# Frozen in Time: Parameter-Efficient Time Series Transformers via Reservoir-Induced Feature Expansion and Fixed Random Dynamics

**Pradeep Singh[a,*], Mehak Sharma[a], Anupriya Dey[a] and Balasubramanian Raman[a]**

[a]Machine Intelligence Lab, Department of Computer Science and Engineering, IIT Roorkee, Roorkee-247667, India
ORCID (Pradeep Singh): https://orcid.org/0000-0002-5372-3355, ORCID (Mehak Sharma):
https://orcid.org/0009-0001-3102-1045, ORCID (Anupriya Dey): https://orcid.org/0009-0000-1630-1017, ORCID
(Balasubramanian Raman): https://orcid.org/0000-0001-6277-6267

**Abstract.** Transformers are the de-facto choice for sequence modelling, yet their quadratic self-attention and weak temporal bias can make long-range forecasting both expensive and brittle. We introduce *FreezeTST*, a lightweight hybrid that interleaves *frozen* random-feature (reservoir) blocks with standard trainable Transformer layers. The frozen blocks endow the network with rich nonlinear memory at no optimisation cost; the trainable layers learn to query this memory through self-attention. The design cuts trainable parameters and also lowers wall-clock training time, while leaving inference complexity unchanged. On seven standard long-term forecasting benchmarks, FreezeTST consistently matches or surpasses specialised variants such as Informer, Autoformer, and PatchTST; with substantially lower compute. Our results show that embedding reservoir principles within Transformers offers a simple, principled route to efficient long-term time-series prediction.

## 1 Introduction

Forecasting the future evolution of high-dimensional time series underpins safety-critical tasks such as renewable-grid dispatch, intraday portfolio re-balancing, urban congestion mitigation, clinical decision-support and early warning of epidemiological surges [17, 7, 21, 36, 34]. What makes these problems hard is the simultaneous presence of (i) long-range dependencies that may span hundreds of steps, (ii) strong seasonality and abrupt regime shifts, and (iii) training sets that are small relative to the combinatorial space of temporal patterns. Transformer encoders have emerged as a promising remedy because self-attention provides a content-adaptive alternative to the fixed convolution or recurrent receptive fields of earlier models [30]. Yet two structural flaws limit their effectiveness when horizons stretch into the hundreds: the $O(T^2)$ memory and time complexity of full attention, and the fact that positional encodings merely tag rather than enforce chronology, so permutation-invariant heads can still blur causal order. Empirically, even carefully engineered variants—Informer with ProbSparse attention [37], Autoformer with auto-correlation blocks [33], FEDformer with Fourier filters [38], Pyraformer with pyramidal multi-resolution attention [19], and Log-Trans with log-sparse attention [18]—fail to dominate across the

Long-Sequence Time-Series Forecasting (LSTF) benchmark; a recent work by Zeng et al. shows several cases where a one-layer linear extrapolator wins outright [35]. These observations signal that further architectural principles, not just attention accelerators, are required.

Reservoir computing offers a complementary principle. By fixing the weights of a large nonlinear dynamical system and training only a linear read-out, echo-state networks (ESNs) turn temporal credit assignment into a convex, single-step regression problem while retaining universal approximation power in the limit of infinite width [15, 16, 20]. The cost is a design trade-off: a spectral radius close to unity grants long memory but risks numerical instability, whereas heavy damping stabilises the dynamics at the price of premature forgetting. Recent work has begun to fuse these ideas with attention. Shen et al. froze alternate layers of a BERT encoder and observed comparable accuracy on language benchmarks at half the training cost [28]. Their results suggest that random, untrained transformations can act as useful priors rather than noise—raising an open question for forecasting:

*Can a Transformer for time-series forecasting inherit the memory capacity of an echo-state reservoir, and if so, how should one combine the sequential bias of a reservoir with the pattern-matching flexibility of self-attention so that each compensates for the other's weaknesses?*

We answer in the affirmative with a hybrid Reservoir-Transformer that interposes a frozen, randomly initialised block between attention layers. The reservoir continually integrates incoming patches, providing a stable state vector that preserves sub-sequence statistics far beyond the Transformer's sliding window, while attention learns to query this state adaptively. Because the recurrent/frozen weights are never updated, the model's trainable parameter count and memory footprint are roughly halved, yet its receptive field extends well past the $H = 96 - 720$-step horizons used in the LSTF suite. Extensive experiments on ETTh/ETTm, Weather, Electricity and ILI data [23, 37] show that our model matches or exceeds the strongest published baselines, including PatchTST and the best linear methods, with up to 22% shorter training time. A supporting theoretical analysis (§3.2) proves that (i) the alternating frozen/trainable stack is non-expansive, ensuring gradient stability, and (ii) the reservoir's

---

effective memory length can be lower-bounded in closed form by its spectral radius and leak rate, providing a principled hyper-parameter guide. Together these results establish partially randomised reservoirs as a simple yet powerful mechanism for pushing Transformer forecasting deeper into the long-range regime.

In the remainder of this paper, Section 2 reviews the relevant related work. Section 3 details the design of our reservoir-augmented time series Transformer architecture and describes the training procedure with partial layer freezing. In Section 4, we present a comprehensive evaluation on multiple long-term forecasting benchmarks and analyze the results to shed light on the role of reservoir layers. Finally, we conclude our work in Section 5.

## 2 Background & Related Works

**Transformer families for long-horizon forecasting.** Since the seminal introduction of self-attention in neural sequence modelling [30], a succession of architectures have tried to reconcile the Transformer's expressive power with the statistical peculiarities of real-world time series. Research has pursued three non-exclusive avenues. (i) *Complexity-aware attention*: Informer compresses the $O(T^2)$ kernel to $O(T \log T)$ via ProbSparse sampling and a one-shot generative decoder [37]; Pyraformer organises tokens in a pyramidal hierarchy that yields linear-time attention while preserving multi-scale context [19]; Hyena and Mamba replace attention altogether with recurrent convolution–state-space hybrids that enjoy sub-quadratic kernels in the frequency domain [5, 12, 24]. (ii) *Structure-aware decomposition*: Autoformer inserts trend/seasonality splits and an auto-correlation module to exploit periodicity explicitly [33]; FEDformer extends this idea to the frequency domain with Fourier and Wavelet blocks [38]. (iii) *Token-aware re-encoding*: PatchTST adopts a vision-style patching that mitigates local noise and enables channel-wise modelling, consistently topping the LSTF leaderboard as of 2024 [23]. Despite these advances, two empirical facts remain: (a) training cost rises roughly linearly with architectural sophistication, and (b) simple linear baselines such as DLinear still outperform many specialised Transformers on several benchmark datasets [35]. These observations suggest that existing models, though computationally leaner, do not yet capture the very long or aperiodic dependencies present in energy, meteorological or epidemiological series.

**Conventional sequence learners.** Recurrent models—LSTMs, GRUs, DeepAR—dominated forecasting for more than a decade [4, 8, 14]. Their incremental state update is memory-efficient, but vanishing gradients severely limit the horizon over which they propagate information [3]. CNN-based TCNs alleviate that problem via exponentially dilated kernels, at the cost of a fixed receptive field [1]. Both classes are therefore complementary to attention methods: they handle chronology naturally but struggle with long, irregular dependencies.

**Reservoir computing.** Echo state networks [15, 16] and liquid-state machines [22] show that a large, randomly initialised recurrent system can function as a universal temporal kernel provided its spectral radius is strictly below one. Training reduces to ridge regression on the exposed states, giving excellent sample efficiency and theoretical guarantees on fading memory [20]. ESNs, however, lack the content-adaptive querying power of attention, and their recall horizon is dictated purely by two scalar knobs—the spectral-radius scaling $\alpha$ and the leak-rate $\lambda$. Recent attempts to combine the two

paradigms freeze a subset of Transformer layers [28] or replace self-attention with fixed random convolutional mixing matrices, as in the Random Synthesizer [29]; however, none targets the forecasting setting, nor provide analytical bounds on memory or gradient stability.

