. However, S2S prediction is particularly challenging due to the interplay between atmospheric initial conditions, essential for short-term and medium-range forecasting accuracy, and boundary conditions dominating seasonal and climate predictions [25, 26]. Traditional numerical weather prediction (NWP) models, built upon solving differential equations of fluid dynamics and thermodynamics, have been instrumental in advancing S2S weather prediction [34, 45, 46]. However, numerical methods incur substantial computational costs due to the complexity of integrating large systems of differential equations, particularly at fine spatial and temporal resolutions. This computational bottleneck also constrains the ensemble size of ensemble systems, which is crucial for achieving accurate S2S predictions.

To overcome the challenges of NWP systems, there has been a growing interest in data-driven approaches based on deep learning for weather forecasting [9, 41, 48]. These approaches involve training deep neural networks on historical datasets, such as ERA5 [13, 14, 37, 38], to learn the underlying weather patterns. Once trained, they can produce forecasts in seconds compared to

the hours required by NWP models. Recent deep learning methods such as PanguWeather [1], Graphcast [20], and Stormer [31] have also shown superior accuracy in medium-range weather 39 forecasting, surpassing operational IFS [47], the state-of-the-art NWP system. However, their 40 application to the S2S timescale has been limited [29]. One possible explanation for this limitation is 41 the rapid error compounding in their autoregressive designs, in which a model learns to forecast the 42 future weather state at a small interval and iteratively feeds its prediction back as input to achieve 43 longer-horizon forecasts. Even though previous works have proposed multi-step finetuning to mitigate this issue, back-propagation through a large number of forward passes required for S2S timescales is computationally prohibitive. Moreover, training a neural network to forecast at a small interval 46 only allows the model to learn the initial conditions problem, ignoring boundary conditions that are 47 critical for prediction at S2S timescales. 48

We propose SeasonCast, a novel latent diffusion model for skillful probabilistic S2S prediction. 49 SeasonCast follows a two-stage training process. First, a VAE model compresses raw weather data 50 into a continuous, lower-dimensional latent space. Second, a transformer is trained to model the distribution of future latent tokens using a masked generative framework [2, 53]. Specifically, during training, we randomly mask a subset of future tokens, and task the transformer to unmask these 53 tokens based on the conditioning tokens and the visible tokens. Since the latent tokens lie in a 54 continuous space, we use a small diffusion network on top of the transformer model to estimate 55 the per-token distribution of unmasked tokens. After training, SeasonCast generates forecasts for 56 the full sequence of future tokens through an iterative process. At inference, SeasonCast iteratively 57 generates forecasts for the full sequence of future tokens by unmasking a subset of tokens in each step until all are generated. This joint generation of future tokens across time and space significantly mitigates the compounding errors issue of an autoregressive approach. Furthermore, training on the 60 full sequence of future frames enables SeasonCast to address both initial condition problems and 61 boundary condition challenges, which are critical for S2S prediction. 62

We evaluate SeasonCast on ChaosBench [29], a recent benchmark for subseasonal-to-seasonal prediction. SeasonCast achieves state-of-the-art performance on key atmospheric variables across various accuracy, physics-based, and probabilistic metrics. Additionally, we carefully study the impact 65 of different design choices, including the auxiliary MSE loss, training sequence lengths, unmasking order, and diffusion sampling temperature, on the forecasting performance of SeasonCast.

Related Work 2

63

64

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

86

87

88

89

Data-driven weather forecasting Deep learning has become a promising approach in the field of weather forecasting. Recent advancements with powerful architectures have achieved significant successes, providing faster inference and superior forecasting accuracy compared to IFS, the goldstandard numerical weather prediction system. Notable methods include FourCastNet [33], which utilizes an adaptive neural operator architecture; Keisler [16]'s, GraphCast [20], and AIFS [22], which leverage graph neural networks; and a series of transformer-based models such as PanguWeather [1], Stormer [31], and others [30, 5, 3, 6]. Beyond deterministic predictions, the field has increasingly focused on probabilistic forecasting to account for forecast uncertainty. Common approaches involve integrating existing architectures with generative frameworks, including diffusion models [35, 28], normalizing flows [6], and latent variable models [32]. Others explore ensemble predictions through initial condition perturbations, exemplified by methods like AIFS-CRPS [22] and NeuralGCM [18].

Data-driven S2S prediction Recent benchmarks have emerged to evaluate data-driven methods at S2S timescales. While many focus on regional forecasts such as the US [15, 27], ChaosBench [29] offers a comprehensive framework for global S2S prediction, providing extensive numerical baselines and physics-based metrics. A key finding from ChaosBench shows that state-of-the-art deep learning methods struggle to extend to S2S timescales. These methods predominantly rely on autoregressive approaches that generate predictions iteratively at short time intervals, leading to error accumulation with increasing lead times. While multi-step finetuning helps mitigate this issue for medium-range forecasts, it becomes computationally prohibitive for S2S predictions due to the extensive number of required forward passes. Moreover, training models with short time intervals fails to capture boundary conditions essential for long-term weather patterns. While Fuxi-S2S [4] was proposed for S2S prediction, it focuses on forecasting daily averaged statistics, which fundamentally alters the underlying weather dynamics and makes it inapplicable to forecasting at instantaneous time steps.

3 Background and Preliminaries

3.1 Weather forecasting

93

105

121

122

123

125

126

127

128

129

The goal of weather forecasting is to forecast future weather conditions $X_T \in \mathbb{R}^{V \times H \times W}$ based on initial conditions $X_0 \in \mathbb{R}^{V \times H \times W}$, where T represents the target lead time, V denotes the number of 94 95 input and output physical variables (e.g., temperature and geopotential), and $H \times W$ corresponds 96 to the spatial resolution of the data, determined by the density of the global grid. In subseasonal-97 to-seasonal (S2S) forecasting, we focus on lead times ranging from 2 to 6 weeks. Autoregressive 98 modeling is a dominant paradigm in data-driven weather forecasting, where a model iteratively produces forecasts $X_{\delta t}$ at a short interval δt to reach the target lead time T. In this work, we propose 100 an alternative approach: training a generative model to estimate the distribution of the entire sequence 101 of future weather states $X_{1:T}$ given initial conditions X_0 . This approach mitigates error accumulation 102 and enables the model to learn both initial and boundary condition dynamics by considering the 103 complete sequence of weather states. 104

