
TS-ICL: A Flexible Time-Indexed Foundation Model for Time Series via In-Context Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Foundation models mark a profound paradigm shift in time series modeling, with
2 task-specific models being superseded by general-purpose zero-shot models. Yet,
3 current approaches primarily focus on forecasting, while real-world time series
4 are often irregularly and partially observed, requiring models that can jointly fore-
5 cast, impute missing values, and handle degraded sampling conditions. To address
6 these challenges, we introduce TS-ICL, a novel probabilistic In-Context Learning
7 encoder–regressor Transformer that unifies forecasting and imputation. TS-ICL
8 formulates time series tasks as timestamp-aligned regression and naturally incor-
9 porates covariates by training on synthetic dependency structures generated from a
10 novel causal data prior. Empirically, TS-ICL achieves a new state-of-the-art in im-
11 putation, while remaining competitive with leading forecasting foundation mod-
12 els across both univariate and covariate-aware benchmarks. It shows particularly
13 strong performance in forecasting with partially observed look-back windows.

14 1 Introduction

15 The recent advent of Time Series Foundation Models (TSFMs) drastically changes time series mod-
16 eling, shifting from task-specific to transferable models that leverage large-scale pretraining and
17 adaptation mechanisms to provide zero-shot inference on unseen data [4, 27, 10, 28]. TSFMs ad-
18 dress limited data regimes and distribution shifts [44], but still face fundamental limitations. First,
19 many tasks cannot be solved in practice without leveraging external information, requiring covariate-
20 aware inference [40]. Second, real-world observations are often incomplete, with missing values and
21 asynchronous measurements [9]. This challenges standard modeling assumptions and motivates uni-
22 fied frameworks that can jointly handle forecasting and imputation in flexible observation settings.

23 To address these challenges, (1) recent TSFMs such as Chronos-2 [3] build on In-Context Learn-
24 ing [5] (ICL) to support covariate-informed inference and handle missing values, while remaining
25 highly efficient. Nevertheless, they do not natively address imputation, thus hindering their practical
26 use. (2) In a remarkably different approach, TabPFN [17] and TabICLv2 [34] have revolutionized the
27 tabular domain with strong few-shot regression capabilities via Transformer-based ICL and synthetic
28 data priors. When adapted to time series (e.g., TabPFN-TS [19]), these Tabular Foundation Models
29 (TFMs) naturally support covariates and enable both zero-shot imputation [26] and forecasting [38].
30 Yet, they lack temporal inductive bias and rely on handcrafted time features. As a result, TFMs lag
31 behind pure TSFMs in forecasting benchmarks, while also incurring high inference cost [38].

32 Consequently, existing approaches fail to jointly provide (i) unified forecasting and imputation,
33 (ii) covariate-aware inference, and (iii) efficient zero-shot performance (see Table 1). In this pa-
34 per, we tackle this challenge and introduce TS-ICL, a unified probabilistic Transformer foundation
35 model for imputation and forecasting. (i) TS-ICL casts time series modeling as an in-context re-
36 gression problem, where observations are represented as timestamp-aligned inputs and encoded into

Table 1: Capabilities of recent time series foundation models. Only TS-ICL supports forecasting and imputation, while enabling efficient covariate-aware inference and supporting irregular sampling.

Method	Handles Forecasting	Handles Imputation	Covariate Integration	Probabilistic Predictions	Designed for Time Series	Irregular Sampling	Fast Inference
TiRex [4], Toto [10], TimesFMv2.5 [11]	✓	✗	✗	✓	✓	✗	✓
Chronos-2 [3]	✗	✗	✓	✓	✓	✗	✓
TabPFNv2.5-TS [19], TabICLv2-TS [34]	✓	✓	✓	✓	✗	✓	✗
TS-ICL (ours)	✓	✓	✓	✓	✓	✓	✓

37 contextual representations that enable forecasting or imputation in a single forward pass. (ii) To
 38 enable effective covariate-aware inference, a structured synthetic prior over target–covariate rela-
 39 tionships is introduced using Directed Acyclic Graphs (DAGs) to define dependency structure, with
 40 node-level mechanisms inspired by structural causal models [32, 17]. (iii) Unlike existing TSFMs,
 41 TS-ICL operates directly on timestamped observations rather than fixed grids, allowing flexible
 42 handling of missing or irregularly sampled data in practice.

43 **Contributions.** The main contributions are as follows:

- 44 • **A unified and flexible TSFM architecture.** We introduce TS-ICL, a novel probabilistic TSFM
 45 that casts time series modeling as a time-indexed in-context regression problem, unifying fore-
 46 casting and imputation with native support for covariates.
- 47 • **Structured synthetic prior.** We design a novel DAG-based causal prior over target–covariate
 48 time series, enabling robust zero-shot generalization to unseen dependency structures.
- 49 • **State-of-the-art imputation performance.** TS-ICL sets a new state-of-the-art on zero-shot impu-
 50 tation benchmarks, consistently outperforming both task-specific models and TFMs, while being
 51 up to $50\times$ faster than TFMs at inference.
- 52 • **Competitive forecasting performance.** On forecasting benchmarks, TS-ICL matches state-of-
 53 the-art TSFMs while supporting covariate-aware inference, and remains particularly robust to
 54 missing observations due to its time-indexed formulation.

55 2 Related Work

56 **Time series foundation models.** Time Series Foundation Models (TSFMs) are pretrained general-
 57 purpose models for time series, typically trained on large mixtures of real and synthetic data and
 58 often based on patch-based architectures [27, 4, 10, 11]. While they achieve strong zero-shot fore-
 59 casting performance on standard benchmarks [1, 33], they exhibit several limitations in practical
 60 settings: they are typically not designed for covariate-aware inference, and primarily focus on fore-
 61 casting without addressing imputation. Chronos-2 [3] partially mitigates the covariate limitation
 62 by enabling inference-time conditioning on exogenous time series, leading to improved performance
 63 when informative covariates are available. However, its patch-based formulation assumes regularly
 64 sampled inputs and does not natively support imputation, which is critical in many real-world time
 65 series applications [36, 9, 6, 13].

66 **Tabular foundation models for time series.** Tabular Foundation Models (TFMs) such as TabPFN
 67 [17] and TabICLv2 [34] leverage in-context learning over synthetic task distributions to achieve
 68 strong few-shot regression performance. Extensions to time series [19] demonstrate competitive re-
 69 sults for both forecasting [38] and imputation [26], while naturally supporting covariates at inference
 70 time. However, these approaches typically rely on handcrafted temporal features, such as Fourier
 71 features operating at predefined frequencies, and incur higher inference costs compared to dedicated
 72 TSFMs [38], limiting their scalability in practice.

73 **Supervised imputation models.** Classical supervised imputation methods such as BRITS [6] and
 74 SAITS [13] achieve strong performance in in-domain settings but require task-specific training and
 75 generalize poorly across domains. Recent large-scale evaluations [26] suggest that TFM-based ap-
 76 proaches can significantly outperform these methods in realistic scenarios.

77 **Time-indexed and Neural Field models.** Continuous-time modeling of time series (also referred
 78 to as time-indexed modeling [42]) has been explored through Neural Ordinary Differential Equa-
 79 tions [8] (ODE) and latent ODE frameworks [35]. More recently, Neural Field-based represen-
 80 tations [25, 37] encode irregular observations into continuous latent functions without requiring
 81 explicit temporal alignment. These approaches provide a principled foundation for flexible time
 82 representations, but are not designed for zero-shot inference in forecasting or imputation settings.

83 Despite these advances, existing approaches remain fragmented across forecasting, imputation, and
 84 representation learning. No framework jointly provides efficient zero-shot inference, covariate-
 85 aware modeling, and a unified treatment of forecasting and imputation within a single architecture.

86 3 TS-ICL Architecture

87 This section introduces the proposed TS-ICL architecture, designed for efficient probabilistic zero-
 88 shot time series forecasting and imputation, while maintaining flexibility in handling covariates and
 89 irregularly sampled observations. A high-level overview of the architecture and its main components
 90 is provided, while detailed descriptions of each module are deferred to Appendix A.

91 **Problem setting.** Let $\mathbf{x} = (x_t)_{t \in \mathcal{T}}$ denote a time series defined over a (possibly irregular) set of
 92 timestamps \mathcal{T} . The index set \mathcal{T} is partitioned into two disjoint subsets: (i) the context timestamps
 93 $\mathcal{T}^{\text{ctxt}}$, corresponding to observed values, and (ii) the target timestamps \mathcal{T}^{tgt} , corresponding to values
 94 to be predicted. Accordingly, $\mathbf{x}^{\text{ctxt}} = (x_t)_{t \in \mathcal{T}^{\text{ctxt}}}$ denotes the observed context values, and $\mathbf{x}^{\text{tgt}} =$
 95 $(x_t)_{t \in \mathcal{T}^{\text{tgt}}}$ the target values.

96 Additionally, an optional set of exogenous covariate time series is considered:

$$\mathbf{X}^{\text{covar}} = \left\{ \mathbf{x}_c^{\text{covar}} = (x_{c,t}^{\text{covar}})_{t \in \mathcal{T}_c^{\text{covar}}} \right\}_{c=1}^{C-1},$$

97 where C refers to the channel dimension, with one channel corresponding to the time series of
 98 interest and the remaining $C - 1$ channels corresponding to exogenous covariates.

99 Each covariate is defined over its own set of timestamps $\mathcal{T}_c^{\text{covar}} \subseteq \mathcal{T}$. Covariates may be observed
 100 over arbitrary subsets of timestamps, including only the context (e.g. the look-back window in
 101 forecasting) or both context and target (e.g. look-back and horizon windows), and may be sparsely
 102 observed. The TS-ICL framework can handle such heterogeneous availability without requiring any
 103 imputation preprocessing.

104 The objective, common to both forecasting and imputation, is to infer the target values conditioned
 105 on the observed context and, when available, the covariates:

$$p(\mathbf{x}^{\text{tgt}} \mid \mathbf{x}^{\text{ctxt}}, \mathbf{X}^{\text{covar}}).$$

106 This formulation unifies forecasting and imputation as conditional inference problems. For clarity,
 107 the next section omits batch and timestamp indices whenever no ambiguity arises.

108 3.1 TS-ICL Overview

109 The key idea behind TS-ICL is to reformulate time series prediction as an in-context regression prob-
 110 lem over learned temporal representations. Thus, the architecture is composed of four successive
 111 modules that transform raw observations into global and local context-aware representations used
 112 for prediction. The pipeline first encodes each time series, then aggregates information across co-
 113 variates, and finally produces timestamp-specific representations that enable in-context regression.
 114 The overall architecture is illustrated in Figure 1.

115 **(i) Time Series Encoder \mathcal{E} .** This module extracts representations from the context time series and
 116 optional $C - 1$ covariates in a channel-independent manner. It follows a Perceiver encoder design
 117 [21, 37] based on a fixed set of M learnable tokens that sequentially cross-attend to timestamp-value
 118 pairs. This allows compressing an arbitrary-length input into a compact latent sequence.

$$(\mathbf{x}^{\text{ctxt}}, \mathcal{T}^{\text{ctxt}}, \mathbf{X}^{\text{covar}}, \mathcal{T}^{\text{covar}}) \xrightarrow{\mathcal{E}} \mathbf{Z}^{\text{val}} \in \mathbb{R}^{C \times M \times d}.$$

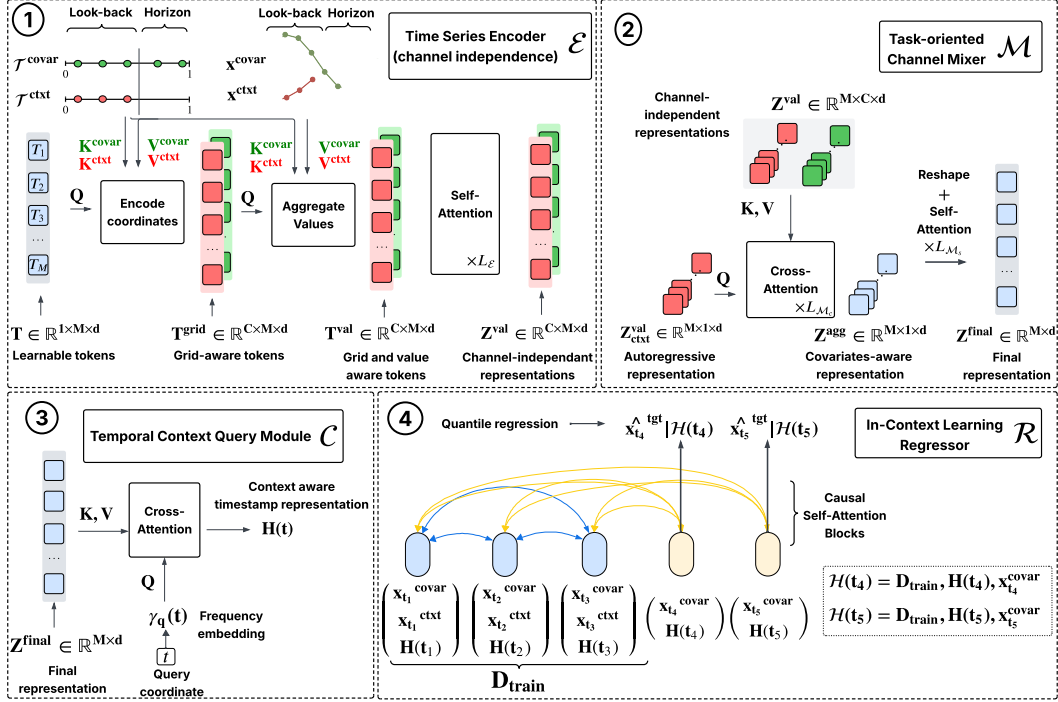


Figure 1: The TS-ICL pipeline. From temporal encoding to in-context regression. The diagram illustrates the four-module transformation for forecasting with one covariate observed on the horizon.

119 **(ii) Channel Mixer \mathcal{M} .** This module aggregates information across the C channels to produce
 120 a unified representation. For each of the M latent positions, the token corresponding to the time
 121 series of interest queries the tokens of the $C - 1$ covariates via cross-attention. This mechanism col-
 122 lapses the channel dimension by selectively integrating exogenous information into the main series’
 123 representation. A subsequent stack of self-attention layers models global dependencies among the
 124 resulting M tokens, yielding a compact task-oriented representation. This step is critical to capture
 125 inter-channel dependencies, which are not modeled by the channel-independent encoder.

$$\mathbf{Z}^{\text{val}} \in \mathbb{R}^{C \times M \times d} \xrightarrow{\mathcal{M}} \mathbf{Z}^{\text{final}} \in \mathbb{R}^{M \times d}.$$

126 **(iii) Temporal Context Query Module \mathcal{C} .** This module aims at bridging the gap between the
 127 discrete latent tokens $\mathbf{Z}^{\text{final}}$ and the continuous time domain. Any timestamp $t \in \mathcal{T}$ is first embed-
 128 ded using a frequency-based positional encoding [29]. This local encoding then cross-attends to the
 129 time series representation $\mathbf{Z}^{\text{final}}$, providing a single local and global context-aware representation
 130 of arbitrary timestamps. This design enables querying at arbitrary timestamps, supporting irregular
 131 sampling and unifying forecasting and imputation.

$$(t, \mathbf{Z}^{\text{final}}) \xrightarrow{\mathcal{C}} \mathbf{H}(t) \in \mathbb{R}^d.$$

132 **(iv) In-Context Learning Regressor \mathcal{R} .** Given the representations $\mathbf{H}(t)$, prediction is formulated
 133 as an in-context regression problem [15, 41], where the observed context defines a training set:

$$\mathcal{D}_{\text{train}} = \{\mathbf{H}(t), \mathbf{x}_t^{\text{ctxt}}, \mathbf{X}_t^{\text{covar}}\}_{t \in \mathcal{T}^{\text{ctxt}}},$$

134 used to infer target values at unseen target timestamps t^{tgt} . The regressor is implemented as a
 135 Transformer that performs causal self-attention over the training (context) set, effectively learning to
 136 map representations to values in an in-context manner. When available, covariates are incorporated
 137 in the training set via cross-attention, allowing the model to condition on exogenous time series
 138 under arbitrary availability patterns (e.g., context-only, context and target, or sparse observations).
 139 TS-ICL outputs a dense set of quantile estimates of the target distribution, trained using a smoothed
 140 pinball loss [39].

$$\mathbf{H}(t^{\text{tgt}}), \mathbf{X}_{t^{\text{tgt}}}^{\text{covar}} \xrightarrow{\mathcal{R}} p(\mathbf{x}^{\text{tgt}} | \mathbf{H}(t^{\text{tgt}}), \mathbf{X}_{t^{\text{tgt}}}^{\text{covar}}, \mathcal{D}_{\text{train}}).$$

141 This four-step formulation unifies time series representation learning and in-context regression, en-
 142 abling TS-ICL to perform forecasting and imputation while flexibly operating over timestamped
 143 observations and optionally observed covariates. Additional architectural details and hyperparame-
 144 ters are provided in Appendix A and Appendix B.2, respectively.

145 4 Data Prior and Training Procedure

146 TS-ICL is trained on a structured data prior combining real-world and synthetic time series, spanning
 147 both univariate signals and multivariate covariate–target structures. This prior is designed to jointly
 148 capture temporal dynamics and inter-variable dependencies within a unified training distribution.
 149 Concretely, training samples consist of either univariate time series or multivariate problems where
 150 target–covariate relationships are generated via structured transformations of base signals.

151 **Univariate time series.** In the univariate setting, a large and highly heterogeneous pretraining
 152 prior is constructed by combining real-world and synthetic time series. This mixture is designed to
 153 expose TS-ICL to a broad spectrum of temporal dynamics. For *real-world data*, the prior leverages
 154 a collection of 31 datasets spanning multiple domains, as detailed in Appendix B.1.1. These datasets
 155 cover diverse temporal phenomena, including trends, seasonal patterns, regime shifts, and varying
 156 levels of non-stationarity. The training distribution is further augmented with *synthetic data* sam-
 157 pled from the TempoPFN univariate time series generator [30], which induces controlled but diverse
 158 stochastic processes. Overall, this yields a large-scale pretraining distribution over univariate time
 159 series, comprising 40 datasets and approximately 2M time series with lengths ranging from ~ 100
 160 to $\sim 600k$ time steps.

161 **Covariates generator.** To enable TS-ICL to
 162 learn under structured covariates, synthetic
 163 multivariate time series are constructed from
 164 base univariate signals (either real or gener-
 165 ated as described above). Specifically, (i) A
 166 Directed Acyclic Graph (DAG) is generated
 167 over univariate signals where nodes correspond
 168 to time series and edges encode causal de-
 169 pendencies. Following [34], each non-root
 170 node is produced by applying a transforma-
 171 tion sampled from a Structural Causal Model
 172 (SCM) registry, including both linear and non-
 173 linear operators (e.g., linear mappings, MLP,
 174 RNN). This yields heterogeneous dependency
 175 structures while preserving temporal coher-
 176 ence. (ii) Given the resulting graph, one node
 177 is selected as the prediction target and a subset
 178 of the remaining nodes is sampled as covariates,
 179 producing multivariate problems with varying
 180 numbers of covariates and controllable depen-
 181 dency strength. This design ensures that cov-
 182 ariates can be causally informative, redundant,
 183 or entirely independent of the target, preventing
 184 reliance on spurious correlations. Figure 2
 185 illustrates the synthetic target–covariate gener-
 186 ation process. Additional details on the data
 187 prior are provided in Appendix B.1.2.

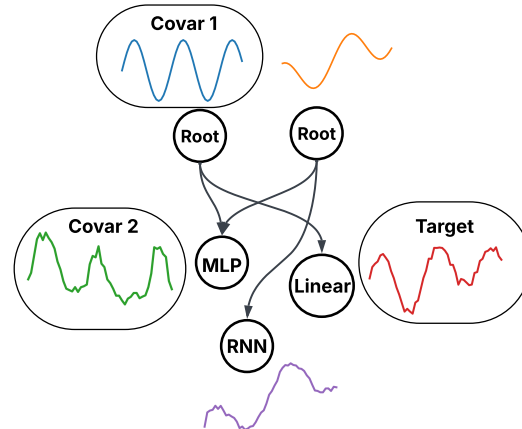


Figure 2: Synthetic target–covariate generation. Multivariate structures are constructed from a sampled DAG, where base signals are transformed via linear and non-linear SCMs. One node is selected as the target, while others serve as informative or redundant covariates.

186 **Whole procedure.** Overall, each training sample is constructed by first sampling base time series
 187 (real or synthetic), which may either be used directly or serve as building blocks for multivariate
 188 structures via the DAG-based generator. This yields a unified training distribution over univariate
 189 and multivariate time series. The training task is then defined dynamically. • For imputation, ob-
 190 servations are masked either point-wise or in contiguous segments with randomly sampled masking
 191 ratios. • For forecasting, a future horizon of random length is masked. The full data generation
 192 pipeline is summarized in Algorithm 1 in Appendix B.1.3, while hyperparameters and training de-
 193 tails are reported in Appendix B.2.

194 **5 Experiments**

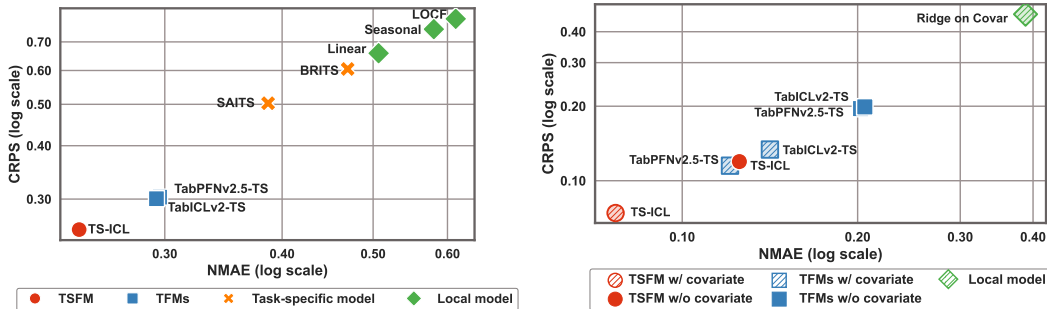
195 TS-ICL is evaluated on zero-shot imputation (Section 5.1) and forecasting (Section 5.2) under two
 196 settings: (i) *univariate time series* and (ii) *time series with known covariates* available at inference
 197 time. All experiments use a 27M-parameter version of TS-ICL (see Appendix B.2 for architectural
 198 details). Imputation experiments rely on `fm-impute-bench` [26], while forecasting is evaluated on
 199 `fev-bench` [38], following recent protocols for time series foundation models. Ablation studies on
 200 TS-ICL components are provided in Appendix C and additional results on the TIME benchmark [33]
 201 are presented in Appendices E and F.