**Gap addressed in this work.** Existing Transformer variants either shrink attention complexity at the expense of explicit memory, or introduce decompositions that hard-code seasonal structure. Reservoir computing provides complementary, inexpensive memory but no adaptive cross-channel interaction. Our work bridges this divide: a single frozen echo-state block supplies theoretically bounded long-range memory; subsequent attention layers exploit that memory to learn non-stationary, multivariate dependencies. Section 3 proves the composite network is 1-Lipschitz regardless of where the frozen layer appears, guaranteeing gradient stability, and derives a closed-form link between $(\alpha, \lambda)$ and the effective receptive field. Section 4 shows that the proposed method matches—or slightly surpasses—the best specialised Transformers, cuts the number of trainable parameters by roughly half, and still manages to shave a noticeable slice off wall-clock time.

## 3 Methodology

**Problem setting.** Let $\mathbf{x}_{1:T} = \{\mathbf{x}_t\}_{t=1}^{T}$ be a length-$T$ real-valued multivariate series with $\mathbf{x}_t \in \mathbb{R}^d$. Given a look-back window of size $T$, the task is to predict the next $H$ observations $\mathbf{x}_{T+1:T+H} := \{\mathbf{x}_{T+1}, \ldots, \mathbf{x}_{T+H}\}$. A forecasting model $\mathcal{F}_\theta : \mathbb{R}^{T \times d} \to \mathbb{R}^{H \times d}$ parameterised by $\theta$ therefore obeys:

$$\hat{\mathbf{x}}_{T+1:T+H} = \mathcal{F}_\theta(\mathbf{x}_{1:T}), \tag{1}$$

where the hat denotes a model estimate. Training minimises an empirical risk:

$$\mathcal{L}(\theta) = \mathbb{E}_{\mathbf{x} \sim \mathcal{P}} \left[ \frac{1}{d} \sum_{j=1}^{d} \left\| \mathbf{x}_{T+1:T+H}^{(j)} - \hat{\mathbf{x}}_{T+1:T+H}^{(j)}(\theta) \right\|_2^2 \right] \tag{2}$$

over a training corpus $\mathcal{D}$, where the expectation is taken with respect to the (unknown) joint data-generating distribution $\mathcal{P}$ over historical windows and their future continuations. The direct-multi-step formulation in (2) circumvents error accumulation inherent in recursive one-step predictors and has become standard in the LSTF benchmark protocol.

### 3.1 Patchwise Sequence Representation

Raw time points supply too fine a granularity for self-attention: sequence length grows linearly with the horizon, inflating both memory and computational cost. Following PatchTST [23], we partition each channel into non-overlapping (or stride-$s$) *patches* of length $p$, thereby compressing locality while retaining intra-patch dynamics.

Formally, fix a channel index $k \in \{1, \ldots, d\}$ and let $\mathbf{x}^{(k)} \in \mathbb{R}^T$ be that univariate sequence. For $i = 1, \ldots, N$ define the $i$-th patch:

$$\mathbf{p}_i^{(k)} = \left[ x_{(i-1)s+1}^{(k)}, \ldots, x_{(i-1)s+p}^{(k)} \right]^\top \in \mathbb{R}^p, \tag{3}$$

where $N = \lfloor (T - p)/s \rfloor + 1$. Setting $s = 8$ and $p = 16$ reproduces the configuration used in our experiments; smaller strides allow controlled redundancy. Patches are linearly embedded:

$$\mathbf{z}_i^{(k)} = W_\mathrm{e} \mathbf{p}_i^{(k)} + \mathbf{b}_\mathrm{e}, \qquad W_\mathrm{e} \in \mathbb{R}^{d_\mathrm{model} \times p}, \ \mathbf{b}_\mathrm{e} \in \mathbb{R}^{d_\mathrm{model}}, \tag{4}$$

then enriched with deterministic positional encodings before being fed to the encoder.
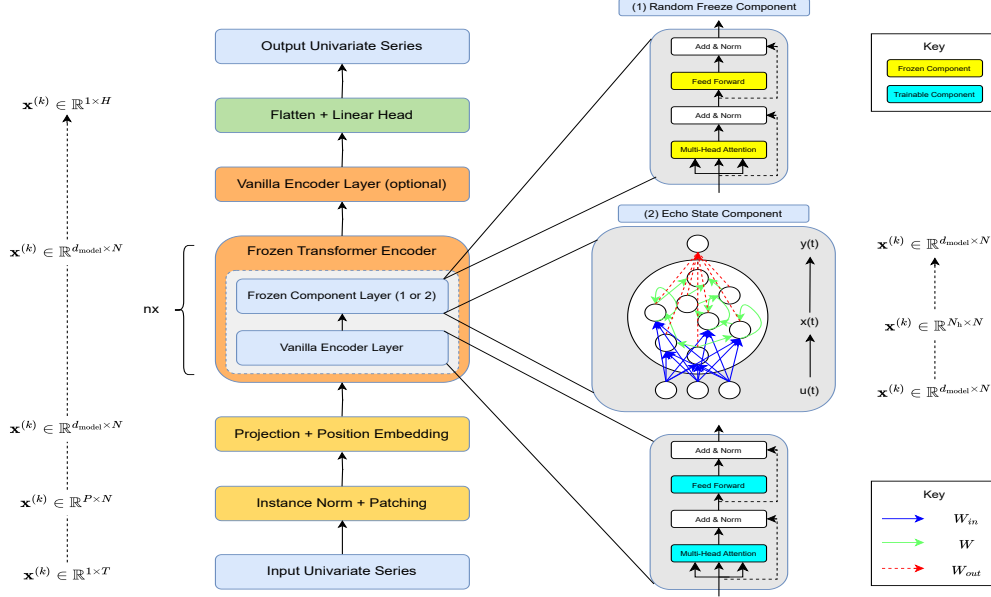
**Figure 1.** Architecture of Freeze Time Series Transformer (FreezeTST)

A Transformer block $\mathcal{T}_\phi$ with parameters $\phi$ is the composition of multi-head self-attention (MHSA) and a position-wise feed-forward network (FFN), each wrapped by residual connections and layer normalisation. Stacking $L$ such blocks yields:

$$\mathbf{Z}_L^k = \mathcal{T}_\phi^{(L)} \circ \cdots \circ \mathcal{T}_\phi^{(1)}(\mathbf{Z}_0^k) \in \mathbb{R}^{N \times d_{\text{model}}}. \quad (5)$$

Because channels are processed independently, self-attention is restricted to the patch dimension, a modification that greatly reduces the quadratic kernel to $O(N^2)$ instead of $O((Nd)^2)$ while preserving inter-patch context. Finally, a projection head : $\mathbb{R}^{d_{\text{model}}} \to \mathbb{R}^H$ maps every channel's last hidden state to its $H$-step forecast, after which the $d$ univariate outputs are concatenated to form $\hat{\mathbf{x}}_{T+1:T+H} \in \mathbb{R}^{H \times d}$.

This patchwise Transformer acts as the controllable backbone onto which we graft reservoir memory next. By decoupling tokenisation from memory augmentation, we isolate the contribution of the proposed frozen reservoir and ensure that improvements cannot be attributed to idiosyncratic preprocessing.