3.2 Masked generative modeling

Masked generative modeling is an efficient and powerful approach for image and video generation in computer vision [2, 53, 23]. In this framework, visual data $X_{1:T} \in \mathbb{R}^{T \times V \times H \times W}$ (T = 1 for images)106 107 is first embedded by a VAE encoder into a sequence of tokens $\mathbf{x} \in \mathbb{R}^{N \times D}$, where N represents 108 the length of the flattened token sequence. During training, we apply a binary mask to randomly 109 select a subset of tokens to be predicted, creating a corrupted sequence. We then train a transformer 110 model to recover the original tokens at masked positions based on both the visible tokens and any 111 112 additional conditioning information such as initial frames. For generation, the framework employs an iterative decoding process that starts with a fully masked sequence of future tokens. In each iteration, 113 the model predicts a random subset of masked tokens in parallel, where the number and positions 114 of the unmasked tokens follow a predefined schedule and order. This process continues until all 115 tokens are unmasked, at which point the generated tokens are decoded back to the original domain 116 through a VAE decoder. This framework offers key advantages for weather forecasting: it allows 117 the model to capture long-range dependencies across the entire sequence while avoiding the error accumulation typical in autoregressive approaches, and the iterative refinement process enables the model to maintain consistency across both spatial and temporal dimensions. 120

3.3 Modeling continuous tokens with diffusion models

In the masked generative modeling framework, a common practice is to embed the raw visual data into a discrete latent space using vector-quantized VAE models [43]. However, discretization is sensitive to gradient approximation strategies [39, 36, 19] and has lower reconstruction quality than continuous VAEs. Recent works [42, 24] have demonstrated that discretization can be eliminated by directly modeling the per-token continuous probability distribution by using diffusion models. Given data $x \in \mathbb{R}^D$ and its conditioning information $z \in \mathbb{R}^D$, we model the conditional distribution $p(x \mid z)$ using a diffusion process that gradually transforms a Gaussian prior into the target distribution. The forward diffusion process progressively adds Gaussian noise to the data x following:

$$x_s = \sqrt{\alpha_s}x + \sqrt{1 - \alpha_s}\epsilon,\tag{1}$$

where s indicates the diffusion step, α_s determines the noise schedule, and $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ represents Gaussian noise. The reverse process employs a denoising network $\epsilon_{\theta}(x_s, s, z)$ parameterized by θ to predict the noise component from the noisy input x_s and condition z:

$$\mathcal{L}_{\text{diff}}(\theta) = \mathbb{E}_{\epsilon,x} \left[\| \epsilon_{\theta}(x_s, s, z) - \epsilon \|^2 \right]. \tag{2}$$

At inference time, conditional sampling begins with a random Gaussian noise $x_S \sim \mathcal{N}(0, \mathbf{I})$ and iteratively applies the reverse diffusion process:

$$x_{s-1} = \frac{1}{\sqrt{\alpha_s}} \left(x_s - \frac{1 - \alpha_s}{\sqrt{1 - \bar{\alpha}_s}} \epsilon_{\theta}(x_s, s, z) \right) + \tau \sigma_s \delta, \tag{3}$$

where $\bar{\alpha}_s = \prod_{k=1}^s \alpha_k$, $\delta \sim \mathcal{N}(0, \mathbf{I})$ and σ_s controls the magnitude of noise added at each step. This iterative process generates samples from the learned conditional distribution $p_{\theta}(x \mid z)$. Following [24], we additionally scale the noise $\sigma_s \delta$ by the temperature τ that controls the sample diversity from the diffusion model.

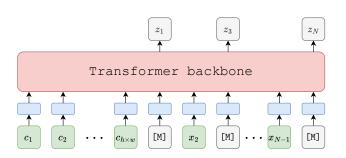


Figure 1: SeasonCast processes the latent tokens through a transformer backbone that outputs a vector z_i for each position i in the sequence.

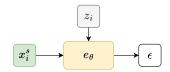


Figure 2: The denoising network e_{θ} predicts the noise ϵ from z_i and x_i^s .



Figure 3: The deterministic network predicts directly x_i from z_i .

139 4 Methodology

162

We present SeasonCast, a novel method for subseasonal-to-seasonal prediction. Similar to previous works in video generation, SeasonCast consists of two components: a VAE model that compresses the raw weather data into a lower-dimensional space, and a masked generative transformer model in this latent space. We present the two components and their key design choices in this section.

4.1 VAE for weather data embedding

A VAE encoder embeds a weather state $X \in \mathbb{R}^{V \times H \times W}$ into a map of $h \times w$ latent tokens, where 145 h < H and w < W. In vector-quantized VAEs, each entry in the latent map is an integer index from a 147 fixed-size vocabulary, representing a discrete latent space. While this discretization is widely adopted in computer vision due to its compatibility with cross-entropy training and straightforward sampling 148 from softmax distributions, it presents significant challenges for weather data. Unlike RGB images 149 with three channels, weather states can contain hundreds of physical variables, resulting in an extreme 150 compression requirement. For instance, consider compressing weather data with 100 variables (32 bits 151 per value) by a factor of 4 in each spatial dimension, using a vocabulary size of $2^{13} = 8192$ (13 bits 152 per latent token). This results in a compression ratio of $(32 \times 100 \times H \times W)/(13 \times (H/4) \times (W/4)) \approx$ 3938. Such aggressive compression leads to substantial reconstruction errors, ultimately degrading 154 the performance of the second-stage generative modeling. 155 Therefore, we adopt a continuous VAE model for SeasonCast, where each token in the $h \times w$ 156 latent map is a continuous vector of D dimensions. With D=16, for example, the compression 157

latent map is a continuous vector of D dimensions. With D=16, for example, the compression ratio becomes $(32\times 100\times H\times W)/(32\times 16\times (H/4)\times (W/4))=100$, substantially lower than the discrete approach. While it is also possible to compress a sequence of weather states $X_{1:T}\in\mathbb{R}^{T\times V\times H\times W}$ in both temporal and spatial dimensions, our preliminary experiments showed no clear benefits from temporal compression, leading us to adopt per-frame embedding.

4.2 Masked generative modeling for S2S prediction

After training the VAE, we embed the initial condition into a sequence of tokens $\mathbf{c}=(c_1,c_2,\ldots,c_{h\times w})$. Similarly, each future weather state is embedded into a sequence of tokens, which are concatenated to form the complete sequence of future tokens $\mathbf{x}=(x_1,x_2,\ldots,x_N)$, where $N=T\times h\times w$ represents the total number of future tokens. Each latent token is a continuous vector of dimension D. Our generative modeling objective is to estimate the conditional distribution $p(\mathbf{x}\mid\mathbf{c})$ from the training data.