202 **5.1 Imputation Experiments**

203 The zero-shot imputation capability of TS-ICL is evaluated on `fm-impute-bench` [26], covering
 204 both *univariate* settings and scenarios with *covariates* available at inference time. The benchmark
 205 spans diverse missingness patterns, sequence lengths, and application domains.

206 **Setting.** (i) In the *univariate setting*, `fm-impute-bench` comprises 33 datasets across multiple
 207 domains (e.g., energy, climate, etc.) with varying lengths and frequencies (see Table 10, Appendix
 208 E.1). Each sample corresponds to a four-week window. TS-ICL is evaluated under four masking
 209 scenarios, namely: • 50% or • 70% pointwise missingness, and • two or • four disjoint one-day gaps.
 210 This results in 132 tasks and about 1.3M windows to impute. (ii) The same four masking scenarios
 211 are applied to the *known-covariates* setting on six datasets providing informative covariates (see
 212 Table 11, Appendix E.1). This results in 24 tasks and approximately 1K windows to impute.

213 **Baselines.** (i) *Univariate comparisons* include Tabular Foundation Models (TFMs) adapted to
 214 time series: TabPFNv2.5-TS [19], TabICLv2-TS [34], along with standard local methods: Linear
 215 Interpolation, Seasonal Naive, and Last Observation Carried Forward (LOCF). In
 216 addition, supervised imputation models SAITS [13] and BRITS [6], trained per dataset, serve as
 217 strong task-specific baselines. (ii) In the *known-covariates* setting, the same foundation models are
 218 considered, together with ridge regression on covariates to quantify the predictive signal of exoge-
 219 nous variables. Variants of foundation models without covariates are also included for comparison.
 220 Following `fm-impute-bench`, pointwise performance is evaluated with Normalized Mean Absolute
 221 Error (NMAE) and probabilistic performance with Continuous Ranked Probability Score (CRPS).
 222 Metric definitions are provided in Appendix D. Figure 3 illustrates the trade-off between pointwise
 223 and probabilistic performance, while Table 2 reports the median inference time.

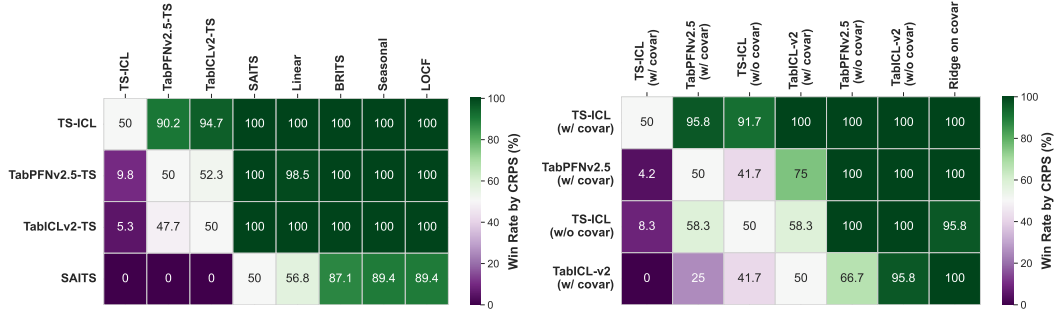


(a) Univariate imputation across 132 tasks. (b) Imputation with *known covariates* across 24 tasks.

Figure 3: NMAE-CRPS (lower is better) on the `fm-impute` benchmark. Each point corresponds to a method, averaged across tasks.

224 **Results.** TS-ICL sets a new state-of-the-art for zero-shot imputation, improving both pointwise
 225 and probabilistic scores over TFMs while being up to two orders of magnitude faster at inference.

226 (i) *Univariate setting.* As shown in Figure 3a, TS-ICL achieves lower NMAE and CRPS than
 227 competing TFMs, improving over TabICLv2-TS by 17% and 15%, respectively, while being $\sim 50\times$
 228 faster at inference (Table 2). TabPFNv2.5-TS and TabICLv2-TS perform similarly and outperform



(a) Univariate imputation across 132 tasks. (b) Imputation with *known covariates* across 24 tasks.

Figure 4: Pairwise win rates of the top-4 models on the `fm-impute` benchmark. Each entry indicates the fraction of tasks where a method outperforms another according to the CRPS.

task-specific and local baselines by a wide margin. Pairwise win rates (Figure 4a) further show that TS-ICL dominates on the vast majority of tasks, indicating robustness across tasks.

(ii) *Covariate-aware setting*. Results in Figure 3b show similar trends when covariates are available. TS-ICL improves over TabPFNv2.5-TS by 36% in NMAE and 35% in CRPS, and gains 39% (NMAE) and 38% (CRPS) over its variant without covariates. While all TFMs benefit from covariates, they remain below TS-ICL. Ridge regression indicates limited predictive power of covariates alone, and pairwise comparisons (Figure 4b) demonstrate the clear dominance of TS-ICL.

Table 2: Median imputation inference time on `fm-impute-bench univar` with a H100 GPU.

	TFSM	Tabular Foundation models		Task-Specific Models		Local models		
	TS-ICL	TabPFNv2.5-TS	TabICLv2-TS	SAITS	BRITS	Linear interpolation	Seasonal Naive	LOCF
Inference time (s per window)	6.51×10^{-3}	2.80×10^{-1}	3.07×10^{-1}	1.33×10^{-2}	4.52×10^{-1}	1.61×10^{-4}	7.54×10^{-4}	5.58×10^{-4}

Overall, TS-ICL outperforms state-of-the-art TFMs for zero-shot imputation in both *univariate* and *covariate-informed* settings, while being two orders of magnitude faster at inference. Additional qualitative results in Appendix E, including Figures 16 and 17, as well as results on the TIME benchmark [33], further support these findings.

5.2 Forecasting Experiments

The zero-shot forecasting capability of TS-ICL is evaluated on `fev-bench` [38], a comprehensive benchmark covering diverse datasets, horizons, and sampling frequencies, with controlled evaluations in both *univariate* and *known-covariates* settings. An additional univariate setting with *missing values in the look-back window* is also considered across the entire `fev-bench` benchmark.

Setting. Evaluation considers two forecasting regimes. (i) In the *univariate setting*, the benchmark comprises 100 tasks and $\sim 235k$ forecasting windows (see Table 18, Appendix F.1). Models rely solely on past observations, with no access to covariates or cross-series information. (ii) In the *known-covariate setting*, we evaluate all methods on the same 100 tasks. Among these, 30 datasets include meaningful exogenous time series, denoted as “known dynamics” covariates in Table 18. For these datasets, methods that support covariates are evaluated both with and without covariate inputs, while methods that do not are kept unchanged. This protocol enables a controlled assessment of the effect of covariate information while preserving comparability across tasks.

Baselines. TS-ICL is compared against a broad set of baselines covering foundation models and local methods. (i) *Univariate comparisons* include state-of-the-art TFSMs (Chronos-2 [3], TiRex [4], Chronos-bolt and Toto [10]), TFMs adapted to time series (TabPFNv2.5-TS, TabICLv2-TS) and local baselines (Seasonal Naive, LOCF). Supervised forecasting models are omitted, as prior large-scale studies indicate that TFSMs generally outperform them on established benchmarks [1]. TimesFM 2.5 [11] and Moirai2 [27] are excluded from evaluations due to substantial data leakage with `fev-bench`. (ii) In the *known-covariate setting*, we evaluate both TFSMs and TFMs un-

260 der the unified protocol described above. Methods supporting covariates (TS-ICL, Chronos-2,
 261 TabPFNv2.5-TS, TabICLv2-TS) are evaluated with and without covariates on the 30 relevant
 262 datasets, while covariate-agnostic TSMs are included unchanged to quantify the benefit of incorpo-
 263 rating exogenous information. Figure 5 reports the results following the *fev-bench* protocol, using
 264 Mean Absolute Scaled Error (MASE) and CRPS (metric definitions provided in Appendix D) while
 265 Table 3 reports the median inference time.

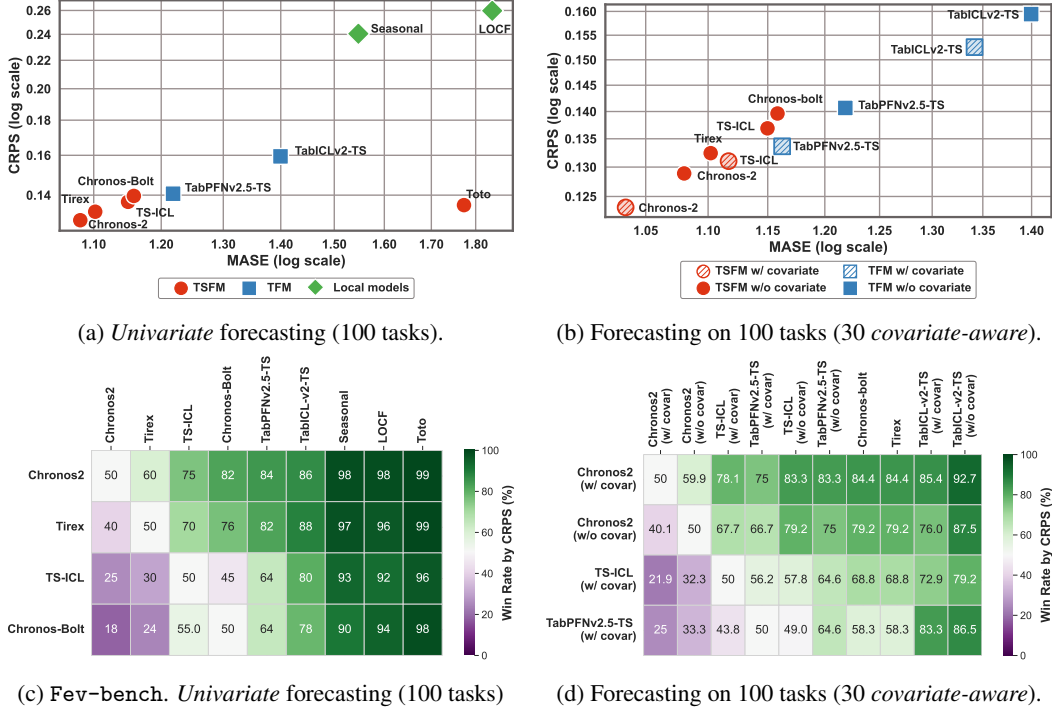


Figure 5: *fev-bench* forecasting benchmark. (a)-(b) MASE-CRPS trade-off (lower is better). Points correspond to per-method scores averaged across tasks. (c)-(d) Pairwise win rates. Each entry indicates the fraction of tasks where a method outperforms another according to the CRPS.

266 **Results.** TS-ICL achieves strong zero-shot performance on *fev-bench*, remaining competitive
 267 with leading TSMs while consistently outperforming TFMs and local baselines on both point and
 268 probabilistic metrics.

269 (i) In the *univariate setting*, TS-ICL ranks among the top-performing methods (Figure 5a), re-
 270 maining within $\sim 6\%$ of Chronos-2 and $\sim 3\%$ of TiRex, while outperforming all other baselines,
 271 including TFMs and local methods. Pairwise comparisons (Figure 5c) confirm consistent majority
 272 wins across tasks against tabular and local approaches. Finally, it offers a strong accuracy-efficiency
 273 trade-off (Table 3), with inference time on the order of 10^{-2} seconds per window, around $4\times$ slower
 274 than Chronos2 but still about $40\times$ faster than TFMs.

275 (ii) In the *known-covariate setting*, performance improves with exogenous information (Figure 5b).
 276 Chronos2 remains the strongest overall method, while TS-ICL benefits from covariates, as illus-
 277 trated qualitatively in Figure 6, and consistently outperforms TFMs under identical inputs. It also
 278 improves its relative ranking among TSMs, surpassing TiRex in CRPS and achieving $\sim 70\%$ pair-
 279 wise win rates (Figure 5d), indicating stable gains on *fev-bench* when leveraging covariates.

Table 3: Median forecasting inference time on *fev-bench* (*univariate setting*) with a H100 GPU.

	Time Series Foundation Models				Tabular Foundation Models		Local Methods	
	TS-ICL	Chronos-2	TiRex	TOTO	TabPFNv2.5-TS	TabICLv2-TS	S-Naive	LOCF
Median Inference time (s / window)	1.54×10^{-2}	3.53×10^{-3}	5.20×10^{-2}	1.28×10^{-1}	4.33×10^{-1}	3.76×10^{-1}	3.09×10^{-4}	1.73×10^{-4}

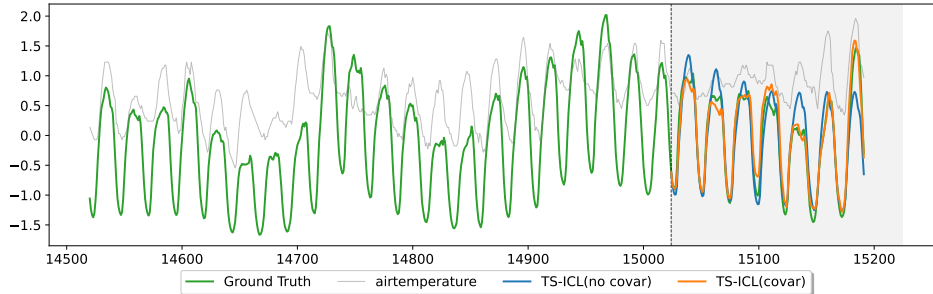


Figure 6: TS-ICL forecast on an horizon of length 168, with one additional covariate, on *GFC17*.

280 **Forecasting with missing values.** Beyond controlled evaluation settings on *fev-bench*, evaluat-
 281 ing TS-ICL robustness in realistic scenarios with partially missing historical observations is crucial.
 282 Table 4 compares TS-ICL and Chronos-2 forecasting performance under increasing levels of look-
 283 back window missingness (30%–90%) across the entire *fev-bench* benchmark. While both models
 284 degrade as the context becomes sparser, TS-ICL consistently outperforms Chronos-2, with signif-
 285 icant gaps at all missingness levels. As a result, TS-ICL remains substantially more robust for
 286 forecasting under missing historical observations.

Table 4: Zero-shot forecasting MASE under increasing missingness in the look-back window for TS-ICL and Chronos-2 on *fev-bench* (100 univariate tasks, look-back = 4092, arithmetic mean). Relative degradation (%) is reported in parentheses. `Seasonal Naive` serves as a simple baseline for evaluating forecasting performance under degraded conditions. Best results are in **bold**.

	0 % missing	30 % missing	50 % missing	70 % missing	90 % missing
Chronos-2	1.62 (0%)	2.16 (-33%)	2.44 (-50%)	2.54 (-56%)	3.97 (-144%)
TS-ICL	1.70 (0%)	1.77 (-4%)	1.89 (-11%)	2.16 (-27%)	3.63 (-113%)
Seasonal(0% missing)	2.48	2.48	2.48	2.48	2.48

287 Overall, TS-ICL is thus a competitive non patch-based TSFM for zero-shot forecasting. It remains
 288 close to state-of-the-art TSFMs such as Chronos-2 and TiRex across standard benchmarks, while
 289 effectively leveraging covariates when available. Its most notable advantage lies in its robustness to
 290 missing history, where it consistently outperforms Chronos-2, highlighting the sensitivity of patch-
 291 based models to incomplete context. In contrast, TS-ICL benefits from its time-indexed formulation
 292 to tackle realistic forecasting settings. Additional analyses and forecasting plots are provided in
 293 Appendix F.1, while Appendix F.2 reports leakage-free comparisons against 12 TSFMs on the TIME
 294 benchmark [33] across 98 zero-shot tasks.

295 6 Conclusion

296 TS-ICL is a flexible probabilistic time series foundation model based on an in-context regression
 297 formulation, unifying forecasting and imputation within a continuous-time framework. It sets new
 298 state-of-the-art performance on zero-shot imputation, while remaining competitive with leading
 299 TSFMs on forecasting benchmarks and serving as a strong alternative to patch-based models. A
 300 key advantage of TS-ICL lies in its robustness to incomplete historical observations: it consistently
 301 outperforms Chronos-2 under partially observed look-back windows, highlighting the benefits of
 302 its time-indexed formulation in realistic forecasting settings. It also enables efficient covariate-aware
 303 inference, further improving performance when exogenous information is available. Despite these
 304 strengths, TS-ICL exhibits higher inference cost than highly optimized models such as Chronos-2
 305 (up to $4\times$ slower despite having $4\times$ fewer parameters), mainly due to its pointwise regression for-
 306 mulation, which increases computational cost during both training and inference. These limitations
 307 may be mitigated through architectural optimizations such as caching or mixed-precision training.
 308 In addition, further increasing the diversity of the training data prior may improve the zero-shot gen-
 309 eralization of TS-ICL, as observed in tabular foundation models [18, 34]. Finally, the flexibility of
 310 the TS-ICL encoder–regressor framework makes it a natural foundation for extending beyond fore-
 311 casting and imputation to tasks such as zero-shot anomaly detection and time series classification.

References

- 312
- 313 [1] Taha Aksu, Gerald Woo, Juncheng Liu, Xu Liu, Chenghao Liu, Silvio Savarese, Caiming
314 Xiong, and Doyen Sahoo. GIFT-eval: A benchmark for general time series forecasting model
315 evaluation. In *NeurIPS Workshop on Time Series in the Age of Large Models*, 2024. URL
316 <https://openreview.net/forum?id=Z2cMO0ANFX>.
- 317 [2] Abdul Fatir Ansari, Lorenzo Stella, Ali Caner Türkmen, Xiyuan Zhang, Pedro Mercado,
318 Huibin Shen, Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda-Arango, Shub-
319 ham Kapoor, Jasper Zschiegner, Danielle C. Maddix, Hao Wang, Michael W. Mahoney, Kari
320 Torkkola, Andrew Gordon Wilson, Michael Bohlke-Schneider, and Bernie Wang. Chronos:
321 Learning the language of time series. *Trans. Mach. Learn. Res.*, 2024.
- 322 [3] Abdul Fatir Ansari, Oleksandr Shchur, Jaris Küken, Andreas Auer, Boran Han, Pedro Mercado,
323 Syama Sundar Rangapuram, Huibin Shen, Lorenzo Stella, Xiyuan Zhang, et al. Chronos-2:
324 From univariate to universal forecasting. *arXiv preprint arXiv:2510.15821*, 2025.
- 325 [4] Andreas Auer, Patrick Podest, Daniel Klotz, Sebastian Böck, Günter Klambauer, and Sepp
326 Hochreiter. Tirez: Zero-shot forecasting across long and short horizons with enhanced in-
327 context learning. In *The Thirty-ninth Annual Conference on Neural Information Processing*
328 *Systems*, 2025. URL <https://openreview.net/forum?id=v7UqniC9pF>.
- 329 [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-
330 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language mod-
331 els are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901,
332 2020.
- 333 [6] Wei Cao, Dong Wang, Jian Li, Hao Zhou, Yitan Li, and Lei Li. BRITS: bidirectional recurrent
334 imputation for time series. In *Advances in Neural Information Processing Systems*, volume 31,
335 2018.
- 336 [7] Mouxiang Chen, Lefei Shen, Zhuo Li, Xiaoyun Joy Wang, Jianling Sun, and Chenghao
337 Liu. VisionTS: Visual masked autoencoders are free-lunch zero-shot time series forecast-
338 ers. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=5DSj3MfWrB>.
- 340 [8] Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary
341 differential equations. In *Advances in Neural Information Processing Systems*, volume 31,
342 2018.
- 343 [9] James S Clark and Ottar N Bjørnstad. Population time series: process variability, observation
344 errors, missing values, lags, and hidden states. *Ecology*, 85(11):3140–3150, 2004.
- 345 [10] Ben Cohen, Emaad Khwaja, Youssef Doubli, Salahidine Lemaachi, Chris Lettieri, Charles
346 Masson, Hugo Miccinilli, Elise Ramé, Qiqi Ren, Afshin Rostamizadeh, Jean Ogier du Ter-
347 rail, Anna-Monica Toon, Kan Wang, Stephan Xie, Zongzhe Xu, Viktoriya Zhukova, David
348 Asker, Ameet Talwalkar, and Othmane Abou-Amal. This time is different: An observability
349 perspective on time series foundation models. In *The Thirty-ninth Annual Conference on Neu-
350 ral Information Processing Systems*, 2025. URL <https://openreview.net/forum?id=1jDAYXfcS2>.
- 352 [11] Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation
353 model for time-series forecasting. In *Forty-first International Conference on Machine Learn-
354 ing, ICML 2024, Vienna, Austria, July 21-27, 2024*, Proceedings of Machine Learning Re-
355 search, 2024.
- 356 [12] Samuel Dooley, Gurnoor Singh Khurana, Chirag Mohapatra, Siddhartha V Naidu, and Colin
357 White. ForecastPFN: Synthetically-trained zero-shot forecasting. In *Advances in Neural In-
358 formation Processing Systems*, volume 36, pp. 2403–2426, 2023.
- 359 [13] Wenjie Du, David Côté, and Yan Liu. SAITS: Self-attention-based imputation for time se-
360 ries. *Expert Systems with Applications*, 219:119619, 2023. doi: <https://doi.org/10.1016/j.eswa.2023.119619>.
- 361

- 362 [14] Wenjie Du, Yiyuan Yang, Linglong Qian, Jun Wang, and Qingsong Wen. PyPOTS: A Python
363 Toolkit for Machine Learning on Partially-Observed Time Series. *arXiv:2305.18811*, 2023.
- 364 [15] Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transform-
365 ers learn in-context? a case study of simple function classes. In S. Koyejo, S. Mohamed,
366 A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Process-*
367 *ing Systems*, volume 35, pp. 30583–30598. Curran Associates, Inc., 2022.
- 368 [16] Tilmann Gneiting, Fadoua Balabdaoui, and Adrian E Raftery. Probabilistic forecasts, calibra-
369 tion and sharpness. *Journal of the Royal Statistical Society Series B: Statistical Methodology*,
370 69(2):243–268, 2007.
- 371 [17] Noah Hollmann, Samuel Müller, Katharina Eggenberger, and Frank Hutter. TabPFN: A trans-
372 former that solves small tabular classification problems in a second. In *The Eleventh Interna-*
373 *tional Conference on Learning Representations*, 2022.
- 374 [18] Noah Hollmann, Samuel Müller, Lennart Purucker, Arjun Krishnakumar, Max Körfer, Shi Bin
375 Hoo, Robin Tibor Schirrmeyer, and Frank Hutter. Accurate predictions on small data with a
376 tabular foundation model. *Nature*, 637(8045):319–326, 2025.
- 377 [19] Shi Bin Hoo, Samuel Müller, David Salinas, and Frank Hutter. From tables to time:
378 How TabPFN-v2 outperforms specialized time series forecasting models. *arXiv preprint*
379 *arXiv:2501.02945*, 2025.
- 380 [20] Rob J Hyndman and George Athanasopoulos. *Forecasting: principles and practice*. OTexts,
381 2018.
- 382 [21] Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Car-
383 reira. Perceiver: General perception with iterative attention. In *International conference on*
384 *machine learning*, pp. 4651–4664. PMLR, 2021.
- 385 [22] Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz Cesista, Laker Newhouse, and
386 Jeremy Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. URL
387 <https://kellerjordan.github.io/posts/muon/>.
- 388 [23] Diederik P Kingma and Jimmy Lei Ba. Adam: A Method for Stochastic Optimization. In
389 *International Conference on Learning Representations*, 2015.
- 390 [24] Roger Koenker and Kevin F Hallock. Quantile Regression. *Journal of Economic Perspectives*,
391 15(4):143–156, 2001. doi: 10.1257/jep.15.4.143.
- 392 [25] Etienne Le Naour, Louis Serrano, Léon Migus, Yuan Yin, Ghislain Agoua, Nicolas Baskiotis,
393 Patrick Gallinari, and Vincent Guigüe. Time series continuous modeling for imputation and
394 forecasting with implicit neural representations. *Transactions on Machine Learning Research*,
395 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=P1vzXDklar>.
- 396 [26] Etienne Le Naour, Tahar Nabil, Adrien Petralia, and Ghislain Agoua. Are time-indexed foun-
397 dation models the future of time series imputation? *Transactions on Machine Learning Re-*
398 *search*, 2026. ISSN 2835-8856. URL <https://openreview.net/forum?id=cTk56KpsP5>.
- 399 [27] Chenghao Liu, Taha Aksu, Juncheng Liu, Xu Liu, Hanshu Yan, Quang Pham, Silvio Savarese,
400 Doyen Sahoo, Caiming Xiong, and Junnan Li. Moirai 2.0: When less is more for time series
401 forecasting. *arXiv preprint arXiv:2511.11698*, 2025.
- 402 [28] Yong Liu, Guo Qin, Zhiyuan Shi, Zhi Chen, Caiyin Yang, Xiangdong Huang, Jianmin
403 Wang, and Mingsheng Long. Sundial: A family of highly capable time series founda-
404 tion models. In *Forty-second International Conference on Machine Learning*, 2025. URL
405 <https://openreview.net/forum?id=L07ciRpjI5>.
- 406 [29] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi,
407 and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commu-*
408 *nications of the ACM*, 65(1):99–106, 2021.