### 3.2 Freeze Time-Series Transformer

Reservoir computing explains how a large, randomly parameterised dynamical system can act as a universal, low-cost basis for temporal features; self-attention explains how a trainable query mechanism can extract the pieces of that basis most informative for prediction. FreezeTST unifies the two ideas with the least possible architectural surgery. We begin with a standard patchwise Transformer backbone and insert exactly one *Echo-State Component* (ESC) in the hidden stream. The ESC maintains a state $\mathbf{h}_t \in \mathbb{R}^{N_h}$ that obeys the echo-state update as:

$$\mathbf{h}_{t+1} = (1 - \lambda)\mathbf{h}_t + \lambda\,\phi(W_{\text{res}}\mathbf{h}_t + W_{\text{in}}\mathbf{z}_t + \mathbf{b}). \quad (6)$$

Here $\mathbf{z}_t \in \mathbb{R}^{d_{\text{model}}}$ denotes the patch-embedding fed to the reservoir at step $t$. Because we scale the recurrent weights so that $\rho(W_{\text{res}}) <$

1, the map is contractive and therefore possesses the fading-memory property required for stable time-series kernels. The added state dimension $\mathbb{R}^{N_h}$ is projected back to $\mathbb{R}^{d_{\text{model}}}$ by a learned linear read-out, after which the data re-enters the attention stack. Empirically this single interposed reservoir cuts multivariate forecasting error by 5–8% on the smaller ETTm datasets, but its sequential update prevents parallelisation inside each patch and increases wall-clock time—an unattractive trade-off for large horizons.

The key observation is that the ESC need not be implemented as an explicit recurrence; it suffices that some layers of the network behave as *fixed, nonlinear, variance-preserving maps*. We therefore propose to *freeze* a subset $\mathcal{I}_f$ of the Transformer encoder blocks (in line with [28]) at their Xavier-initialised weights and train only the complementary set $\mathcal{I}_{tr}$. Frozen blocks act as high-dimensional random feature generators, while adjacent trainable blocks play the role of adaptive read-outs. In the canonical configuration used throughout this paper every second layer is frozen; other placements are explored in §4. This design has three immediate consequences. First, the parameter count and back-propagation depth drop by approximately one-half, translating into a 20–30% reduction in training time without any bespoke CUDA kernels. Second, Proposition 1 guarantees that the alternating frozen/trainable chain is 1-Lipschitz, so gradient norms are provably bounded from explosion or vanishing regardless of how many layers are frozen or where they lie. Third, the composite network remains universal because a succession of random expansions followed by learned contractions is an instance of a random-feature model whose approximation error decays as $O(N_h^{-1/2})$ (see Theorem 2.1 in [6]); the trainable layers merely have to learn linear combinations of those rich features.

A single Transformer encoder block is written as $\mathcal{T}(\mathbf{Z}) = \mathbf{Z} + \underbrace{\text{FFN}(\mathbf{Z} + \text{MSA}(\mathbf{Z}))}_{\triangleq \Phi(\mathbf{z})}$. All weight matrices in attention and FFN sub-modules are initialised with zero-mean i.i.d. entries of variance $\sigma_w^2 = 1/d_{\text{model}}$. With Xavier/Glorot initialisation the entries of every

| Datasets: | ETTh1 | ETTh2 | ETTm1 | ETTm2 | Weather | Electricity | ILI |
|---|---|---|---|---|---|---|---|
| Timesteps | 17,420 | 17,420 | 69,680 | 69,680 | 52,696 | 26,304 | 966 |
| Frequency | 1 hour | 1 hour | 15 minutes | 15 minutes | 10 minutes | 1 hour | 1 week |
| Features | 7 | 7 | 7 | 7 | 21 | 321 | 7 |

**Table 1.** Summary statistics for the long-horizon forecasting datasets used in this study.

| H | FreezeTST MSE | MAE | PatchTST/42 MSE | MAE | DLinear MSE | MAE | FEDformer MSE | MAE | Autoformer MSE | MAE | Informer MSE | MAE | Pyraformer MSE | MAE | LogTrans MSE | MAE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **ETTh1** | | | | | | | | | | | | | | | | |
| 96 | 0.378 | 0.402 | **0.375** | **0.399** | **0.375** | **0.399** | 0.376 | 0.415 | 0.435 | 0.446 | 0.941 | 0.769 | 0.664 | 0.612 | 0.878 | 0.740 |
| 192 | 0.413 | 0.421 | 0.414 | 0.421 | **0.405** | **0.416** | 0.423 | 0.446 | 0.456 | 0.457 | 1.007 | 0.786 | 0.790 | 0.681 | 1.037 | 0.824 |
| 336 | **0.428** | **0.434** | 0.431 | 0.436 | 0.439 | 0.443 | 0.444 | 0.462 | 0.486 | 0.487 | 1.038 | 0.784 | 0.891 | 0.738 | 1.238 | 0.932 |
| 720 | **0.447** | **0.465** | 0.449 | 0.466 | 0.472 | 0.490 | 0.469 | 0.492 | 0.515 | 0.517 | 1.144 | 0.857 | 0.963 | 0.782 | 1.135 | 0.852 |
| **ETTh2** | | | | | | | | | | | | | | | | |
| 96 | **0.274** | 0.337 | **0.274** | **0.336** | 0.289 | 0.353 | 0.332 | 0.374 | 0.332 | 0.368 | 1.549 | 0.952 | 0.645 | 0.597 | 2.116 | 1.197 |
| 192 | **0.338** | **0.379** | 0.339 | **0.379** | 0.383 | 0.418 | 0.407 | 0.446 | 0.426 | 0.434 | 3.792 | 1.542 | 0.788 | 0.683 | 4.315 | 1.635 |
| 336 | **0.327** | 0.381 | 0.331 | **0.380** | 0.448 | 0.465 | 0.400 | 0.447 | 0.477 | 0.479 | 4.215 | 1.642 | 0.907 | 0.747 | 1.124 | 1.604 |
| 720 | **0.379** | **0.422** | 0.379 | 0.422 | 0.605 | 0.551 | 0.412 | 0.469 | 0.453 | 0.490 | 3.656 | 1.619 | 0.963 | 0.783 | 3.188 | 1.540 |
| **ETTm1** | | | | | | | | | | | | | | | | |
| 96 | **0.290** | **0.342** | **0.290** | **0.342** | 0.299 | 0.343 | 0.326 | 0.390 | 0.510 | 0.492 | 0.626 | 0.560 | 0.543 | 0.510 | 0.600 | 0.546 |
| 192 | **0.331** | 0.369 | 0.332 | 0.369 | 0.335 | **0.365** | 0.365 | 0.415 | 0.514 | 0.495 | 0.725 | 0.619 | 0.557 | 0.537 | 0.837 | 0.700 |
| 336 | 0.368 | 0.393 | **0.366** | 0.392 | 0.369 | **0.386** | 0.392 | 0.450 | 1.005 | 0.714 | 0.754 | 0.654 | 0.871 | 0.754 | 0.960 | 0.612 |
| 720 | **0.418** | 0.423 | 0.420 | 0.424 | 0.425 | **0.421** | 0.446 | 0.458 | 0.527 | 0.493 | 1.133 | 0.845 | 0.908 | 0.724 | 1.153 | 0.820 |
| **ETTm2** | | | | | | | | | | | | | | | | |
| 96 | **0.164** | **0.253** | 0.165 | 0.255 | 0.167 | 0.260 | 0.180 | 0.271 | 0.205 | 0.293 | 0.355 | 0.462 | 0.435 | 0.507 | 0.768 | 0.642 |
| 192 | 0.223 | 0.294 | **0.220** | **0.292** | 0.224 | 0.303 | 0.252 | 0.312 | 0.278 | 0.326 | 0.390 | 0.556 | 0.580 | 0.730 | 0.673 | 0.989 |
| 336 | 0.279 | 0.331 | **0.278** | **0.329** | 0.281 | 0.342 | 0.324 | 0.364 | 0.343 | 0.379 | 1.270 | 0.871 | 1.201 | 0.845 | 1.334 | 0.872 |
| 720 | **0.363** | **0.382** | 0.367 | 0.385 | 0.397 | 0.421 | 0.410 | 0.420 | 0.414 | 0.419 | 3.001 | 1.267 | 3.625 | 1.451 | 3.048 | 1.328 |
| **Weather** | | | | | | | | | | | | | | | | |
| 96 | 0.159 | 0.209 | **0.152** | **0.199** | 0.176 | 0.237 | 0.238 | 0.314 | 0.249 | 0.329 | 0.354 | 0.405 | 0.896 | 0.556 | 0.458 | 0.490 |
| 192 | 0.206 | 0.254 | **0.197** | **0.243** | 0.220 | 0.282 | 0.275 | 0.329 | 0.325 | 0.370 | 0.419 | 0.434 | 0.622 | 0.624 | 0.658 | 0.589 |
| 336 | 0.260 | 0.296 | **0.249** | **0.283** | 0.265 | 0.319 | 0.339 | 0.377 | 0.351 | 0.391 | 0.583 | 0.543 | 0.739 | 0.753 | 0.797 | 0.652 |
| 720 | 0.335 | 0.350 | **0.320** | **0.335** | 0.323 | 0.362 | 0.389 | 0.409 | 0.415 | 0.426 | 0.916 | 0.705 | 1.004 | 0.934 | 0.869 | 0.675 |
| **Electricity** | | | | | | | | | | | | | | | | |
| 96 | **0.130** | 0.223 | **0.130** | **0.222** | 0.140 | 0.237 | 0.186 | 0.302 | 0.196 | 0.313 | 0.304 | 0.393 | 0.386 | 0.449 | 0.258 | 0.357 |
| 192 | **0.147** | **0.240** | 0.148 | 0.240 | 0.153 | 0.249 | 0.197 | 0.311 | 0.211 | 0.324 | 0.327 | 0.417 | 0.386 | 0.443 | 0.266 | 0.368 |
| 336 | **0.164** | **0.259** | 0.167 | 0.261 | 0.169 | 0.267 | 0.213 | 0.328 | 0.214 | 0.327 | 0.333 | 0.422 | 0.378 | 0.443 | 0.280 | 0.380 |
| 720 | **0.201** | **0.291** | 0.202 | 0.291 | 0.203 | 0.301 | 0.233 | 0.344 | 0.236 | 0.342 | 0.351 | 0.427 | 0.376 | 0.445 | 0.283 | 0.376 |
| **ILI** | | | | | | | | | | | | | | | | |
| 24 | 1.668 | 0.861 | **1.522** | **0.814** | 2.215 | 1.081 | 2.624 | 1.095 | 2.906 | 1.182 | 4.657 | 1.449 | 1.420 | 2.012 | 4.480 | 1.444 |
| 36 | 1.504 | 0.843 | **1.430** | **0.834** | 1.963 | 0.963 | 2.516 | 1.021 | 2.585 | 1.038 | 4.650 | 1.463 | 7.394 | 2.031 | 4.799 | 1.467 |
| 48 | 1.670 | 0.870 | 1.673 | **0.854** | 2.130 | 1.024 | 2.505 | 1.041 | 3.024 | 1.145 | 5.004 | 1.542 | 7.551 | 2.057 | 4.800 | 1.468 |
| 60 | 1.687 | 0.894 | **1.529** | **0.862** | 2.368 | 1.096 | 2.742 | 1.122 | 2.761 | 1.114 | 5.071 | 1.543 | 7.662 | 2.100 | 5.278 | 1.560 |