We achieve this using a masked generative framework, as illustrated in Figure 1. During training, we sample a binary mask $\mathbf{m} = [m_i]_{i=1}^N \sim p_{\mathcal{U}}$ and replace tokens x_i with a learnable, continuous [MASK] token where $m_i = 1$, creating a corrupted sequence $\overline{\mathbf{x}} = \mathbf{m}(\mathbf{x})$. The generative objective is to estimate the distribution of masked tokens conditioned on the visible and conditioning tokens:

$$\mathcal{L}_{\text{gen}}(\theta) = \mathbb{E}_{\mathbf{m} \sim p_{\mathcal{U}}} \left[\sum_{i \text{ s.t. } m_i = 1} -\log p_{\theta}(x_i \mid \mathbf{c}, \overline{\mathbf{x}}) \right]. \tag{4}$$

The model processes the input by concatenating the conditioning tokens c with the corrupted future tokens $\overline{\mathbf{x}}$, adding positional encodings to the embedded sequence, and passing it through a bi-directional transformer backbone to obtain vectors z_i for each masked position. Given these vectors, the per-token objective $\log p_{\theta}(x_i \mid \mathbf{c}, \overline{\mathbf{x}})$ in Equation 4 simplifies to $\log p_{\theta}(x_i \mid z_i)$. To model this continuous distribution, we employ a diffusion model where z_i serves as conditional information for a denoising network – implemented as a small MLP on top of the transformer (Figure 2). We train the denoising network and the transformer backbone jointly using the diffusion loss specified in Equation 2. Conceptually, this diffusion objective encourages the model to produce representations z_i that facilitate effective denoising.

Auxiliary deterministic objective To encourage accurate predictions of near-term future tokens, we incorporate an auxiliary mean-squared error loss in the latent space. We implement this through a separate MLP head that produces deterministic predictions \hat{x}_i from z_i , training it jointly with the transformer backbone. Since weather dynamics become increasingly chaotic beyond day 10, making deterministic predictions progressively less meaningful, we apply this loss only to the first 10 future frames. Furthermore, we employ an exponentially decreasing weighting scheme to emphasize the importance of accurate predictions for earlier frames. The deterministic objective is thus:

$$\mathcal{L}_{\text{deter}}(\theta) = \underset{\mathbf{m} \sim p_{\mathcal{U}}}{\mathbb{E}} \left[\sum_{m_i = 1} w(i) ||x_i - \hat{x}_i||_2^2 \right]. \tag{5}$$

Appendix A.2 presents the details of this objective. The complete training objective combines both losses: $\mathcal{L}(\theta) = \mathcal{L}_{\text{gen}}(\theta) + \mathcal{L}_{\text{deter}}(\theta)$.

Sampling from SeasonCast At inference time, we generate samples from $p(\mathbf{x} \mid \mathbf{c})$ through an iterative decoding process, starting from a sequence of fully masked future tokens. Each iteration consists of three steps: first, the transformer backbone processes the conditioning tokens and corrupted future tokens to produce vectors z_i for each masked position; second, a subset of masked positions is randomly selected according to a predefined schedule for unmasking; third, for each selected position, the diffusion model generates token x_i by conditioning on z_i and performing a fixed number of diffusion steps. This process iterates until all future tokens are revealed, at which point the VAE decoder maps the generated tokens back to the weather domain. To generate an ensemble of forecasts, we simply replicate the initial tokens and perform independent sampling for each copy. Four hyperparameters affect the sampling procedure: the number of unmasking iterations, the unmasking order, the number of diffusion steps, and the diffusion temperature.

4.3 Implementation details

Architectural details For the transformer backbone, we adopt the encoder-decoder architecture from Masked Autoencoder (MAE) [12]. The model processes an input sequence in two stages: first, the encoder processes the conditioning and visible tokens; second, the encoded sequence is augmented with learnable [MASK] tokens at appropriate positions and passed through the decoder to produce z_i for each position i. Both the encoder and decoder are bidirectional, employing full attention. Before feeding to either the encoder or decoder, we add the input sequences with positional embeddings that combine two components: temporal embeddings to distinguish different frames, and spatial embeddings to differentiate tokens within each frame. The encoder and decoder follow the Transformer [44] implementation in ViT [8], each having 16 layers with 16 attention heads, a hidden dimension of 1024, and a dropout rate of 0.1.

Mask sampling During training, we sample a masking ratio $\gamma \sim \mathcal{U}[0.5, 1.0]$ and generate a corresponding binary mask \mathbf{m} , where $\gamma = 0.75$ indicates that 75% of entries in \mathbf{m} are 1. For inference, we start with full masking ($\gamma = 1.0$) and gradually reduce it to 0.0 with a cosine schedule [2]. We set the number of unmasking iterations to match the number of future weather states T. We employ random masking orders across both spatial and temporal dimensions for training and inference.

Diffusion loss details We use a linear noise schedule with 1000 steps at training time that are resampled to 100 steps at inference. The denoising network ϵ_{θ} is implemented as a small MLP following Li et al. [24]. Specifically, the network consists of six residual blocks, each comprising a LayerNorm (LN), a linear layer, a SiLU activation, and another linear layer, with a residual connection around the block. Each block maintains a width of 2048 channels. The network takes the vector z_i from the transformer as conditioning information, which is combined with the time embedding of the diffusion step s through adaptive layer normalization (AdaLN) in each block's LN layers.

25 5 Experiments

We compare SeasonCast with state-of-the-art deep learning and numerical methods on medium-range 226 weather forecasting and S2S prediction, using WeatherBench2 [38] (WB2) and ChaosBench [29] as 227 benchmarks, respectively. We also conduct extensive ablation studies to assess the contribution of 228 each component in SeasonCast, and evaluate its scalability under varying inference compute budgets. 229 Across both tasks, we train and evaluate SeasonCast on 69 variables from the ERA5 reanalysis dataset [14], including four surface-level variables – 2-meter temperature (T2m), 10-meter U and V wind components (U10, V10), and mean sea-level pressure (MSLP), as well as five atmospheric 232 233 variables – geopotential (Z), temperature (T), U and V wind components, and specific humidity (Q), each at 13 pressure levels {50, 100, 150, 200, 250, 300, 400, 500, 600, 700, 850, 925, 1000} hPa. 234 For medium-range forecasting, we use native 0.25° resolution (721×1440 grids) and follow WB2 235 to train on years 1979–2018, validate on 2019, and test on 2020 using initial conditions at 00UTC 236 and 12UTC. For S2S prediction, we downsample the data to 1.40625° (128×256 grids) and follow 237 ChaosBench to train on 1979–2020, validate on 2021, and test on 2022 using 00UTC initializations. 238

5.1 SeasonCast for S2S prediction

239

249

250

251

252

253

254

255 256

257

258

259

260 261

263

264

265

Training and inference details We train a VAE that embeds each weather state of shape $69 \times 128 \times 256$ into a latent map of shape $1024 \times 8 \times 16$, reducing spatial dimensions by a factor of 16. The architectural details and training process of the VAE are described in Appendix A.1. We train SeasonCast to forecast a sequence of T=44 future weather states at 24hr intervals, covering lead times from 1 to 44 days. Each training example consists of $45 \times 8 \times 16 = 5760$ latent tokens, including the initial condition. During inference, we generate the complete future sequence in 44 iterations (1 iteration per frame) using a diffusion temperature of $\tau=1.3$. We produce an ensemble of 50 forecast sequences for each initial condition.