- 409 [30] Vladyslav Moroshan, Julien Siems, Arber Zela, Timur Carstensen, and Frank Hutter. TempoPFN: Towards synthetic pre-training of linear RNNs for zero-shot time series forecasting. In *EurIPS 2025 Workshop: AI for Tabular Data*, 2025. URL <https://openreview.net/forum?id=lqex1gfnvc>.
410
411
412
- 413 [31] Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *International Conference on Learning Representations, ICLR*, 2023.
414
415
- 416 [32] Jonas Peters, Dominik Janzing, and Bernhard Scholkopf. *Elements of causal inference: foundations and learning algorithms*. MIT press, 2017.
417
- 418 [33] Zhongzheng Qiao, Sheng Pan, Anni Wang, Viktoriya Zhukova, Yong Liu, Xudong Jiang, Qingsong Wen, Mingsheng Long, Ming Jin, and Chenghao Liu. It’s TIME: Towards the next generation of time series forecasting benchmarks. *arXiv preprint arXiv:2602.12147*, 2026.
419
420
- 421 [34] Jingang Qu, David Holzmüller, Gaël Varoquaux, and Marine Le Morvan. TabICLv2: A better, faster, scalable, and open tabular foundation model. *arXiv preprint arXiv:2602.11139*, 2026.
422
- 423 [35] Yulia Rubanova, Ricky T. Q. Chen, and David Duvenaud. Latent odes for irregularly-sampled time series. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2019. Curran Associates Inc.
424
425
- 426 [36] Michael Schulz and Karl Stattegger. Spectrum: Spectral analysis of unevenly spaced paleoclimatic time series. *Computers & Geosciences*, 23(9):929–945, 1997.
427
- 428 [37] Louis Serrano, Thomas X Wang, Etienne Le Naour, Jean-Noël Vittaut, and Patrick Gallinari. Aroma: Preserving spatial structure for latent pde modeling with local neural fields. *Advances in Neural Information Processing Systems*, 37:13489–13521, 2024.
429
430
- 431 [38] Oleksandr Shchur, Abdul Fatir Ansari, Caner Turkmen, Lorenzo Stella, Nick Erickson, Pablo Guerron, Michael Bohlke-Schneider, and Yuyang Wang. fev-bench: A realistic benchmark for time series forecasting. *arXiv preprint arXiv:2509.26468*, 2025.
432
433
- 434 [39] Ingo Steinwart and Andreas Christmann. Estimating conditional quantiles with the help of the pinball loss. *Bernoulli*, 17(1), 2011. doi: 10.3150/10-BEJ267.
435
- 436 [40] Sean J Taylor and Benjamin Letham. Forecasting at scale. *The American Statistician*, 72(1): 37–45, 2018.
437
- 438 [41] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *Proceedings of the 40th International Conference on Machine Learning*, 2023.
439
440
441
- 442 [42] Gerald Woo, Chenghao Liu, Doyen Sahoo, Akshat Kumar, and Steven Hoi. Learning deep time-index models for time series forecasting. In *International Conference on Machine Learning*, pp. 37217–37237. PMLR, 2023.
443
444
- 445 [43] Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen Sahoo. Unified training of universal time series forecasting transformers. In *Forty-first International Conference on Machine Learning*, 2024.
446
447
- 448 [44] Xin Wu, Fei Teng, Xingwang Li, Ji Zhang, Qiang Duan, and Tianrui Li. Out-of-distribution generalization in time series: A survey. *Information Fusion*, pp. 104336, 2026.
449
- 450 [45] Shifeng Xie, Vasilii Feofanov, Jianfeng Zhang, Themis Palpanas, and Ievgen Redko. Cauker: Classification time series foundation models can be pretrained on synthetic data. In *The Fourteenth International Conference on Learning Representations*, 2026. URL <https://openreview.net/forum?id=xBW2FifswU>.
451
452
453

454 Appendices

455	A TS-ICL Architecture Details	14
456	A.1 The Time Series Encoder \mathcal{E}	14
457	A.2 The Channel Mixer \mathcal{M}	15
458	A.3 The Temporal Context Query Module \mathcal{C}	16
459	A.4 The In-Context Learning Regressor Module \mathcal{R}	16
460	B Training Details	19
461	B.1 Training Prior	19
462	B.2 Architecture Hyperparameters and Implementation Details	25
463	C Ablation Studies	28
464	C.1 Architecture Scaling	28
465	C.2 Component Ablations: Synergy between Encoder and Regressor	28
466	C.3 Covariate Management Strategies	29
467	D Evaluation Metrics	30
468	D.1 Metrics Definition	30
469	E Extended Imputation Experiments	32
470	E.1 Fm-impute-bench Benchmark	32
471	E.2 TIME Benchmark	38
472	F Extended Forecasting Experiments	40
473	F.1 Fev-bench Benchmark	40
474	F.2 TIME Benchmark	47

475 A TS-ICL Architecture Details

476 This section provides a detailed breakdown of the TS-ICL framework introduced in Section 3. Cast-
 477 ing both imputation and forecasting tasks as in-context regression problems over learned temporal
 478 representation, TS-ICL consists in four successive modules that transform raw observations into
 479 global and local context-aware representations used for prediction.

480 A.1 The Time Series Encoder \mathcal{E}

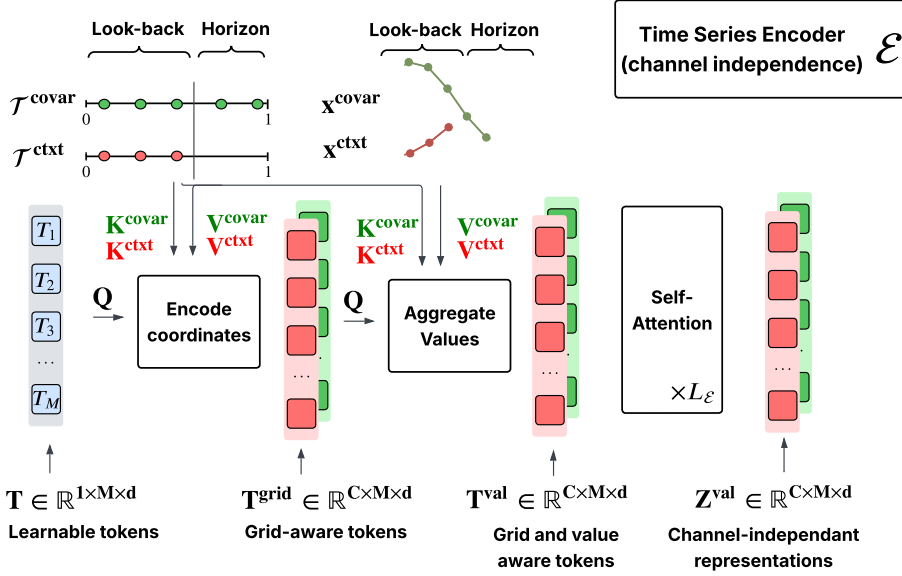


Figure 7: Overview of the Time Series Encoder \mathcal{E} . Forecasting task shown for illustration.

481 **Encoder Overview.** The encoder \mathcal{E} maps the observed context $(\mathcal{T}^{\text{ctxt}}, \mathbf{x}^{\text{ctxt}})$ and $C - 1$ optional
 482 covariates $(\mathcal{T}^{\text{covar}}, \mathbf{X}^{\text{covar}})$, where $C \geq 1$, jointly into a channel-independent latent representation
 483 $\mathbf{Z}^{\text{val}} \in \mathbb{R}^{C \times M \times d}$. Shown in Figure 7, the module operates through the following steps:

- 484 (i) **Temporal Encoding.** The timestamps from both the context $\mathcal{T}^{\text{ctxt}}$ and covariates $\mathcal{T}^{\text{covar}}$
 485 are independently mapped into higher-dimensional representations using Fourier features
 486 [29] followed by a linear projection:

$$\mathcal{T} \xrightarrow{\text{Fourier} + \text{Linear}} \gamma(\mathcal{T}) \in \mathbb{R}^{T \times d}.$$

- 487 (ii) **Coordinate Encoding.** A set of M learnable latent tokens $\mathbf{T} \in \mathbb{R}^{1 \times M \times d}$ serves as a
 488 query (Q) and attends to the temporal embeddings of all channels (context and covariates)
 489 through cross-attention:

$$\mathbf{T}^{\text{grid}} = \text{CrossAttn}(Q = \mathbf{T}, K = V = \gamma(\mathcal{T})) \in \mathbb{R}^{C \times M \times d}.$$

490 This step produces *grid-aware tokens* that capture the geometric structure of the sampling
 491 grid for each channel.

- 492 (iii) **Value Aggregation.** The observed values $(\mathbf{x}^{\text{ctxt}}, \mathbf{X}^{\text{covar}})$ are first projected into the latent
 493 dimension (*value lifting*). The grid-aware tokens \mathbf{T}^{grid} then attend to these value embed-
 494 dings:

$$\mathbf{T}^{\text{val}} = \text{CrossAttn}(Q = \mathbf{T}^{\text{grid}}, K = \gamma(\mathcal{T}), V = \text{lift}(\mathbf{X})) \in \mathbb{R}^{C \times M \times d}.$$

495 This step integrates the specific observed values into the latent representations, resulting in
 496 *grid and value aware tokens*.

497 (iv) **Latent Refinement.** Finally, $L_{\mathcal{E}}$ layers of channel-independent self-attention are applied
 498 to refine the representations:

$$\mathbf{Z}^{\text{val}} = \text{Transformer}_L(\mathbf{T}^{\text{val}}) \in \mathbb{R}^{C \times M \times d}.$$

499 This produces the final latent representations \mathbf{Z}^{val} , where each channel has been com-
 500 pressed into M informative tokens.

501 A.2 The Channel Mixer \mathcal{M}

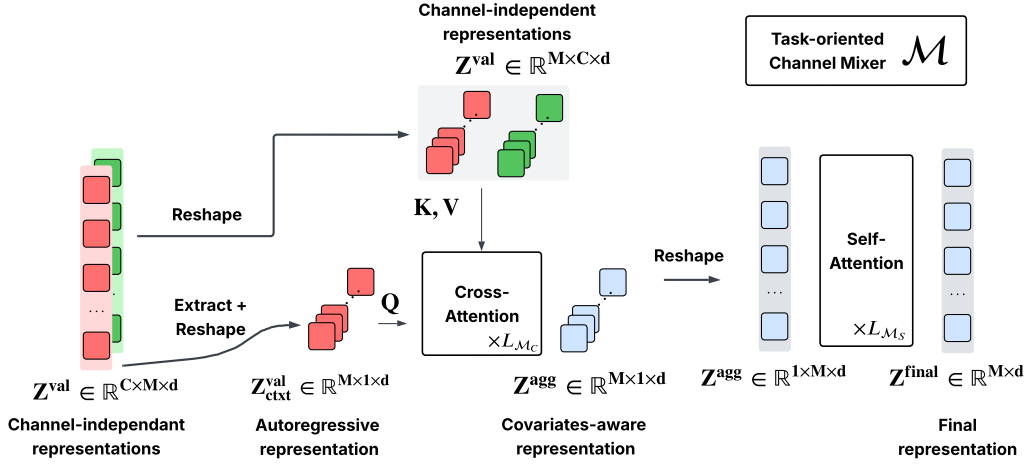


Figure 8: Overview of the Channel Mixer \mathcal{M} .

502 **Module \mathcal{M} Overview.** The Channel Mixer \mathcal{M} is designed to aggregate information across mul-
 503 tiple channels (time series of interest and covariates representations) by conditioning the channel-
 504 independent features on the target series context. It transforms a set of independent representations
 505 \mathbf{Z}^{val} into a unified, covariate-aware representation $\mathbf{Z}^{\text{final}}$.

506 Shown in Figure 8, the module follows a three-step process:

507 (i) **Cross-Channel Attention.** The autoregressive context representation $\mathbf{Z}^{\text{val}}_{\text{ctxt}} \in \mathbb{R}^{M \times 1 \times d}$,
 508 which represents the specific temporal dynamics of the target time series, acts as a query
 509 (Q). It attends to the channel-independent representations $\mathbf{Z}^{\text{val}} \in \mathbb{R}^{M \times C \times d}$ (where C is
 510 the number of channels), which serve as keys (K) and values (V):

$$\mathbf{Z}^{\text{agg}} = \text{CrossAttn}_L(Q = \mathbf{Z}^{\text{val}}_{\text{ctxt}}, K = V = \mathbf{Z}^{\text{val}}) \in \mathbb{R}^{M \times 1 \times d}.$$

511 This operation, repeated over $L_{\mathcal{M}_c}$ layers, compresses the multi-channel information into
 512 a single "covariates-aware" latent representation, effectively selecting the most relevant
 513 features from the covariates for the given context.

514 (ii) **Latent Reshaping.** To prepare the aggregated representation for global sequence process-
 515 ing, the tensor is reshaped to treat the M latent tokens as a sequence:

$$\mathbf{Z}^{\text{agg}} \in \mathbb{R}^{M \times 1 \times d} \xrightarrow{\text{reshape}} \mathbf{Z}^{\text{agg}} \in \mathbb{R}^{1 \times M \times d}.$$

516 (iii) **Global Latent Refinement.** Finally, $L_{\mathcal{M}_s}$ self-attention blocks are applied to the sequence
 517 of tokens. This allows the model to capture global dependencies across the aggregated
 518 latent space:

$$\mathbf{Z}^{\text{final}} = \text{Transformer}_L(\mathbf{Z}^{\text{agg}}) \in \mathbb{R}^{1 \times M \times d}.$$

519 The resulting $\mathbf{Z}^{\text{final}}$ is the final time series representation, integrating both the local channel
 520 information and the global context necessary for the downstream task.

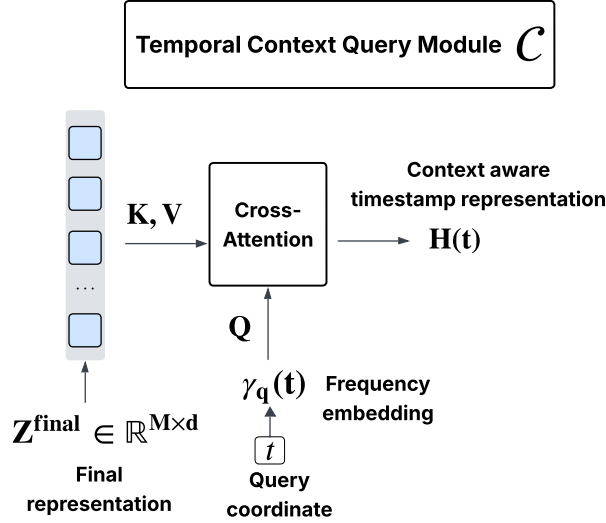


Figure 9: Overview of the Temporal Context Query Module \mathcal{C} .

521 A.3 The Temporal Context Query Module \mathcal{C}

522 **Module \mathcal{C} Overview.** The Temporal Context Query Module \mathcal{C} maps the final latent representation
 523 $\mathbf{Z}^{\text{final}} \in \mathbb{R}^{M \times d}$ and a query coordinate t to a time-series-aware representation $H(t) \in \mathbb{R}^d$. This
 524 module enables continuous-time querying of the encoded context by bridging the gap between the discrete latent tokens and the continuous time domain.
 525

526 Shown in Figure 9, the module operates as follows:

527 (i) **Frequency Encoding.** A target timestamp $t \in \mathbb{R}$ is mapped into a higher-dimensional
 528 frequency embedding $\gamma_q(t)$ [29]. This encoding uses sinusoidal functions at multiple scales
 529 to capture both coarse and fine-grained temporal patterns:

$$t \xrightarrow{\text{Frequency Encoding}} \gamma_q(t) \in \mathbb{R}^d.$$

530 (ii) **Contextual Querying via Cross-Attention.** The frequency embedding $\gamma_q(t)$ serves as the
 531 query (Q) in a cross-attention mechanism. It attends to the final time series representation
 532 $\mathbf{Z}^{\text{final}}$, which provides the keys (K) and values (V):

$$H(t) = \text{CrossAttn}\left(Q = \gamma_q(t), K = V = \mathbf{Z}^{\text{final}}\right).$$

533 This operation extracts a localized, time-series-aware representation from the global latent
 534 context, specifically conditioned on the query coordinate t .

535 (iii) **Representation Output.** The resulting vector $H(t)$ constitutes the "time-series-aware
 536 timestamp representation". It integrates the global context stored in the latent tokens with
 537 the specific temporal information of the query, serving as the primary input for the down-
 538 stream in-context regressor.

539 A.4 The In-Context Learning Regressor Module \mathcal{R}

540 **Module \mathcal{R} Overview.** The In-Context Learning Regressor \mathcal{R} is the final component of the archi-
 541 tecture. It treats both forecasting and imputation as in-context regression tasks, where the model
 542 learns to map target representations to values by conditioning on observed "input-output" pairs
 543 [5, 15, 41]. \mathcal{R} leverages a specific token construction mechanism to align context-aware embed-
 544 dings with raw covariates.

545 **Input Projection and Token Construction via Cross-Attention.** As depicted in Figure 10, prior
 546 to token construction, all raw inputs (including the observed values $\mathbf{x}_t^{\text{txt}}$ and covariates $\mathbf{x}_t^{\text{covar}}$) are

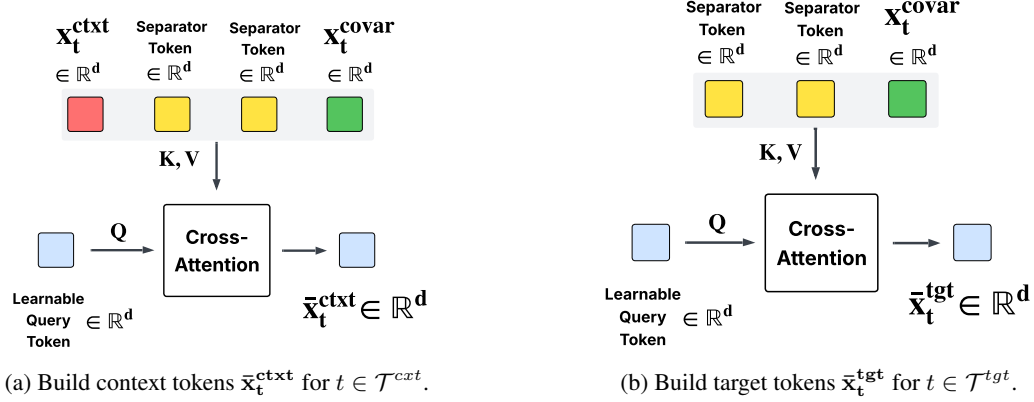


Figure 10: Cross-Attention Mechanism for Context and Target Token Construction.

547 linearly projected into a common d -dimensional latent space \mathbb{R}^d . A Cross-Attention mechanism
 548 then fuses these projected representations with a learnable query token $Q \in \mathbb{R}^d$.

549 Notably, the covariate input x_t^{covar} is strictly optional and can consist of zero, one, or multiple distinct
 550 covariates. The Cross-Attention mechanism accommodates this variability: since attention operates
 551 over sets, varying the number of covariates simply changes the sequence length of the Keys and
 552 Values, requiring no architectural modifications. This allows the regressor \mathcal{R} , and thus TS-ICL, to
 553 operate on covariate grids $\mathcal{T}_c^{\text{covar}}$ unaligned with the context grid $\mathcal{T}^{\text{ctxt}}$.

- 554 • **Context Tokens (\bar{x}_t^{ctxt}):** For $t \in \mathcal{T}^{\text{ctxt}}$, the model groups the projected observed value x_t^{ctxt}
 555 with a set of learned separator tokens and the projected covariates x_t^{covar} (if any). These act
 556 as Keys (K) and Values (V) for the learnable Query token Q , resulting in a unified context
 557 representation $\bar{x}_t^{\text{ctxt}} \in \mathbb{R}^d$.
- 558 • **Target Tokens (\bar{x}_t^{tgt}):** For $t \in \mathcal{T}^{\text{tgt}}$, the ground truth value is unknown. The target token
 559 $\bar{x}_t^{\text{tgt}} \in \mathbb{R}^d$ is similarly constructed by applying Cross-Attention between the learnable Query
 560 Q and the available information: the covariate embeddings x_t^{covar} and the separator tokens.