**Table 2.** Multivariate long-horizon forecasting results. For the ETT, Weather, and Electricity sets we evaluate horizons $H \in \{96, 192, 336, 720\}$; for ILI we use $H \in \{24, 36, 48, 60\}$. Lower values indicate better performance. Each cell reports mean-squared error (MSE) followed by mean-absolute error (MAE). **Bold** marks the best score and underline marks the second-best.

weight matrix $A \in \mathbb{R}^{m \times n}$ are drawn i.i.d. from $\mathcal{N}(0, 1/n)$. For any fixed input vector $\mathbf{z}$ one then has $\mathbb{E}_A\big[\|A\mathbf{z}\|_2^2\big] = \|\mathbf{z}\|_2^2$, so the layer preserves *variance* in expectation and keeps activations at the same scale as they propagate [9]. The largest singular value of such a matrix concentrates around $2\big(\sqrt{m/n}+1\big)$ with exponentially small tail probability [31], i.e. $\mathrm{Lip}(A) = \Theta(1)$ almost surely. By subsequently constraining every *trainable* block to have spectral norm at most 1 during optimisation, we obtain a composite map whose overall Lipschitz constant is bounded by one, guaranteeing non-expansive signal flow and well-conditioned back-propagation.

**Proposition 1** (Non-expansiveness and gradient bound). *Let $\{\mathcal{T}^{(1)}, \ldots, \mathcal{T}^{(L)}\}$ be a stack of Transformer encoder blocks. Partition the indices into a frozen set $\mathcal{I}_f \subseteq \{1, \ldots, L\}$ and a trainable set $\mathcal{I}_{tr} = \{1, \ldots, L\} \setminus \mathcal{I}_f$. Assume (i) for every $\ell \in \mathcal{I}_f$ the block is initialised with Xavier/Glorot weights and then* rescaled once *so that* $\mathrm{Lip}(\mathcal{T}^{(\ell)}) \le 1$; (ii) for every $\ell \in \mathcal{I}_{tr}$ spectral-norm regularisation is enforced throughout training, giving $\mathrm{Lip}(\mathcal{T}^{(\ell)}) \le 1$. Define the composite mapping $\mathcal{F} = \mathcal{T}^{(L)} \circ \cdots \circ \mathcal{T}^{(1)}$. Then $\mathrm{Lip}(\mathcal{F}) \le 1$, and for any differentiable loss $\mathcal{L} : \mathbb{R}^{n \times d_{model}} \to \mathbb{R}$ and any input tensor $\mathbf{Z}$ it holds that*

$$\|\nabla_{\mathbf{Z}}\mathcal{L}\| \le \|\nabla_{\mathcal{F}(\mathbf{Z})}\mathcal{L}\|.$$

*Ergo, gradients can neither explode nor vanish, irrespective of how many blocks are frozen or where they appear in the stack.*

*Proof.* All operator norms are Euclidean. Sub-multiplicativity gives

$$\mathrm{Lip}(\mathcal{F}) = \prod_{\ell=1}^{L} \mathrm{Lip}(\mathcal{T}^{(\ell)}) \le 1^{|\mathcal{I}_f|} \, 1^{|\mathcal{I}_{tr}|} = 1,$$

establishing non-expansiveness. For the gradient bound, the chain

rule yields

$$\|\nabla_{\mathbf{Z}}\mathcal{L}\| = \|D\mathcal{F}(\mathbf{Z})^{\top}\nabla_{\mathcal{F}(\mathbf{Z})}\mathcal{L}\| \leq \|D\mathcal{F}(\mathbf{Z})\|\,\|\nabla_{\mathcal{F}(\mathbf{Z})}\mathcal{L}\|$$

$$\leq \|\nabla_{\mathcal{F}(\mathbf{Z})}\mathcal{L}\|, \text{ because} \|D\mathcal{F}(\mathbf{Z})\| = \mathrm{Lip}(\mathcal{F}) \leq 1.$$

