

SUPPLEMENTARY: THE LOGARITHM TRICK: ACHIEVE BETTER LONG TERM FORECAST VIA MEAN LOGARITHM SQUARE LOSS

Anonymous authors

Paper under double-blind review

1 PHENOMENOLOGICAL THEORY

A typical Autoregressive training task (ATT) consists of a model $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and a sequence of unanimous data

$$\vec{X}_t, \vec{X}_{t+1}, \dots, \vec{X}_{t+N},$$

where each $\vec{X}_t = [\mathbf{X}_{t,1}, \mathbf{X}_{t,2}, \dots, \mathbf{X}_{t,s}]$ represents a time snapshot of the state at step t with s degrees of freedom. For example, a full resolution WeatherBenchsnapshot \vec{X}_t is a large 4D tensor with dimensions $s = 5 \times 37 \times 720 \times 1440 \approx 2 \times 10^8$; A typical time series prediction task such as ETTh1 consist of series of 1×7 vector. When the input is a sequence of 96 time stamp, the freedom is viewed as $s = 96 \times 7 = 672$.

An ideal model f is expected to reproduce any \vec{X}_{t+n+1} by $f(\vec{X}_{t+n})$

$$\vec{X}_t \xrightarrow{f} \vec{X}_{t+1} \xrightarrow{f} \vec{X}_{t+2} \xrightarrow{f} \vec{X}_{t+3} \xrightarrow{f} \dots \xrightarrow{f} \vec{X}_{t+N}$$

However, in most tasks, it is impossible to find the ideal model. Therefore, we need to use approximation functions such as neural networks to simulate this mapping, and the ideal model degenerates to a parameterized function $f \rightarrow f_\theta$. The task becomes an optimization problem of order N :

$$\min \sum_t |\vec{X}_{t+N} - \underbrace{f \dots f}_{N}(\vec{X}_t)| = \min \sum_t \mathcal{E}_t^N = \min \sum_t \mathcal{E}_{t,s}^N$$

Researchers usually back-propagate on the first low-order errors. For example, in the FourCastNetPathak et al. (2022) forwards twice on the "fine-tune" phase and Pangu only deal with order-1 errors. Several factors are concerned: Firstly, achieving the ideal model on the order-1 loss $\mathcal{E}_{\sqrt{t}}^1 \rightarrow 0$ directly implies that $\mathcal{E}_{\sqrt{t}}^N \rightarrow 0$. Secondly, training on \mathcal{E}_t^N requires N times forward-prediction and back-propagation, making it quite unaffordable in large system simulations compared to single loss optimization. Thirdly, the marginal benefit of increasing N decays rapidly in experience, and there must be a trade-off between accuracy and speed.

Thus, a comprehensive ATT optimization task of order N is formalized as

$$\min \mathcal{E} = \sum_t (\mathcal{E}_t^1 + \mathcal{E}_t^2 + \dots + \mathcal{E}_t^N) = \sum_N \sum_t \mathcal{E}_t^N = \sum_N \sum_t \sum_s \mathcal{E}_{t,s}^N$$

We expect a model trained on order N to be capable of forecasting a longer future beyond N , as the physical world is believed to have localized correlations, and a machine learning model trained on neighboring relationships should be able to extrapolate to a further distance. This pipeline has been validated in many tasksBi et al. (2022); Lam et al. (2022); Chen et al. (2023), which have significantly outperformed traditional methods.

For any step t , its time snapshot \mathbf{X}_t and a differentiable and smooth parameter neural function f_θ , we will not only obtain a snapshot chain

$$\begin{array}{ccccccc} \mathbf{X}_t^O & \xrightarrow{f_\theta} & \mathbf{X}_{t+1}^I & \xrightarrow{f_\theta} & \mathbf{X}_{t+2}^{II} & \xrightarrow{f_\theta} & \mathbf{X}_{t+3}^{III} & \xrightarrow{f_\theta} & \dots & \xrightarrow{f_\theta} & \mathbf{X}_{t+N}^N \\ \Downarrow & & \Downarrow & & \Downarrow & & \Downarrow & & \dots & & \Downarrow \\ \mathbf{X}_t^O & & \mathbf{X}_{t+1}^O & & \mathbf{X}_{t+2}^O & & \mathbf{X}_{t+3}^O & & \dots & & \mathbf{X}_{t+N}^O \end{array}$$

but also a snapshot tree. lets take $N = 5$ as an example:

$$\left(\begin{array}{cccccc} \mathbf{X}_t^O & \mathbf{X}_{t+1}^I & \mathbf{X}_{t+2}^{II} & \mathbf{X}_{t+3}^{III} & \mathbf{X}_{t+4}^{IV} & \mathbf{X}_{t+5}^V \\ & \mathbf{X}_{t+1}^O & \mathbf{X}_{t+2}^I & \mathbf{X}_{t+3}^{II} & \mathbf{X}_{t+4}^{III} & \mathbf{X}_{t+5}^{IV} \\ & & \mathbf{X}_{t+2}^O & \mathbf{X}_{t+3}^I & \mathbf{X}_{t+4}^{II} & \mathbf{X}_{t+5}^{III} \\ & & & \mathbf{X}_{t+3}^O & \mathbf{X}_{t+4}^I & \mathbf{X}_{t+5}^{II} \\ & & & & \mathbf{X}_{t+4}^O & \mathbf{X}_{t+5}^I \\ & & & & & \mathbf{X}_{t+5}^O \end{array} \right) \quad (1)$$

In our notation, the superscript represents the order order. For instance, O refers to the real data, while I represents the data generated by performing one forward pass from a ground truth through the function f . The subscript indicates the time stamp of the prediction. It is important to note that only tensors with the same time stamp can be compared with each other. As a result, we can categorize all errors according to their order by follow identity:

$$\begin{aligned} \mathbf{X}_{t+1}^O &= \mathbf{X}_{t+1}^I + \boldsymbol{\varepsilon}_{t+1}^I \\ \mathbf{X}_{t+2}^O &= \mathbf{X}_{t+2}^I + \boldsymbol{\varepsilon}_{t+2}^I = \mathbf{X}_{t+2}^{II} + \boldsymbol{\varepsilon}_{t+2}^{II} \\ \mathbf{X}_{t+3}^O &= \mathbf{X}_{t+3}^I + \boldsymbol{\varepsilon}_{t+3}^I = \mathbf{X}_{t+3}^{II} + \boldsymbol{\varepsilon}_{t+3}^{II} = \mathbf{X}_{t+3}^{III} + \boldsymbol{\varepsilon}_{t+3}^{III} \\ &\dots \end{aligned}$$

when the order order meets the same time stamp, the error is just the original error we want to minimized. $\mathcal{E}_t^1 = \|\boldsymbol{\varepsilon}_{t+1}^I\|, \mathcal{E}_t^2 = \|\boldsymbol{\varepsilon}_{t+2}^{II}\|, \dots$