In-Context Input Sequence. The regressor processes a sequence \mathbf{S} organized to facilitate relational learning (Figure 11). The sequence consists of paired context tokens followed by target queries:

$$\mathbf{S} = \left[\underbrace{(\bar{x}_{t_1}^{\text{ctxt}}, \mathbf{H}(t_1)), \dots, (\bar{x}_{t_n}^{\text{ctxt}}, \mathbf{H}(t_n))}_{\mathcal{D}_{\text{train}}}, (\bar{x}_{t_{n+1}}^{\text{tgt}}, \mathbf{H}(t_{n+1})), \dots \right]$$

561 where $\mathbf{H}(t)$ represents the context-aware temporal embedding. All $(\bar{x}_t^{\text{ctxt}}, \mathbf{H}(t))$ pairs are summed to
 562 form the final regressor input sequence in the d -dimensional latent space.

563 **Causal In-Context Regression.** The sequence \mathbf{S} is processed by L layers of causal self-attention
 564 blocks. The causal mask is critical as it ensures a specific information flow:

- 565 • Each target token $(\bar{x}_{t_j}^{\text{tgt}}, \mathbf{H}(t_j))$ attends to all previous pairs in $\mathcal{D}_{\text{train}}$ to infer the underlying
 566 mapping $\mathbf{H}(t) \rightarrow \mathbf{x}(t)$.
- 567 • The attention mechanism allows the model to dynamically weigh past observations based
 568 on their similarity to the current target query in the representation space, without attending
 569 to future target values.

Quantile Prediction and Loss. To capture predictive uncertainty, for each target timestamp $t_j \in \mathcal{T}^{\text{tgt}}$, the model outputs 99 quantiles $\hat{\mathbf{q}}(t_j) = (\hat{q}_{\alpha_k}(t_j))_{k=1}^{99}$ via a linear projection of the final hidden states. The model is trained by minimizing a Smooth Pinball Loss:

$$\mathcal{L} = \sum_{t_j \in \mathcal{T}^{\text{tgt}}} \sum_{k=1}^{99} \left[\alpha_k e_{j,k} + \beta \log \left(1 + \exp(-e_{j,k}/\beta) \right) \right]$$

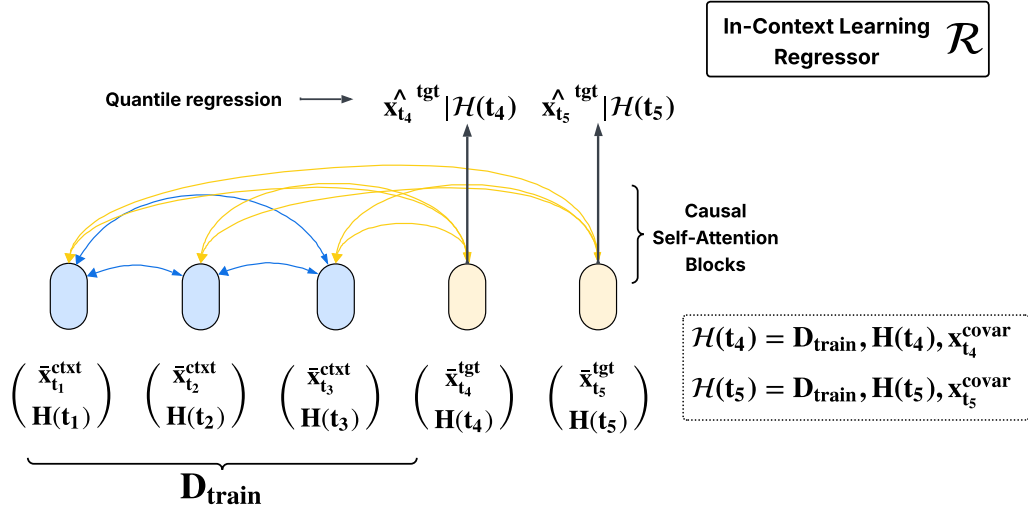


Figure 11: Overview of the In-Context Learning Regressor Module \mathcal{R} . Forecasting task shown for illustration.

570 where $e_{j,k} = x(t_j) - \hat{q}_{\alpha_k}(t_j)$, $\alpha_k \in (0, 1)$ is the quantile level, and $\beta > 0$ is a smoothing parameter.
 571 In practice, we set $\beta = 0.01$. During training, gradients are only backpropagated through the target
 572 predictions; the context set $\mathcal{D}_{\text{train}}$ is treated strictly as conditioning data and does not contribute to
 573 the loss.

574 **B Training Details**

575 This appendix provides detailed information about the training procedure of TS-ICL. Section B.1
 576 covers the data aspect, from pretraining corpus to prior generation, whereas Section B.2 describes
 577 the set of hyperparameters used to instantiate and train TS-ICL.

578 **B.1 Training Prior**

579 **B.1.1 Univariate Time Series Datasets.**

580 The univariate pretraining datasets of TS-ICL originate from three main sources, namely: (i) LOTSA
 581 [43]; (ii) Chronos training data [2] and (iii) TempoPFN synthetic data [30]. The latter includes in
 582 particular the synthetic generators from ForecastPFN [12], Chronos [2] and CauKer [45]. Overall,
 583 the pretraining corpus comprises 40 datasets, listed in Table 5 with their key features.

Table 5: All 40 univariate time series datasets used to pretrain TS-ICL and their key properties. The weight column reports the down- or upsampling coefficient applied dynamically at each epoch. †: offline downsampling from the original dataset.

Dataset	Release Platform	Domain	Freq	Num. Series	Num. Variates	Max Length	Weight
Australian Electricity	Chronos	Energy	30T	5	1	232,272	220
BDG-2 Bull	LOTSAs	Energy	H	41	1	12,280	25
BDG-2 Fox	LOTSAs	Energy	H	135	1	12,280	5
BDG-2 Panther	LOTSAs	Energy	H	105	1	6,132	2.5
BuildingsBench900k	LOTSAs	Energy	H	100,000†	1	8,759	0.02048
Residential Load Power	LOTSAs	Energy	1T	271	3	614,880	1.2
Residential PV Power	LOTSAs	Energy	1T	233	3	614,880	1.5
Wind Farms H	Chronos	Energy	H	337	1	6,148	4
Wind Farms D	Chronos	Energy	D	337	1	366	2
China Air Quality	LOTSAs	Climate	H	437	6	397,335	0.3
CMIP6 2000	LOTSAs	Climate	6H	8,192	22†	7,300	0.057
ERA5 1989	LOTSAs	Climate	H	8,192	15†	8,736	0.085
ERA5 1990	LOTSAs	Climate	H	8,192	15†	8,736	0.085
ERA5 1991	LOTSAs	Climate	H	8,192	15†	8,736	0.085
Spanish Weather	Kaggle	Climate	H	5	3	24,544	105
Subseasonal	LOTSAs	Climate	1D	862	4	16,470	0.3
Subseasonal Precipitation	LOTSAs	Climate	1D	862	1	11,323	1.2
Weatherbench daily	Chronos	Climate	1D	10,000†	1	14,610	0.1024
Mexico City Bikes	Chronos	Traffic	H	494	1	104,449	2.5
PEMS04	LOTSAs	Traffic	5T	307	3	16,992	1.2
PEMS07	LOTSAs	Traffic	5T	883	1	28,224	1.2
PEMS08	LOTSAs	Traffic	5T	170	3	17,856	2.1
Q-TRAFFIC	LOTSAs	Traffic	15T	45,148	1	5,856	0.024
Taxi (30 Min.)	Chronos	Traffic	30T	2,428	1	1,488	0.88
Taxi (Hourly)	Chronos	Traffic	H	2,428	1	744	0.88
Uber TLC (Hourly)	Chronos	Traffic	H	262	1	4,344	4
Alibaba Cluster Trace 2018	LOTSAs	Cloud	5T	58,409	2	1,728	0.009
Wiki Daily	Chronos	Web	D	100,000	1	2,741	0.00512
Monash M3 Monthly	LOTSAs	Econ./Fin.	M	1,428	1	126	0.72
NN5 Weekly	LOTSAs	Econ./Fin.	W	111	1	105	5
Project Tycho	LOTSAs	Health	W	1,258	1	3,854	0.21
Anomaly	TempoPFN	Synthetic	-	5,000	1	10,000	0.0256
ForecastPFN	TempoPFN	Synthetic	-	5,000	1	10,000	1
GP	TempoPFN	Synthetic	-	5,000	1	10,000	0.4096
Kernel Synth 1M	Chronos	Synthetic	-	1,000,000	1	1,024	0.001024
Sawtooth	TempoPFN	Synthetic	-	5,000	1	10,000	0.0512
Sinewave	TempoPFN	Synthetic	-	5,000	1	10,000	0.1024
Spikes	TempoPFN	Synthetic	-	5,000	1	10,000	0.0256
Step	TempoPFN	Synthetic	-	5,000	1	10,000	0.0512
OU	TempoPFN	Synthetic	-	5,000	1	10,000	0.4096

584 Table 5 highlights that our selection (i) spans *multiple domains*, including energy, nature/climate,
 585 transport, cloud, health/economics; (ii) covers a broad range of *frequencies*, from minutely- to
 586 weekly- and monthly-sampled time series; (iii) includes series of *varying context length*, from 126
 587 timesteps to 614k. In total, this forms a corpus of about 2M series, strictly non-overlapping with the
 588 different benchmarks on which TS-ICL is evaluated (fm-imputation, fev-bench and TIME).

589 **Sampling strategy.** We adopt a simple three-step stratified sampling strategy to encourage training
 590 on maximally diverse patterns. (i) Very large hourly datasets are downsampled offline, thereby
 591 reducing memory footprint: we use 100k random samples from *BuildingsBench900k*, 10k random
 592 samples from *Weatherbench daily* and 15 (resp., 22) representative channels out of 45 (resp. 53) for
 593 the *ERA5* (resp., *CMIP6*) datasets. (ii) Similarly to Moirai [43], a random subsample or upsample is
 594 drawn from each dataset at each training epoch, to avoid biases towards hourly datasets in the energy
 595 and climate domains. The sampling coefficients are reported in Table 5. (iii) For each sample, we
 596 select a single window drawn at random from the available series length.

597 B.1.2 Covariates Problem Generators.

598 In this section, we provide the technical specifications for the synthetic target-covariate(s) genera-
 599 tion process described in Section 4.

600 **Graph Construction Logic** The generator constructs a fixed Directed Acyclic Graph (DAG) for
 601 each time series batch. The construction follows a topological ordering to ensure acyclicity:

- 602 (i) **Node Initialization.** We start with R root nodes, where R is randomly drawn from a
 603 geometric distribution (favoring smaller values), and capped at a maximum of $R_{\max} = 3$.
 604 Each root corresponds to one of the base univariate time series in the pretraining corpus in
 605 Table 5 (e.g., signals from real-world datasets or TempoPFN).
- 606 (ii) **Graph Size.** The total number of generated series C (channels) corresponds to the num-
 607 ber of time series that will define the final covariates-target problem. Specifically, the
 608 constructed problem will consist of $C - 1$ covariate time series and 1 target time series.
 609 The value of C is sampled from a shifted exponential distribution to favor smaller, more
 610 manageable graphs while allowing for high complexity:

$$C = \min(C_{\max}, \lfloor \text{Exp}(\lambda) \rfloor + C_{\min}),$$

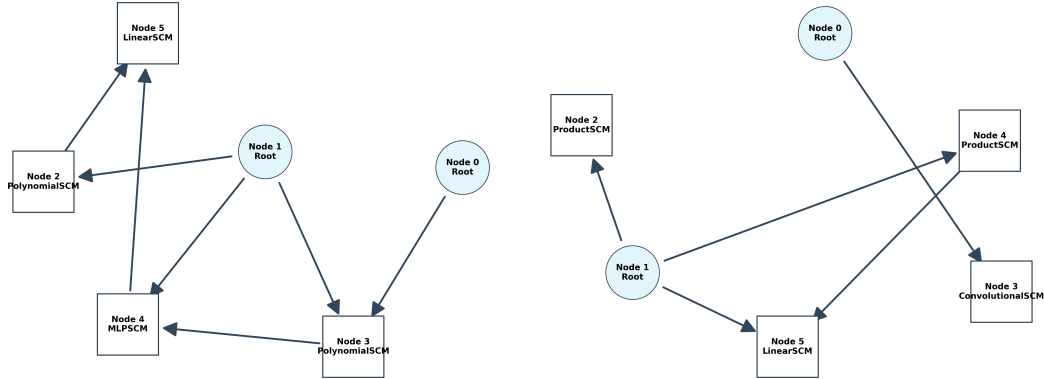
611 where we set $C_{\min} = 2$, $C_{\max} = 20$, and $\lambda = 0.8$ by default. The total number of nodes in
 612 the graph is $V = 2C + R$.

- 613 (iii) **Edge Sampling.** For each non-root node $i \in \{R, \dots, V - 1\}$, the number of parents k_i
 614 is sampled from a geometric distribution $k_i \sim \text{Geometric}(p)$, clipped to the number of
 615 available ancestors. Parents are then sampled uniformly without replacement from the set
 616 of all preceding nodes $\{0, \dots, i - 1\}$.
- 617 (iv) **Operator Assignment.** Each non-root node is assigned a Structural Causal Model (SCM)
 618 sampled from the Structural Causal Model registry.

619 This DAG structure allows for the generation of both dependent and independent child nodes, which
 620 is essential for constructing realistic target-covariate time series relationships. This design also
 621 ensures that covariates are not always informative, thereby encouraging models to learn to ignore
 622 irrelevant inputs when appropriate. As an illustration, Figure 12 displays two randomly generated
 623 graphs with 6 nodes each, including 2 channels (for one target and one covariate).

624 **Structural Causal Model (SCM) Registry.** To ensure a wide variety of functional relationships,
 625 we implement a set of diverse SCMs inspired by [34]. Let $\mathbf{Xpa}(i)$ denote the collection of parent
 626 time series for node i . The child node Y_i is generated as $Y_i = \text{Normalize}(f(\mathbf{Xpa}(i)))$, where
 627 $f \in \text{Registry}$. The registry includes:

- 628 • **LinearSCM.** A simple linear combination $Y = \mathbf{WX} + b$.
- 629 • **MLPSCM.** A multi-layer perceptron with random depth (2–10 layers) and hidden dimen-
 630 sions (8–128). Activations are randomly selected for each layer (ReLU, Tanh, ELU, etc.).
 631 We use a sparsity-inducing "block-wise dropout" initialization to create specific feature-
 632 group dependencies.



(a) Fully dependent structure (all nodes are connected). (b) Partially independent structure (not all nodes are connected).

Figure 12: Examples of randomly generated DAG structures with six nodes. The left graph enforces strong dependencies between nodes, while the right graph allows for partial independence, leading to more diverse causal structures.

- 633 • **ConvolutionalSCM.** Models local temporal dependencies using 1D convolutions with random kernel sizes (3–8) and random channel depths.
- 634
- 635 • **RNNSCM.** Captures deep temporal dependencies using a GRU architecture. These are strictly causal, ensuring the value at time t only depends on $t' \leq t$.
- 636
- 637 • **PolynomialSCM.** Each input is raised to a random power $d \in \{1, 2, 3, 4\}$ before being linearly combined, inducing symmetric non-linearities.
- 638
- 639 • **DiscretizeSCM.** Simulates quantization effects by mapping a linear mixture of inputs into a fixed number of discrete bins (2 to 15).
- 640
- 641 • **ProductSCM.** Computes the element-wise product of all parent signals, representing multiplicative interactions.
- 642

643 For visual illustrations of the types of dependencies induced by each SCM, we refer the reader to
644 Figure 13.

645 **Normalization and Stability.** To prevent numerical instability and exploding values across deep
646 DAGs, every node’s output is z-normalized. This ensures that every generated time series Y_i main-
647 tains a mean of 0 and a standard deviation of 1 before it is used as an input for further child nodes.

648 **Problem Formulation (Target and Covariates).** Once the DAG is computed, we finalize the
649 target–covariate problem by:

- 650 (i) Sampling a subset of C nodes from the graph to be “visible” to the model.
- 651 (ii) Randomly designating one of these nodes as the Target (y).
- 652 (iii) Designating the remaining $C - 1$ nodes as Covariates (\mathbf{x}).

653 To encourage learning of informative covariate-target relationships, we first seek to fulfill steps (ii)
654 and (iii) by looking for connected components of the DAG. We thus sample the C channels from the
655 subset of nodes that have all, or at least one, root node as common ancestor.

- 656 • This approach ensures that covariates are not merely “noise” but share a common underlying
657 causal structure with the target, sometimes acting as direct causes, sometimes as effects,
658 and sometimes as siblings sharing a latent root cause.
- 659 • If this subset is empty or contains less than C nodes, the remaining nodes are drawn uni-
660 formly at random within the entire graph, allowing for independent or unrelated (covariate,
661 target) pairs.

662 Examples of such multivariate problems drawn from the data prior are shown in Figure 14.

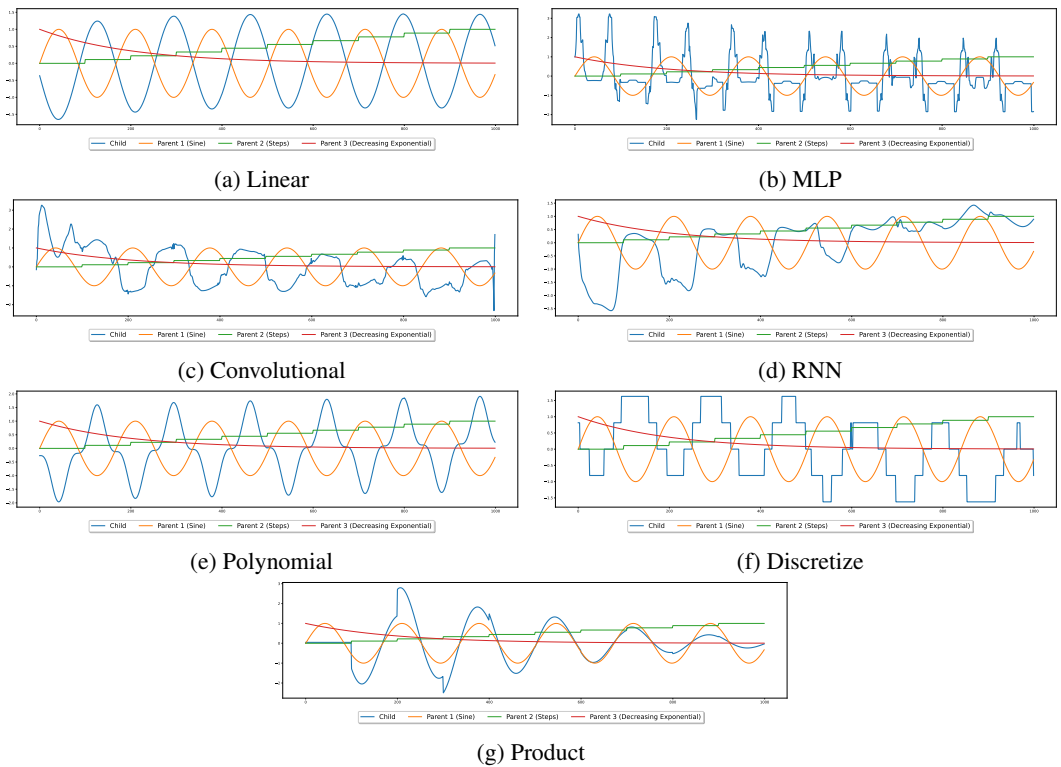


Figure 13: Examples of time series generated by different Structural Causal Models (SCMs) from the registry. Each plot shows the transformation of three root signals (sinusoidal, linear trend, and exponential decay) into a child time series through the corresponding operator.



Figure 14: Examples of different covariate-target time series sampled from the data prior.

663 **B.1.3 The whole procedure**

664 We summarize the full data generation pipeline in Algorithm 1. The procedure samples het-
 665 erogeneous training tasks combining real-world time series, synthetic time series and target-
 666 covariates time series problems. This enables the model to learn both standard univariate forecast-
 667 ing/imputation tasks and more complex multivariate settings generated via structural causal models.

668 A Bernoulli switch controlled by probability π determines whether a task is sampled from the uni-
 669 variate or covariate (SCM-based) mode.

Algorithm 1 Data Prior Sampling Procedure for TS-ICL

Require: Collection of base univariate time series \mathcal{D} (real + synthetic)

Require: SCM registry Υ

Require: Probability π of sampling a univariate task

```

1: Sample  $u \sim \mathcal{U}(0, 1)$ 
2: if  $u < \pi$  then
3:   Univariate mode:
4:   Sample time series  $x \sim \mathcal{D}$ 
5:   Apply task masking (imputation or forecasting)
6:   return  $x$ 
7: else
8:   Covariate mode:
9:   Sample number of task channels:
10:     $C \sim \text{Exp}(\lambda)$ , clipped to  $[C_{\min}, C_{\max}]$ 
11:   Sample number of root signals:
12:     $R \sim \text{Geometric}(q)$ , clipped to  $R_{\max}$ 
13:   Define latent DAG size:
14:     $V = 2C + R$ 
15:   Sample root time series:
16:     $S_1, \dots, S_R \sim \mathcal{D}$ 
17:   Initialize DAG nodes:
18:     $\{X_i\}_{i=1}^R \leftarrow \{S_i\}_{i=1}^R$ 
19:   DAG construction
20:   for  $i = R + 1, \dots, V$  do
21:     Sample number of parents  $k_i \sim \text{Geometric}(p)$ 
22:     Sample parent set  $\mathcal{P}_i \subset \{1, \dots, i - 1\}$ 
23:     Sample SCM  $f_i \sim \Upsilon$ 
24:     Compute:
25:      $X_i \leftarrow \text{Normalize}(f_i(\mathbf{X}_{\mathcal{P}_i}))$ 
26:   end for
27:   Observation step
28:   Sample observed set  $\mathcal{O} \subset \{1, \dots, V\}$  such that  $|\mathcal{O}| = C$ 
29:   Sample target node  $y \sim \mathcal{O}$ 
30:   Define covariates  $\mathbf{x} = \mathcal{O} \setminus \{y\}$ 
31:   Task masking
32:     - imputation: random point/block masking
33:     - forecasting: future horizon masking
34:   return  $(\mathbf{x}, y)$ 
35: end if

```

670 **B.2 Architecture Hyperparameters and Implementation Details**

671 The TS-ICL architecture is governed by a set of hyperparameters controlling model capacity, tok-
672 enization granularity, and depth across modules. All latent vectors across the four modules share a
673 common dimensionality d , ensuring seamless information flow.

674 **B.2.1 Architectural Hyperparameters.**

675 **Time Series Encoder \mathcal{E} hyperparameters.** The encoder serves as the primary feature extractor,
676 compressing multi-channel time series into a fixed-size latent bottleneck - see Figure 7.