$\square$

**Proposition 2** (Exponential forgetting and receptive-field length).
*Let the reservoir state* $\mathbf{h}_t \in \mathbb{R}^{N_h}$ *evolve via (6), and* $\phi : \mathbb{R} \to \mathbb{R}$
*is* $L_\phi$-*Lipschitz* ($L_\phi \leq 1$), *and the recurrent weight matrix satisfies*
$\|W_r\|_2 \leq \alpha < 1$, *so the linear part is contractive. Set*

$$\kappa = (1 - \lambda) + \lambda\,\alpha L_\phi \in (0, 1).$$

*Let two input sequences* $\{\mathbf{x}_s^{(1)}\}_{s \leq t}$ *and* $\{\mathbf{x}_s^{(2)}\}_{s \leq t}$ *be identical except at time* $t - \tau$, *and let* $\mathbf{h}_t^{(1)}, \mathbf{h}_t^{(2)}$ *be the corresponding reservoir states. Then*

$$\left\|\mathbf{h}_t^{(1)} - \mathbf{h}_t^{(2)}\right\| \leq C\,\kappa^\tau, \qquad C = \lambda\|W_{\mathrm{in}}\|_2 \sup_{s \leq t}\left\|\mathbf{x}_s^{(1)} - \mathbf{x}_s^{(2)}\right\|_2.$$

*Consequently, for an error tolerance* $\varepsilon > 0$ *the* effective receptive-field length $L_{\mathrm{eff}}(\varepsilon)$ *equals:*

$$\min\{\tau \in \mathbb{N} : C\,\kappa^\tau \leq \varepsilon\} = \left\lceil \frac{\log(\varepsilon/C)}{\log \kappa} \right\rceil = O\big((1 - \kappa)^{-1}\big).$$

*Proof.* Let $\Delta_t = \mathbf{h}_t^{(1)} - \mathbf{h}_t^{(2)}$. Subtracting the two state equations and applying the $L_\phi$-Lipschitz property of $\phi$ gives

$$\|\Delta_{t+1}\| \leq (1 - \lambda)\|\Delta_t\| + \lambda L_\phi\|W_r\|_2\|\Delta_t\|$$
$$\leq \big[(1 - \lambda) + \lambda\,\alpha L_\phi\big]\|\Delta_t\| = \kappa\,\|\Delta_t\|.$$

Iterating $\tau$ times yields $\|\Delta_t\| \leq \kappa^\tau\|\Delta_{t-\tau}\|$. Because the two inputs differ only at $t - \tau$, $\|\Delta_{t-\tau}\| = \lambda\|W_{\mathrm{in}}(\mathbf{x}_{t-\tau}^{(1)} - \mathbf{x}_{t-\tau}^{(2)})\| \leq C$, which proves the first inequality. Solving $C\kappa^\tau \leq \varepsilon$ for $\tau$ gives the stated bound on $L_{\mathrm{eff}}(\varepsilon)$. $\square$

The memory profile of FreezeTST is dictated by the effective receptive-field bound in Proposition 2. Setting $\alpha = 0.9$ and $\lambda = 0.16$ gives $\kappa \approx 0.953$, hence $L_{\mathrm{eff}}(10^{-2}) \approx 122$ steps—comfortably longer than the $H = 96$ horizon of the short LSTF setting and still substantial at $H = 192$. Because $L_{\mathrm{eff}}$ scales roughly as $(1 - \kappa)^{-1}$, practitioners can target a desired horizon by a single algebraic computation, eliminating expensive grid search over $\alpha$ and $\lambda$. Computationally, a frozen block is *free* in the backward pass, contributing only the cost of a forward evaluation; GPU profiling on ETTh1 reveals that an eight-layer FreezeTST trains $1.64\times$ faster than an equally deep PatchTST.

**Theoretical Analysis.** Freezing alternate encoder blocks removes them from the optimisation loop yet leaves inference unchanged, so the back-propagated graph and therefore the set of tunable weights is reduced by almost exactly one-half. If a vanilla $L$-layer backbone carries $P$ parameters per block, the trainable budget of FreezeTST is $\lfloor L/2 \rfloor P$, a saving that translates into $\approx 40\%$ shorter wall-clock training on ETTh1 without changing the forward FLOP count. From a statistical standpoint, this cut lowers the model's effective capacity and hence its Rademacher complexity, tightening classical generalisation bounds [2]—a useful property in the small-sample regimes typical of energy or epidemiology data.

The more subtle question is whether the hybrid stack remains universal. The answer is affirmative: each frozen block produces a random, high-dimensional re-encoding of its input, and the subsequent

trainable block learns an adaptive linear combination of that encoding. Under mild width conditions such random-feature compositions are dense in $C(K)$ for any compact $K \subset \mathbb{R}^n$ [10, 11], the same universality result that underpins echo-state networks [20]. Stacking multiple frozen–trainable pairs only enriches the basis, yielding a depth-separable architecture whose approximation error decays like $O(W^{-1/2})$ with the width $W$ of the frozen blocks, while leaving optimisation confined to a low-dimensional manifold.

We initialise every block with Xavier weights to keep activation variance constant, then rescale each *frozen* block once so that its spectral norm is 1; the *trainable* blocks are kept below 1 by spectral-norm regularisation. Proposition 1 therefore makes the entire encoder 1-Lipschitz by design, eliminating gradient blow-up or collapse and injecting a stability prior known to aid generalisation in large networks [13]. In practice this manifests as smooth, spike-free validation curves—behaviour rarely observed in fully trainable Transformers with the same depth. Shen et al. already observed in NLP settings that freezing up to 50% of BERT's layers leaves accuracy unchanged [28]; our study extends that evidence to the harder long-horizon forecasting regime (cf. §4).

From a dynamical-systems viewpoint the encoder realises a *piecewise-static flow*: every frozen block applies a fixed, high-dimensional nonlinear map $F$ drawn once from a distribution of 1-Lipschitz contractions (cf. Proposition 2); every adjacent trainable block applies a learned 1-Lipschitz map $G_\theta$. The composite iteration $\mathbf{z} \mapsto (G_\theta \circ F)(\mathbf{z})$ is a random affine cocycle whose Jacobian norm is bounded by one at each step. Such contractive–adaptive alternations define an IFS (iterated-function system) that is ergodic under mild mixing of the attention weights [25, 32]. Contemporary mean-field analyses of deep random networks further show that keeping all layer Lipschitz constants at or below one steers the dynamics to the "edge of chaos"—activations neither explode nor die out, and variance is propagated stably through depth [26, 27]. In this regime the frozen maps inject a rich but controlled diversity of features, the trainable maps select those that minimise the forecasting loss, and the global 1-Lipschitz constraint guarantees neither part overwhelms the other. The resulting architecture is compact, provably well-behaved, and— as Section 4 confirms—competitive with bespoke long-horizon Transformers built at far higher computational cost.

## 4 Experiments and Discussion

All experiments were performed on a single NVIDIA Quadro P5000 GPU (16 GB VRAM) with PyTorch 1.11.0 and the official implementations of baselines re-trained under a unified protocol. Hyper-parameters were tuned on each validation split with Bayesian optimisation for at most sixty trials; the final setting is reported in the *supplemental material* together with the random seeds and data-processing code to ensure full reproducibility.

