

Time-History Wind Response Analysis of High-Rise Buildings Using Latent Space Time-Step Operators

Chan Ho Kim^a, Thomas Kang^b

^a *Interdisciplinary Program in Artificial Intelligence, Seoul National University, Seoul, Korea, tntndi001@snu.ac.kr*

^b *Department of Architecture and Architectural Engineering & Interdisciplinary Program in Artificial Intelligence, Seoul National University, Seoul, Korea, tkang@snu.ac.kr*

SUMMARY

One-shot neural operators are strong baselines for fixed-grid structural wind-response prediction, but they require the full input history and do not support causal step-by-step rollout. This paper studies LSR-AR, a causal latent-step autoregressive surrogate based on a learned residual state update conditioned on the current wind forcing. On the 37-story, 997-step along-wind benchmark of Goswami et al. (2025), LSR-AR attains 6.4% relative L^2 error (RelL2) with 236K parameters. As the primary causal comparison, LSR-AR reduces the test RelL2 from 16.8% (NARX-MLP, retrained) to 6.4%. In a separate continuation study using a time-embedding LSR-AR variant, stable beyond-horizon rollout is observed from sufficiently long prefixes.

Keywords: *Aerodynamic response, Autoregressive surrogate, Causal rollout, Temporal extrapolation, Latent-space dynamics, High-rise buildings*

1. MOTIVATION

Neural operators map wind-load histories to full response trajectories in a single forward pass, but they cannot be queried causally, step by step, as new loading data arrive. Causal rollout is relevant to applications such as active damper control, real-time monitoring, and temporal extrapolation beyond the training horizon. In wind-engineering practice, different design codes and wind-environment characteristics require time-series analyses of varying length, yet most neural-operator surrogates are tied to a fixed temporal grid and can only infer on sequences of the length seen during training. A new model must therefore be trained whenever the required duration changes, and an already-trained model cannot be reused even in the same wind environment to examine longer-duration behavior. Autoregressive surrogates, in principle, lift this restriction, but they face compounding prediction errors over hundreds of sequential steps. In addition, gradient vanishing and increased training cost for long time-series data make causal models more difficult to train than one-shot models.

This paper presents Latent-Step Residual Autoregression (LSR-AR), a latent-step autoregressive surrogate that achieves 6.4% RelL2 on a 37-story, 997-step along-wind benchmark with only 236K parameters. The manuscript reports a primary full-horizon causal benchmark together with a separate secondary continuation study.

35 2. RELATED WORK

36 2.1. One-shot neural operators

37 One-shot operator-learning surrogates map an observed loading history to the response trajectory
38 in a single non-autoregressive pass. Within this class, Fourier Neural Operators (FNO) provide a
39 widely used spectral baseline, while Deep Operator Networks (DeepONet) learns operator
40 mappings through a branch-trunk decomposition. Recent wind-response work has also adapted
41 this paradigm through Self-Adaptive FNO (SA-FNO) (Goswami et al., 2025). In the wind
42 benchmark considered here, these full-input baselines require the available forcing history and
43 therefore do not natively support causal rollout beyond it.

44 2.2. Autoregressive surrogates

45 Autoregressive models enable causal rollout but face error accumulation over long horizons. A
46 common practical formulation is window-based autoregression, in which the next response is
47 predicted from a finite history of past responses and exogenous inputs. This finite-window view
48 underlies nonlinear autoregressive models with exogenous inputs (NARX) and provides a simple
49 strictly causal baseline, but the flattened input dimension grows with both window length and
50 response dimension. The NARX-MLP baseline considered here follows this formulation through
51 a windowed multilayer perceptron (MLP) adaptation inspired by Thedy et al. (2025). By
52 contrast, the present latent residual rollout reuses a shared latent update map at every step,
53 maintaining a fixed latent-state dimension during long-horizon rollout.