Baselines We compare SeasonCast with PanguWeather (PW) [1] and GraphCast (GC) [20], two leading open-sourced deep learning methods, and ensemble systems of four numerical models from different national agencies: UKMO-ENS (UK) [51], NCEP-ENS (US) [40], CMA-ENS (China) [52], and ECMWF-ENS (Europe) [10]. We refer to ChaosBench for details about these baselines. Following ChaosBench, we report results on T850, Z500, and Q700 at lead times from 1 to 44 days. We additionally compare SeasonCast with ClimaX [30] and Stormer [31] in Appendix B.2. We do not compare against Fuxi-S2S [4] as Fuxi-S2S forecasts daily average values from past daily averages, making it incomparable with SeasonCast and the rest of the methods, which perform point-in-time weather forecasting based on an initial condition. We are also not able to run Gencast [35] and NeuralGCM [18] for S2S due to their significant computational demands.

Results Figure 4 compares different methods on three deterministic metrics: Root Mean-Squared Error (RMSE), Absolute Bias (ABS BIAS), and Multi-scale Structural Similarity (SSIM). At shorter lead times, SeasonCast shows slightly worse performance on RMSE and SSIM than other baselines, which is expected since we train SeasonCast to model a full sequence of future weather states rather than optimizing for short- and medium-range predictions. However, SeasonCast's relative performance improves with increasing lead time, ultimately matching ECMWF-ENS as one of the top two performing methods beyond day 10. Notably, SeasonCast demonstrates the lowest bias among all baselines, maintaining near-zero bias across all three target variables.

Physical consistency also plays a crucial role in S2S prediction, particularly for ensemble systems.

We evaluate this aspect using two physics-based metrics: Spectral Divergence (SDIV) and Spectral Residual (SRES), which measure how closely the power spectra of predictions match those of ground-truths. As shown in Figure 5, SeasonCast achieves substantially better physical consistency than other deep learning methods, and often outperforms all baselines on these metrics. These results demonstrate how SeasonCast effectively preserves signals across the frequency spectrum.

Finally, we compare SeasonCast with the four numerical ensemble systems on two probabilistic metrics: Continuous Ranked Probability Score (CRPS) and Spread/Skill Ratio (SSR) (closer to 1 is better). Figure 6 shows that SeasonCast and ECMWF-ENS are the two leading methods across variables and lead times. Similar to deterministic results, SeasonCast performs worse than ECMWF-ENS at shorter lead times but outperforms this baseline beyond day 15.

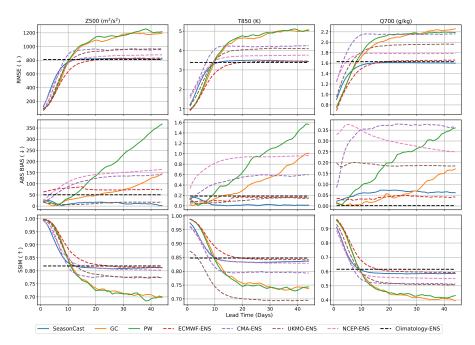


Figure 4: Deterministic performance of different methods at lead times from 1 to 44 days across three key variables. Solid curves are deep learning methods and dashed curves are numerical methods.

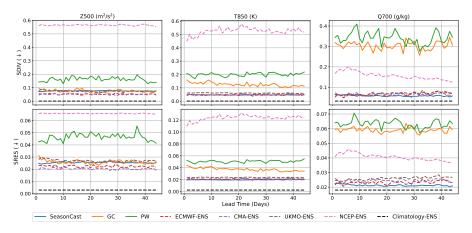


Figure 5: Physics-based metrics of different methods at lead times from 1 to 44 days across three key variables. Solid curves are deep learning methods and dashed curves are numerical methods.

5.2 SeasonCast for medium-range forecasting

In addition to its strong performance on the S2S task, we demonstrate that SeasonCast also performs competitively at the medium-range timescale. We train a VAE model with a spatial downsampling ratio of 16, compressing each weather state of shape $69 \times 721 \times 1440$ into a latent representation of size $256 \times 45 \times 90$. We then train SeasonCast to predict two steps ahead at 12-hour intervals, following the setup of Gencast [35]. During inference, we use autoregressive sampling, recursively feeding the most recent predicted frame as the new initial condition until the target lead time is reached. We generate forecasts using a single sampling iteration per frame with a diffusion temperature $\tau=1.0$, and produce an ensemble of 50 members.

We compare SeasonCast against Gencast [35], a leading deep learning method for probabilistic fore-casting, and IFS-ENS [21], the gold-standard numerical ensemble system. Following WeatherBench2, we use ensemble RMSE, CRPS, and spread-skill ratio (SSR) as evaluation metrics. Figure 7 shows that SeasonCast performs comparably with IFS-ENS across all variables and metrics, and is only

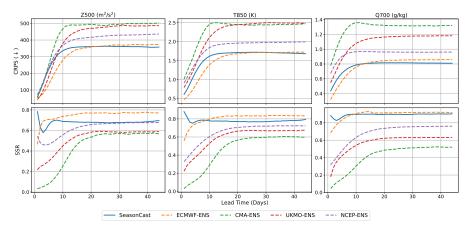


Figure 6: Probabilistic performance of different methods at lead times from 1 to 44 days across three key variables. Solid curves are deep learning methods and dashed curves are numerical methods.

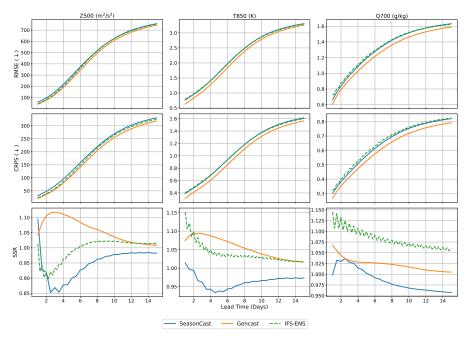


Figure 7: Probabilistic performance of different methods in medium-range forecasting. Solid curves are deep learning methods and dashed curves are numerical methods.

slightly behind Gencast. Moreover, our analysis in Appendix B.1 further shows that SeasonCast is $10 \times$ to $20 \times$ faster than all baselines. These results indicate strong performance across both medium-range and S2S timescales of SeasonCast.

6 Conclusion

We present SeasonCast, a novel latent diffusion model for S2S prediction. By combining the masked generative framework with a diffusion objective, our approach enables direct modeling of long sequences of future weather states while avoiding error accumulation inherent in autoregressive methods. SeasonCast achieves state-of-the-art performance in deterministic and probabilistic metrics while maintaining exceptional physical consistency. In medium-range forecasting, SeasonCast performs competitively with existing methods while being significantly more efficient. Future work could study the fundamental trade-off between VAE reconstruction quality and transformer modeling capacity, and explore more sophisticated generative frameworks to enhance the diffusion objective.