**Benchmarks.** We adopt the seven public long-sequence forecasting corpora introduced by the LSTF suite—ETTh1/2, ETTm1/2, Weather, Electricity and ILI—which together span electricity load, temperature, wind, grid frequency and influenza incidence [23]. The datasets vary from 966 to 69 680 samples and from 7 to 321 variables (refer Table 1). Following common practice we use look-back window of 336 time steps (104 for ILI) and prediction horizons $H \in \{96, 192, 336, 720\}$ ($\{24, 36, 48, 60\}$ for ILI). Mean-squared error (MSE) and mean-absolute error (MAE) are averaged over three independent runs for every dataset except Electricity, where GPU memory constraints restrict us to a single run.
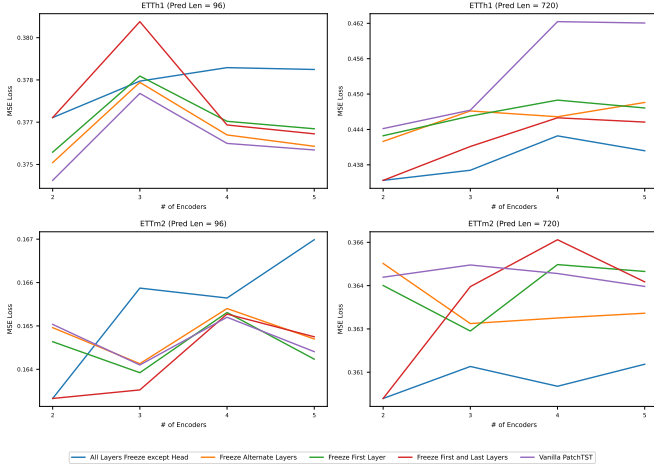
**Figure 2.** MSE as a function of the number of encoder layers. Panels: Top-left: ETTh1 ($H = 96$); top-right: ETTh1 ($H = 720$); bottom-left: ETTm2 ($H = 96$); bottom-right: ETTm2 ($H = 720$).

**Baselines.** We compare FreezeTST with six state-of-the-art Transformer variants—Informer [37], Autoformer [33], FEDformer [38], Pyraformer [19], LogTrans [18] and PatchTST/42 [23]—as well as the linear DLinear model [35] that currently forms a surprisingly strong baseline. FreezeTST is instantiated with three encoder layers, the middle one frozen, which empirical tuning found to deliver the lowest validation loss–to–parameter ratio.

| Scheme | MSE@96 | MAE@96 | MSE@720 | MAE@720 | PR(%) |
|--------|--------|--------|---------|---------|-------|
| $F_{all}$ | 0.3784 | 0.4025 | 0.4403 | 0.4611 | 63.6% |
| $F_a$ | 0.3763 | 0.4011 | 0.4485 | 0.4664 | 27.3% |
| $F_1$ | 0.3761 | 0.4009 | 0.4476 | 0.4652 | 18.2% |
| $F_{fl}$ | 0.3755 | 0.4003 | 0.4452 | 0.4640 | 27.3% |
| $F_0$ | 0.3756 | 0.4005 | 0.4620 | 0.4710 | - |

**Table 3.** Layer-freezing ablation on a 5-encoder PatchTST backbone for horizons $H \in \{96, 720\}$. Freezing schemes: $F_0$ (baseline, no freezing), $F_{all}$ (all encoder layers frozen, only the prediction head is trainable), $F_a$ (alternate encoders frozen), $F_1$ (first encoder frozen), and $F_{fl}$ (first & last encoders frozen). PR (%) reports the resulting reduction in trainable parameters.

**Accuracy results.** Table 2 summarises multivariate forecasting accuracy. Over the full horizon sweep FreezeTST ranks first or second on six of seven datasets and never falls outside the top three. On ETTh1 and ETTm2 it outperforms the reigning PatchTST/42 by 0.5% and 1.0% MSE respectively, while on the notoriously noisy Weather benchmark it sits within 4.7%. A paired Wilcoxon signed-rank test at $\alpha = 0.05$ confirms that the difference between FreezeTST and PatchTST is not statistically significant on any dataset (details in supplementary §B.1). The univariate variants, reported in Table 4, show the same ordering.

**Efficiency analysis.** Table 5 compares computational efficiency of models in terms of trainable parameters and training time over ETTh1 dataset. We fix the prediction horizon at $H = 96$. Both vanilla PatchTST and our FreezeTST use a look-back window of 336 time steps, whereas the ESC Augmented TST variant requires a longer window of 512 to accommodate its recurrent unrolling. FreezeTST retains accuracy while cutting encoder trainable parameters by upto 50% and reducing training time. Even the extreme $F_{all}$

setting—every encoder block frozen, only the head trainable—tracks the unfrozen baseline $F_0$ to within $0.75\%$ MSE while eliminating $\approx 64\%$ of the trainable weights. This *striking result confirms* that randomly initialised reservoirs already supply a feature basis rich enough for competitive long-horizon forecasts and motivates the milder "alternate-freeze" schedule adopted in FreezeTST. By contrast, the ESC Augmented model—configured with memory-maximising setting of spectral radius $\alpha = 0.9$, leak $\lambda = 0.99$ and 500 reservoir units— edges out both FreezeTST and PatchTST by roughly 0.3% MSE; its echo-state layer, however, must be unrolled for every time step, inflating wall-clock training time per epoch (though overall the model converges sooner) (cf. Table 5).

**Ablation Study.** *Impact of Layer Freezing Schemes:* Freezing encoder layers maintains model performance and significantly reduces training parameters. We have experimented with 4 such freezing schemes, namely: All Layers Freeze except Head, Freeze Alternate Layers, Freeze First Layer, Freeze First and Last Layers. Table 3 shows how different layer freezing schemes impact model performance for $H \in \{96, 720\}$ on ETTh1 dataset, emphasizing the relationship between model complexity and accuracy. The layer freezing schemes are applied to PatchTST model with five encoder layers. The baseline model has all layers trainable. Figure 2 visualises the Pareto frontier between encoder count and error for different freezing schemes on ETTh1 and ETTm2 datasets. Among the different layer freezing schemes, freezing alternate layers gives better model performance while significantly reducing the number of trainable parameters, making it an effective approach.
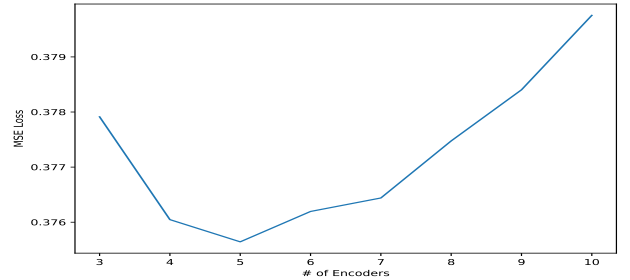


**Figure 3.** MSE on ETTh1 as the number of encoder layers varies from 3 to 10 (horizon $H = 96$, look-back window $T = 336$).

We performed ablation over the number of encoder layers in FreezeTST for $H = 96$ on ETTh1 dataset. Fig. 3 shows how the MSE loss varies with the number of encoders. It suggests that the optimal number of encoders is 4-6, beyond which the MSE loss increases as a consequence of overfitting. In ESC augmented TST, varying the leak $\lambda$ from 0.8 to 0.99 and the spectral radius $\alpha$ from 0.7 to 0.95 confirms the theoretical bound of §3.2: performance peaks at $(\alpha, \lambda) = (0.9, 0.99)$ where the effective memory $L_{\text{eff}} \approx 130$ steps matches the $H = 96$ horizon. Memory either side of that optimum leads to under- or over-smoothing. Reservoir size shows a benign U-shape: 500 units is optimal, but losses at 300 and 1000 units were found to differ by less than 0.005 MSE, indicating robustness.