- 677 • **Latent dimension (d):** The feature dimensionality used throughout the encoder is set to
678 $d = 256$.
- 679 • **Temporal Encoding:** Timestamps are mapped using Fourier features in logarithmic scale
680 with base 2, with $L_{\text{fourier}} = 128$ frequencies up to maximum frequency 2^{10} , followed by a
681 linear projection layer to \mathbb{R}^d .
- 682 • **Number of latent tokens (M):** The number of learnable tokens per channel is $M = 32$.
683 This parameter controls the resolution of the latent representation.
- 684 • **Refinement block:** The refinement stage consists of $L_{\mathcal{E}} = 3$ channel-independent Trans-
685 former blocks, each with $n_h^{\mathcal{E}} = 8$ self-attention heads of dimension 64.

686 **Channel Mixer Module \mathcal{M} hyperparameters.** The mixer facilitates task-oriented channel mix-
687 ing by conditioning the target series on covariates - see Figure 8.

- 688 • **Latent dimension (d):** The feature dimensionality used throughout the channel mixer is
689 set to $d = 256$.
- 690 • **Cross-Channel Attention:** The number of cross-attention layers used to aggregate
691 channel-independent tokens into a covariate-aware representation is $L_{\mathcal{M}}^{\text{cross}} = 3$. Each layer
692 contains $n_h^{\mathcal{M}} = 8$ attention heads of dimension 64.
- 693 • **Global Latent Refinement:** After reshaping, the latent sequence is processed by $L_{\mathcal{M}}^{\text{self}} = 3$
694 self-attention layers to capture dependencies across the aggregated latent space. Each layer
695 contains $n_h^{\mathcal{M}} = 8$ self-attention heads of dimension 64.

696 **Temporal Context Query Module \mathcal{C} hyperparameters.** This module acts as a continuous-time
697 interface between the latent tokens and the regressor - see Figure 9.

- 698 • **Latent dimension (d):** The feature dimensionality used throughout the context query mod-
699 ular is set to $d = 256$.
- 700 • **Frequency Encoding:** Query coordinates t are encoded using high-frequency sinusoidal
701 features (NeRF-style) to preserve fine-grained temporal localizations before being pro-
702 jected to \mathbb{R}^d . We use three frequency bands each with 128 frequencies and respective
703 maximum frequencies $2^6, 2^7, 2^{10}$.
- 704 • **Query mechanism:** A single 8×64 cross-attention layer is used to extract $H(t)$ from
705 $\mathbf{Z}^{\text{final}}$. This layer operates in parallel on the three frequency bands. After concatenating,
706 the context-aware time representation $H(t)$ has dimension $3 \times d = 768$.

707 **In-Context Learning Regressor \mathcal{R} hyperparameters.** The regressor performs the final predic-
708 tive task using a causal Transformer architecture - see Figures 10 & 11.

- 709 • **Latent dimension (d):** The feature dimensionality used throughout the in-context learning
710 regressor is set to $d = 512$.
- 711 • **Input Value Projection:** Before entering the cross-attention block for token construction,
712 all available observed values x_t are projected to \mathbb{R}^d via a linear layer. Similarly, if available,
713 covariates X_t are projected to \mathbb{R}^d via another linear layer, shared by all covariates.
- 714 • **ICL Transformer:** The causal sequence processing is performed by $L_{\mathcal{R}} = 12$ Transformer
715 layers, each with $n_h^{\mathcal{R}} = 8$ heads of dimension 64.

716 • **Quantile Head:** A final linear layer maps the d -dimensional hidden states of target tokens
 717 to a 99-dimensional vector representing the equidistant predicted quantiles.

718 **B.2.2 Training hyperparameters and task-specific strategies.**

719 This section provides additional details about the training strategy.

720 **Input scaling.** Following common practice in time series forecasting, raw time series are pre-
 721 processed by an instance-normalization layer, scaling each sample to zero mean and unit variance.
 722 We then apply a pointwise \sinh^{-1} transform to the standardized inputs, to stabilize training against
 723 outlier values [3].

724 **Task mixing and zero-shot robustness.** We employ a task mixing probability $\pi = 0.8$: 20%
 725 of the training batches are drawn from the univariate pretraining corpus in Table 5, whereas the
 726 remaining 80% are further processed by the SCM prior in Algorithm 1 to create target-covariates
 727 tasks (with $C - 1$ covariates, $C \geq 2$) or new univariate tasks ($C = 1$). This mixture ensures the
 728 model remains a robust univariate predictor while learning to leverage exogenous signals.

729 • **Covariate Complexity Sampling:** For covariate tasks, the number of channels C is (i)
 730 either $C = 1$ (univariate task) with probability 0.2 (ii) or sampled from a truncated expo-
 731 nential distribution:

$$P(K = k) \propto e^{-\lambda k}, \quad k \in \{2, \dots, 20\}.$$

732 We set $\lambda = 0.5$ to favor simpler tasks with few covariates while maintaining significant
 733 exposure to high-dimensional inputs (up to 20 covariates).

734 **Specialized Checkpoints.** While the architecture \mathcal{R} is identical for all tasks, we provide two spe-
 735 cialized checkpoints. The *Forecasting* checkpoint is trained strictly with right causal masking,
 736 whereas the *Imputation* checkpoint is trained to reconstruct missing values using both preceding
 737 and succeeding context.

738 **Imputation Training.** We enforce task diversity while training the imputation checkpoint with a
 739 two-step procedure.

740 1. **Window sampling:** a per-batch context length T sampled from multiple regimes to handle
 741 varying context lengths:

$$T \sim \begin{cases} \mathcal{U}(128, 336) & \text{with probability } p_1 = 0.15 \\ \mathcal{U}(512, 1024) & \text{with probability } p_2 = 0.6 \\ \mathcal{U}(1300, 1400) & \text{with probability } p_3 = 0.05, \\ \mathcal{U}(2000, 2100) & \text{with probability } p_4 = 0.1 \\ \mathcal{U}(4000, 4096) & \text{with probability } p_5 = 0.1 \end{cases}$$

742 where p_1, \dots, p_5 are dataset balancing probabilities. The maximum supported window
 743 length for imputation is $T = 4096$.

744 2. **Masking Procedure:** once a context window has been sampled, we then apply a random
 745 masking strategy. A fraction ρ of the observations are held out as target queries $\mathcal{T}^{\text{test}}$. The
 746 remaining points form $\mathcal{D}_{\text{train}}$.

- 747 • ρ is either a pointwise missingness rate, sampled randomly from
 748 $\{0.05, 0.06, \dots, 0.95\}$;
- 749 • or corresponds to up to four missing blocks of random length between 12 and 168
 750 (adjusted depending on available context length) and at most 50% pointwise missing
 751 values.

752 **Forecasting training.** A similar procedure is applied to the forecasting checkpoint;

753 1. **Window sampling:** the lookback and horizon pairs are sampled jointly:

$$(L, H) \sim \begin{cases} (\mathcal{U}(50, 200), \mathcal{U}(8, 20)) & \text{with probability } p_1 = 0.2 \\ (\mathcal{U}(200, 672), \mathcal{U}(18, 36)) & \text{with probability } p_2 = 0.15 \\ (\mathcal{U}(1000, 1100), \mathcal{U}(48, 336)) & \text{with probability } p_3 = 0.5 \\ (\mathcal{U}(2000, 2100), \mathcal{U}(48, 336)) & \text{with probability } p_4 = 0.1 \\ (\mathcal{U}(3062, 4096), \mathcal{U}(48, 672)) & \text{with probability } p_5 = 0.05 \end{cases} .$$

754 The maximum supported configuration is:

$$T_{\text{look-back}} \leq 4096, \quad T_{\text{horizon}} \leq 672.$$

755 2. **Masking Procedure:** We enforce robustness to irregularly sampled time series by remov-
 756 ing a fraction ρ of the observations in the lookback window, with probability 0.15. We
 757 draw ρ from the set $\{0.1, 0.25, 0.5, 0.75\}$.

758 **Quantile head.** The model predicts a set of 99 quantiles $\{\alpha_k\}_{k=1}^{99}$, uniformly spaced in $(0, 1)$,
 759 allowing for full density estimation. We set the smoothing coefficient of the pinball loss to $\beta = 0.01$
 760 to ensure differentiability.

761 B.2.3 Optimization Hyperparameters.

762 For both imputation and forecasting checkpoints, the optimization is configured as follows:

- 763 • **Optimizer:** *Muon* [22] for 2D parameters and AdamW for 1D parameters (embeddings,
 764 scales, biases).
- 765 • **Learning rate:** Max $lr = 4e - 4$.
- 766 • **Scheduler:** Cosine decay down to $5e - 4$.
- 767 • **Hardware:** $4 \times$ Nvidia H100 GPUs (92GB).
- 768 • **Batch size:** the global batch size is $B = 256$, obtained through a mini-batch size of 32 and
 769 two steps of gradient accumulation.
- 770 • **Training Budget:** The imputation checkpoint is trained for $500k$ optimization steps, rep-
 771 resenting approximately 5 training days. The forecasting checkpoint is trained for $650k$
 772 optimization steps, representing approximately 9 training days.

773 **C Ablation Studies**

774 In this section, we provide a comprehensive analysis of the architectural choices and scaling prop-
 775 erties of TS-ICL. All ablation models were trained for a fixed budget of 100k steps on a H100
 776 GPU to ensure fair comparison. Evaluation is performed on a subset of 11 datasets (44 tasks) of
 777 fm-impute-bench, namely fm-impute-mini, commonly used for ablations [26] (see Table 6).

Table 6: fm-impute-mini subset for zero-shot imputation ablation studies.

Dataset	Domain	Freq	Num. Series	Series Length	Num. Test Windows	Window Size
BDG2-Bear	Energy	1H	91	17544	7522	672
BDG2-Rat	Energy	1H	280	17544	24915	672
Covid19 Energy	Energy	1H	1	31912	195	672
GFC12 Load	Energy	1H	20	39414	4960	672
Hog	Energy	1H	24	17544	2310	672
Jena Weather 10T	Climate	10min	21	52704	1428	4032
Jena Weather 1H	Climate	1H	21	8784	1344	672
Oikolab Weather	Climate	1H	8	100057	5288	672
PDB	Energy	1H	1	17520	96	672
Pedestrian Counts	Transport	1H	66	96400	7733	672
Weather	Climate	1H	11	35064	2398	672

778 **C.1 Architecture Scaling**

779 We evaluate the impact of model capacity by varying (i) the number of attention heads
 780 $n_h = n_h^{\mathcal{E}} = n_h^{\mathcal{M}} = n_h^{\mathcal{C}}$ in the three Transformers of the encoder; (ii) the number $\mathcal{L}_{\mathcal{R}}$ of self-at-
 781 tention layers in the in-context regressor \mathcal{R} ; (iii) the number $n_h^{\mathcal{R}}$ of self-attention heads in \mathcal{R} . We
 782 define three model configurations: *Small*, *Medium*, and *Large*, with respectively **8.5M**, **12M** and
 783 **27M** parameters.

Table 7: Model Capacity Ablation. Average CRPS over 44 univariate imputation tasks from the fm-impute-mini benchmark (lower is better).

Model	n_h	$\mathcal{L}_{\mathcal{R}}$	$n_h^{\mathcal{R}}$	Params	fm-impute-mini
Small	4	8	4	~8.5M	0.201
Medium	8	12	4	~12M	0.197
Large	8	12	8	~27M	0.194

784 **Results.** Table 7 shows that increasing model capacity consistently improves imputation accuracy,
 785 with the *Large* configuration obtaining the best average CRPS. In particular, moving from *Small* to
 786 *Large* reduces the average CRPS from 0.201 to 0.194, suggesting that additional attention capacity
 787 in both the encoder and the in-context regressor improves the model’s ability to exploit contextual
 788 information across tasks. The gains, however, are relatively smooth and exhibit diminishing returns:
 789 the *Medium* model already closes a substantial fraction of the gap to the *Large* model, while using
 790 less than half the number of parameters. Moreover, the *Small* model remains competitive despite
 791 its lower capacity. Overall, this suggests a favorable accuracy–efficiency trade-off, with smaller
 792 variants remaining attractive under computational constraints.

793 **C.2 Component Ablations: Synergy between Encoder and Regressor**

794 To justify the hybrid structure of TS-ICL, we compare the full architecture against two baseline
 795 configurations:

- 796 • **Encoder-Only.** We remove the In-Context Regressor \mathcal{R} . The context-aware representation
 797 $H(t)$ from module \mathcal{C} is passed directly through a dense MLP network - with 5×256 hidden
 798 layers - to project onto the target values.
- 799 • **Regressor-Only (Pure ICL).** We remove the Encoder \mathcal{E} and the Mixer \mathcal{M} . The Regressor
 800 \mathcal{R} receives only raw temporal Fourier features [29] as representation $H(t)$, performing pure
 801 in-context regression without the benefit of the refined latent context.

Table 8: **Model Capacity Ablation.** Average CRPS over 44 univariate imputation tasks from the `fm-impute-mini` benchmark (lower is better).

Configuration	Architecture Change	fm-impute-mini
Encoder-Only	$H(t) \rightarrow$ MLP Head	0.297
Regressor-Only	No \mathcal{E} , Raw Fourier Features	0.204
Full TS-ICL	Encoder \mathcal{E} + Regressor \mathcal{R}	0.194

802 **Results.** Table 8 highlights the complementarity between the encoder and the in-context regressor.
 803 The *Encoder-Only* variant performs substantially worse, indicating that the latent context representa-
 804 tion alone is not sufficient without an expressive regression mechanism. Conversely, the *Regressor-*
 805 *Only* baseline is much stronger, but still underperforms the full model, showing that raw Fourier
 806 features provide a competitive ICL baseline but lack the refined context produced by the encoder.
 807 The full TS-ICL architecture achieves the best CRPS, suggesting that the encoder and regressor act
 808 synergistically: the encoder builds informative context-aware representations, while \mathcal{R} effectively
 809 uses them for in-context prediction.

810 C.3 Covariate Management Strategies

811 We investigate the optimal placement of the covariate mixing mechanism. Since TS-ICL allows for
 812 multi-stage conditioning, we compare three strategies for handling exogenous signals:

- 813 • **Early Mixing (Encoder Only):** Covariates are processed in the Encoder \mathcal{E} and mixed in
 814 the Mixer \mathcal{M} . The Regressor \mathcal{R} only receives $H(t)$ and context-target pairs of the main
 815 series, with no cross-attention on covariates during token construction.
- 816 • **Late Mixing (Regressor Only):** The Encoder \mathcal{E} is univariate (no covariates). All covariate
 817 information is provided directly to the Regressor \mathcal{R} via the Cross-Attention mechanism
 818 during input token construction.
- 819 • **Dual Mixing (Full Architecture):** Covariates are leveraged both in the global representa-
 820 tion (Encoder/Mixer) and for local token conditioning (Regressor).

Table 9: Covariate Mixing Strategy. Average CRPS over 6 covariate-aware imputation tasks from `fm-impute-covars` (lower is better).

Mixing Strategy	Mixing Stage	TMLR (Covar)
Early Mixing	$\mathcal{E} + \mathcal{M}$ only	0.123
Late Mixing	\mathcal{R} only	0.085
Dual Mixing	$\mathcal{E} + \mathcal{M}$ and \mathcal{R}	0.085

821 **Results.** The results in Table 9 show that late covariate mixing is already sufficient to obtain strong
 822 performance, with *Late Mixing* matching the *Dual Mixing* variant. This suggests that injecting co-
 823 variate information directly in the regressor \mathcal{R} provides effective pointwise conditioning at predic-
 824 tion time. In contrast, *Early Mixing* alone performs worse, indicating that global covariate infor-
 825 mation in the encoder is not sufficient without late-stage conditioning. We nevertheless retain the
 826 dual strategy, as early mixing may provide useful global context about covariate structure in harder
 827 settings, while late mixing supplies local, pointwise covariate information to the regressor.

828 **D Evaluation Metrics**

829 This section provides formal definitions for the evaluation metrics employed across the various
 830 experimental benchmarks presented in Section 5 and Appendices E, F. These metrics are designed
 831 to provide a comprehensive assessment of model performance, covering (i) computational efficiency
 832 during inference, (ii) the calibration and quality of the predicted probability distributions, and (iii /
 833 iv) the accuracy of point predictions via scale-independent error measures.

834 **D.1 Metrics Definition**

835 **(i) Inference efficiency score definition.** The computational efficiency is measured by the median
 836 inference time per window (as in [38]), expressed in milliseconds (ms). The calculation follows
 837 a two-step aggregation: first, the mean inference time is computed for each individual dataset to
 838 account for domain-specific variations; second, the median of these mean values is taken. For a
 839 collection of D datasets, where each dataset d contains M_d windows with individual inference times
 840 $\Delta t_{m,d}$, the efficiency score is defined as:

$$\mu_d = \frac{1}{M_d} \sum_{m=1}^{M_d} \Delta t_{m,d},$$

841 Efficiency = median($\{\mu_1, \dots, \mu_D\}$).

842 This procedure ensures that the final metric is representative of typical performance while remaining
 843 robust to outliers across different data distributions.

844 **(ii) Weighted Quantile Loss (WQL) and Continuous Ranked Probability Score (CRPS) defi-**
 845 **nitions.** To evaluate the quality of the predicted distribution, the Continuous Ranked Probability
 846 Score (CRPS) is employed, which measures the compatibility between the predicted cumulative dis-
 847 tribution function \hat{F} and the observed ground truth x . The CRPS can be expressed in its integral
 848 form as:

$$\text{CRPS}(\hat{F}, x) = \int_0^1 2 \cdot \text{QL}_\alpha(\hat{F}^{-1}(\alpha), x) d\alpha,$$

849 where $\text{QL}_\alpha(q, x)$ represents the quantile loss (or pinball loss) at level α :

$$\text{QL}_\alpha(q, x) = (\alpha - \mathbb{I}_{\{x < q\}})(x - q).$$

850 To ensure computational tractability and provide a normalized metric for cross-dataset comparison,
 851 a normalized discrete approximation of the CRPS is utilized, known as the Weighted Quantile Loss
 852 (WQL) [24, 16]. For a set of K discrete quantiles $\{\alpha_1, \dots, \alpha_K\}$, the WQL is calculated as:

$$\text{WQL} = \frac{1}{K} \sum_{j=1}^K \text{WQL}_{\alpha_j},$$

853 where each individual WQL_α is normalized by the absolute scale of the targets:

$$\text{WQL}_\alpha = \frac{2 \sum_{i,t \in \mathcal{T}^{\text{tgt}}} \text{QL}_\alpha(q_{i,t}^{(\alpha)}, x_{i,t})}{\sum_{i,t \in \mathcal{T}^{\text{tgt}}} |x_{i,t}|}.$$

854 In the evaluation, $K = 9$ equidistant quantiles $\alpha \in \{0.1, 0.2, \dots, 0.9\}$ are used, following standard
 855 practice, e.g. [38, 33]. This formulation allows the WQL to serve as a robust, scale-invariant proxy
 856 for the CRPS, capturing the accuracy of the entire predicted distribution.

857 *Important note.* For deterministic baselines such as Linear, Seasonal, or LOCF, the same point
 858 prediction is used across all quantile levels when computing the WQL.

859 **(iii) Normalized Mean Absolute Error (NMAE) definition.** To assess point prediction accuracy
 860 while accounting for differing scales across series, the Normalized Mean Absolute Error (NMAE)
 861 is used (as in [31]). This metric rescales the standard Mean Absolute Error (MAE) by the standard
 862 deviation of the observations in the context set.

863 For a series i and a target horizon \mathcal{T}^{tgt} , let $x_{i,t}$ be the ground truth and $\hat{x}_{i,t}$ the predicted median
 864 (quantile 0.5 for TS-ICL). The NMAE for series i is defined as:

$$\text{NMAE}_i = \frac{\frac{1}{|\mathcal{T}^{\text{tgt}}|} \sum_{t \in \mathcal{T}^{\text{tgt}}} |x_{i,t} - \hat{x}_{i,t}|}{\sigma_{i,\text{ctxt}}},$$

865 where $\sigma_{i,\text{ctxt}}$ is the standard deviation of the series i calculated over the observed context set $\mathcal{T}^{\text{ctxt}}$:

$$\sigma_{i,\text{ctxt}} = \sqrt{\frac{1}{|\mathcal{T}^{\text{ctxt}}|} \sum_{t \in \mathcal{T}^{\text{ctxt}}} (x_{i,t} - \bar{x}_{i,\text{ctxt}})^2}.$$

866 This normalization provides a scale-independent measure of the error relative to the inherent volatili-
 867 ty of the series. The global NMAE is obtained by averaging across all N series:

$$\text{NMAE} = \frac{1}{N} \sum_{i=1}^N \text{NMAE}_i.$$

868 **(iv) Mean Absolute Scaled Error (MASE) definition.** To evaluate point prediction accuracy
 869 across datasets with varying scales, the Mean Absolute Scaled Error (MASE) is used [20]. MASE
 870 normalizes the Mean Absolute Error (MAE) of the model by the mean absolute error of a seasonal
 871 naïve baseline.

872 For a series i , let $x_{i,t}$ be the ground truth and $\hat{x}_{i,t}$ the predicted median. The MASE for series i is
 873 defined as:

$$\text{MASE}_i = \frac{1}{|\mathcal{T}^{\text{tgt}}|} \sum_{t \in \mathcal{T}^{\text{tgt}}} \frac{|x_{i,t} - \hat{x}_{i,t}|}{a_i},$$

874 where a_i is the seasonal normalization factor calculated over the target set \mathcal{T}^{tgt} :

$$a_i = \frac{1}{|\mathcal{T}^{\text{tgt}}| - s} \sum_{t \in \mathcal{T}^{\text{tgt}}, t > s} |x_{i,t} - x_{i,t-s}|,$$

875 and s is the seasonal periodicity. This normalization ensures that the metric is scale-independent by
 876 comparing the model's error to the typical seasonal variations within the same target horizon.

877 The global metric is obtained by averaging across all N series:

$$\text{MASE} = \frac{1}{N} \sum_{i=1}^N \text{MASE}_i.$$

878 **E Extended Imputation Experiments**

879 This section provides broader insights into TS-ICL imputation performances. A detailed description
 880 of the `fm-impute-bench` benchmark used in the main text (Section 5.1) is given in Section E.1,
 881 together with complementary results and qualitative visualizations. Section E.2 further extends the
 882 evaluation to a second benchmark, TIME [33], which we adapt to the univariate imputation setting.

883 **E.1 Fm-impute-bench Benchmark**

884 **E.1.1 Datasets and Baselines**

885 **Univariate inference datasets.** Table 10 details the *univariate datasets* used for the zero-shot
 886 imputation experiments in Section 5.1. These datasets cover a diverse range of domains, including
 887 energy, transport, and climate science, with sampling frequencies varying from 5 minutes to 1 hour
 888 (specifically 5, 10, 15, 30, and 60 minutes). The imputation tasks are performed on four-week
 889 windows. When considering the four distinct missingness scenarios, this benchmark represents a
 890 large-scale evaluation involving approximately 1.3 million windows to be imputed.

Table 10: All datasets used for zero-shot imputation in the `fm-impute-bench` benchmark under the *univariate setting*.