Those identity reveal the relation between low order error and its higher order partner. Start from $N = 2$, we can achieve

$$\begin{aligned} \mathbf{X}_{t+2}^O &= \mathbf{X}_{t+2}^I + \boldsymbol{\varepsilon}_{t+2}^I = \mathbf{X}_{t+2}^{II} + \boldsymbol{\varepsilon}_{t+2}^{II} \\ &\rightarrow f(\mathbf{X}_{t+1}^O) + \boldsymbol{\varepsilon}_{t+2}^I = f(\mathbf{X}_{t+1}^I) + \boldsymbol{\varepsilon}_{t+2}^{II} \\ &\rightarrow f(\mathbf{X}_{t+1}^O) - f(\mathbf{X}_{t+1}^I) + \boldsymbol{\varepsilon}_{t+2}^I = \boldsymbol{\varepsilon}_{t+2}^{II} \\ &\rightarrow \nabla f(\mathbf{X}_\delta)(\mathbf{X}_{t+1}^O - \mathbf{X}_{t+1}^I) + \boldsymbol{\varepsilon}_{t+2}^I = \boldsymbol{\varepsilon}_{t+2}^{II} \\ &\rightarrow \nabla f(\mathbf{X}_\delta)(\boldsymbol{\varepsilon}_{t+1}^I) + \boldsymbol{\varepsilon}_{t+2}^I = \boldsymbol{\varepsilon}_{t+2}^{II} \\ &\rightarrow M_{t+1}^I(\boldsymbol{\varepsilon}_{t+1}^I) + \boldsymbol{\varepsilon}_{t+2}^I = \boldsymbol{\varepsilon}_{t+2}^{II} \end{aligned}$$

Here, we utilize the Mean Value Theorem of multi-dimensional calculus, where $\mathbf{X}_\delta \in [\mathbf{X}_{t+1}^O, \mathbf{X}_{t+1}^I]$ and ∇f represents the Jacobian operator. This operator produces a matrix $M_{t+1}^I(\delta)$ for \mathbf{X}_δ , which is abbreviated as M_{t+1}^I . The metric δ is dependent on the model parameter θ and fluctuate between the predicted value \mathbf{X}_{t+1}^I and the ground truth \mathbf{X}_{t+1}^O . As $\delta \rightarrow 0$, the model f_θ approaches the ideal model $f_\theta \rightarrow f$ and can accurately estimate any $\mathbf{X}_\delta = \mathbf{X}_{t+1}^O = \mathbf{X}_{t+1}^I = f(\mathbf{X}_t^O)$.

Similarly, any order propagation can be expressed as follows:

$$\boldsymbol{\varepsilon}_{t+N}^N = \boldsymbol{\varepsilon}_{t+N}^I + M_{t+N-1}^{N-1} \boldsymbol{\varepsilon}_{t+N-1}^{N-1}$$

where ϵ_{t+N}^N is the N -order error at time $t+N$, ϵ_{t+N-1}^1 is the first-order error at time $t+N-1$, and M_{t+N-1}^{N-1} is the propagation matrix. We can observe that this is a recursive formulation that goes from ϵ_{t+N-1}^{N-1} to ϵ_{t+N}^N .

Since all ϵ terms represent the prediction's error, they are expected to behave like high-dimensional, non-correlated, and time-independent noise fluctuation. Therefore, when the number of freedoms in ϵ is large enough, the elements from different time stamps can be considered "almost" orthogonal: $\langle \epsilon_{t+N-1}^1 | M_{t+N-1}^{N-1} \epsilon_{t+N-1}^{N-1} \rangle = 0$. Then the size of error for each sample is measured as

$$\|\epsilon_{t+N}^N\| \approx \|\epsilon_{t+N}^1\| + \|M_{t+N-1}^{N-1} \epsilon_{t+N-1}^{N-1}\|$$

And the average error is measured as

$$\begin{aligned} \mathcal{E}^N &= \frac{1}{n} \sum_t \|\epsilon_{t+N}^N\| \\ &= E[\|\epsilon_{t+N}^1\|] + E[\|M_{t+N-1}^{N-1} \epsilon_{t+N-1}^{N-1}\|] \\ &= \mathcal{E}^1 + E[\alpha_t^N \|\epsilon_{t+N-1}^{N-1}\|] \\ &= \mathcal{E}^1 + E[\alpha_t^N] E[\|\epsilon_{t+N-1}^{N-1}\|] \\ &= \mathcal{E}^1 + \alpha^N \mathcal{E}^N \end{aligned}$$

It is important to note that $E[\alpha_t^N] E[\|\epsilon_{t+N-1}^{N-1}\|]$ may not be equal to $E[\alpha_t^N \|\epsilon_{t+N-1}^{N-1}\|]$ since there is no guarantee that the random variables α_t^N and $\|\epsilon_{t+N-1}^{N-1}\|$ are independent. Geometrically, α_t^N is a factor that amplifies any potential fluctuations of the input $\mathbf{X}_{t+N-1} + \epsilon_{t+N-1}^{N-1}$ on the result. Since different t corresponds to different orientations in variable space, and our model has no anisotropic prior, such an amplifier should also be "almost" independent of the orientation. Therefore, we assume that $\alpha_t^N = \|M_{t+N-1}^{N-1} \epsilon_{t+N-1}^{N-1}\| / \|\epsilon_{t+N-1}^{N-1}\|$ is a number that only depends on the order N and the model θ .

Finally, we get the error propagation law of ATT problem when 1-order prediction error δ as small as enough $\delta \rightarrow 0$:

$$\mathcal{E}^N = (1 + \alpha^N + \alpha^N \alpha^{N-1} + \alpha^N \alpha^{N-1} \alpha^{N-2} + \dots) \mathcal{E}^1 \quad (2)$$

Furthermore, experiments in Sec.2 reveal that $\alpha^{i:1 \rightarrow N}$ follow an exponential decay law:

$$\alpha^N = 1 - \beta_1 \exp(-\beta_2 * N)$$

This finding implies that optimize α^1 and α^2 will spontaneously suppress all $\alpha^{i>2}$. Moreover, we can estimate long-range performance by computing only the first three errors \mathcal{E}^{III} , \mathcal{E}^{II} , and \mathcal{E}^{I} , using them to compute α^1 and α^2 , and fitting all α^N via this empirical formula. This allows us to quickly estimate the long-range prediction error \mathcal{E}^N .

According to Equ.1, to optimize a model with minimum N -order error \mathcal{E}^N , we need to optimize both the 1-order error \mathcal{E}^1 and the error amplifier $\alpha_{i:1 \rightarrow N}$. The traditional ATT loss directly optimizes $\mathcal{E}^1 + \mathcal{E}^2$, which is actually a coupled version of this view. For example, when $N = 2$, the traditional loss configuration treats the 2-order error and 1-order error equally, leading to:

$$\mathcal{E} = \mathcal{E}^{\text{II}} + \mathcal{E}^{\text{I}} = \sum_t \sum_s (\epsilon_{t+2,s}^{\text{II}})^2 + \sum_t \sum_s (\epsilon_{t+1,s}^{\text{I}})^2$$

This enforces that both errors are minimized rather than optimizing the amplifier $\alpha_2 = \frac{\mathcal{E}^{\text{II}}}{\mathcal{E}^{\text{I}}} - 1$. However, this approach risks amplifying $\frac{\mathcal{E}^{\text{II}}}{\mathcal{E}^{\text{I}}}$ when both \mathcal{E}^{II} and \mathcal{E}^{I} are minimized, which can lead to very poor long-term forecasts.