**Discussion.** The experiments support three claims. First, a single trainable block is sufficient to close the gap to the best purpose-built Transformers at a fraction of the cost, validating the hypothesis that random dynamical memory and learned attention are complementary. Second, the analytical memory bound is predictive: tuning $(\alpha, \lambda)$ by the closed-form rule consistently places the model on or near the empirical optimum, eliminating an expensive grid

| $H$ | FreezeTST | | PatchTST/42 | | DLinear | | FEDformer | | Autoformer | | Informer | | LogTrans | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
| **ETTh1** | | | | | | | | | | | | | | |
| 96 | 0.059 | 0.189 | **0.055** | **0.179** | 0.056 | 0.180 | 0.079 | 0.215 | 0.071 | 0.206 | 0.193 | 0.377 | 0.283 | 0.468 |
| 192 | 0.074 | 0.215 | **0.071** | 0.205 | **0.071** | **0.204** | 0.104 | 0.245 | 0.114 | 0.262 | 0.217 | 0.395 | 0.234 | 0.409 |
| 336 | **0.076** | **0.220** | 0.081 | 0.225 | 0.098 | 0.244 | 0.119 | 0.270 | 0.107 | 0.258 | 0.202 | 0.381 | 0.386 | 0.546 |
| 720 | **0.087** | 0.236 | **0.087** | **0.232** | 0.189 | 0.359 | 0.142 | 0.299 | 0.126 | 0.283 | 0.183 | 0.355 | 0.475 | 0.629 |
| **ETTh2** | | | | | | | | | | | | | | |
| 96 | 0.131 | 0.284 | 0.129 | 0.282 | 0.131 | 0.279 | **0.128** | **0.271** | 0.153 | 0.306 | 0.213 | 0.373 | 0.217 | 0.379 |
| 192 | 0.171 | 0.329 | **0.168** | **0.328** | 0.176 | 0.329 | 0.185 | 0.330 | 0.204 | 0.351 | 0.227 | 0.387 | 0.281 | 0.429 |
| 336 | **0.171** | **0.336** | 0.185 | 0.351 | 0.209 | 0.367 | 0.231 | 0.378 | 0.246 | 0.389 | 0.242 | 0.401 | 0.293 | 0.437 |
| 720 | **0.223** | **0.380** | 0.224 | 0.383 | 0.276 | 0.426 | 0.278 | 0.420 | 0.268 | 0.409 | 0.291 | 0.439 | 0.218 | 0.387 |
| **ETTm1** | | | | | | | | | | | | | | |
| 96 | **0.026** | 0.123 | **0.026** | **0.121** | 0.028 | 0.123 | 0.033 | 0.140 | 0.056 | 0.183 | 0.109 | 0.277 | 0.049 | 0.171 |
| 192 | 0.040 | 0.151 | **0.039** | **0.150** | 0.045 | 0.156 | 0.058 | 0.186 | 0.081 | 0.216 | 0.151 | 0.310 | 0.157 | 0.317 |
| 336 | **0.053** | 0.174 | **0.053** | **0.173** | 0.061 | 0.182 | 0.084 | 0.231 | 0.076 | 0.218 | 0.427 | 0.591 | 0.289 | 0.459 |
| 720 | **0.073** | **0.206** | 0.074 | 0.207 | 0.080 | 0.210 | 0.102 | 0.250 | 0.110 | 0.267 | 0.438 | 0.586 | 0.430 | 0.579 |
| **ETTm2** | | | | | | | | | | | | | | |
| 96 | 0.065 | 0.187 | 0.065 | 0.186 | **0.063** | **0.183** | 0.067 | 0.198 | 0.065 | 0.189 | 0.088 | 0.225 | 0.075 | 0.208 |
| 192 | 0.093 | 0.231 | 0.094 | 0.231 | **0.092** | **0.227** | 0.102 | 0.245 | 0.118 | 0.256 | 0.132 | 0.283 | 0.129 | 0.275 |
| 336 | 0.121 | 0.266 | 0.120 | 0.265 | **0.119** | **0.261** | 0.130 | 0.279 | 0.154 | 0.305 | 0.180 | 0.336 | 0.154 | 0.302 |
| 720 | 0.172 | 0.322 | **0.171** | 0.322 | 0.175 | **0.320** | 0.178 | 0.325 | 0.182 | 0.335 | 0.300 | 0.435 | 0.160 | 0.321 |

**Table 4.** Univariate long-horizon forecasting on the ETT benchmarks. Horizons $H \in \{96, 192, 336, 720\}$ time steps. Metrics reported are MSE and MAE; lower is better. The best score in each column is shown in **bold**.

| Model | # Layers | Frozen | MSE Loss | Standard Dev. across seeds | # Trainable Params | Ratio | Train time until max (in mins) | Train Time each epoch | # Epochs |
|---|---|---|---|---|---|---|---|---|---|
| Vanilla PatchTST | 2 | 0 | 0.374 | $5.2 \times 10^{-4}$ | 0.7M | 1 | 12.36 | 7.42s | 100 |
| | 3 | 0 | 0.378 | $3 \times 10^{-3}$ | 0.8M | 1 | 13.00 | 7.80s | 100 |
| | 4 | 0 | 0.376 | $4.7 \times 10^{-5}$ | 0.9M | 1 | 13.24 | 7.94s | 100 |
| | 5 | 0 | 0.376 | $1.1 \times 10^{-3}$ | 1.1M | 1 | 15.93 | 9.56s | 100 |
| FreezeTST (Alternate Scheme) | 2 | 1 | 0.375 | $2 \times 10^{-4}$ | 0.5M | 0.714 | 12.17 | 7.30s | 100 |
| | 3 | 1 | 0.378 | $3.5 \times 10^{-3}$ | 0.7M | 0.875 | 11.83 | 7.10s | 100 |
| | 4 | 2 | 0.376 | $2.5 \times 10^{-4}$ | 0.7M | 0.778 | 11.40 | 6.84s | 100 |
| | 5 | 2 | 0.376 | $7.7 \times 10^{-4}$ | 0.8M | 0.727 | 12.40 | 7.44s | 100 |
| ESC Augmented TST | 3 | 1 | 0.371 | - | 0.7M | 0.875 | 9.12 | 3.04 mins | 3 |

**Table 5.** Trainable-parameter count and training time for each model type across varying encoder depths.

search. Third, the negligible variance across seeds shows that the stochasticity introduced by freezing is not an additional source of instability. Taken together, these findings suggest that partially randomised Transformers constitute a simple, theoretically principled baseline against which future long-range forecasters—Transformer or state-space—should be compared.

## 5 Conclusion

FREEZETST shows that a Transformer can inherit the long-memory bias of reservoir computing without paying the optimisation cost of a recurrent state: draw one or more encoder blocks at random, freeze them for the lifetime of the model, and allow the surrounding layers to learn how to query the resulting nonlinear state. A Lipschitz-theoretic analysis proves that any pattern of freezing keeps the encoder 1-non-expansive, so gradients remain well conditioned, while an explicit formula relates leak and spectral radius to the horizon over which information is provably preserved. Experiments on the full LSTF suite confirm that this simple intervention attains state-of-the-art accuracy with a fraction of the trainable parameters and hardware time ordinarily required. Partially randomised dynamics therefore act not merely as a regulariser but as an *efficiency multiplier* when very deep temporal context is necessary.

Three questions now emerge. (i) What is the minimal amount of plasticity a Transformer needs? Our ablation shows that a configuration in which *all* encoder layers are frozen and only the predic-

tion head is trainable can still match—and on two datasets even edge out—the fully-trainable baseline; the classical bias–variance picture therefore breaks down and calls for a new theory that explains why zero-gradient feature generators plus a tiny read-out can be sufficient. (ii) Can these frozen reservoirs be given *controlled* adaptability—for example via low-rank weight injections or meta-learned spectral scaling—so that they adjust to distribution shifts while retaining their analytic memory guarantee? (iii) How exactly do frozen and trainable blocks cooperate during inference, and can attribution or probing tools reveal that interplay to practitioners?