2 References

- [1] Kaifeng Bi, Lingxi Xie, Hengheng Zhang, Xin Chen, Xiaotao Gu, and Qi Tian. Accurate medium-range global weather forecasting with 3D neural networks. *Nature*, 619(7970):533–538, 2023.
- Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked
 generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision* and Pattern Recognition, pages 11315–11325, 2022.
- [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] Lei Chen, Xiaohui Zhong, Jie Wu, Deliang Chen, Shangping Xie, Qingchen Chao, Chensen Lin, Zixin Hu, Bo Lu, Hao Li, et al. Fuxi-s2s: An accurate machine learning model for global subseasonal forecasts. *arXiv preprint arXiv:2312.09926*, 2023.
- [5] Lei Chen, Xiaohui Zhong, Feng Zhang, Yuan Cheng, Yinghui Xu, Yuan Qi, and Hao Li. FuXi:
 A cascade machine learning forecasting system for 15-day global weather forecast. *arXiv* preprint arXiv:2306.12873, 2023.
- [6] Guillaume Couairon, Christian Lessig, Anastase Charantonis, and Claire Monteleoni. Archesweather: An efficient ai weather forecasting model at 1.5 {\deg} resolution. arXiv preprint arXiv:2405.14527, 2024.
- [7] Daniela IV Domeisen, Christopher J White, Hilla Afargan-Gerstman, Ángel G Muñoz, Matthew A Janiga, Frédéric Vitart, C Ole Wulff, Salomé Antoine, Constantin Ardilouze, Lauriane Batté, et al. Advances in the subseasonal prediction of extreme events: Relevant case studies across the globe. *Bulletin of the American Meteorological Society*, 103(6):E1473–E1501, 2022.
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image
 recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
- [9] P. D. Dueben and P. Bauer. Challenges and design choices for global weather and climate models based on machine learning. *Geoscientific Model Development*, 11(10):3999–4009, 2018. doi: 10.5194/gmd-11-3999-2018. URL https://gmd.copernicus.org/articles/11/3999/2018/.
- [10] ECMWF. IFS Documentation CY41R1 Part V: The Ensemble Prediction System. Number 5. ECMWF, 2015 2015. doi: 10.21957/eow1lonc. URL https://www.ecmwf.int/node/9212. operational implementation 12 May 2015.
- Jayesh K Gupta and Johannes Brandstetter. Towards multi-spatiotemporal-scale generalized pde modeling. *arXiv preprint arXiv:2209.15616*, 2022.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022.
- [13] Hans Hersbach, Bill Bell, Paul Berrisford, Gionata Biavati, András Horányi, Joaquín Muñoz Sabater, Julien Nicolas, Carole Peubey, Raluca Radu, Iryna Rozum, Dinand Schepers, Adrian Simmons, Cornel Soci, Dick Dee, and Jean-Noël Thépaut. ERA5 hourly data on single levels from 1979 to present. Copernicus Climate Change Service (C3S) Climate Data Dtore (CDS), 10(10.24381), 2018.
- [14] Hans Hersbach, Bill Bell, Paul Berrisford, Shoji Hirahara, András Horányi, Joaquín Muñoz Sabater, Julien Nicolas, Carole Peubey, Raluca Radu, Dinand Schepers, , Adrian Simmons,
 Cornel Soci, Saleh Abdalla, Xavier Abellan, Gianpaolo Balsamo, Peter Bechtold, Gionata
 Biavati, Jean Bidlot, Massimo Bonavita, Giovanna De Chiara, Per Dahlgren, Dick Dee, Michail

- Diamantakis, Rossana Dragani, Johannes Flemming, Richard Forbes, Manuel Fuentes, Alan Geer, Leo Haimberger, Sean Healy, Robin J. Hogan, Elías Hólm, Marta Janisková, Sarah Keeley, Patrick Laloyaux, Philippe Lopez, Cristina Lupu, Gabor Radnoti, Patricia de Rosnay, Iryna Rozum, Freja Vamborg, Sebastien Villaume, and Jean-Noël Thépaut. The ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730):1999–2049, 2020.
- In Jessica Hwang, Paulo Orenstein, Judah Cohen, Karl Pfeiffer, and Lester Mackey. Improving subseasonal forecasting in the western us with machine learning. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2325–2335, 2019.
- 360 [16] Ryan Keisler. Forecasting global weather with graph neural networks. *arXiv preprint* 361 *arXiv*:2202.07575, 2022.
- [17] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint
 arXiv:1412.6980, 2014.
- [18] Dmitrii Kochkov, Janni Yuval, Ian Langmore, Peter Norgaard, Jamie Smith, Griffin Mooers,
 Milan Klöwer, James Lottes, Stephan Rasp, Peter Düben, et al. Neural general circulation
 models for weather and climate. *Nature*, 632(8027):1060–1066, 2024.
- [19] Alexander Kolesnikov, André Susano Pinto, Lucas Beyer, Xiaohua Zhai, Jeremiah Harmsen,
 and Neil Houlsby. Uvim: A unified modeling approach for vision with learned guiding codes.
 Advances in Neural Information Processing Systems, 35:26295–26308, 2022.
- [20] Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirnsberger, Meire Fortunato,
 Ferran Alet, Suman Ravuri, Timo Ewalds, Zach Eaton-Rosen, Weihua Hu, Alexander Merose,
 Stephan Hoyer, George Holland, Oriol Vinyals, Jacklynn Stott, Alexander Pritzel, Shakir
 Mohamed, and Peter Battaglia. Learning skillful medium-range global weather forecasting.
 Science, 0(0):eadi2336, 2023. doi: 10.1126/science.adi2336. URL https://www.science.org/doi/abs/10.1126/science.adi2336.
- Simon Lang, Mark Rodwell, and Dinand Schepers. Ifs upgrade brings many improvements and unifies medium-range resolutions. *ECMWF Newsletter*, 176:21–28, 2023.
- Simon Lang, Mihai Alexe, Matthew Chantry, Jesper Dramsch, Florian Pinault, Baudouin Raoult, Mariana CA Clare, Christian Lessig, Michael Maier-Gerber, Linus Magnusson, et al. Aifs-ecmwf's data-driven forecasting system. *arXiv preprint arXiv:2406.01465*, 2024.
- [23] Tianhong Li, Huiwen Chang, Shlok Mishra, Han Zhang, Dina Katabi, and Dilip Krishnan.
 Mage: Masked generative encoder to unify representation learning and image synthesis. In
 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages
 2142–2152, 2023.
- Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. *arXiv preprint arXiv:2406.11838*, 2024.
- Edward N Lorenz. Forced and free variations of weather and climate. *Journal of Atmospheric Sciences*, 36(8):1367–1376, 1979.
- [26] Annarita Mariotti, Paolo M Ruti, and Michel Rixen. Progress in subseasonal to seasonal prediction through a joint weather and climate community effort. *Npj Climate and Atmospheric Science*, 1(1):4, 2018.
- Soukayna Mouatadid, Paulo Orenstein, Genevieve Flaspohler, Miruna Oprescu, Judah Cohen,
 Franklyn Wang, Sean Knight, Maria Geogdzhayeva, Sam Levang, Ernest Fraenkel, et al.
 Subseasonalclimateusa: a dataset for subseasonal forecasting and benchmarking. Advances in
 Neural Information Processing Systems, 36, 2024.
- ³⁹⁶ [28] Congyi Nai, Xi Chen, Shangshang Yang, Yuan Liang, Ziniu Xiao, and Baoxiang Pan. Boosting weather forecast via generative superensemble. *arXiv preprint arXiv:2412.08377*, 2024.