Dataset	Release Platform	Domain	Freq	Num. Series	Series Length	Num. Test Windows	Window Size
BDG2-Bear	LOTSAs	Energy	1H	91	17544	7522	672
BDG2-Rat	LOTSAs	Energy	1H	280	17544	24915	672
Borealis	LOTSAs	Energy	1H	15	7447	77	672
Covid19 Energy	LOTSAs	Energy	1H	1	31912	195	672
GFC12 Load	LOTSAs	Energy	1H	20	39414	4960	672
Hog	LOTSAs	Energy	1H	24	17544	2310	672
Ideal	LOTSAs	Energy	1H	217	16167	156	672
PDB	LOTSAs	Energy	1H	1	17520	96	672
KDD Cup2022	LOTSAs	Energy	10min	134	35279	2546	4032
ERA5 geopotential	LOTSAs	Climate	1H	500	8736	19000	672
ERA5 humidity	LOTSAs	Climate	1H	500	8736	19000	672
ERA5 temperature	LOTSAs	Climate	1H	500	8736	19000	672
ERA5 wind speed	LOTSAs	Climate	1H	500	8736	19000	672
Oikolab Weather	LOTSAs	Climate	1H	8	100057	5288	672
Pedestrian Counts	LOTSAs	Transport	1H	66	96400	7733	672
Traffic	LOTSAs	Transport	1H	861	17544	83479	672
PEMS BAY	LOTSAs	Transport	5min	325	52128	2275	8064
PEMS 03	LOTSAs	Transport	5min	358	26208	358	8064
SHMETRO	LOTSAs	Transport	15min	576	8809	576	2688
ETT1-15T	GIFT-eval	Energy	15min	7	69680	1050	2688
ETT1-1H	GIFT-eval	Energy	1H	7	17420	1092	672
ETT2-15T	GIFT-eval	Energy	15min	7	69680	1050	2688
ETT2-1H	GIFT-eval	Energy	1H	7	17420	1092	672
Solar-1H	GIFT-eval	Energy	1H	137	8760	8768	672
Jena Weather 10T	GIFT-eval	Climate	10min	21	52704	1428	4032
Jena Weather 1H	GIFT-eval	Climate	1H	21	8784	1344	672
Loop Seattle 5T	GIFT-eval	Transport	5min	323	105120	21964	8064
Loop Seattle 1H	GIFT-eval	Transport	1H	323	8760	20672	672
MDense	GIFT-eval	Transport	1H	30	17520	4710	672
Enedis LDM Small	Zenodo	Energy	30min	500	17424	20500	1344
London Smart Meters Small	Chronos	Energy	30min	500	22000	25779	1344
Spanish Energy	Kaggle	Energy	1H	9	35064	1962	672
Weather	Informer	Climate	1H	11	35064	2398	672

891 **Inference datasets with covariates.** Table 11 details the six datasets used to evaluate zero-shot
 892 imputation with *exogenous covariates*. Following the protocol in Table 10, four-week windows are
 893 generated for these experiments. As described in [26], the PV and Wind datasets map regional re-
 894 newable energy production in 2021 to solar irradiance and wind speed, respectively. In contrast,
 895 Load-France tracks national electricity demand using average temperature as the primary covari-

896 ate to model consumption patterns. When considering the four distinct missingness scenarios, the
 897 covariate benchmark represents an evaluation involving approximately 1k windows to be imputed.

Table 11: All datasets used for zero-shot imputation in the `fm-impute-bench` benchmark under the *known-covariate setting*.

Dataset	Release Platform	Domain	Freq	Target / Covariate	Series Length	Num. Test Windows	Window Size
PV-OCC	RTE / Meteo	Energy	1H	1 / 1	8760	38	672
PV-PACA	RTE / Meteo	Energy	1H	1 / 1	8760	38	672
Wind-HDF	RTE / Meteo	Energy	1H	1 / 1	8760	38	672
Wind-GE	RTE / Meteo	Energy	1H	1 / 1	8760	38	672
Load-France 21	RTE / Enedis	Energy	30min	1 / 1	17520	41	1344
Load-France 22	RTE / Enedis	Energy	30min	1 / 1	17520	41	1344

898 **Baselines details.** A brief description of the baselines used in the benchmark is provided below.

- 899 • TabPFNv2.5-TS [19] is a time series foundation model that adapts the tabular founda-
 900 tion model TabPFN - a transformer-based model pretrained on synthetic supervised-
 901 learning tasks for in-context prediction of unseen tabular datasets [17] - to temporal data.
 902 TabPFN-TS leverages a TabPFN regression backbone and reformulates time-series predic-
 903 tion as an in-context tabular regression problem. Originally proposed for zero-shot fore-
 904 casting, TabPFN-TS still naturally applies to imputation: observed target values form the
 905 in-context training set, while missing timestamps are treated as query points, enabling the
 906 model to impute gaps using all available non-missing observations. For consistency with
 907 TabICLv2-TS, we use the same feature-generation pipeline as in the TabICLv2-TS pack-
 908 age: each timestamp is converted into a tabular row with temporal features, automatically
 909 extracted seasonal features, and, when available, exogenous covariates [34]. The pretrained
 910 TabPFN regressor is then queried in zero-shot to obtain point predictions. Our experiments
 911 use the TabPFNv2.5 regression checkpoint as the backbone, together with the official im-
 912 plementation: <https://github.com/PriorLabs/tabpfn-time-series>.
- 913 • Similarly to TabPFNv2.5-TS, TabICLv2-TS is a foundation model for time series analysis
 914 adapted from the TabICLv2 [34] tabular foundation model, designed for scalable in-context
 915 learning on regression and classification tasks. We use the package-provided time-series
 916 transformation pipeline to construct the tabular representation. TabICLv2 then performs
 917 regression by conditioning on the resulting table in a single in-context inference procedure.
 918 In our experiments, we use the TabICLv2 implementation and forecasting utilities from the
 919 official release: <https://github.com/soda-inria/tabicl>.
- 920 • Linear imputes missing values by linearly interpolating between the closest observed
 921 neighbors surrounding the gap. If a gap has no future (resp., past) anchor, it falls back to
 922 NOCB, next-observation-carried-backward (resp., LOCF, last-observation-carried-forward).
- 923 • Seasonal Naive imputes a missing value at timestamp t by repeating the observation
 924 from the previous seasonal period, i.e., the value at $t - S$. S is pre-defined for each dataset
 925 based on its dominant frequency (e.g., daily). If the value at $t - S$ is also missing, the
 926 method sequentially searches for an available observation at other seasonal timestamps
 927 (e.g., $t + S$, then $t - 2S$, etc.). The method falls back to LOCF in case this search fails.
- 928 • LOCF imputes a missing value by copying the most recent available past value.
- 929 • SAITS [13] is a supervised Transformer-based imputation model designed for partially ob-
 930 served multivariate time series. It uses two diagonally-masked self-attention blocks to
 931 capture temporal and cross-variable dependencies, whose outputs are combined through
 932 a learned gating mechanism to reconstruct missing values.
- 933 • BRITS [6] is a supervised recurrent imputation model based on bidirectional RNNs with
 934 learned temporal decay. It processes each window forward and backward to account for
 935 irregular gaps, jointly estimating hidden states and missing values while encouraging con-
 936 sistency between directions.

937 **Task specific baseline detailed training.** The two supervised baselines, SAITS and BRITS, are
 938 trained on the training split of `fm-impute-bench` (see [26] for more details). Both implementa-
 939 tions are taken from the PyPOTS toolbox [14] and trained on fixed-length windows with a masked-
 940 reconstruction objective: a subset of observed entries is randomly masked, and the model recon-
 941 structs these values from the remaining observations and the corresponding binary observation mask.
 942 Inputs are z-score normalized per variable, training minimizes MAE on the artificially masked po-
 943 sitions, and model selection is performed using validation MSE with early stopping. We use the de-
 944 fault hyperparameter configurations recommended by the original authors, the Adam optimizer [23],
 945 a batch size of 64, at most 50 training epochs, and early stopping with patience 5. Since SAITS and
 946 BRITS produce pointwise imputations, we adapt them to quantile-based evaluation by replicating
 947 each point prediction across all requested quantile levels.

948 E.1.2 Extended Results

949 This section extends the empirical evaluation in Section 5.1 with a more detailed analysis of impu-
 950 tation performance across all experimental settings. Specifically, we provide:

- 951 • **Aggregated detailed performance tables:** We report the average NMAE and CRPS (met-
 952 rics definition in Section D) across the 132 *univariate* tasks and 24 *covariate-aware* tasks
 953 of `fm-impute-bench`. These results, detailed in Table 12 and Table 13, provide a detailed
 954 view of both point-wise and probabilistic performance. It is important to note that metrics
 955 are aggregated across tasks using the arithmetic mean, following the evaluation protocol
 956 established in `fm-impute-bench` [26].
- 957 • **NMAE pairwise win rates:** To complement the CRPS-based win rate diagrams presented
 958 in the main text Figure 4, we include the corresponding pairwise win rate visualizations
 959 in terms of NMAE for both *univariate* and *known-covariate* experiments in Figure 15.
 960 The NMAE-based pairwise comparisons provide additional evidence of the robustness of
 961 TS-ICL, showing consistent superiority regardless of the chosen accuracy metric.

Table 12: Aggregated imputation metrics on the 132 tasks of the *univariate setting* in `fm-impute-bench` (mean \pm std). Best in **bold**.

	TSM	Tabular Foundation Models			Task Specific Models		Local Models		
	TS-ICL	TabPFNv2.5-TS	TabICLv2-TS	SAITS	BRITS	Linear interp.	Seasonal Naive	LOCF	
NMAE (\downarrow)	0.243 \pm 0.118	0.296 \pm 0.145	0.294 \pm 0.136	0.386 \pm 0.140	0.470 \pm 0.181	0.507 \pm 0.287	0.580 \pm 0.177	0.612 \pm 0.255	
CRPS (\downarrow)	0.255 \pm 0.137	0.303 \pm 0.156	0.301 \pm 0.148	0.503 \pm 0.193	0.605 \pm 0.227	0.658 \pm 0.377	0.750 \pm 0.230	0.793 \pm 0.335	

Table 13: Aggregated imputation performance metrics across the 24 tasks of the *known covariates setting* in `fm-impute-bench` (mean \pm std). Best in **bold**.

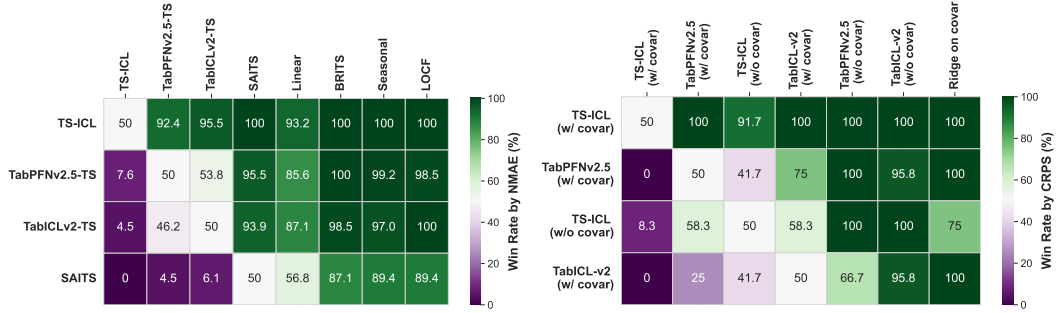
	TSM		Tabular Foundation Models				Local Model
	TS-ICL		TabPFNv2.5-TS		TabICLv2-TS		Ridge on Covar
	w/ covar	w/o covar	w/ covar	w/o covar	w/ covar	w/o covar	w/ covar
NMAE (\downarrow)	0.077 \pm 0.053	0.125 \pm 0.113	0.121 \pm 0.081	0.202 \pm 0.158	0.141 \pm 0.099	0.206 \pm 0.136	0.388 \pm 0.253
CRPS (\downarrow)	0.074 \pm 0.047	0.119 \pm 0.100	0.115 \pm 0.074	0.196 \pm 0.147	0.134 \pm 0.092	0.199 \pm 0.125	0.471 \pm 0.306

962 E.1.3 Qualitative Analysis and Visualizations

963 This section presents visual examples of TS-ICL imputations for both *univariate* and *known-*
 964 *covariate* settings. We illustrate in Figure 16 and Figure 17 the model imputation capabilities across
 965 various missingness patterns covering pointwise and blockwise scenarios.

966 **Results.** Several observations emerge from these plots.

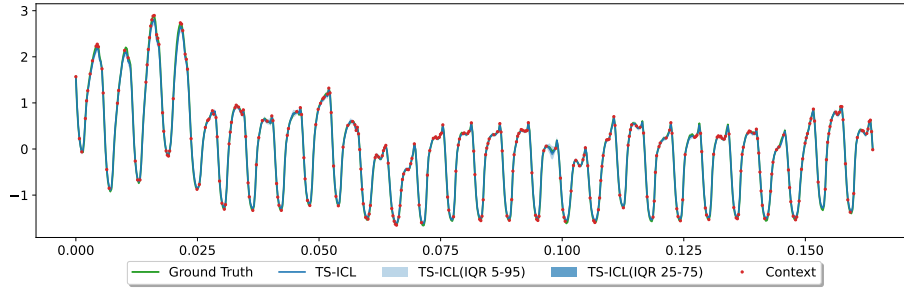
- 967 (i) With rich context, TS-ICL provides accurate reconstructions of smooth patterns, with tight
 968 inter-quantile ranges (Figure 16a). On the contrary, TS-ICL adjusts its uncertainty estima-
 969 tion of sparsely observed yet regular patterns (Figures 17b and 17e).
- 970 (ii) TS-ICL adapts well to distribution shifts (Figures 16b and 17c).



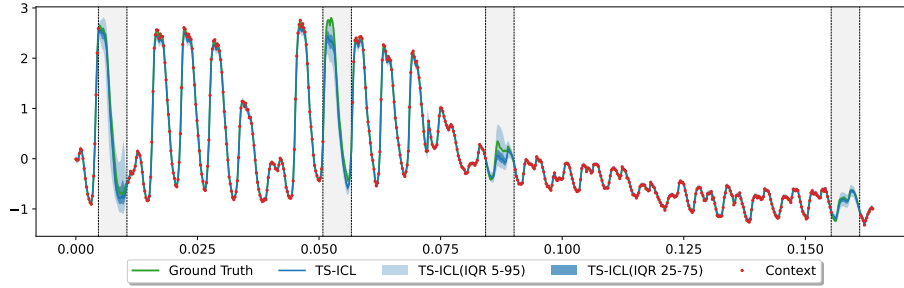
(a) Univariate imputation across 132 tasks. (b) Imputation with known covariates across 24 tasks.

Figure 15: Pairwise win rates of the top-4 models for imputation on the fm-impute benchmark. Each entry indicates the fraction of tasks where a method outperforms another according to the NMAE.

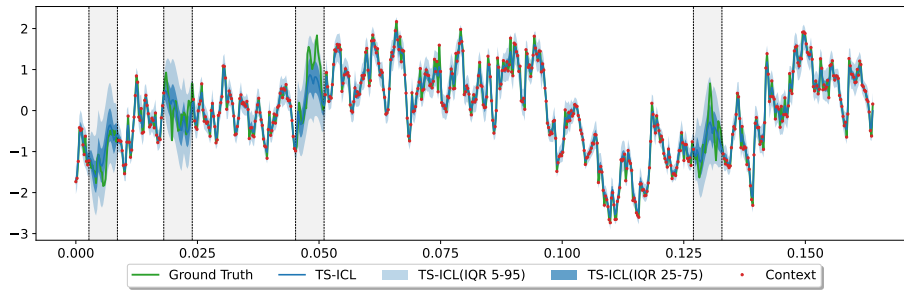
- 971 (iii) TS-ICL tends to provide smooth reconstructions of sparse high-frequency signals, with
 972 higher interquantile ranges accomodating their stronger variability (Figures 16c to 16e).
- 973 (iv) Figure 17a (third block) suggests that TS-ICL can serve as a counterfactual estimator, re-
 974 placing unusual sequences with expected ones under regular conditions.
- 975 (v) Finally, the ability of TS-ICL to incorporate covariate information at inference is key to
 976 produce meaningful imputations in challenging scenarios (Figure 17d).



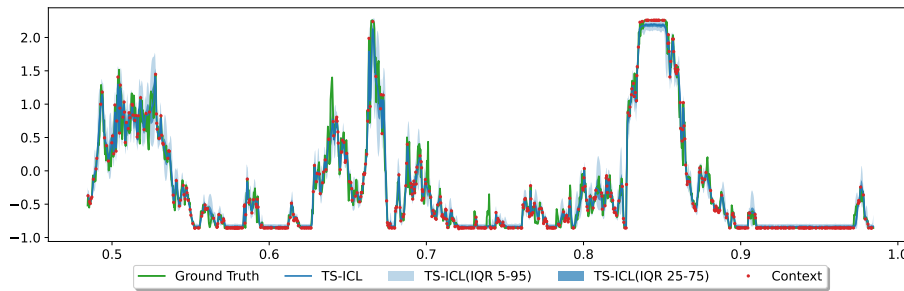
(a) *PDB*, 50% missing values.



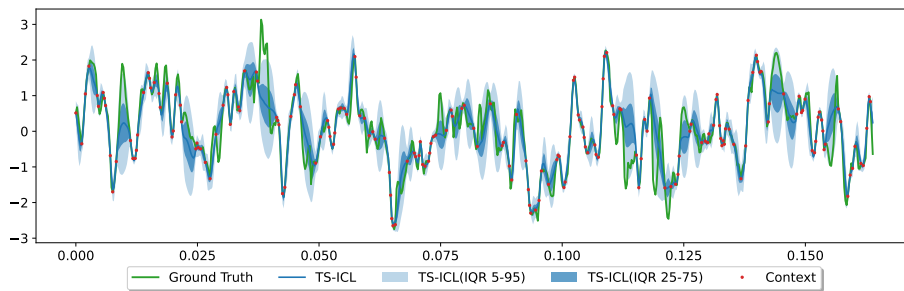
(b) *Covid19 Energy*, four one-day missing blocks.



(c) *BDG2-Rat*, four one-day missing blocks.

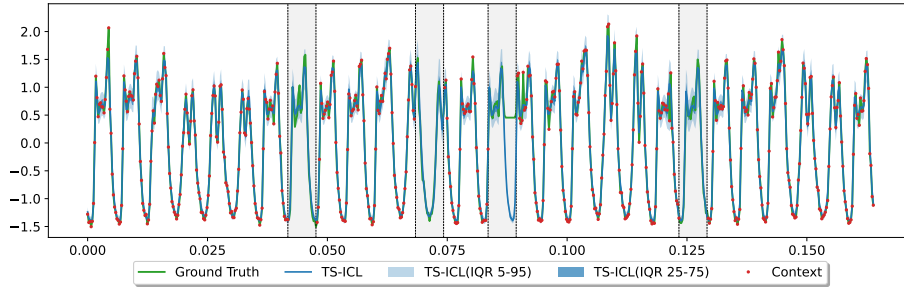


(d) *KDD Cup2022*, 70% missing values.

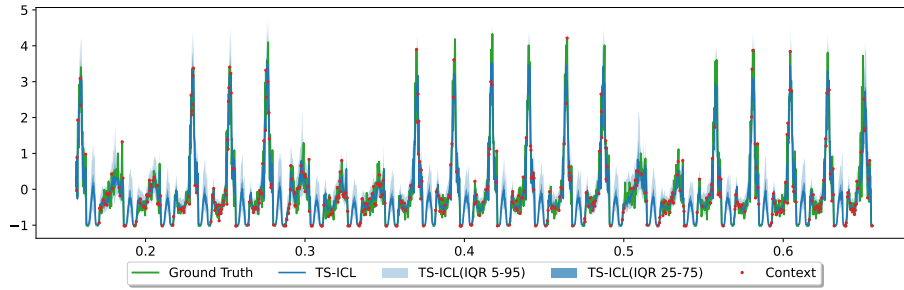


(e) *ERA5 wind speed*, 70% missing values.

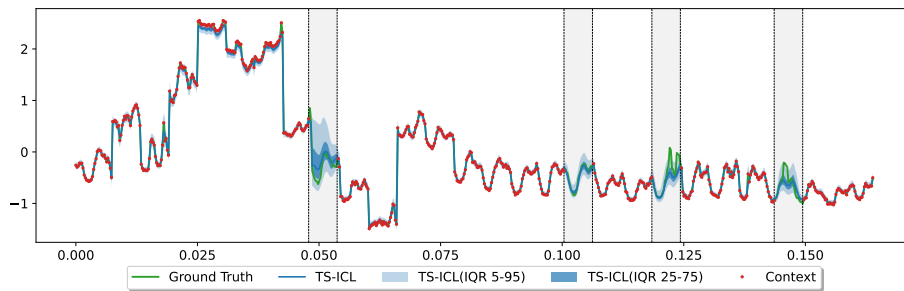
Figure 16: Qualitative assessment of TS-ICL imputations on the *fm-impute-bench* benchmark.



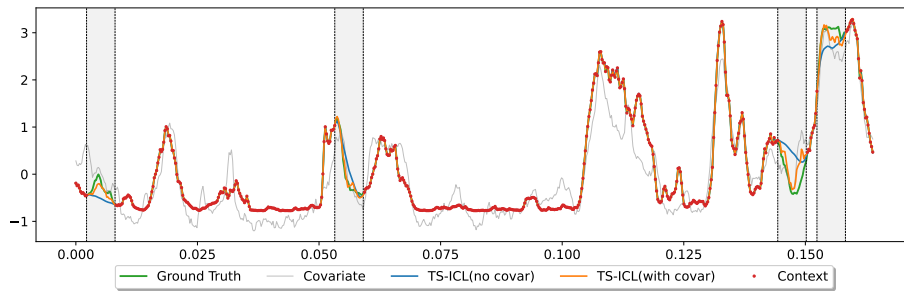
(a) *MDense*, four one-day missing blocks.



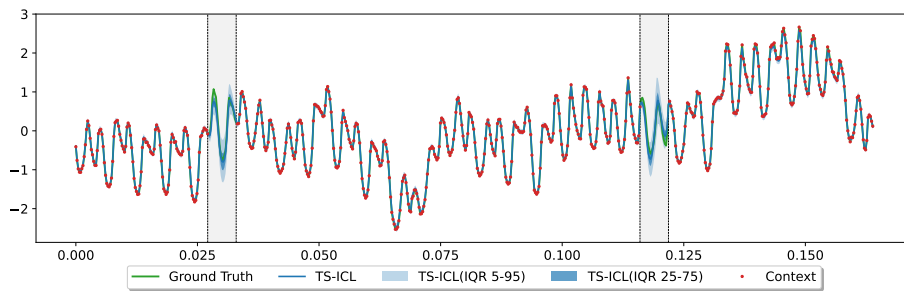
(b) *SHMETRO*, 70% missing values.