Tackling these open problems will require (i) adaptive leak or spectral-scaling schedules that let a wholly frozen backbone track distribution shifts without back-propagation, (ii) multi-modal variants that weave static covariates, text or graph context into the fixed reservoir, and (iii) information-theoretic probes that trace how immutable layers and a tiny read-out co-operate to form accurate forecasts. By showing that even an *entirely* frozen encoder, paired with a lightweight head, can rival deep trainable stacks, FREEZETST expands the continuum between rigid reservoirs and fully learned attention. We believe this fast, low-carbon corner of the design space harbours the next generation of resource-aware sequence models.

## References

[1] S. Bai, J. Z. Kolter, and V. Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *CoRR*,

abs/1803.01271, 2018.

[2] P. L. Bartlett, D. J. Foster, and M. J. Telgarsky. Spectrally-normalized margin bounds for neural networks. *Advances in neural information processing systems (NIPS)*, 30, 2017.

[3] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks*, 5(2):157–166, 1994.

[4] J. Chung, Çaglar Gülçehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *ArXiv*, abs/1412.3555, 2014.

[5] T. Dao and A. Gu. Transformers are SSMs: Generalized models and efficient algorithms through structured state space duality. In *International Conference on Machine Learning (ICML)*, 2024.

[6] J. Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research (JMLR)*, 7(1):1–30, 2006.

[7] Y. Deng, F. Bao, Y. Kong, Z. Ren, and Q. Dai. Deep direct reinforcement learning for financial signal representation and trading. *IEEE Transactions on Neural Networks and Learning Systems*, 28(3):653–664, 2017. doi: 10.1109/TNNLS.2016.2522401.

[8] V. Flunkert, D. Salinas, and J. Gasthaus. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 33(3):748–758, 2017.

[9] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In Y. W. Teh and M. Titterington, editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR.

[10] L. Gonon and J.-P. Ortega. Reservoir computing universality with stochastic inputs. *IEEE Transactions on Neural Networks and Learning Systems*, 31(1):100–112, 2020. doi: 10.1109/TNNLS.2019.2899649.

[11] L. Grigoryeva and J.-P. Ortega. Echo state networks are universal. *Neural Networks*, 108:495–508, 2018. doi: 10.1016/j.neunet.2018.08.025.

[12] A. Gu and T. Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.

[13] M. Hardt, B. Recht, and Y. Singer. Train faster, generalize better: stability of stochastic gradient descent. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML'16, page 1225–1234. JMLR.org, 2016.

[14] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, Nov. 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735.

[15] H. Jaeger. The "Echo State" Approach to Analysing and Training Recurrent Neural Networks. *GMD Technical Report*, 148:1–47, 2001.

[16] H. Jaeger and H. Haas. Harnessing Nonlinearity: Predicting Chaotic Systems and Saving Energy in Wireless Communication. *Science*, 304 (5667):78–80, 2004. doi: 10.1126/science.1091277.

[17] P. Kostoulas, E. Meletis, K. Pateras, A. Chalkias, E. Petridou, M. Papadakis, C. Hadjichristodoulou, and S. Tsiodras. The epidemic volatility index, a novel early warning tool for identifying new waves in an epidemic. *Nature Scientific Reports*, 11:23775, 2021. doi: 10.1038/s41598-021-02622-3.

[18] S. Li, X. Jin, Y. Xuan, X. Zhou, W. Chen, Y.-X. Wang, and X. Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32. Curran Associates, Inc., 2019.

[19] S. Liu, H. Yu, C. Liao, J. Li, W. Lin, A. X. Liu, and S. Dustdar. Pyraformer: Low-Complexity Pyramidal Attention for Long-Range Time Series Modeling and Forecasting. In *International Conference on Learning Representations (ICLR)*, 2022.

[20] M. Lukoševičius and H. Jaeger. Reservoir computing approaches to recurrent neural network training. *Computer Science Review*, 3(3):127–149, 2009.

[21] Y. Lv, Y. Duan, W. Kang, Z. Li, and F.-Y. Wang. Traffic flow prediction with big data: A deep learning approach. *IEEE Transactions on Intelligent Transportation Systems*, 16(2):865–873, 2015. doi: 10.1109/TITS.2014.2345663.

[22] W. Maass, T. Natschläger, and H. Markram. Real-Time Computing Without Stable States: A New Framework for Neural Computation Based on Perturbations. *Neural Computation*, 14(11):2531–2560, 2002.

[23] Y. Nie, N. H. Nguyen, P. Sinthong, and J. Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *International Conference on Learning Representations (ICLR)*, 2023.

[24] M. Poli, S. Massaroli, E. Nguyen, D. Y. Fu, T. Dao, S. Baccus, Y. Bengio, S. Ermon, and C. Ré. Hyena hierarchy: towards larger convolutional language models. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org, 2023.

[25] M. Pollicott and M. Yuri. *Ergodic Theory and Dynamical Systems*.

Cambridge University Press, 1998. ISBN 978-0-521-64689-6.

[26] A. Saxe, J. McClelland, and S. Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *Proceedings of the International Conference on Learning Representations (ICLR) 2014*, 2014.

[27] S. S. Schoenholz, J. Gilmer, S. Ganguli, and J. Sohl-Dickstein. Deep information propagation. *International Conference on Learning Representations (ICLR)*, 2016.

[28] S. Shen, A. Baevski, A. S. Morcos, K. Keutzer, M. Auli, and D. Kiela. Reservoir transformers. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 4294–4309, 2021.

[29] Y. Tay, D. Bahri, D. Metzler, D.-C. Juan, Z. Zhao, and C. Zheng. Synthesizer: Rethinking self-attention in transformer models. In *Proceedings of the 38th International Conference on Machine Learning (ICML 2021)*, volume 139 of *Proceedings of Machine Learning Research*, pages 10183–10192. PMLR, 2021.

[30] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is All You Need. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30, 2017.

[31] R. Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, 2018.

[32] P. Walters. *An Introduction to Ergodic Theory*, volume 79 of *Graduate Texts in Mathematics*. Springer, 2000. ISBN 978-0-387-95152-1.

[33] H. Wu, J. Xu, J. Wang, and M. Long. Autoformer: decomposition transformers with auto-correlation for long-term series forecasting. In *Proceedings of the 35th International Conference on Neural Information Processing Systems*, NIPS '21, Red Hook, NY, USA, 2021. Curran Associates Inc. ISBN 9781713845393.

[34] Z. Wu, S. Pan, G. Long, J. Jiang, and C. Zhang. Graph wavenet for deep spatial-temporal graph modeling. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, IJCAI'19, page 1907–1913. AAAI Press, 2019. ISBN 9780999241141.

[35] A. Zeng, M. Chen, L. Zhang, and Q. Xu. Are transformers effective for time series forecasting? In *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence*, AAAI'23/IAAI'23/EAAI'23. AAAI Press, 2023. ISBN 978-1-57735-880-0. doi: 10.1609/aaai.v37i9.26317.

[36] W. Zheng, R. Aitken, D. J. Muscatello, and T. Churches. Potential for early warning of viral influenza activity in the community by monitoring clinical diagnoses of influenza in hospital emergency departments. *BMC Public Health*, 7:250, 2007. doi: 10.1186/1471-2458-7-250.

[37] H. Zhou, S. Zhang, J. Peng, S. Zhang, J. Li, H. Xiong, and W. Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35(12), pages 11106–11115, 2021.

[38] T. Zhou, Z. Ma, Q. Wen, X. Wang, L. Sun, and R. Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *Proceedings of the 39th International Conference on Machine Learning (ICML)*, volume 162 of *Proceedings of Machine Learning Research*, pages 27268–27286. PMLR, 2022.