A good long-term forecast trainer must decouple the errors $\mathcal{E}^N, \mathcal{E}^{N-1}, \dots, \mathcal{E}^1$. One approach to achieve this is to directly optimize $\mathcal{E} = c_1 \frac{\mathcal{E}^{\text{II}}}{\mathcal{E}^{\text{I}}} + c_2 \mathcal{E}^{\text{I}}$. However, the reciprocal loss makes training unstable. An alternative is to use the logarithm of the mean square error, $\text{Log}\mathcal{E} = \ln \mathcal{E}^{\text{II}} + \ln \mathcal{E}^{\text{I}}$. Unfortunately, such a loss is not additive during mini-batch training, especially. It is equivalent between $\mathcal{E} = \sum_b^{\text{batch}} \mathcal{E}_b$ and $\text{Log}\mathcal{E} \neq \sum_b^{\text{batch}} \text{Log}\mathcal{E}_b$ with dynamic coefficients, as shown below:

$$d\text{Log}\mathcal{E}_b = d \ln \mathcal{E}_b^{\text{II}} + d \ln \mathcal{E}_b^{\text{I}} = \frac{1}{\mathcal{E}_b^{\text{II}}} d\mathcal{E}_b^{\text{II}} + \frac{1}{\mathcal{E}_b^{\text{I}}} d\mathcal{E}_b^{\text{I}}$$

Take $N = 2$ as an example, the goal is to minimize

$$\min \mathcal{E}^{\text{I}} \ \& \ \min \alpha^{\text{II}} = \frac{\mathcal{E}^{\text{II}}}{\mathcal{E}^{\text{I}}} - 1$$

Notice minimize α^{II} equals to minimize $\frac{\mathcal{E}^{\text{II}}}{\mathcal{E}^{\text{I}}}$ which can be converted to log format as minimized $\log(\mathcal{E}^{\text{II}}) - \log(\mathcal{E}^{\text{I}})$. Remember we still need minimize \mathcal{E}^{I} which is same as minimize $\log(\mathcal{E}^{\text{I}})$. These two part are independent implying that we need to assign two coefficients to combine them like:

$$\text{final loss} = a * [\log(\mathcal{E}^{\text{II}}) - \log(\mathcal{E}^{\text{I}})] + b * \log(\mathcal{E}^{\text{I}})$$

They are hyper-parameter and we just set $a = 1, b = 2$ in this paper. This results in:

$$\text{final loss} = \log(\mathcal{E}^{\text{II}}) + \log(\mathcal{E}^{\text{I}}) = \text{MLSE}$$

Machine learning training typically involves batch updating. Thus, the key point here is to decide whether to "average after computing" or "compute after averaging". We leave this to the experiment.

For three-timestamps optimization MLSE, the derivation is:

$$\min \mathcal{E}^{\text{I}} \ \& \ \min \alpha^{\text{II}} = \frac{\mathcal{E}^{\text{II}}}{\mathcal{E}^{\text{I}}} - 1 \ \& \ \min \alpha^{\text{III}} = \frac{\mathcal{E}^{\text{III}}}{\mathcal{E}^{\text{II}}} - \frac{\mathcal{E}^{\text{I}}}{\mathcal{E}^{\text{II}}}$$

which is equivalent to:

$$\min\{\mathcal{E}^{\text{I}}, \frac{\mathcal{E}^{\text{II}}}{\mathcal{E}^{\text{I}}}, \frac{\mathcal{E}^{\text{III}}}{\mathcal{E}^{\text{II}}}\} \ \text{and} \ \max\{\frac{\mathcal{E}^{\text{I}}}{\mathcal{E}^{\text{II}}}\}.$$

The last term is same as $\min\{\frac{\mathcal{E}^{\text{II}}}{\mathcal{E}^{\text{I}}}\}$. Using the hyper-parameter coefficient trick, we can finally obtain:

$$\text{MLSE}_{\text{order3}} = \log(\mathcal{E}^{\text{III}}) + \log(\mathcal{E}^{\text{II}}) + \log(\mathcal{E}^{\text{I}})$$

In Equation 1, there are many other useful error information that can be utilized. For instance, we can calculate the error between $\mathbf{X}_{t+2}^{\text{II}}$ and $\mathbf{X}_{t+2}^{\text{I}}$ instead of $\|\mathcal{E}_{t+2}^{\text{II}} = \mathbf{X}_{t+2}^{\text{O}} - \mathbf{X}_{t+2}^{\text{II}}\|$. Notice the equation $\ln[\alpha_t^{\text{II}}] = \ln[\|\mathbf{X}_{t+2}^{\text{II}} - \mathbf{X}_{t+2}^{\text{I}}\|] - \ln[\|\mathbf{X}_{t+1}^{\text{I}} - \mathbf{X}_{t+1}^{\text{O}}\|] = \ln[\varepsilon_{t+2}^{\text{II,I}}] - \ln[\varepsilon_{t+1}^{\text{I}}]$ (as shown in Equation 2), thus a loss simultaneously optimizing $\ln \alpha_t^{\text{II}}$ and $\ln \varepsilon_{t+1}^{\text{I}}$ is equivalent to optimizing $\ln \varepsilon_{t+2}^{\text{II,I}}$ and $\ln \varepsilon_{t+1}^{\text{I}}$.

$$\left(\begin{array}{cccccc} \mathbf{X}_t^{\text{O}} & \rightarrow & \mathbf{X}_{t+1}^{\text{I}} & \rightarrow & \mathbf{X}_{t+2}^{\text{II}} & \rightarrow & \mathbf{X}_{t+3}^{\text{III}} & \rightarrow & \mathbf{X}_{t+4}^{\text{IV}} \\ \mathbf{X}_{t+1}^{\text{O}} & \rightarrow & \mathbf{X}_{t+2}^{\text{I}} & & & & & & \\ \mathbf{X}_{t+2}^{\text{O}} & \rightarrow & \mathbf{X}_{t+3}^{\text{I}} & & & & & & \\ \mathbf{X}_{t+3}^{\text{O}} & \rightarrow & \mathbf{X}_{t+4}^{\text{I}} & & & & & & \end{array} \right) \quad (3)$$

The ideas presented here have inspired us to develop the MASE

$$\text{MASE} = \sum_N \sum_t [\|\mathbf{X}_{t+N}^{\text{N}} - \mathbf{X}_{t+N}^{\text{I}}\| + \|\mathbf{X}_{t+N-1}^{\text{N-1}} - \mathbf{X}_{t+N-1}^{\text{O}}\|]$$

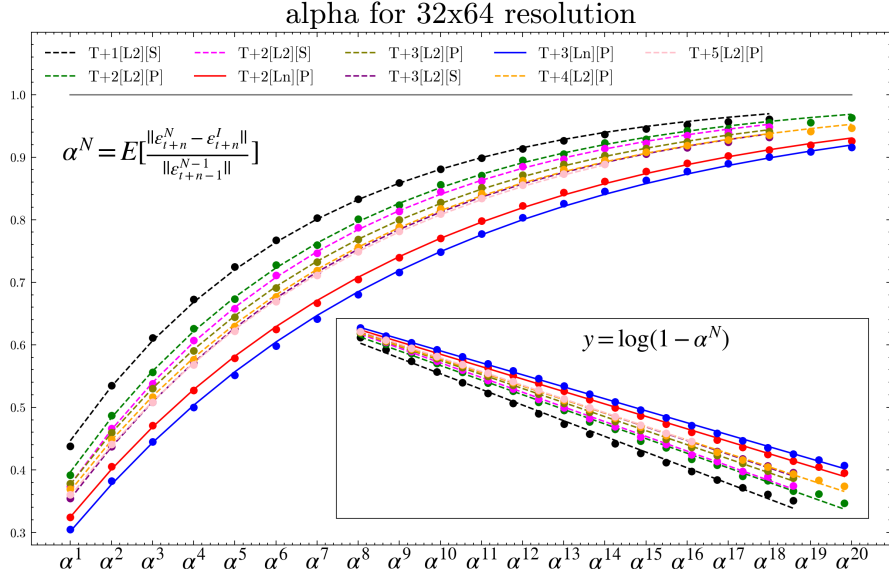


Figure 1: The figure shows the amplifier computed using Equation 2 for different convergence points resulting from various training strategies. The [S] label indicates that the model was trained from a random initialization, while the [P] label indicates that the model was trained from pre-trained weights of model T+1[L2][L2]. The [L2] label indicates that the mean squared error (MSE) loss was used, while the [Ln] label indicates that the mean logarithmic squared error (MLSE) loss was used. The T+n notation indicates that the error for the attention task (ATT) was optimized from 0 to N-th orders.