(c) *Spanish Energy*, four one-day missing blocks.



(d) *Wind-GE*, four one-day missing blocks.



(e) *GFC12 Load*, two one-day missing blocks.

Figure 17: Qualitative assessment of TS-ICL imputations on the `fm-impute-bench` benchmark (continued).

977 **E.2 TIME Benchmark**

978 In this section, we evaluate the zero-shot imputation capability of TS-ICL on TIME [33] a recently
 979 introduced benchmark originally designed for TSFM forecasting. We adapt it to cover imputation
 980 for the *univariate setting*. Performance is assessed across diverse missingness patterns, sequence
 981 lengths, and application domains (details in Table 14).

Table 14: All datasets used for zero-shot imputation in the TIME benchmark.

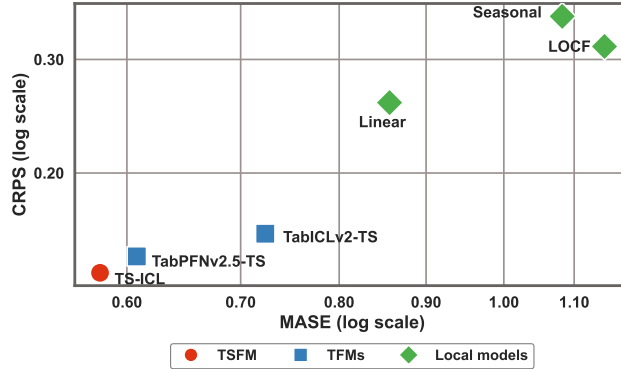
Dataset	Release Platform	Domain	Freq	Num. Series	Num. Variate	Avg Series Length	Short-term		Med-term		Long-term	
							Window Size	Num. Test Windows	Window Size	Num. Test Windows	Window Size	Num. Test Windows
Water Quality-Darwin	IMOS	Nature	15T	7	6	15,229	256	3,780	1024	630	4096	210
Current Velocity	IMOS	Nature	5T	1	6	26,486	256	720	1024	90	4096	30
Current Velocity	IMOS	Nature	10T	10	6	20,669	256	7,200	1024	900	4096	300
Current Velocity	IMOS	Nature	15T	5	6	8,503	256	3,600	1024	450	4096	150
Current Velocity	IMOS	Nature	20T	27	6	6,460	256	19,440	1024	2,430	4096	810
Current Velocity	IMOS	Nature	H	21	6	3,502	256	3,528	1024	504	4096	252
CPHL	IMOS	Nature	15T	2	1	10,831	256	240	1024	30	4096	10
CPHL	IMOS	Nature	30T	2	1	14,687	256	240	1024	60	4096	20
CPHL	IMOS	Nature	H	4	1	4,971	256	112	1024	16	4096	8
Coastal T-S	IMOS	Nature	5T	18	3	68,604	256	6,480	1024	810	4096	270
Coastal T-S	IMOS	Nature	15T	5	3	20,870	256	1,800	1024	225	4096	75
Coastal T-S	IMOS	Nature	20T	1	3	8,198	256	360	1024	45	4096	15
Coastal T-S	IMOS	Nature	H	24	3	5,489	256	2,016	1024	288	4096	144
SG Weather	data.gov.sg	Nature	D	6	4	2,953	256	2,928	1024	1,272	4096	648
SG PM 2.5	data.gov.sg	Nature	H	1	5	38,688	256	460	1024	150	4096	65
NE China Wind	GitHub	Nature	H	1	4	8,764	256	120	1024	40	4096	16
Australia Solar	Pvoutput	Energy	H	1	3	35,064	256	315	1024	105	4096	45
EPF Electricity	Academic	Energy	H	5	1	52,416	256	525	1024	175	4096	75
OpenElectricity	OpenElec	Energy	5T	1	10	43,488	256	1,680	1024	420	4096	140
EWELD Load	Academic	Energy	15T	1	10	20,544	256	560	1024	140	4096	20
SG Carpark	data.gov.sg	Transport	15T	354	1	14,332	256	14,868	1024	2,478	4096	354
Finland Traffic	Digitraffic	Transport	15T	1	1	35,136	256	186	1024	31	4096	4
Port Activity	Competition	Transport	D	99	2	2,127	256	2,376				
Port Activity	Competition	Transport	W	99	2	304	256	792				
ECDC COVID	ECDC	Healthcare	D	9	1	1,117	256	45				
ECDC COVID	ECDC	Healthcare	W	16	1	165	256	64				
Global Influenza	WHO	Healthcare	W	15	4	205	256	240				
Crypto	FRED	Finance	D	1	4	2,842	256	36				
US Term Structure	FRED	Finance	B	1	40	9,327	256	1,400				
Oil Price	FRED	Finance	B	1	12	5,035	256	420				
Job Claims	FRED	Finance	W	1	2	196	256	8				
Uncertainty-1M	FRED	Economics	M	1	3	780	256	21				
Housing Inventory	FRED	Economics	M	1	4	114	256	12				
JOLTS	FRED	Economics	M	1	6	297	256	30				
US Labor	FRED	Economics	M	1	14	380	256	70				
Vehicle Supply	FRED	Economics	M	1	6	391	256	30				
Auto Production-SF	FRED	Economics	M	1	1	367	256	5				
Commodity Prod.	FRED	Economics	M	32	1	325	256	160				
Commodity Import	FRED	Economics	M	8	1	697	256	40				
WUI-Global	FRED	Economics	Q	1	15	294	256	75				
Global Price	FRED	Economics	Q	1	60	142	256	300				
Vehicle Sales	FRED	Sales	M	1	10	596	256	50				
Online Retail II	Competition	Sales	D	1	1	739	256	6				
Supply Chain-Cust.	Competition	Sales	D	1	36	2,007	256	432				
Supply Chain-Loc.	Competition	Sales	D	1	51	2,007	256	612				
Azure2019-D	GitHub	CloudOPS	5T	989	3	8,627	256	8,901				
Azure2019-I	GitHub	CloudOPS	5T	492	3	8,630	256	4,428				
Azure2019-U	GitHub	CloudOPS	5T	78	3	1,406	256	1,404				
Smart Mfg.	Competition	Industry	H	34	5	1,666	256	2,380	1024	340	4096	170
MetroPT-3	Competition	Industry	5T	1	6	17,809	256	216	1024	36	4096	18

982 **Setting.** We reuse the datasets and windowing protocol originally designed for forecasting, treat-
 983 ing each lookback window as a partially observed sequence with synthetically introduced missing
 984 values. The benchmark spans multiple domains (e.g., energy, finance, healthcare, transportation)
 985 and covers diverse lengths and frequencies. Considering short, medium, and long window size
 986 configurations yields 98 imputation datasets, each composed of multiple samples to impute. The
 987 per-task window sizes can be found in Table 14. To evaluate robustness under different missingness
 988 patterns, we define four masking scenarios, namely:

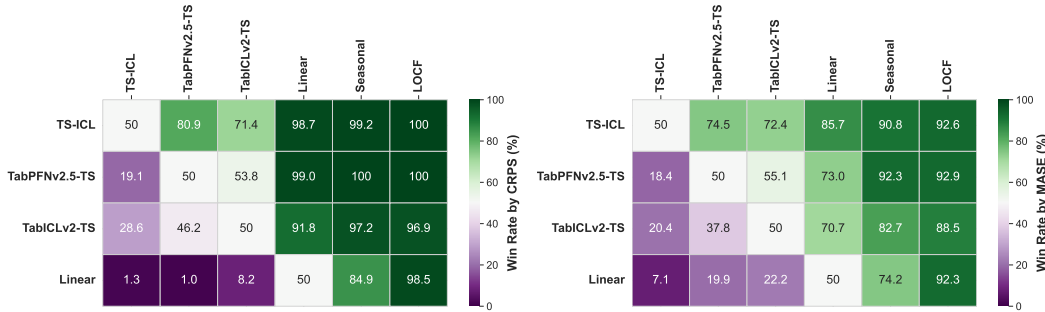
- 989 • pointwise masking with 50% missing values;
- 990 • pointwise masking with 70% missing values;
- 991 • a single contiguous missing block of size 1/15 of the window; and
- 992 • two disjoint missing blocks, each of size 1/15 of the window.

993 This setup captures both random and structured missingness, reflecting realistic data corruption
 994 patterns. Overall, this results in 98×4 imputation tasks and approximately 440k windows to recon-
 995 struct.

996 **Baselines.** TS-ICL is evaluated against the same tabular foundation model and local baselines as
 997 those described in Section 5.1. Final aggregated results (arithmetic mean) are reported in Figure 18
 998 (and Table 15), using both probabilistic (CRPS) and scale-normalized pointwise metrics (MASE;
 999 see Section D for a definition).



(a) Aggregated scores across 392 tasks - TIME benchmark.



(b) Win rates across 392 tasks (CRPS).

(c) Win rates across 392 tasks (MASE).

Figure 18: Aggregated metrics on the TIME univariate time series imputation benchmark. (a) MASE-CRPS (lower is better). Each point corresponds to a method, averaged across 392 tasks. (b-c) Pairwise win rates of the top-4 models. Each entry indicates the fraction of tasks where a method outperforms another according to the (b) CRPS or (c) the MASE.

1000 **Results.** Results consistently highlight the strong performance of TS-ICL. As illustrated in Fig-
 1001 ure 18a, TS-ICL achieves the lowest overall error, yielding relative improvements of 5.4% in
 1002 CRPS and 4.8% in MASE over TabPFNv2.5-TS, the current state-of-the-art tabular foundation
 1003 model (TFM). Compared to TabICLv2-TS, these gains extend to 13.0% (CRPS) and 20.0%
 1004 (MASE). Beyond aggregate metrics, TS-ICL demonstrates a dominant pairwise win rate against
 1005 TabPFNv2.5-TS, outperforming it on 80.9% of tasks for CRPS and 74.5% for MASE (Figure 18c).
 1006 Notably, TS-ICL achieves this peak accuracy with significant efficiency gains: it maintains an aver-
 1007 age inference runtime two orders of magnitude faster compared to the most competitive TFMs. This
 combination positions TS-ICL as a highly scalable solution for large-scale time series imputation.

Table 15: Detailed performance metrics (mean \pm std) for the univariate time series on the TIME imputation benchmark (392 tasks). Best in **bold**.

	TFSM	Tabular FMs		Local models		
	TS-ICL	TabPFNv2.5	TabICLv2	Linear	Seasonal	LOCF
MASE (\downarrow)	0.579 \pm 0.323	0.608 \pm 0.281	0.724 \pm 0.662	0.857 \pm 0.539	1.082 \pm 0.243	1.146 \pm 0.700
CRPS (\downarrow)	0.140 \pm 0.118	0.148 \pm 0.123	0.161 \pm 0.134	0.257 \pm 0.228	0.350 \pm 0.398	0.314 \pm 0.256

1008

1009 F Extended Forecasting Experiments

1010 This section provides broader insights into TS-ICL forecasting performances. A detailed description
1011 of the `fev-bench` datasets used in the main benchmark in Section 5.2 is given in Section F.1,
1012 together with complementary results and qualitative visualizations. Section F.2 further extends the
1013 zero-shot evaluation in the *univariate setting* to a second benchmark, TIME [33], across 98 tasks and
1014 against 12 foundation models.

1015 F.1 Fev-bench Benchmark

1016 **Inference datasets.** Table 18 details the datasets used for zero-shot forecasting in Section 5.2. As
1017 described by [38], these datasets cover a diverse range of domains, with frequencies ranging from
1018 5 minutes to quarterly data. Each forecasting task has its own prediction horizon. In Section 5.2,
1019 we distinguish two scenarios: *univariate zero-shot forecasting* and *zero-shot forecasting with known*
1020 *covariates* when available.

1021 **Baseline details.** We consider the strongest time-series foundation model baselines reported in
1022 `fev-bench` at the time of writing, while excluding models with substantial training-data overlap
1023 with the benchmark. In particular, we do not include `Moirai-2.0` and `TimesFm2.5` among main
1024 baselines due to their high reported leakage rates of 28% and 10%, respectively, on `fev-bench`.
1025 All baselines are evaluated in the zero-shot setting, without task-specific fine-tuning. We also report
1026 the leakage indicator from [38], defined as the fraction of model-task pairs for which the model
1027 pretraining data overlaps with the benchmark data.

- 1028 • `Chronos-2` [3] is a 120M-parameter, patch-based encoder-only transformer closely follow-
1029 ing the T5 encoder design, with alternating time and group attention layers for in-context
1030 learning across related series and covariates. It is the only TSFM baseline in `fev-bench`
1031 that natively supports known-future covariates and handles missing values in the look-back
1032 window. Its reported leakage rate is 0%, making it a clean zero-shot baseline.
- 1033 • `TiRex` [4] is a 35M-parameter decoder-only xLSTM model for zero-shot probabilistic fore-
1034 casting. It predicts quantiles directly and does not use covariates in the `fev-bench` setup.
1035 Its reported leakage rate is 1%.
- 1036 • `TimesFM-2.5` [11] is a 200M-parameter patched decoder-only transformer designed for
1037 long-context forecasting and direct quantile prediction. It is evaluated as a univariate fore-
1038 caster in `fev-bench`. Its reported leakage rate is 10%, so its aggregate score should be
1039 interpreted with some caution.
- 1040 • `Toto-1.0` [10] is a 151M-parameter decoder-only transformer optimized for multivariate
1041 observability time series. It supports multivariate inputs, but does not use known future co-
1042 variates in the `fev-bench` setting. Its reported leakage rate is 8%, which is non-negligible
1043 but substantially lower than that of `Moirai-2.0`.
- 1044 • `Chronos-Bolt` [2] is a 205M-parameter T5 encoder-decoder model and a patch-based
1045 variant of Chronos. It chunks the historical context into patches and produces multi-step
1046 quantile forecasts. It is evaluated as a univariate forecaster in `fev-bench`. Its reported
1047 leakage rate is 0%.

1048 F.1.1 Extended Results

1049 This section extends the empirical evaluation in Section 5.2 with a more detailed analysis of fore-
1050 casting performance across all experimental settings. Specifically, we provide:

- 1051 • **Aggregated detailed performance tables:** We report the average MASE and CRPS (met-
1052 rics definition in Section D) across the *univariate* tasks and *covariates-aware* tasks of
1053 `fm-impute-bench`. These results, detailed in Table 16 and Table 17, providing a detailed
1054 view of both point-wise and probabilistic performance. Note that metrics are aggregated
1055 across tasks using the geometric mean, following the evaluation protocol established in
1056 `fev-bench`.

1057
1058
1059
1060

- **MASE pairwise win rates:** To complement the CRPS-based win rate diagrams presented in the main text (Figures 5c and 5d), we include the corresponding pairwise win rate visualizations in terms of MASE for both *univariate* and *known-covariate* experiments in Figure 19.

Table 16: Fev-bench *univariate* forecasting (100 tasks). Performance metrics aggregated (geometric mean \pm geometric std). Best in **bold**.

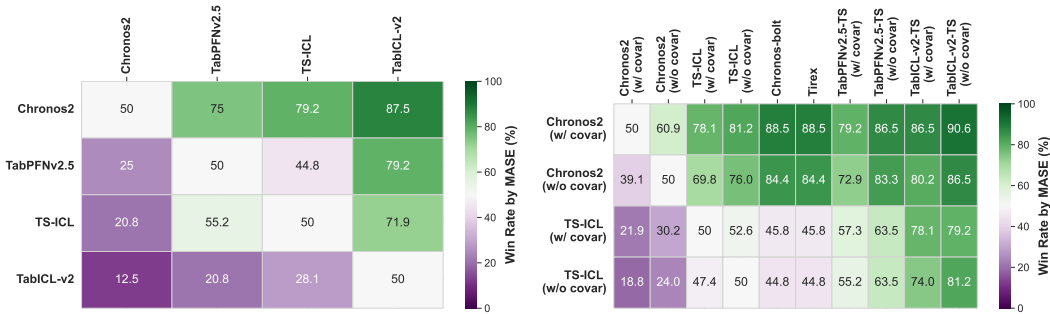
	TSFM				
	TS-ICL	Chronos-2	Chronos-Bolt	TiRex	Toto
MASE (\downarrow)	1.150 \pm 2.175	1.081 \pm 2.201	1.158 \pm 2.154	1.102 \pm 2.134	1.773 \pm 2.148
CRPS (\downarrow)	0.137 \pm 2.995	0.129 \pm 3.009	0.140 \pm 3.027	0.132 \pm 3.064	0.135 \pm 2.958

	Tabular Foundation models		Local models	
	TabPFNv2.5	TabICLv2	Seasonal	LOCF
MASE (\downarrow)	1.218 \pm 2.173	1.400 \pm 2.129	1.547 \pm 2.062	1.839 \pm 2.243
CRPS (\downarrow)	0.141 \pm 3.028	0.159 \pm 3.252	0.241 \pm 2.954	0.260 \pm 2.682

Table 17: Fev-bench forecasting on 100 tasks (30 *covariate-aware*) forecasting (100 tasks). Performance metrics aggregated (geometric mean \pm geometric std). Best in **bold**.

	TSFM			
	TS-ICL		Chronos-2	
	w/ covar	w/o covar	w/ covar	w/o covar
MASE (\downarrow)	1.117 \pm 2.232	1.150 \pm 2.175	1.034 \pm 2.250	1.081 \pm 2.202
CRPS (\downarrow)	0.131 \pm 3.000	0.137 \pm 2.995	0.123 \pm 3.062	0.129 \pm 3.018

	Tabular Foundation Models		Univariate FMs	
	TabPFNv2.5-TS (w/ covar)	TabICLv2-TS (w/ covar)	Chronos-Bolt (univar.)	TiRex (univar.)
MASE (\downarrow)	1.162 \pm 2.235	1.342 \pm 2.194	1.158 \pm 2.154	1.102 \pm 2.134
CRPS (\downarrow)	0.134 \pm 3.060	0.153 \pm 3.279	0.140 \pm 3.027	0.132 \pm 3.064



(a) *Univariate* forecasting across 100 tasks. (b) Forecasting with *known covariates* across 30 tasks.

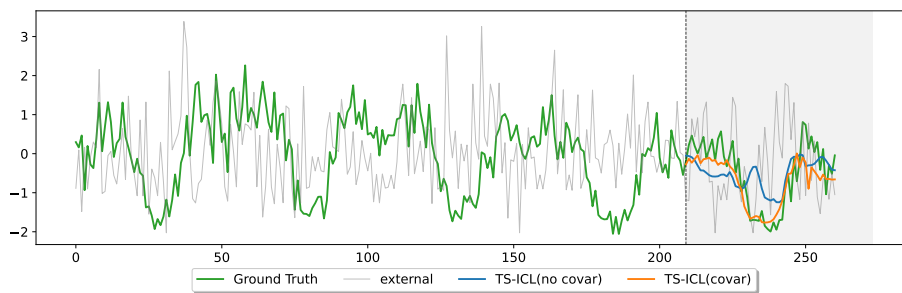
Figure 19: Pairwise win rates for forecasting on the fev-bench benchmark. Each entry indicates the fraction of tasks where a method outperforms another according to the MASE.

1061 F.1.2 Qualitative Analysis and Visualizations

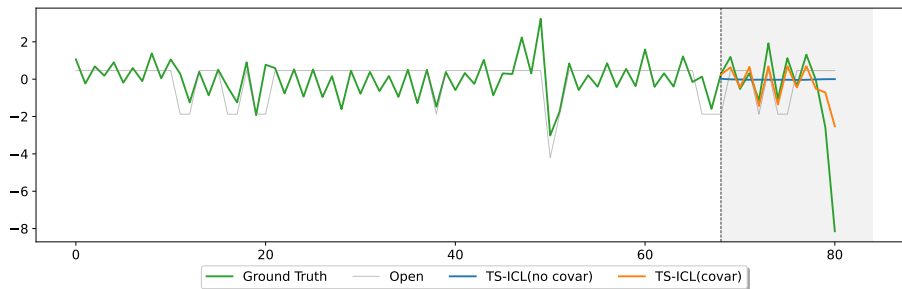
1062 This section presents visual examples of TS-ICL imputations for both *univariate* and *known-*
1063 *covariates* settings. We illustrate model forecasting capabilities across various *fev-bench* tasks.

1064 **Results.** Several observations emerge from the forecasting plots in Figure 20 (*known covariate*
1065 *setting*) and Figures 21 to 23 (*univariate setting*), where the median forecast is shown together with
1066 the corresponding 25-75 and 5-95 inter-quantile ranges.

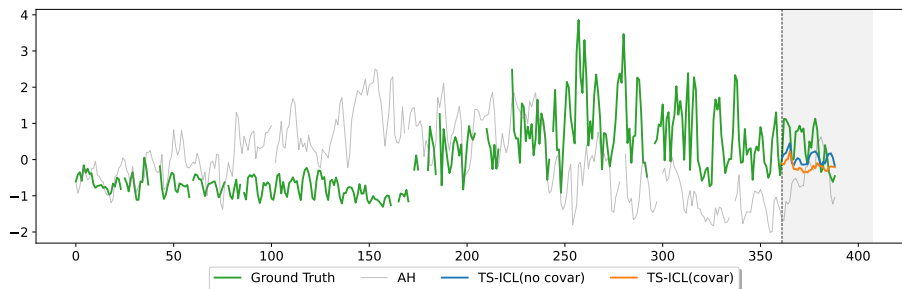
- 1067 (i) The plots highlight the general ability of TS-ICL to extrapolate from long context windows
 1068 (with a maximum lookback length of 4096) and regular patterns in heterogeneous sampling
 1069 rates, domains and seasonalities (e.g. Figures 21a, 21b, 22a and 23c).
 1070 (ii) Similarly to the imputation setting, TS-ICL tends to provide smooth forecasts of high-
 1071 frequency phenomena and adjusts its inter-quantile range accordingly (e.g. Figures 21d,
 1072 22c and 23e).
 1073 (iii) In the *univariate setting*, extrapolating from very short contexts is particularly challenging.
 1074 TS-ICL compensates with wider inter-quantile ranges, with mixed success depending on
 1075 the regularity of the underlying phenomena (Figures 21e, 22e, 23a and 23b).
 1076 (iv) In the *known covariate setting*, TS-ICL manages to leverage additional covariate, when the
 1077 latter informs about the target, while mostly ignoring it otherwise (Figure 20).
 1078 (v) Figure 23d gives an example of forecasting with missing values, with TS-ICL providing
 1079 adequate uncertainty estimates.



(a) *Hermes/W* - $H = 52$.