- [29] Juan Nathaniel, Yongquan Qu, Tung Nguyen, Sungduk Yu, Julius Busecke, Aditya Grover, and
 Pierre Gentine. Chaosbench: A multi-channel, physics-based benchmark for subseasonal-to-seasonal climate prediction. arXiv preprint arXiv:2402.00712, 2024.
- 401 [30] Tung Nguyen, Johannes Brandstetter, Ashish Kapoor, Jayesh K Gupta, and Aditya Grover. ClimaX: A foundation model for weather and climate. *arXiv preprint arXiv:2301.10343*, 2023.
- Tung Nguyen, Rohan Shah, Hritik Bansal, Troy Arcomano, Romit Maulik, Veerabhadra Kotamarthi, Ian Foster, Sandeep Madireddy, and Aditya Grover. Scaling transformer neural networks for skillful and reliable medium-range weather forecasting. arXiv preprint arXiv:2312.03876, 2023.
- Joel Oskarsson, Tomas Landelius, Marc Peter Deisenroth, and Fredrik Lindsten. Probabilistic
 weather forecasting with hierarchical graph neural networks. arXiv preprint arXiv:2406.04759,
 2024.
- [33] Jaideep Pathak, Shashank Subramanian, Peter Harrington, Sanjeev Raja, Ashesh Chattopadhyay,
 Morteza Mardani, Thorsten Kurth, David Hall, Zongyi Li, Kamyar Azizzadenesheli, Pedram
 Hassanzadeh, Karthik Kashinath, and Animashree Anandkumar. FourCastNet: A global data driven high-resolution weather model using adaptive Fourier neural operators. arXiv preprint
 arXiv:2202.11214, 2022.
- Kathy Pegion, Ben P Kirtman, Emily Becker, Dan C Collins, Emerson LaJoie, Robert Burgman, Ray Bell, Timothy DelSole, Dughong Min, Yuejian Zhu, et al. The subseasonal experiment (subx): A multimodel subseasonal prediction experiment. *Bulletin of the American Meteorological Society*, 100(10):2043–2060, 2019.
- [35] Ilan Price, Alvaro Sanchez-Gonzalez, Ferran Alet, Timo Ewalds, Andrew El-Kadi, Jacklynn
 Stott, Shakir Mohamed, Peter Battaglia, Remi Lam, and Matthew Willson. Gencast: Diffusion-based ensemble forecasting for medium-range weather. arXiv preprint arXiv:2312.15796,
 2023.
- 423 [36] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark
 424 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on*425 *machine learning*, pages 8821–8831. Pmlr, 2021.
- [37] Stephan Rasp, Peter D Dueben, Sebastian Scher, Jonathan A Weyn, Soukayna Mouatadid, and
 Nils Thuerey. WeatherBench: a benchmark data set for data-driven weather forecasting. *Journal* of Advances in Modeling Earth Systems, 12(11):e2020MS002203, 2020.
- 429 [38] Stephan Rasp, Stephan Hoyer, Alexander Merose, Ian Langmore, Peter Battaglia, Tyler Russel,
 430 Alvaro Sanchez-Gonzalez, Vivian Yang, Rob Carver, Shreya Agrawal, Matthew Chantry,
 431 Zied Ben Bouallegue, Peter Dueben, Carla Bromberg, Jared Sisk, Luke Barrington, Aaron Bell,
 432 and Fei Sha. WeatherBench 2: A benchmark for the next generation of data-driven global
 433 weather models. arXiv preprint arXiv:2308.15560, 2023.
- 434 [39] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32, 2019.
- [40] Suranjana Saha, Shrinivas Moorthi, Xingren Wu, Jiande Wang, Sudhir Nadiga, Patrick Tripp,
 David Behringer, Yu-Tai Hou, Hui-ya Chuang, Mark Iredell, et al. The ncep climate forecast
 system version 2. *Journal of climate*, 27(6):2185–2208, 2014.
- 439 [41] Sebastian Scher. Toward data-driven weather and climate forecasting: Approximating a simple
 440 general circulation model with deep learning. *Geophysical Research Letters*, 45(22):12–616,
 441 2018.
- 442 [42] Michael Tschannen, Cian Eastwood, and Fabian Mentzer. Givt: Generative infinite-vocabulary transformers. In *European Conference on Computer Vision*, pages 292–309. Springer, 2025.
- 444 [43] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.

- [44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in Neural Information
 Processing Systems, 30, 2017.
- 449 [45] Frédéric Vitart. Evolution of ecmwf sub-seasonal forecast skill scores. *Quarterly Journal of the Royal Meteorological Society*, 140(683):1889–1899, 2014.
- [46] Frederic Vitart, Constantin Ardilouze, Axel Bonet, Anca Brookshaw, M Chen, C Codorean,
 M Déqué, L Ferranti, E Fucile, M Fuentes, et al. The subseasonal to seasonal (s2s) prediction
 project database. *Bulletin of the American Meteorological Society*, 98(1):163–173, 2017.
- [47] NP Wedi, P Bauer, W Denoninck, M Diamantakis, M Hamrud, C Kuhnlein, S Malardel,
 K Mogensen, G Mozdzynski, and PK Smolarkiewicz. The modelling infrastructure of the
 Integrated Forecasting System: Recent advances and future challenges. European Centre for
 Medium-Range Weather Forecasts, 2015.
- [48] Jonathan A Weyn, Dale R Durran, and Rich Caruana. Can machines learn to predict weather?
 Using deep learning to predict gridded 500-hPa geopotential height from historical weather
 data. *Journal of Advances in Modeling Earth Systems*, 11(8):2680–2693, 2019.
- [49] Christopher J White, Henrik Carlsen, Andrew W Robertson, Richard JT Klein, Jeffrey K Lazo,
 Arun Kumar, Frederic Vitart, Erin Coughlan de Perez, Andrea J Ray, Virginia Murray, et al.
 Potential applications of subseasonal-to-seasonal (s2s) predictions. *Meteorological applications*,
 24(3):315–325, 2017.
- [50] Christopher J White, Daniela IV Domeisen, Nachiketa Acharya, Elijah A Adefisan, Michael L
 Anderson, Stella Aura, Ahmed A Balogun, Douglas Bertram, Sonia Bluhm, David J Brayshaw,
 et al. Advances in the application and utility of subseasonal-to-seasonal predictions. *Bulletin of the American Meteorological Society*, 103(6):E1448–E1472, 2022.
- KD Williams, CM Harris, A Bodas-Salcedo, J Camp, RE Comer, D Copsey, D Fereday,
 T Graham, R Hill, T Hinton, et al. The met office global coupled model 2.0 (gc2) configuration.
 Geoscientific Model Development, 88(55):1509–1524, 2015.
- Tongwen Wu, Yixiong Lu, Yongjie Fang, Xiaoge Xin, Laurent Li, Weiping Li, Weihua Jie, Jie
 Zhang, Yiming Liu, Li Zhang, et al. The beijing climate center climate system model (bcc-csm):
 The main progress from cmip5 to cmip6. Geoscientific Model Development, 12(4):1573–1600,
 2019.
- Lijun Yu, Yong Cheng, Kihyuk Sohn, José Lezama, Han Zhang, Huiwen Chang, Alexander G
 Hauptmann, Ming-Hsuan Yang, Yuan Hao, Irfan Essa, et al. Magvit: Masked generative video transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10459–10469, 2023.