2 OBSERVATION

We monitor the amplifier α^N on the WeatherBench32x64 dataset for a well-trained FourCastNet model. We employ different training strategies to encourage the model to converge to different points. The test dataset for this evaluation consists of a $1480 \times (70 \times 32 \times 64)$ vector sequence. For each starting point \mathbf{X}_t , we calculate the model’s forecast results up to 19 time steps, i.e., $\hat{X}_{t+1}, \dots, \hat{X}_{t+19}$. Finally, we collect 1461×18 different order error \mathcal{E}^N and calculate the amplifier α^N by

$$\alpha^N = E\left[\frac{\|\boldsymbol{\varepsilon}_{t+N}^N - \boldsymbol{\varepsilon}_{t+N}^1\|}{\|\boldsymbol{\varepsilon}_{t+N-1}^{N-1}\|}\right] = E\left[\frac{\|\mathbf{X}_{t+N}^N - \mathbf{X}_{t+N}^1\|}{\|\mathbf{X}_{t+N-1}^{N-1} - \mathbf{X}_{t+N-1}^0\|}\right] \quad (4)$$

The results, depicted in Fig.1, are based on many different training strategies: The [S] label indicates that the model was trained from a random initialization, while the [P] label indicates that the model was trained from pre-trained weights of model T+1[L2][L2]. The [L2] label indicates that the mean squared error (MSE) loss was used, while the [Ln] label indicates that the mean logarithmic squared error (MLSE) loss was used. The T+n notation indicates that the error for the attention task (ATT) was optimized from 0 to N-th orders.

The findings depicted in Fig. 1 provide clear evidence that, once a deep learning model has been trained and converged on a dataset, its long-range prediction errors follow an exponential amplifier law. Specifically, the amplifier α^N between N -order errors and $N - 1$ errors can be approximated by a geometric fitting curve, such as $\alpha^N = 1 - \beta_1 \exp(-\beta_2 * N)$. We validate this law for different models using varying training strategies, highlighting the fact that the long-term prediction error can be effectively optimized by

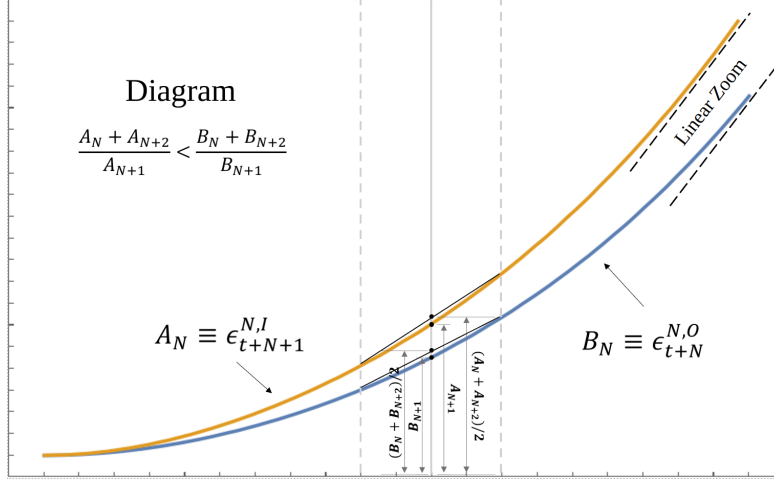


Figure 2: The diagram to show the convex of error propagation. Note it is only a diagram, to vivid demonstrate the convex of amplifier $\alpha^N = \frac{A_N}{B_N}$.

appropriately configuring the short-term error. These findings provide important insights into the relationship between short-term performance and long-term forecasting accuracy, particularly in cases where the model is not exposed to the forecasted data during training.

The geometric formulation exhibits an exponential behavior experimentally. However, a simple analysis reveals that the convex behavior can be obtained. Let us denote the series $\|\epsilon_{t+N+1}^{N+1,I}\|$ as A_N and the series $\|\epsilon_{t+N}^{N,O}\|$ as B_N . We know that both A_N and B_N increase and achieve linear growth at infinity N ; Before that, they are strictly convex. Since A_N is the higher error incresement, we can assume $B_{N+2} \geq A_N \geq B_{N+1}$ for any N . This leads to $\frac{A_{N+2}+A_N}{A_{N+1}} \leq \frac{B_{N+2}+B_N}{B_{N+1}}$, as the error accumulation speed should be approximately consistent for both series A_N and B_N as shown in Diagram.2. Therefore, the second differences of the sequence $\alpha^N = \frac{A_N}{B_N}$ becomes

$$(\alpha^{n+2} - \alpha^{n+1}) - (\alpha^{n+1} - \alpha^n) = \frac{A_{n+2}B_{n+1} - A_{n+1}B_{n+2} - A_{n+1}B_n + A_nB_{n+1}}{B_nB_{n+1}B_{n+2}} \leq 0$$

. As the model is a black-box, it is difficult to obtain further insights into its behavior beyond the geometric formulation. However, we are able to analyze the asymptotic behavior as $N \rightarrow \infty$. Due to the fact that time-correlation is not infinite, the prediction at $N - 1$ step \mathbf{X}_{t+N}^N is no longer related to the first step \mathbf{X}_{t+N-1}^O , as well as \mathbf{X}_{t+N}^N to \mathbf{X}_{t+N}^I . Consequently, if we assume each normlized element $X_{t,s}^N \sim N(0, 1)$ follow the normal distribution, then the α^N degenerates to the division between two χ_k^2 distributions, which corresponds to an F -distribution. The expectation is $\alpha^{N \rightarrow \infty} = \frac{S}{S-2} \rightarrow 1$, where S is the freedom of X . This finding implies that the asymptotic behavior of error propagation is linear, which means the difference between ϵ_N and ϵ_{N-1} is a constant as $N \rightarrow \infty$. By combining this with $\alpha^{N \rightarrow \infty} = 1$, we can qualitatively understand the behavior of alpha.