54 3. METHODOLOGY

55 3.1. Dataset

56 The study uses the wind force–response benchmark of Goswami et al. (2025). The target
57 structure is a 37-story steel frame, and the wind loading is synthesized from the Tokyo
58 Polytechnic University (TPU) aerodynamic database (Tamura, 2012). Displacement responses
59 are obtained by nonlinear time-history analysis in OpenSees; full structural and simulation
60 details are given in Goswami et al. (2025). Table 1 summarizes the benchmark configuration
61 used throughout the paper. Figure 1 shows the target structure and a representative forcing–
62 response pair from the dataset. Figure 1 (a) depicts the 37-story steel frame, and Figure 1 (b)
63 shows selected-floor forcing and along-wind displacement time histories from a representative
64 test scenario.

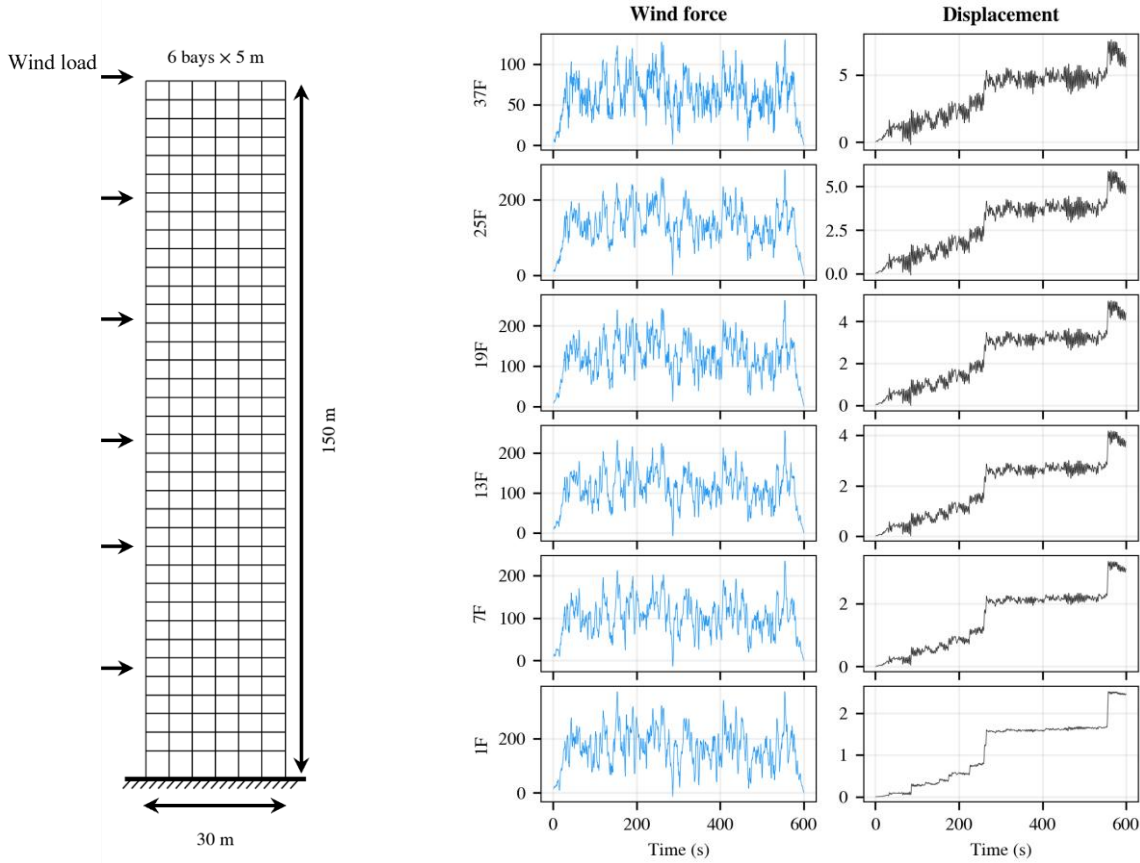
65 **Table 1:** Benchmark summary of the 37-story along-wind response dataset

Item	Value
Structure	37-story steel frame; 150 m height; 60 m × 30 m plan
Wind condition	Along-wind loading; $\alpha = 90^\circ$; reference wind speed $V_H = 65$ m/s
Load synthesis	Spectral proper orthogonal decomposition (POD) from TPU aerodynamic pressure data
Time discretization	Raw: 600 s at $\Delta t_{\text{raw}} = 0.02$ s ; benchmark: crop first 100 steps and stride by 30, giving effective $\Delta t = 0.6$ s and $T = 997$
Data split	720 train / 80 validation / 200 test
Input / output	37 floor-level wind forcing channels → 37 along-wind displacement channels