480 A Implementation details

481 A.1 VAE details

Our VAE model follows the UNet implementation from PDEArena [11]. We use the following hyperparameters for UNet in our experiments.

Table 1: Default hyperparameters of UNet

Lyparnaramatar	Meaning	Value
Hyperparameter	Meaning	value
Padding size	Padding size of each convolution layer	1
Kernel size	Kernel size of each convolution layer	3
Stride	Stride of each convolution layer	1
Input channels	The number of channels of the input	69
Input channels	The number of channels of the output	69
Base channels	The base hidden dimension of the UNet	256
Channel multiplications	Determine the number of output channels for Down and Up blocks	[1, 2, 4, 4, 8]
Dimension of z	The dimension of the latent space	1024
Blocks	Number of blocks	2
Use attention	If use attention in Down and Up blocks	False
Dropout	Dropout rate	0.0

483

485

486

487

488

503

The VAE encoder embeds each weather state of shape $69 \times 128 \times 256$ to a latent map of shape $1024 \times 8 \times 16$, reducing the spatial dimensions by 16. We use a KL weight of 5e-5 and optimize the VAE model with Adam [17] for 200 epochs with a batch size of 32, a base learning rate of 2e-4, parameters ($\beta_1 = 0.9, \beta_2 = 0.95$), and weight decay of 1e-5. The learning rate follows a linear warmup for the first 20 epochs, followed by a cosine decay schedule for the remaining 180 epochs.

489 A.2 Weighted deterministic objective

In SeasonCast, we employ a weighted MSE objective to encourage accurate deterministic predictions for near-term frames. The objective is formulated as:

$$\mathcal{L}_{\text{deter}}(\theta) = \underset{\mathbf{m} \sim p_{\mathcal{U}}}{\mathbb{E}} \left[\sum_{m_i = 1} w(i) ||x_i - \hat{x}_i||_2^2 \right], \tag{6}$$

where w(i) is an exponentially decreasing weighting function. We compute this weight in three steps. First, for each token i, we determine its corresponding frame index $k = \lfloor \frac{i}{h \times w} \rfloor$, where $h \times w$ represents the spatial dimensions of each frame's latent map. Second, we assign weights to tokens based on their frame index: $w(i) = e^{-k} = e^{-\lfloor \frac{i}{h \times w} \rfloor}$, ensuring all tokens from the same frame receive equal weight. Third, we set w(i) = 0 for tokens beyond frame 10 and normalize the remaining weights to sum to one.

498 A.3 Optimization details

We optimize SeasonCast with AdamW [17] for 100 epochs with a batch size of 32, a base learning rate of 2e-4, parameters ($\beta_1=0.9,\beta_2=0.95$), and weight decay of 1e-5. The learning rate follows a linear warmup for the first 10 epochs, followed by a cosine decay schedule for the remaining 90 epochs.

B Additional experiments

504 B.1 Efficiency of SeasonCast

Beyond its empirical performance, SeasonCast offers substantial efficiency gains over existing
 methods. We train SeasonCast for 4 days using 32 NVIDIA A100 GPUs. In comparison, Geneast

requires 5 days of training on 32 TPUv5e devices – hardware significantly more powerful than A100s, and NeuralGCM [18] requires 10 days on 128 TPUv5e devices. Additionally, Gencast employs a two-stage training pipeline, first pretraining on 1.0° resolution and then finetuning on 0.25°, while SeasonCast is trained in a single stage.

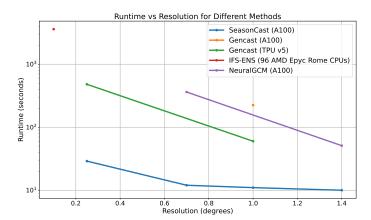


Figure 8: Runtime vs resolution of different methods to produce one forecast at 15-day lead time.

At inference time, SeasonCast is orders of magnitude faster than Gencast, NeuralGCM, and IFS-ENS. Figure 8 compares the runtime (in seconds) required to generate a 15-day forecast across different resolutions. At 0.25° resolution, Gencast requires 480 seconds on TPUv5, whereas SeasonCast achieves the same forecast in just 29 seconds on an A100. At 1.0°, SeasonCast completes inference in only 11 seconds, compared to 224 seconds for Gencast on the same hardware. These results highlight the scalability and practicality of SeasonCast for operational forecasting.

The efficiency of SeasonCast stems from two key architectural innovations. First, SeasonCast operates in a much lower-dimensional latent space $(45 \times 90 \text{ latent grid vs } 721 \times 1440 \text{ original grid})$, significantly reducing the computational cost of training and inference. Second, SeasonCast employs a highly efficient sampling mechanism. Unlike Gencast, which performs 50 full forward passes through the entire network for 50 diffusion steps, SeasonCast requires only a single forward passes through the transformer backbone. The subsequent diffusion steps involve only lightweight forward passes through a compact MLP diffusion head, resulting in orders-of-magnitude lower inference time. Together, these design choices enable SeasonCast to deliver fast and scalable forecasts.

B.2 Comparison with more deep learning baselines

In addition to PanguWeather and GraphCast, we compare SeasonCast with two advanced transformer-based methods: ClimaX [30] and Stormer [31]. Figure 9 shows that Stormer achieves superior accuracy in short-to-medium timescales, consistent with its reported results. However, as an autoregressive method, its performance degrades more rapidly than SeasonCast, eventually falling below Climatology, albeit at a slower rate than PanguWeather and GraphCast. ClimaX takes a different approach as a direct forecasting method, where a model trained on large-scale climate data is finetuned specifically for individual lead times. This approach avoids error accumulation and achieves comparable performance with SeasonCast at S2S scales. However, ClimaX requires fine-tuning separate models for each target lead time, while a single SeasonCast model can simultaneously generate the complete sequence of future weather states.