The geometric behavior can aid in quickly estimating the long-term performance of a model based on a few first-order errors. For example, Fig. 3 shows the estimated performance of model [T+3][P][L2]. Due to the entanglement problem mentioned previously, α_t^N and $\|\epsilon_{t+N-1}^{N-1}\|$ typically have the relationship $E[\alpha_t^N]E[\|\epsilon_{t+N-1}^{N-1}\|] \leq E[\alpha_t^N]E[\|\epsilon_{t+N-1}^{N-1}\|]$, which represents an upper bound on the true error. Fig. 3

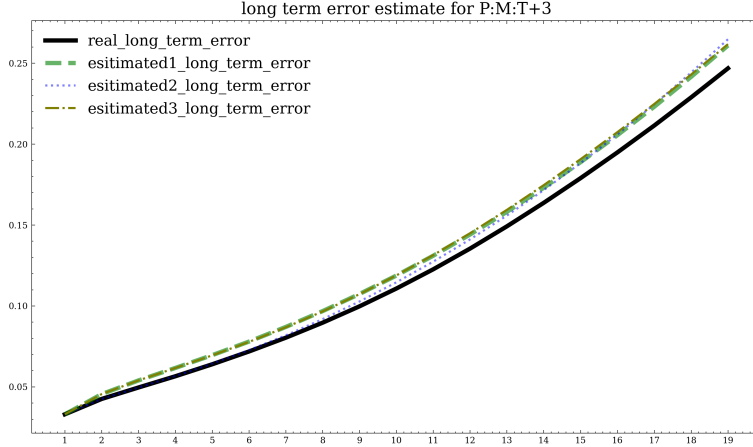


Figure 3: The estimated long-term performance of model [T+3][P][L2] from its first three low order errors. The three esitamed method is low-order expansion, full-amplifier computing, and statistic-fitting. The result demonstrates that estimation can achieve a considerable approximation of the true long-term error by only measuring a few short-term errors. In this example, the 19-step error is smaller than 6%.

depicts the results obtained from three different estimation methods that were used to test the influence of different statistical routes: (1) The estimated1 method uses the statistical means α^0 and α^1 by only accessing the \mathcal{E}^I , \mathcal{E}^{II} , and \mathcal{E}^{III} , then calculates all α values, and finally accumulates the amplifier using the Equ.2. (2) The estimated2 method uses all statistical means $\alpha^{n:1 \rightarrow 19}$. This is not a practical method for estimation since it requires the computation of all $\mathcal{E}^{n:1 \rightarrow 19}$. However, we plot the results to emphasize the feasibility of the exponential law for α . (3) The estimated3 method uses the losses \mathcal{E}^I , \mathcal{E}^{II} , and \mathcal{E}^{III} directly to compute the "mean statistical alpha" $\bar{\alpha}^1 = \frac{\mathcal{E}^{II} - \mathcal{E}^I}{\mathcal{E}^I}$ and $\bar{\alpha}^2 = \frac{\mathcal{E}^{III} - \mathcal{E}^I}{\mathcal{E}^{II}}$. The example of model [T+3][P][L2] demonstrates that estimation can achieve a considerable approximation of the true long-term error by only measuring a few short-term errors. In this example, the 19-step error is smaller than 6%.

In this section, we have shown that the error amplifier in autoregressive prediction tasks follows an obvious geometric law. This implies that constraints on one amplifier will spontaneously distribute to all amplifiers. Therefore, the correct descent direction for minimizing long-term prediction error is the same as the direction for minimizing the amplifier directly. When the model θ is near the convergence point, it may be better to constrain the amplifier α and the 1-order error \mathcal{E}^I rather than directly constraining high-order errors \mathcal{E}^I , \mathcal{E}^{II} , \mathcal{E}^{III} . We have also demonstrated that this phenomenological theory can help estimate long-range forecast performance by only measuring a few short-term errors. Unfortunately, even if we can analytically write the relationship between long-term and short-term errors, we cannot use this analytical expression for descent updates directly for several reasons: (1) It is merely an estimate and is only effective near the convergence point. (2) It is only effective in a statistical sense and is incompatible with the mini-batch training mode. (3) The gradient of the analytical expression for \mathcal{E}^{II} is always negative.

3 THE ALPHA MONITOR FOR TIME SERIES DATASET

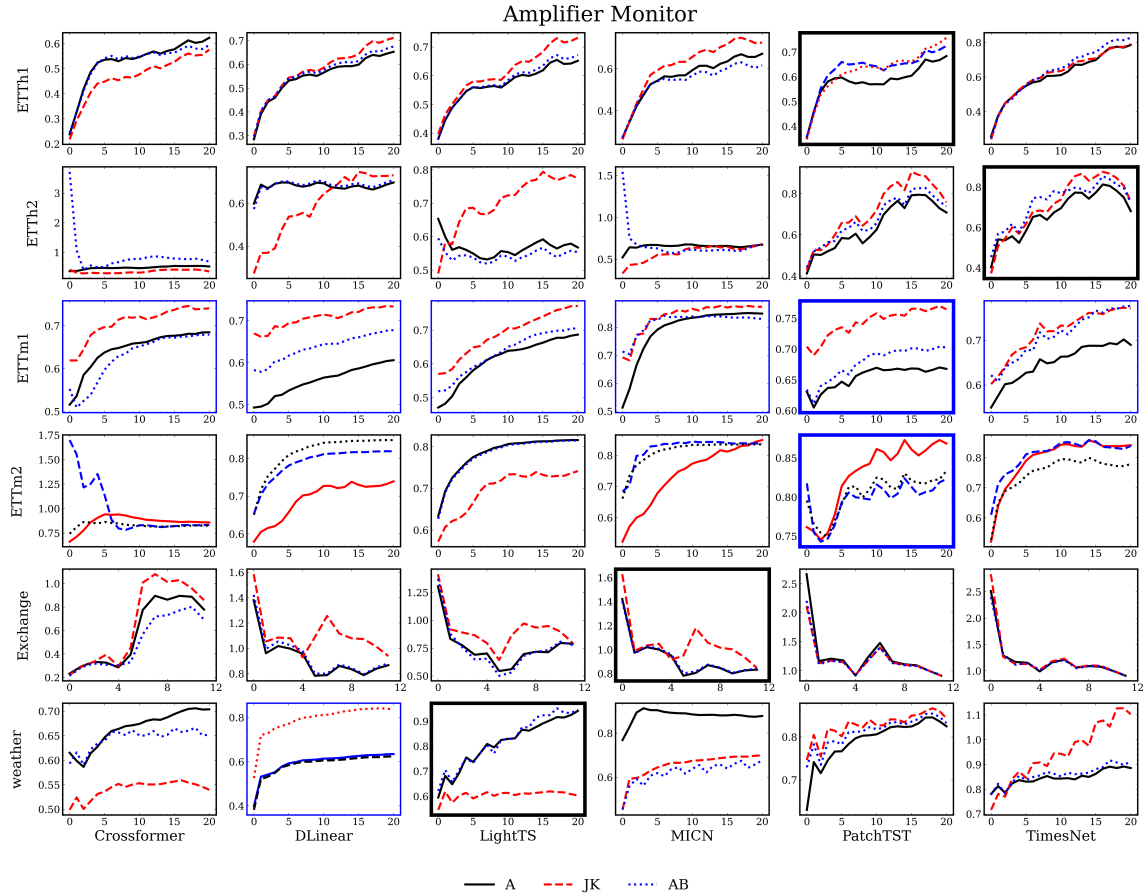


Figure 4: The alpha monitor for the case presented in Table 2 follows the same procedure outlined in Section 2. We only plot the alpha on the prediction sequence equal to 96. The frame colour and size of the 720 rows in Table 2 correspond. A black board indicates that the MLSE+ATT model performs better than the end-to-end model, while a blue board is the opposite. The bold boundary represents the best model for each dataset, which is the case with background shading in Table 2. From this figure, we can see that, except for the Exchange dataset, all alpha relations are approximately monotonically increasing and converge near 1. There are several abnormal examples, such as Crossformer under ETTh2 and ETTm2 datasets. By observing its long-term error, we can see that the end-to-end performance of this model is extremely poor and does not converge at all. However, using MLSE can correct this behavior and bring the long-term prediction into the same convergence basin. The exchange dataset performs differently from other time series datasets, implying that it is not a predictable system since economic systems are often filled with complex external factors.