66 The task is to learn a causal input-output map from trajectory data alone:

$$\mathcal{D} = \left\{ \left(u_{0:T}^{(j)}, x_{0:T}^{(j)} \right) \right\}_{j=1}^N \quad (1)$$

67 where $u_k \in \mathbb{R}^{37}$ is the floor-level wind forcing at discrete time k and $x_k \in \mathbb{R}^{37}$ is the along-wind
 68 displacement response



69 **Figure 1:** (a) Target structure: 37-story steel frame (150 m, 6 bays). (b) Representative test scenario: wind-force
 70 (kN) (left) and along-wind displacement (m) (right) time histories for selected floors.

71 The raw simulation output spans 600 s at $\Delta t_{\text{raw}} = 0.02$ s (30,000 steps). For the benchmark, the
 72 first 100 steps are discarded and every 30th sample is retained, yielding an effective $\Delta t = 0.6$ s
 73 and a sequence length of $T = 997$ steps per scenario. This matches the original dataset setting of
 74 Goswami et al. (2025) and is retained for valid comparison. The benchmark split is 720 training /
 75 80 validation / 200 test. The 200-scenario test split matches the published test cases, while the
 76 80-scenario validation split is held out from the remaining training cases. Checkpoint selection
 77 uses validation metrics only; the test split is reserved for final reporting.

78 3.2. Experimental setup

79 3.2.1. Baseline models

80 Three comparison architectures are considered:

- 81 • **FNO** : one-shot Fourier Neural Operator (7.40M params).
- 82 • **SA-FNO** : Self-Adaptive FNO with adversarial frequency adjustment.
- 83 • **NARX-MLP**: adapted windowed NARX-MLP autoregressive baseline, inspired by the
- 84 autoregressive formulation of Thedy et al. (2025). At each step, the next response is
- 85 predicted from a sliding window of past responses ($n_{wr} = 50$) and forcing ($n_{wg} = 5$).

86 FNO and SA-FNO are one-shot operators; NARX-MLP is the only causal baseline in the present
87 comparison. FNO and SA-FNO were evaluated using the original trained checkpoints published
88 by Goswami et al. (2025) on the identical 200-scenario test set; no retraining was performed.
89 Their reference metrics were recomputed from the released predictions using the definitions in
90 this section. For roof-response diagnostics only, we additionally compare against the DeepONet
91 and DeepFNONet outputs provided with the Goswami dataset release. NARX-MLP was
92 retrained on the same split. At rollout start, unavailable response and forcing history in NARX-
93 MLP are zero-padded, while models that require an observed initial state are conditioned on the
94 first displacement sample x_0 of each scenario.

95 3.2.2. Evaluation metrics

96 All headline ReL2 values are reported as the aggregate/global relative L2 error over the full
97 evaluated set:

$$\text{ReL2} = \frac{1}{N} \sum_{i=1}^N \frac{\|\hat{x}^{(i)} - x^{(i)}\|_2}{\|x^{(i)}\|_2} \quad (2)$$

98 where, for each sample, the norm spans all time steps and output channels.

99 Pearson correlation is computed over the flattened prediction–target vector across all evaluated
100 samples, time steps, and floors:

$$\text{Corr} = \frac{\sum_{m=1}^M (y_m - \bar{y})(\hat{y}_m - \bar{\hat{y}})}{\sqrt{\sum_{m=1}^M (y_m - \bar{y})^2} \sqrt{\sum_{m=1}^M (\hat{y}_m - \bar{\hat{y}})^2}} \quad (3)$$

101 where m indexes the flattened sample–time–floor entries. Peak displacement error is the relative
102 error of per-sample, per-floor peak absolute displacement, averaged over all samples and floors:

$$\text{PeakErr} = \frac{1}{NC} \sum_{i=1}^N \sum_{c=1}^C \frac{|p_c^{(i)} - \hat{p}_c^{(i)}|}{p_c^{(i)} + \varepsilon}, \quad p_c^{(i)} = \max_t |x_{t,c}^{(i)}|, \quad \hat{p}_c^{(i)} = \max_t |\hat{x}_{t,c}^{(i)}| \quad (4)$$