B.3 Impact of IC perturbations

Initial condition (IC) perturbations—adding random noise to initial conditions X_0 – are a standard technique in numerical methods for generating ensemble forecasts. This approach complements our generative framework. Figure 10 evaluates SeasonCast's performance across different noise levels, varying the standard deviation of the Gaussian distribution used for generating perturbations. The results demonstrate SeasonCast's robustness to input noise, maintaining consistent RMSE and CRPS scores across noise levels from 0.0 to 0.2, with only minor variations in SSR scores at short lead times.

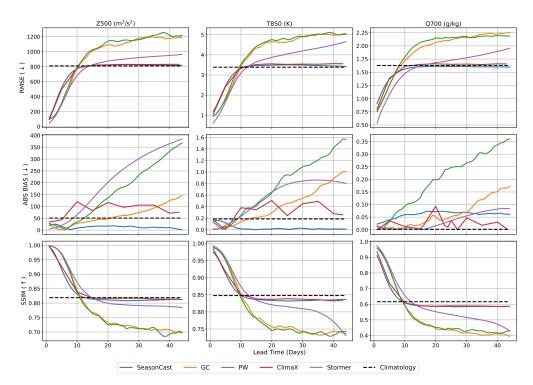


Figure 9: Comparison of deterministic performance of SeasonCast with more deep learning methods.

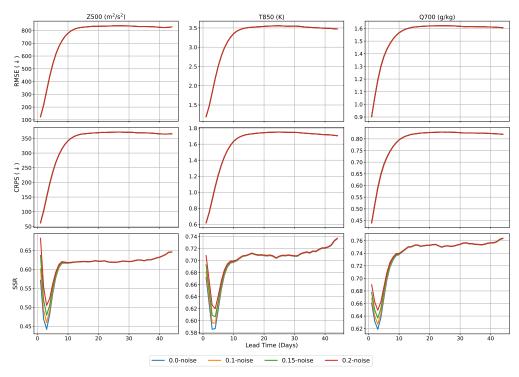


Figure 10: Performance of SeasonCast with different levels of IC noise.

44 B.4 Ablation studies

We analyze four key factors that influence SeasonCast's performance: the auxiliary deterministic objective, training sequence length T, unmasking order during sampling, and diffusion sampling

temperature τ . We present results for T850 on RMSE, CRPS, and SSR. We additionally study the impact of IC perturbations in Appendix B.3.

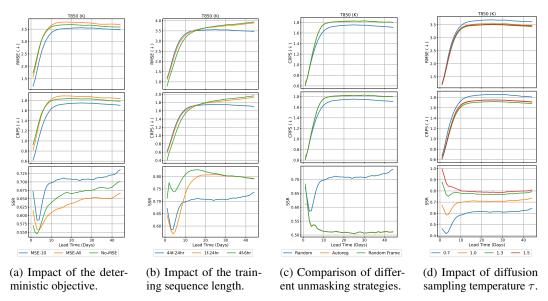


Figure 11: Ablation studies showing the impact of different components in SeasonCast.

Impact of the deterministic objective Figure 11a demonstrates the important role of the deterministic loss in SeasonCast's performance. Removing the MSE objective (No-MSE) degrades both RMSE and CRPS scores, with particularly noticeable impact at short lead times. However, naively applying MSE to all future frames (MSE-All-Frames) also proves counterproductive, as it forces deterministic predictions even for S2S timescales where weather systems become inherently chaotic. Our approach of applying MSE only to the first 10 frames achieves the best RMSE and CRPS scores across medium-range and S2S timescales.

Impact of training sequence length In our main experiments, we train SeasonCast to generate 44 future weather states at 24 hour intervals. One could alternatively train the model on shorter sequences and/or smaller intervals, then apply multiple roll-outs during inference to reach longer horizons, similar to autoregressive approaches. Figure 11b shows that models trained on shorter sequences or smaller intervals excel at short- and medium-range forecasting but underperform at S2S timescales. This trade-off emerges because shorter sequences allow models to specialize in near-term predictions, leading to better performance at shorter lead times. However, these models suffer from error accumulation at longer horizons, ultimately performing worse than the model trained on full sequences.

Impact of unmasking orders While our approach randomly masks tokens across both space and time during training, one may try more structured masking strategies at inference. We evaluate two such alternatives: an autoregressive strategy that unmasks entire frames sequentially, and a random framewise approach that unmasks complete frames in random order. Figure 11c shows that our fully randomized strategy achieves the best SSR scores, while both alternatives produce under-dispersive ensemble predictions. The superior performance of the fully randomized approach stems from its introduction of additional randomness through the fully random unmasking order, generating more diverse ensemble forecasts. This greater diversity consequently leads to better performance across other metrics.

Impact of diffusion sampling temperature The temperature τ controls the generation diversity, with higher values producing more diverse forecasts. Figure 11d demonstrates this relationship empirically. Setting $\tau < 1$ produces under-dispersive ensembles, degrading performance across other metrics. Increasing τ boosts sample diversity, improving SSR scores and overall better performance. However, pushing τ too high (e.g., $\tau=1.5$) causes samples to deviate from the mean prediction, compromising RMSE and CRPS performance. We identify $\tau=1.3$ as the optimal value, providing the best balance between ensemble diversity and forecast quality, which we adopt for our main experiments.

B.5 Scaling inference compute

Finally, we examine how increasing inference compute affects SeasonCast's performance through two hyperparameters: the number of ensemble forecasts and the average number of unmasking iterations per frame, i.e., 1-iter means a total of 44 iterations for 44 frames. Figure 12 shows that generating more ensemble forecasts improves both system diversity (higher SSR) and mean prediction accuracy (lower RMSE). Interestingly, while increasing the number of unmasking iterations shows minimal impact on RMSE, it yields slight improvements in SSR. This improvement likely stems from the increased randomness in unmasking order with more iterations, leading to greater ensemble diversity.

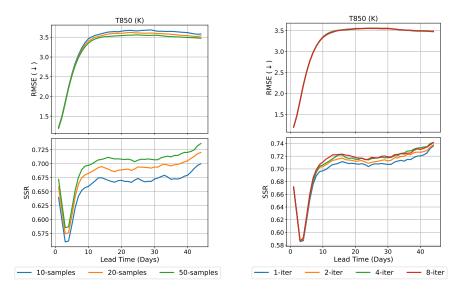


Figure 12: Performance of SeasonCast as we vary the number of ensemble forecasts (left) and the number of unmasking iterations.