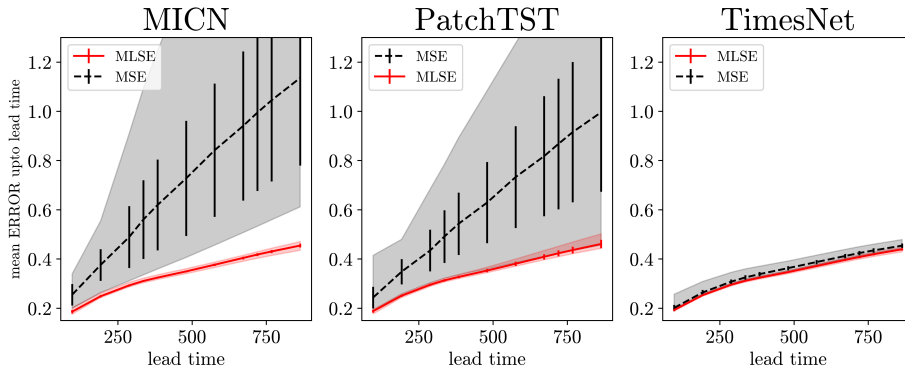


Figure 5: We evaluated the fine-tuned performance of TimesNet, MICN, and PatchTST on the ETTm2 datasets using MSE/MLSE metrics. The learning rate was adjusted from $1e-5$ to $1e-3$, and the random seed was selected from the following options: [19940928, 19950929, 20130901, 20230901]. On the y-axis, we plotted the mean error up to the lead time, which is the average error across the entire sequence. The shaded background represents the min-max region of error. The error bars denote the standard deviation, and the central line indicates the mean value. We conducted 20 samples for both the MLE and MLSE hyperparameter searches.

dataset	pred	model	Crossformer				DLinear				LightS				MICN				PatchTST				TimesNet			
		metric	A	D	R	A	D	R	A	D	R	A	D	R	A	D	R	A	D	R	A	D	R	A	D	R
ETTth1	096		0.414	0.403	0.372	0.396	0.365	0.366	0.435	0.399	0.400	0.394	0.407	0.372	0.378	0.348	0.349	0.389	0.363	0.364						
	192		0.450	0.480	0.420	0.445	0.414	0.411	0.494	0.441	0.439	0.454	0.474	0.434	0.437	0.382	0.399	0.437	0.396	0.403						
	336		0.731	0.554	0.474	0.490	0.443	0.434	0.552	0.469	0.456	0.598	0.565	0.460	0.466	0.403	0.403	0.490	0.413	0.427						
	720		0.609	0.736	0.613	0.509	0.525	0.485	0.613	0.561	0.550	0.696	0.784	0.690	0.507	0.456	0.405	0.520	0.461	0.489						
	096		0.617	1.046	0.504	0.348	0.257	0.211	0.428	0.360	0.295	0.344	0.344	0.211	0.306	0.218	0.209	0.321	0.261	0.235						

Figure 6: A table similar to Table 2, but with the learning rate initialized at $2e-5$ (original value is $1e-5$).

dataset	pred	model	Crossformer				DLinear				LightS				MICN				PatchTST				TimesNet			
		metric	A	D	R	A	D	R	A	D	R	A	D	R	A	D	R	A	D	R	A	D	R	A	D	R
ETTth1	096		0.414	0.384	0.363	0.396	0.366	0.366	0.435	0.402	0.402	0.394	0.389	0.372	0.378	0.348	0.349	0.389	0.361	0.362						
	192		0.450	0.444	0.411	0.445	0.416	0.413	0.494	0.443	0.439	0.454	0.447	0.401	0.437	0.382	0.387	0.437	0.396	0.400						
	336		0.731	0.504	0.451	0.490	0.446	0.437	0.552	0.470	0.457	0.598	0.523	0.490	0.466	0.401	0.401	0.490	0.415	0.424						
	720		0.609	0.676	0.570	0.509	0.529	0.499	0.613	0.560	0.523	0.696	0.736	0.699	0.507	0.455	0.427	0.520	0.464	0.487						
	096		0.617	1.031	0.562	0.348	0.260	0.215	0.428	0.359	0.268	0.344	0.304	0.212	0.306	0.213	0.207	0.321	0.255	0.234						

Figure 7: A table similar to Table 2, but with the random seed initialized at 2023. (original value is 2021)

6 THE ERROR BAR (STD) FOR TIME SERIES DATASET

We fix the random seed and repeat the finetune experiment 5 times to get the standard deviation of the results. (These fluctuations arise from machine and build-in code feature). The results are shown in Table 9.