103 where C is the number of output channels and ε is a small constant. In addition, *Shape-Based*
104 *Distance* (SBD), suggested by Paparrizos and Gravano (2015), is used here as a coarse global

105 diagnostic and computed over the flattened and z-normalized signal across the full evaluated set.
 106 SBD measures waveform similarity through the normalized cross-correlation after z-normalization:

$$\text{SBD}(x, \hat{x}) = 1 - \max_s \text{NCC}_s(\tilde{x}, \tilde{\hat{x}}) \quad (5)$$

107 where x denotes the z-normalized signal and s is the lag index. SBD is invariant to vertical offset,
 108 amplitude scaling, and temporal shift; it ranges from 0 (identical shape) to 2 (maximally dissimilar).
 109 While ReL2 penalizes all pointwise deviations equally, SBD isolates shape fidelity from
 110 amplitude and phase errors, providing a complementary diagnostic of waveform quality.

111 3.2.3. Benchmark tasks

112 Causal models are evaluated on two tasks:

- 113 • **Full-horizon prediction.** The model is trained and evaluated on the complete 997-step
 114 sequence. This is the primary comparison task.
- 115 • **Temporal extrapolation.** The model is trained on a truncated prefix ($L < 997$ steps) and
 116 rolled out over the full 997-step test sequence. The in-distribution region ($k \leq L$) and
 117 extrapolation region ($k > L$) are scored separately, testing whether the learned dynamics
 118 generalize beyond the training horizon.

119 One-shot operators are evaluated on the full-horizon task only, as they cannot extrapolate beyond
 120 their fixed input length by construction.

121 3.2.4. Proposed model: LSR-AR

122 LSR-AR is a latent-step autoregressive surrogate that predicts the next-step structural response via
 123 autoregression in a learned latent space (Figure 2).

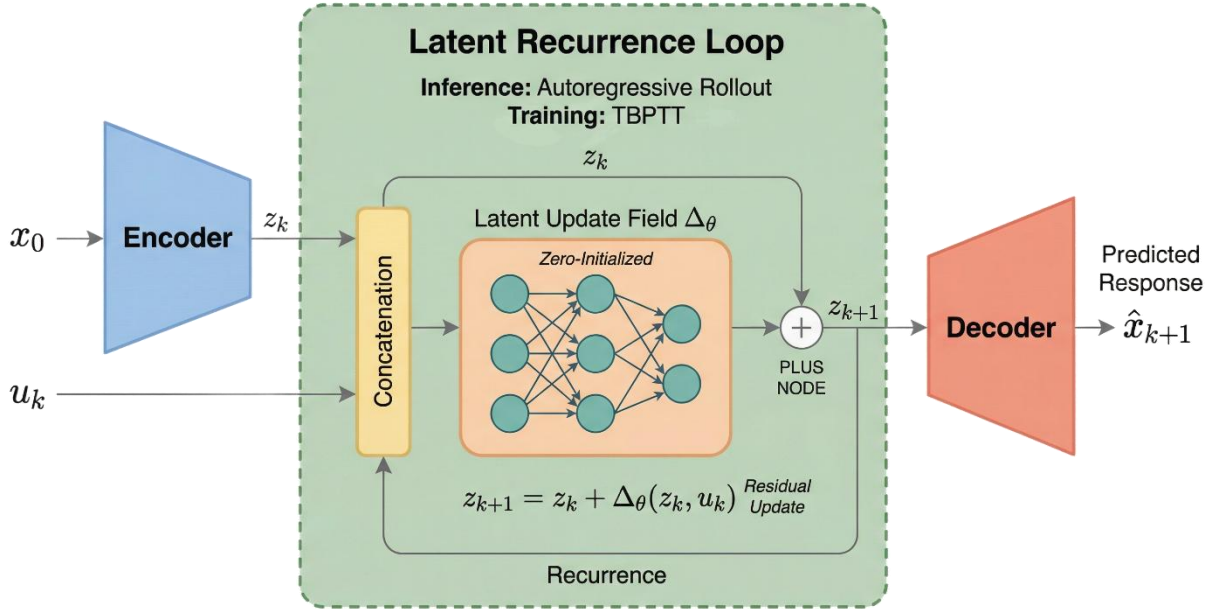
124 The dynamic response of a structure may not be fully observable from displacements alone—
 125 unobserved internal states such as velocities, internal forces, and damping forces govern the
 126 dynamics. LSR-AR addresses this by extracting a latent vector $z_0 \in \mathbb{R}^{d_z}$ from the observed
 127 displacement x_0 through an encoder. The latent vector is intended to capture not only the
 128 observable displacement but also supplementary information learned from data, serving as a
 129 surrogate for the unobserved states.

130 In this latent space, a latent increment field $\Delta_\theta(z_k, u_k)$ conditioned on the current forcing u_k is
 131 learned as an MLP. The increment field can be interpreted as a learned state-update field over the
 132 latent space, and the state update itself takes the form of a forward Euler step:

$$z_{k+1} = z_k + \Delta_\theta(z_k, u_k) \quad (6)$$

133 A decoder recovers the physical-space displacement \hat{x}_{k+1} from the updated latent state z_{k+1} . No
 134 explicit time embedding is used; temporal information enters only through the forcing sequence.

135 The model is trained end-to-end on full-trajectory rollout loss with truncated backpropagation
 136 through time (TBPTT) (Williams and Peng, 1990). The training hyperparameters of the final
 137 selected models are listed in Table 2. Figure 2 summarizes this encoder–increment–decoder
 138 structure.



139
 140 **Figure 2:** LSR-AR architecture: an autoregressive rollout composed of an encoder, latent increment field, and decoder.

141 **Table 2:** Hyperparameter summary for the compared evaluation setups.

	FNO	SA-FNO	NARX-MLP	LSR-AR
Parameters	7.40M	7.40M	111K	236K
Optimizer	Adam	Adam	AdamW	AdamW
Epochs	500	500	300	200
Batch size	200	200	50	32
TBPTT length	-	-	200	200

142 **4. RESULTS**

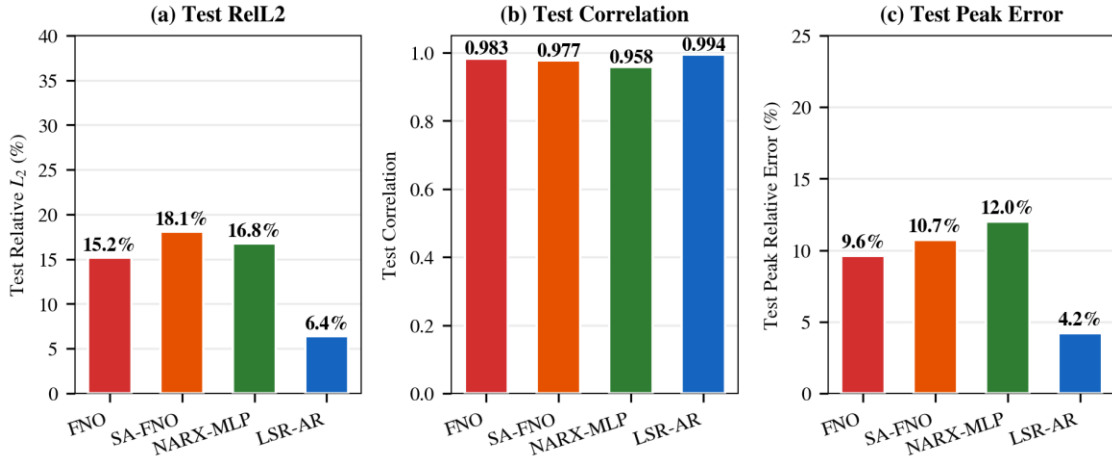
143 **4.1. Full-horizon inference comparison**