dataset	pred	model Autoformer		Crossformer		DLinear		Informer		LightTS		MICN		PatchTST		Pyraformer		Reformer		Stationary		TimesNet		Transformer			
		metric	D	R	D	R	D	R	D	R	D	R	D	R	D	R	D	R	D	R	D	R	D	R	D	R	
ETTh1	096	2e-03	2e-03	3e-03	1e-03	1e-05	1e-05	4e-03	4e-03	2e-04	2e-04	1e-03	2e-03	3e-04	0e+00	4e-03	0e+00	7e-03	5e-03	4e-03	2e-03	5e-04	2e-05	6e-03	4e-03		
	192	3e-03	6e-03	5e-03	2e-03	2e-05	1e-05	3e-03	3e-03	6e-04	5e-04	3e-03	5e-03	3e-04	4e-04	1e-02	8e-03	8e-03	5e-03	2e-03	3e-03	3e-04	4e-04	1e-02	9e-03		
	336	6e-03	9e-03	9e-03	4e-03	5e-05	3e-05	5e-03	3e-03	1e-03	1e-03	6e-03	8e-03	3e-04	2e-04	2e-02	9e-03	7e-03	5e-03	3e-04	4e-03	6e-04	8e-04	1e-02	1e-02		
	720	6e-03	1e-02	1e-02	1e-02	2e-04	2e-04	1e-02	6e-03	3e-03	3e-03	1e-02	2e-02	2e-03	4e-04	2e-02	5e-03	2e-03	1e-03	1e-03	9e-03	1e-03	2e-03	2e-02	8e-03		
ETTh2	096	4e-04	2e-04	1e-02	7e-03	5e-04	2e-04	7e-01	5e-01	4e-03	5e-04	1e-03	4e-04	7e-04	6e-04	4e-02	2e-02	7e-02	6e-02	2e-03	3e-03	3e-05	4e-04	3e-02	7e-02		
	192	2e-03	4e-04	3e-03	1e-02	1e-03	2e-04	6e-01	4e-01	1e-02	8e-04	1e-02	1e-04	7e-04	5e-04	8e-01	5e-01	8e-03	1e-02	3e-03	2e-03	4e-05	6e-04	3e-02	5e-02		
	336	2e-02	7e-04	2e-02	1e-02	2e-03	2e-04	5e-01	4e-01	2e-02	2e-03	3e-02	3e-04	7e-04	6e-04	7e-01	5e-01	1e-01	1e-02	2e-03	2e-03	1e-03	6e-04	2e-02	4e-02		
	720	7e-02	4e-03	6e-02	1e-02	3e-03	1e-04	4e-01	3e-01	2e-02	4e-03	8e-02	1e-03	1e-03	4e-04	4e-01	3e-01	1e-01	2e-02	2e-03	2e-03	2e-03	5e-04	2e-02	3e-02		
ETTm1	096	4e-03	9e-04	5e-03	1e-03	6e-06	4e-05	2e-02	4e-03	4e-05	8e-05	4e-03	2e-04	5e-04	4e-04	2e-02	0e+00	2e-02	1e-02	6e-03	2e-03	4e-04	5e-04	4e-03	5e-03		
	192	1e-02	2e-03	6e-03	2e-03	2e-05	2e-05	2e-02	4e-03	2e-04	1e-04	1e-02	2e-04	3e-04	5e-04	1e-02	0e+00	2e-02	1e-02	9e-03	3e-03	3e-03	6e-04	6e-03	6e-03		
	336	2e-02	1e-02	1e-02	5e-03	4e-05	4e-05	2e-02	3e-03	6e-04	2e-04	2e-02	4e-04	1e-03	3e-04	1e-02	0e+00	2e-02	1e-02	1e-02	5e-03	6e-03	1e-03	7e-03	8e-03		
	720	5e-02	4e-02	3e-02	1e-02	8e-05	9e-05	4e-03	7e-03	1e-03	3e-04	6e-02	2e-03	2e-03	7e-04	2e-02	2e-03	2e-02	1e-02	2e-02	1e-02	1e-02	2e-03	1e-02	1e-02		
ETTm2	096	7e-03	4e-04	9e-04	9e-04	3e-17	3e-06	2e-02	1e-16	1e-03	4e-06	1e-07	7e-05	1e-03	2e-05	8e-03	1e-03	3e-02	3e-02	5e-03	4e-03	4e-04	2e-04	3e-03	2e-03		
	192	7e-03	5e-04	1e-02	2e-03	6e-17	1e-05	2e-02	0e+00	3e-03	1e-05	8e-03	1e-04	3e-04	5e-04	3e-02	5e-03	3e-02	3e-02	4e-03	4e-03	2e-04	3e-04	4e-03	3e-03		
	336	6e-03	7e-04	2e-02	3e-03	0e+00	2e-05	2e-02	1e-16	6e-03	3e-05	2e-02	5e-04	2e-04	7e-04	9e-02	1e-02	3e-02	3e-02	4e-03	4e-03	3e-04	4e-04	5e-03	4e-03		
	720	6e-04	1e-03	2e-02	1e-02	1e-16	3e-05	2e-02	8e-03	2e-02	4e-05	7e-02	1e-03	1e-03	6e-04	1e-01	2e-02	3e-02	3e-02	2e-02	2e-02	9e-04	3e-04	2e-03	5e-03		
Exchange	096	4e-03	3e-03	2e-02	4e-03	1e-04	0e+00	2e-02	3e-02	0e+00	6e-04	5e-04	9e-04	2e-04	5e-04	3e-02	5e-02	2e-01	1e-01	7e-04	1e-03	3e-05	6e-04	6e-03	1e-02		
	192	8e-03	5e-03	4e-02	6e-03	9e-04	2e-03	9e-03	2e-02	0e+00	2e-03	5e-03	1e-03	7e-04	8e-04	9e-02	1e-01	2e-01	1e-01	1e-03	3e-03	2e-04	9e-05	5e-03	7e-03		
	336	2e-02	7e-03	5e-02	9e-03	2e-03	3e-03	9e-03	2e-02	0e+00	6e-03	4e-02	2e-03	9e-04	2e-03	7e-02	9e-02	2e-01	1e-01	2e-03	4e-03	3e-04	3e-04	5e-03	1e-02		
	720	5e-02	3e-02	4e-02	1e-02	2e-03	2e-02	9e-03	2e-02	1e-16	2e-03	3e-02	2e-02	9e-03	1e-02	2e-02	2e-02	2e-01	8e-02	2e-02	2e-02	6e-04	4e-05	6e-03	1e-02		
weather	096	2e-02	8e-03	1e-03	3e-04	0e+00	0e+00	2e-02	1e-03	0e+00	2e-04	3e-03	5e-03	3e-04	2e-03	3e-03	4e-04	8e-03	2e-03	2e-03	1e-03	5e-05	3e-04	2e-03	8e-03		
	192	2e-02	6e-04	4e-04	3e-04	2e-04	3e-17	6e-03	1e-03	2e-04	4e-05	3e-03	2e-03	3e-04	1e-03	4e-03	2e-04	8e-03	3e-03	3e-03	3e-03	9e-05	2e-04	6e-04	1e-02		
	336	4e-02	3e-03	8e-04	3e-04	4e-04	2e-04	7e-03	2e-03	6e-04	8e-05	6e-03	1e-03	6e-05	7e-04	6e-03	7e-04	6e-03	3e-03	5e-03	4e-03	1e-04	1e-04	6e-03	1e-02		
	720	5e-02	2e-02	1e-04	4e-04	9e-04	6e-04	7e-03	3e-03	5e-03	3e-04	7e-03	5e-03	1e-04	5e-05	1e-02	4e-03	3e-03	4e-03	2e-02	2e-02	4e-04	3e-04	2e-02	2e-02		

Figure 9: The table shows the standard derivation for repeated experience among all the time series experiment. We only repeat experiment for finetune task.

7 USE MLSE START FROM SKETCH DOESN'T PERFORM BETTER

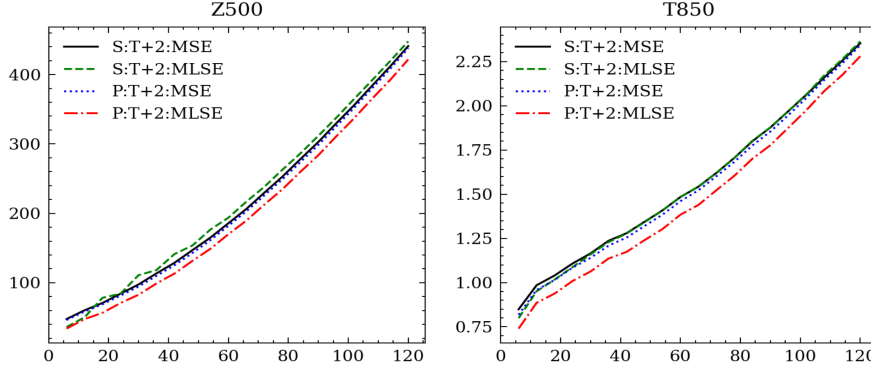


Figure 10: Directly use MLSE and train a model from sketch has no any benefit. It is only a finetune technology.