144 Table 3 compares the main LSR-AR configuration without explicit time embedding against the
 145 baselines. Within the retrained causal comparison, this configuration achieves 6.4% test ReL2,
 146 compared with 16.8% for NARX-MLP. Only task-aligned 37-floor comparisons are included in
 147 Table 3. Published one-shot checkpoints from Goswami et al. (2025) are shown only as
 148 contextual references; they yield 15.2% for FNO and 18.1% for SA-FNO on all 37 floors. Top-
 149 floor-only references such as DeepFNONet and DeepONet are omitted because they are not task-
 150 aligned 37-floor comparisons. The SBD column in Table 3 shows a lower value for LSR-AR
 151 (SBD = 0.006) than for the compared baselines (0.058 - 0.080). Because SBD is invariant to
 152 amplitude and temporal shift, this is consistent with better overall waveform-shape agreement
 153 under the present aggregation. Figure 3 visualizes the same all-floor comparison across ReL2,
 154 correlation, and peak displacement error.

155 **Table 3:** Main model comparison on the full 997-step along-wind benchmark using ReL2, PeakErr, Corr., and SBD

Model	Type	Params	ReL2 (%)	Peak err	Corr.	SBD
LSR-AR (no time embedding)	Causal	236K	6.4	4.2%	0.994	0.006
NARX-MLP	Causal	111K	16.8	12.0%	0.958	0.058
FNO [†]	One-shot	7.40M	15.2	9.6%	0.983	0.063
SA-FNO [†]	One-shot	7.40M	18.1	10.7%	0.977	0.080

[†]Published checkpoints from Goswami et al. (2025) were used.



156
157 **Figure 3:** Model-wise comparison on 200 unseen test scenarios: (a) ReL2, (b) correlation, and (c) peak
158 displacement error

159 **4.2. Ablation study**

160 Table 4 ablates the key design choices and reports test aggregate ReL2 for the selected variants.
161 Here d_z denotes the latent dimension and d_t the time-embedding dimension.

162 **Table 4:** Selected LSR-AR ablations on latent dimension and time embedding, reported with aggregate test ReL2

Variant	Change	Params	Epochs	Best ReL2
Default ($d_t = 0$)	No time embed	236K	200	6.4%
<i>Time embedding:</i>				
$d_t = 32$	With time embed	244K	200	7.3%
<i>Latent dimension (200 epochs):</i>				
$d_z = 4$	Bottleneck	223K	200	8.3%
$d_z = 8$	Reduced	225K	200	7.6%

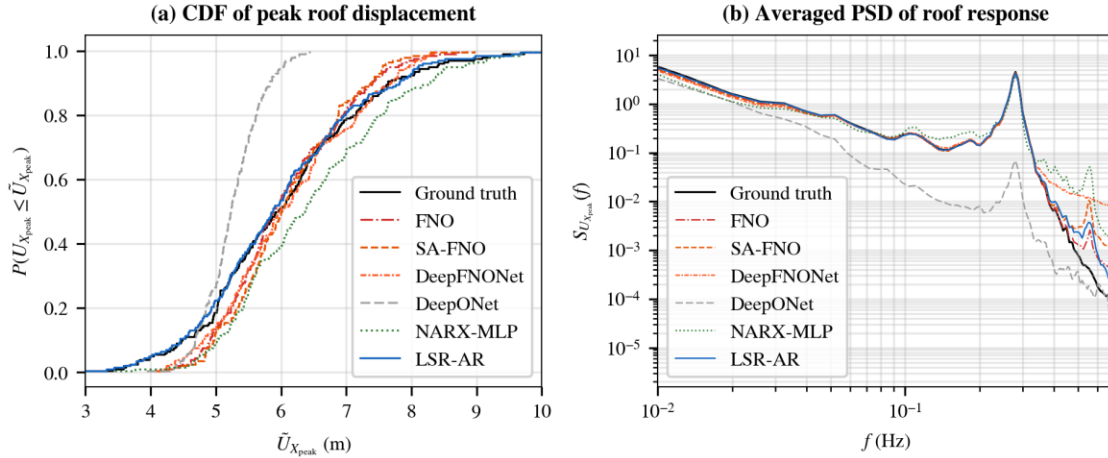
163

164 Table 4 summarizes the completed ablations selected by validation performance; the test metrics
165 are shown only for final reporting. Among these variants, the default configuration without explicit
166 time embedding yields the lowest ReL2 on this benchmark.

167 **4.3. Peak displacement and PSD analysis**

168 For this roof-level diagnostic, the released roof-only DeepONet and DeepFNONet references are
169 also included. Figure 4 (a) compares empirical cumulative distribution functions (CDFs) of peak
170 roof displacement across 200 test scenarios. LSR-AR most closely follows the ground-truth

171 distribution among the compared models. Figure 4 (b) shows the averaged power spectral density
 172 (PSD) at the roof. Most models, including LSR-AR, reproduce the dominant first-mode resonant
 173 peak with good fidelity. However, all models share a common artifact: a spurious secondary
 174 peak at higher frequencies that is absent or much weaker in the ground truth. One possible
 175 explanation is that the compared surrogates fit the dominant low-frequency content more readily,
 176 while high-frequency response contributes less to the overall training loss and may therefore be
 177 learned less faithfully. This behavior may be especially problematic for resonance-sensitive
 178 responses, and it should be considered explicitly in wind-response studies where high-frequency
 179 behavior is important.



180

181 **Figure 4:** Roof-level diagnostics: (a) empirical CDF of peak displacement and (b) averaged PSD of the response

182 4.4. Separate temporal extrapolation study

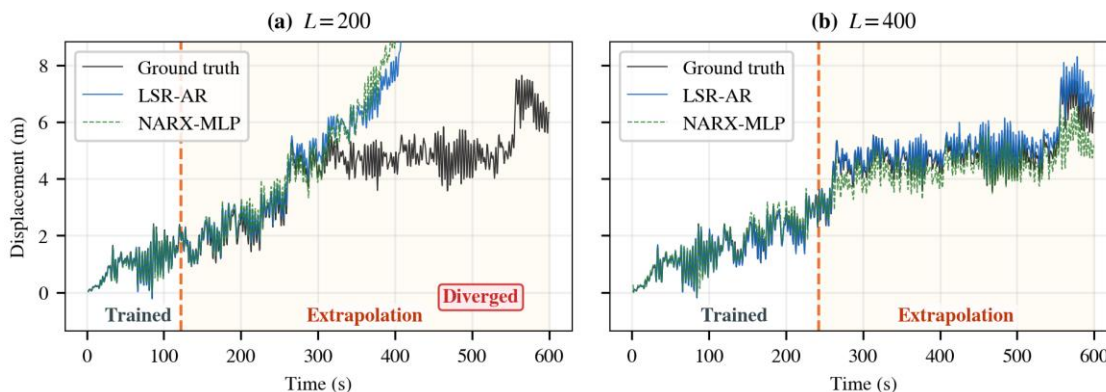
183 To examine continuation beyond the seen training horizon, we conducted a separate rollout study
 184 using a time-embedding LSR-AR variant and NARX-MLP. Each model was trained on truncated
 185 prefixes with $L = 200$ and $L = 400$, then evaluated by rolling out the full 997-step test sequence
 186 on unseen scenarios. No additional hyperparameter tuning was performed for the $L = 200$ and L
 187 $= 400$ runs; each cropped-length experiment reused the corresponding architecture and
 188 optimization settings from the full-length setup. Unlike Table 3, Table 5 reports mean per-
 189 sample ReL2 in order to summarize scenario-wise continuation behavior, including divergent
 190 cases. For that reason, the full-train reference values in Table 5 are not directly numerically
 191 identical to the aggregate/global ReL2 values reported in Table 3.

192 Under this setup, both models diverged at $L = 200$, whereas both remained stable at $L = 400$. The
 193 400-step NARX-MLP run achieved 9.9% full-rollout ReL2, lower than its own full-train result
 194 (15.2%) under the same mean per-sample summary. We do not attempt to interpret that difference
 195 further here. The time-embedding LSR-AR variant also remained stable at $L = 400$, but its full-
 196 rollout error increased to 14.7%; its in-distribution error (4.7%) nevertheless stayed below its own
 197 full-train reference (6.7%). Figure 5 shows a representative roof-floor rollout for the $L = 200$ and
 198 $L = 400$ settings.

199 **Table 5:** Temporal extrapolation results for the separate two-model continuation study, reported as mean per-sample
 200 ReL2 over the seen, extrapolation, and full regions.