8 EULER EQUATION

The latent inductive bias in physical simulation tasks is that all evolution must follow a physics formula. For weather simulation, this can be an explicit formula such as the Euler equation (Equ.8) Laprise (1992) or an unexplicit neural network simulator such as FourCast. This prior implies that the relation between time frames must be "Space to Time", meaning that next state variables such as velocity V_{t+1} , geometry ϕ_{t+1} and temperature T_{t+1} can be derived from the current state information V_t , ϕ_t and T_t .

$$\begin{aligned}\partial_t V &= F - V \cdot \nabla V - \omega \partial_p V - \nabla \phi \\ \partial_t T &= Q/C_v + \frac{RT}{C_p} \omega - V \cdot \nabla T - \omega \partial_p T \\ \partial_t \phi &= wg - V \cdot \nabla V - \omega \partial_p \phi\end{aligned}$$

where V is the horizontal velocity, ω is the vertical velocity driven by $\frac{\partial \omega}{\partial p} = -\nabla \cdot V$, ϕ is the geopotential, T is the temperature, F is the force, all of which are continuous functions of the cartesian coordinates (x, y) and pressure p . Additionally, Q is the external heat, C_v is the heat capacity and R is the gas constant.

9 DATASET INFORMATION

In this paper, we have two major dataset:

- **WeatherBenchRasp** et al. (2020) is a large, high-resolution global atmospheric dataset created to facilitate the development and benchmarking of machine learning models for weather forecasting. It is widely used in the AI for weather predicting community, such as PanguBi et al. (2022), GraphCastLam et al. (2022) and FengwuChen et al. (2023). The dataset includes a range of atmospheric variables, such as temperature (T), pressure (ϕ), humidity (H) and wind speed (V) at multiple levels of the atmosphere, and covers the entire globe with a resolution of approximately 0.25 degrees (or about 28 km). The data is available at hourly intervals and covers a period from 1979 to 2018. This paper only involves the 32×64 and 64×128 resolutions. We use the horizontal wind speed information, temperature, pressure and humidity data, but without the constant variables; these features consist of the 68 channels of one earth snapshot. Following the same habit as previous works Bi et al. (2022); Lam et al. (2022); Chen et al. (2023), we sequentially divide the train/valid/test datasets. More specifically, we take data from years 1979-2016 as the train dataset, year 2017 as the valid dataset, and year 2018 as the test dataset. The WeatherBench-300k dataset task hourly snapshot thus has around three hundred thousands rawdata. The WeatherBench-50k dataset downsamples hourly data into 6-hourly data and holds around 55k rawdata.

dataset	WeatherBench32x64-6hour			WeatherBench32x64-1hour			WeatherBench64x128-1hour		
shape	$68 \times 32 \times 64$			$68 \times 32 \times 64$			$68 \times 32 \times 128$		
split	train	valid	test	train	valid	test	train	valid	test
years	1979 - 2016	2017	2018	1979 - 2015	2016 - 2017	2018	1979 - 2015	2016 - 2017	2018
number	55514	1458	1458	324301	17533	8749	324301	17533	8749

Table 2: The information of weatherbench dataset used in this paper.

- **Time series** use the long-term forecasting setting in TimesNet Wu et al. (2022) and same hyperparameters for benchmark. The datasets consist of ETT Zhou et al. (2021), Electricity Trindade (2016), Traffic PeM, Weather Wet, and Exchange Lai et al. (2018). Each part is a single continuous time series, and we sample rawdata by a sliding window.

dataset	input feature	train	valid	test	categorize
ETTM1, ETTM2	$7 \times 96, 192, 336, 720$	34465	11521	11521	Electricity(15 mins)
ETTh1, ETTh2	$7 \times 96, 192, 336, 720$	8545	2881	2881	Electricity (15 mins)
Electricity	$321 \times 96, 192, 336, 720$	18317	2633	5261	Electricity (Hourly)
Traffic	$862 \times 96, 192, 336, 720$	12185	1757	3509	Transportation (Hourly)
Weather	$21 \times 96, 192, 336, 720$	36792	5271	10540	Weather (10 mins)
Exchange	$8 \times 96, 192, 336, 720$	5120	665	1422	Exchange rate (Daily)

Table 3: The information of time series dataset used in this paper.

10 HYPER-PARAMETER AND TRAINING DETAIL

For the WeatherBench dataset training, the same hyperparameters are used for the settings "MSE", "MASE" or "MLSE". We usually train the model on 4 or 8 A100-80G GPUs using a DataParallel pipeline. The batch size is calculated as the total number and the random seed is fixed at 73001. We use a cosine learning rate scheduler which will anneal the learning rate from $1e-6$ to the set learning rate in the warmup epoch and then fall back to $1e-5$. The optimizer is AdamW with parameters $\beta = (0.9, 0.95)$ and weight decay 0.05. More details can be found in Table 10. Due to computing resource limitations, we only train the model 64×128 with a small number of epochs. It can be seen that the model quickly overfits in the 64×128 resolution during the finetune procedure. We do not use an earlystop strategy, so all models run for the

	dataset	WB32x64-50k		WB32x64-300k				WB64x128-300k	
	model	AFNONet		AFNONet		Lgnet		Lgnet	
	phase	pre-T	fine-T	pre-T	fine-T	pre-T	fine-T	pre-T	fine-T
	epoch	100	100	40	40	20	20	20	4
	batch_size	64	64	64	16	64	16	64	16
optimz	type	AdamW							
	beta	(0.9, 0.95)							
sched	type	Cosine							
	warmup	5							
	value	$1e-6 \rightarrow lr \rightarrow 1e-5$							
	lr	8e-4	8e-4	8e-4	8e-4	2e-4	2e-4	8e-4	1e-6

Table 4: The training information of weatherbench in this paper.

For the timeseries dataset, the hyperparameters for both the "MSE" and "MLSE" settings are identical. The only difference between the pretrain and finetune phases is the learning rate; the finetune learning rate is set to $1e-5$, which is 10x smaller than the pretrain learning rate of $1e-4$. All architecture settings are aligned with the TimesNet Library Wu et al. (2022). The scheduler is cosine with 0 warmup, so the learning rate will decay to $1e-6$ from its set value. The batch size is listed in Table 10. All experiments were run on a single A100-80G GPU.

REFERENCES

Pems. <http://pems.dot.ca.gov/>.

Wetterstation. weather. <https://www.bgc-jena.mpg.de/wetter/>.

Kaifeng Bi, Lingxi Xie, Hengheng Zhang, Xin Chen, Xiaotao Gu, and Qi Tian. Pangu-weather: A 3d high-resolution model for fast and accurate global weather forecast. *arXiv preprint arXiv:2211.02556*, 2022.

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.

dataset	ETTh1	ETTh2	ETTm1	ETTm2	Exchange	traffic	weather
Autoformer	32	32	32	32	32	32	32
Crossformer	32	32	32	32	32	16	32
DLinear	32	32	32	32	32	32	32
Informer	32	32	32	32	32	32	32
LightTS	32	32	32	32	32	32	32
MICN	32	32	32	32	32	32	32
PatchTST	32	32	32	32	32	4	32
Pyraformer	32	32	32	32	32	32	32
Reformer	32	32	32	32	32	32	32
Stationary	32	32	32	32	32	32	32
TimesNet	32	32	32	32	32	32	32
Transformer	32	32	32	32	32	32	32

Table 5: The batch size of time series in this paper.

Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long-and short-term temporal patterns with deep neural networks. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pp. 95–104, 2018.

Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirsberger, 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.

René Laprise. The euler equations of motion with hydrostatic pressure as an independent variable. *Monthly weather review*, 120(1):197–207, 1992.

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.

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.

Artur Trindade. Uci maching learning repository-electricityloadaddiagrams20112014 data set, 2016.

Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. *arXiv preprint arXiv:2210.02186*, 2022.

Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 11106–11115, 2021.