Model	Train length	Seen mean per-sample ReL2	Extrapolation mean per-sample ReL2	Full mean per-sample ReL2
LSR-AR	200	3.8 %	diverge	diverge
LSR-AR	400	4.7 %	15.2 %	14.7 %
LSR-AR	997 (full)	-	-	6.7 %
NARX-MLP	200	7.5 %	diverge	diverge
NARX-MLP	400	8.8 %	9.9 %	9.9 %
NARX-MLP	997 (full)	-	-	15.2 %

201



202

203 **Figure 5:** Representative roof-floor temporal extrapolation trajectories for (a) $L = 200$ and (b) $L = 400$

204 **5. DISCUSSION AND CONCLUSION**

205 This paper presented LSR-AR, a latent-step autoregressive surrogate for along-wind response
 206 prediction of a 37-story high-rise building. As summarized in Table 3 and Figure 3, the main LSR-
 207 AR configuration without explicit time embedding achieves 6.4% test ReL2 with 236K
 208 parameters. Within the full-horizon comparison of Table 3, this improves over NARX-MLP. The
 209 selected ablations in Table 4 further suggest that moderate latent compression and a simple
 210 forward-Euler residual update remain viable on this benchmark, while leaving room for higher-
 211 order integrators and larger-system studies.

212 A secondary continuation study of the separate time-embedding variant, together with Figure 5,
 213 showed stable beyond-horizon rollout at $L = 400$, whereas the $L = 200$ setting diverged. At the
 214 same time, the divergence observed under shorter training horizons indicates that the model cannot
 215 yet be interpreted as fully capturing the governing physics. Likewise, Figure 4 shows common
 216 high-frequency spectral artifacts across models, motivating follow-up work on cross-wind and
 217 torsional response, larger-system extrapolation, sampling sensitivity, and stronger physical
 218 constraints or regularization.

219 **ACKNOWLEDGEMENTS**

220 This work was supported by Institute of Information & communications Technology Planning & Evaluation (IITP)
 221 grant funded by the Korea government (MSIT) [No. RS-2021-II211343, Artificial Intelligence Graduate School
 222 Program (Seoul National University)].

223 **REFERENCES**

- 224 Goswami, S., Giovanis, D.G., Li, B., Spence, S.M.J., Shields, M.D., 2025. Neural operators for stochastic modeling
225 of nonlinear structural system response to natural hazards. *Eng. Struct.* 345, 121284.
226 <https://doi.org/10.1016/j.engstruct.2025.121284>
- 227 Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Bhattacharya, K., Stuart, A., Anandkumar, A., 2021. Fourier neural
228 operator for parametric partial differential equations, in: *International Conference on Learning Representations*
229 (ICLR 2021).
- 230 Lu, L., Jin, P., Pang, G., Zhang, Z., Karniadakis, G.E., 2021. Learning nonlinear operators via DeepONet based on the
231 universal approximation theorem of operators. *Nat. Mach. Intell.* 3, 218–229. [https://doi.org/10.1038/s42256-](https://doi.org/10.1038/s42256-021-00302-5)
232 [021-00302-5](https://doi.org/10.1038/s42256-021-00302-5)
- 233 Paparrizos, J., Gravano, L., 2015. k-Shape: Efficient and accurate clustering of time series, in: *Proceedings of the 2015*
234 *ACM SIGMOD International Conference on Management of Data*, pp. 1855–1870.
235 <https://doi.org/10.1145/2723372.2737793>
- 236 Tamura, Y., 2012. Aerodynamic database for high-rise buildings. Tokyo Polytechnic University, Global Center of
237 Excellence Program. <http://www.wind.arch.t-kougei.ac.jp/system/eng/contents/code/tpu>
- 238 Thedy, J., Liao, K.W., Kim, T., 2025. Auto regressive neural network-driven reliability optimization in base-isolated
239 building design. *Results Eng.* 26, 104713. <https://doi.org/10.1016/j.rineng.2025.104713>
- 240 Williams, R.J., Peng, J., 1990. An efficient gradient-based algorithm for on-line training of recurrent network
241 trajectories. *Neural Comput.* 2, 490–501. <https://doi.org/10.1162/neco.1990.2.4.490>