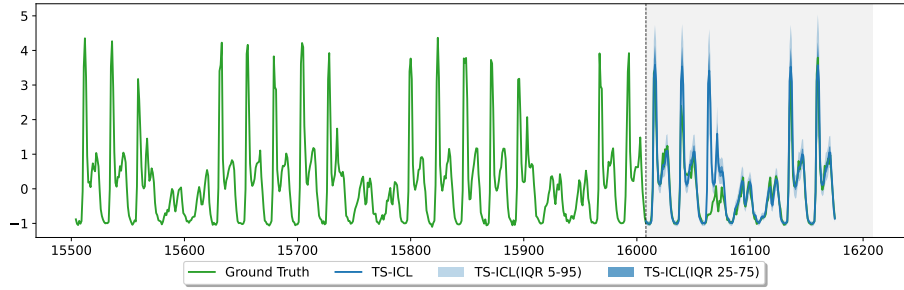


(b) *Rossmann/W* - $H = 13$.

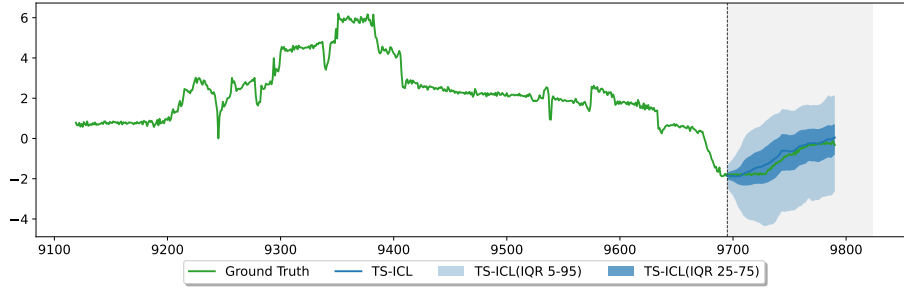


(c) *UCI Air Quality/D* - $H = 28$.

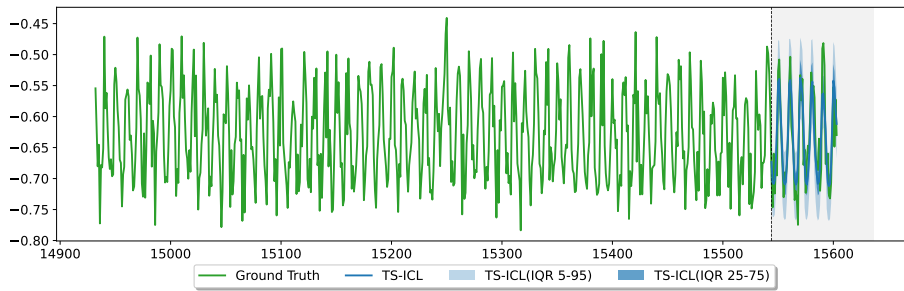
Figure 20: Qualitative assessment of TS-ICL forecasts on the *fev-bench* benchmark, in the *known covariate setting*. Covariates are shown in light gray.



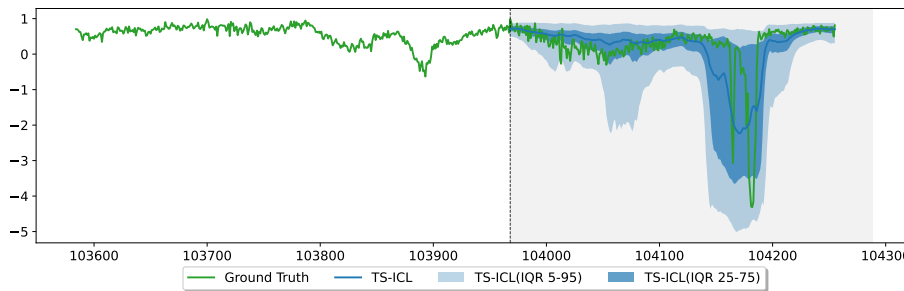
(a) *M-DENSE/IH* - $H = 168$.



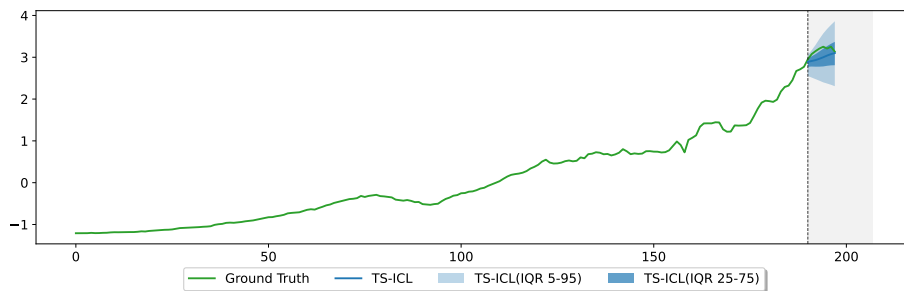
(b) *BOOMLET - 1631/30T* - $H = 96$.



(c) *BOOMLET - 1225/1T* - $H = 96$.

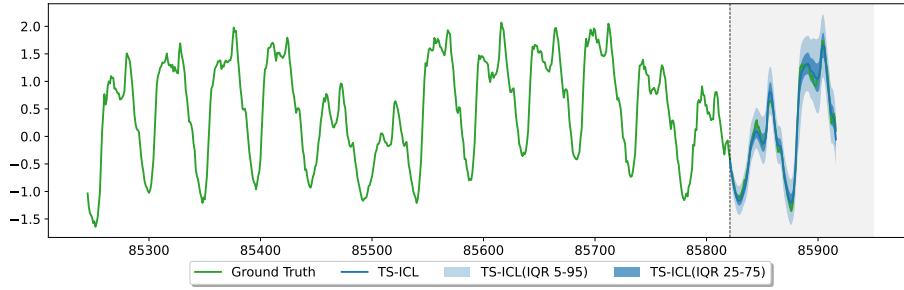


(d) *Loop Seattle/5T* - $H = 288$.

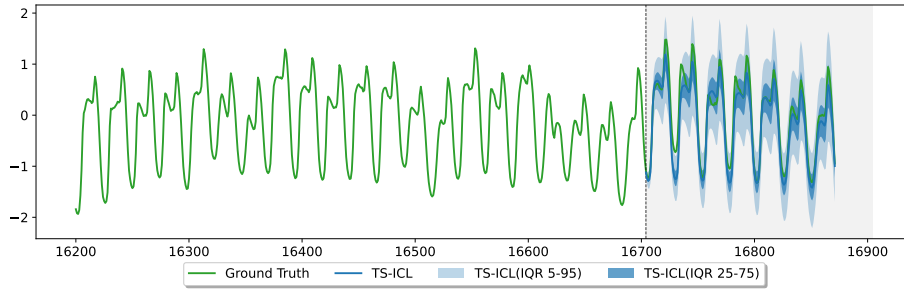


(e) *US Consumption/1Q* - $H = 8$.

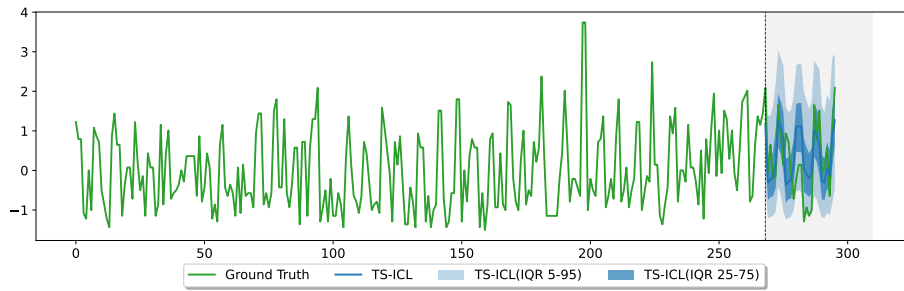
Figure 21: Qualitative assessment of TS-ICL forecasts on the fev-bench benchmark.



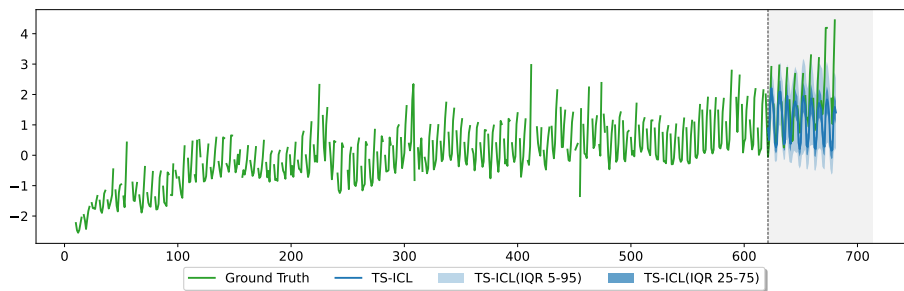
(a) *ENTSOE-e Load/30T* - $H = 96$.



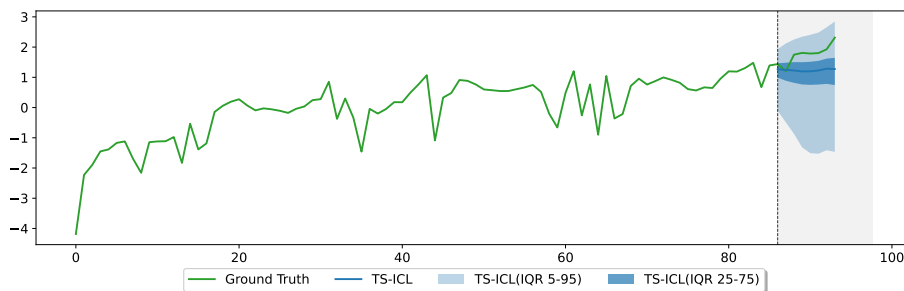
(b) *GFC17/1H* - $H = 168$.



(c) *Restaurant/1D* - $H = 28$.

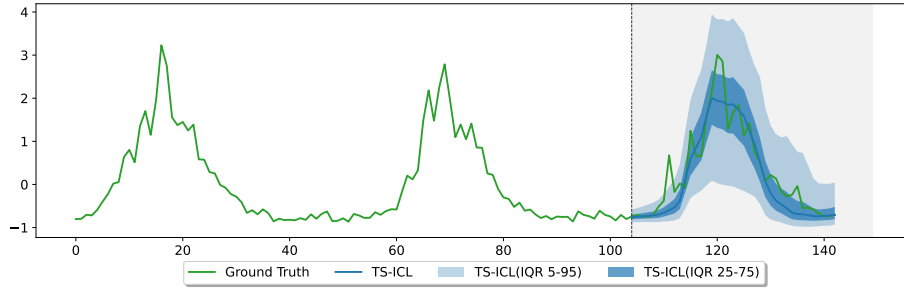


(d) *Rohlik Orders/1D* - $H = 61$.

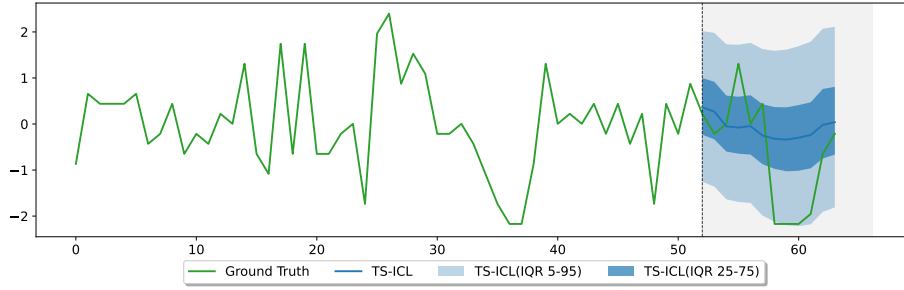


(e) *Rohlik Orders/1W* - $H = 8$.

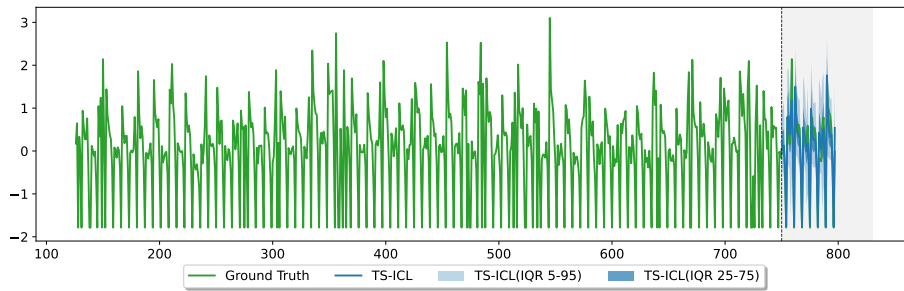
Figure 22: Qualitative assessment of TS-ICL forecasts on the fev-bench benchmark (continued).



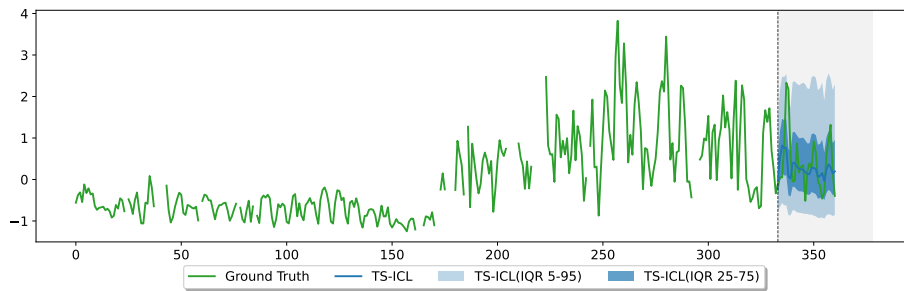
(a) *Walmart/1W* - $H = 39$.



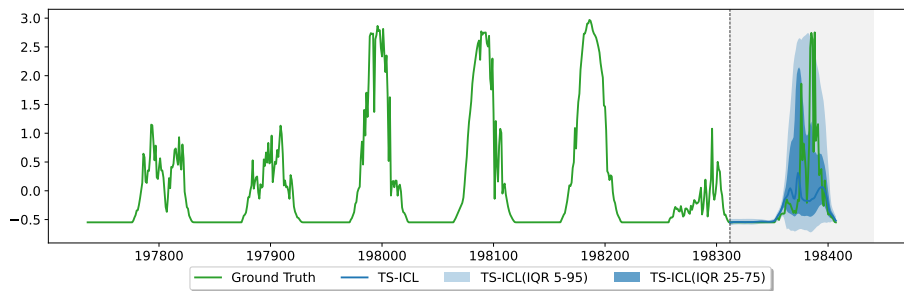
(b) *M5/1M* - $H = 12$.



(c) *Rossmann/1D* - $H = 48$.



(d) *UCI Air Quality/1D* - $H = 28$.



(e) *Solar with Weather/15T* - $H = 96$.

Figure 23: Qualitative assessment of TS-ICL forecasts on the fev-bench benchmark (continued).

Table 18: Full statistics of fev-bench tasks with data sources. Covariates: P (Past), K (Known), S (Static).

Task / Dataset	Source	Domain	Freq	Horizons	Num. Series	Num. Target	Median Length	Covariates P K S	Num. Test Windows
BizTObs - L2C	LOTSAs	cloud	5T	288	1	7	31,968	0 0 0	140
BizTObs - L2C	LOTSAs	cloud	H	24	1	7	2,664	0 0 0	140
ETT	GitHub	energy	15T	96	2	7	69,680	0 0 0	280
ETT	GitHub	energy	H	168	2	7	17,420	0 0 0	280
ETT	GitHub	energy	D	28	2	7	724	0 0 0	280
ETT	GitHub	energy	W	13	2	7	103	0 0 0	70
Hierarchical Sales	LOTSAs	retail	D	28	118	1	1,825	0 0 0	1,180
Hierarchical Sales	LOTSAs	retail	W	13	118	1	260	0 0 0	1,180
Hospital	LOTSAs	healthcare	M	12	767	1	84	0 0 0	3,068
Jena Weather	MPI Jena	nature	10T	144	1	21	52,704	0 0 0	420
Jena Weather	MPI Jena	nature	D	28	1	21	366	0 0 0	231
Jena Weather	MPI Jena	nature	H	24	1	21	8,784	0 0 0	420
Loop Seattle	LOTSAs	mobility	D	28	323	1	365	0 0 0	3,230
Loop Seattle	LOTSAs	mobility	5T	288	323	1	105,120	0 0 0	3,230
Loop Seattle	LOTSAs	mobility	H	168	323	1	8,760	0 0 0	3,230
M-DENSE	LOTSAs	mobility	D	28	30	1	730	0 0 0	300
M-DENSE	LOTSAs	mobility	H	168	30	1	17,520	0 0 0	300
SZ Taxi	LOTSAs	mobility	15T	96	156	1	2,976	0 0 0	1,560
SZ Taxi	LOTSAs	mobility	H	168	156	1	744	0 0 0	312
Solar	LOTSAs	energy	W	13	137	1	52	0 0 0	137
Solar	LOTSAs	energy	D	28	137	1	365	0 0 0	1,370
Australian Tourism	Monash	econ	Q	8	89	1	36	0 0 0	178
FRED-MD - CEE	Fed	econ	M	12	1	3	798	4 0 0	60
FRED-MD - Macro	Fed	econ	M	12	1	51	798	0 0 0	1,020
FRED-QD - CEE	Fed	econ	Q	8	1	3	266	4 0 0	60
FRED-QD - Macro	Fed	econ	Q	8	1	51	266	0 0 0	1,020
GVAR	Mohaddes	econ	Q	8	33	6	178	3 0 0	1,980
US Consumption	FPP3	econ	M	12	31	1	792	0 0 0	310
US Consumption	FPP3	econ	Q	8	31	1	262	0 0 0	310
US Consumption	FPP3	econ	Y	5	31	1	64	0 0 0	310
World CO2 Emissions	WorldBank	econ	Y	5	191	1	60	0 0 0	1,719
World Life Expectancy	WorldBank	econ	Y	5	237	1	74	0 0 0	2,370
World Tourism	WorldBank	econ	Y	5	178	1	21	0 0 0	356
ENTSO-e Load	ENTSO-E	energy	15T	96	6	1	175,292	0 3 0	120
ENTSO-e Load	ENTSO-E	energy	30T	96	6	1	87,645	0 3 0	120
ENTSO-e Load	ENTSO-E	energy	H	168	6	1	43,822	0 3 0	120
EPF-BE	GitHub	energy	H	24	1	1	52,416	0 2 0	20
EPF-DE	GitHub	energy	H	24	1	1	52,416	0 2 0	20
EPF-FR	GitHub	energy	H	24	1	1	52,416	0 2 0	20
EPF-NP	GitHub	energy	H	24	1	1	52,416	0 2 0	20
EPF-PJM	GitHub	energy	H	24	1	1	52,416	0 2 0	20
ERCOT	ERCOT	energy	D	28	8	1	6,452	0 0 0	160
ERCOT	ERCOT	energy	H	168	8	1	154,872	0 0 0	160
ERCOT	ERCOT	energy	M	12	8	1	211	0 0 0	120
ERCOT	ERCOT	energy	W	13	8	1	921	0 0 0	160
GFC12	LOTSAs	energy	H	168	11	1	39,414	0 1 0	110
GFC14	LOTSAs	energy	H	168	1	1	17,520	0 1 0	20
GFC17	LOTSAs	energy	H	168	8	1	17,544	0 1 0	160
Solar with Weather	Kaggle	energy	15T	96	1	1	198,600	2 7 0	20
Solar with Weather	Kaggle	energy	H	24	1	1	49,648	2 7 0	20
BOOMLET - 1062	BOOM	cloud	5T	288	1	21	16,384	0 0 0	420
BOOMLET - 1209	BOOM	cloud	5T	288	1	53	16,384	0 0 0	1,060
BOOMLET - 1225	BOOM	cloud	T	60	1	49	16,384	0 0 0	980
BOOMLET - 1230	BOOM	cloud	5T	288	1	23	16,384	0 0 0	460
BOOMLET - 1282	BOOM	cloud	T	60	1	35	16,384	0 0 0	700
BOOMLET - 1487	BOOM	cloud	5T	288	1	54	16,384	0 0 0	1,080
BOOMLET - 1631	BOOM	cloud	30T	96	1	40	10,463	0 0 0	800
BOOMLET - 1676	BOOM	cloud	30T	96	1	100	10,463	0 0 0	2,000
BOOMLET - 1855	BOOM	cloud	H	24	1	52	5,231	0 0 0	1,040
BOOMLET - 1975	BOOM	cloud	H	24	1	75	5,231	0 0 0	1,500
BOOMLET - 2187	BOOM	cloud	H	24	1	100	5,231	0 0 0	2,000
BOOMLET - 285	BOOM	cloud	T	60	1	75	16,384	0 0 0	1,500
BOOMLET - 619	BOOM	cloud	T	60	1	52	16,384	0 0 0	1,040
BOOMLET - 772	BOOM	cloud	T	60	1	67	16,384	0 0 0	1,340
BOOMLET - 963	BOOM	cloud	T	60	1	28	16,384	0 0 0	560
Favorita Store Sales	Kaggle	retail	M	12	1,579	1	54	1 1 6	3,158
Favorita Store Sales	Kaggle	retail	W	13	1,579	1	240	1 1 6	15,790
Favorita Store Sales	Kaggle	retail	D	28	1,579	1	1,688	1 2 6	15,790
Favorita Transactions	Kaggle	retail	M	12	51	1	54	1 0 5	102
Favorita Transactions	Kaggle	retail	W	13	51	1	240	1 0 5	510
Favorita Transactions	Kaggle	retail	D	28	51	1	1,688	1 1 5	510
KDD Cup 2022	Kaggle	energy	D	14	134	1	243	9 0 0	1,340
KDD Cup 2022	Kaggle	energy	10T	288	134	1	35,279	9 0 0	1,340
KDD Cup 2022	Kaggle	energy	30T	96	134	1	11,758	9 0 0	1,340
M5	Kaggle	retail	M	12	30,490	1	58	0 8 5	30,490
M5	Kaggle	retail	W	13	30,490	1	257	0 8 5	30,490
M5	Kaggle	retail	D	28	30,490	1	1,810	0 8 5	30,490
Restaurant	Kaggle	retail	D	28	817	1	296	0 0 4	6,536
Rohlik Orders	Kaggle	retail	W	8	7	1	170	9 4 0	35
Rohlik Orders	Kaggle	retail	D	61	7	1	1,197	9 4 0	35
Rohlik Sales	Kaggle	retail	W	8	5,243	1	150	1 13 7	5,243
Rohlik Sales	Kaggle	retail	D	14	5,390	1	1,046	1 13 7	5,390
Rossmann	Kaggle	retail	W	13	1,115	1	133	1 4 10	8,920
Rossmann	Kaggle	retail	D	48	1,115	1	942	1 5 10	11,150
Walmart	Kaggle	retail	W	39	2,936	1	143	0 10 4	2,936
ECDC ILI	ECDC	healthcare	W	13	25	1	201	0 0 0	250
Hermes	LOTSAs	retail	W	52	10,000	1	261	0 1 2	10,000
Hospital Admissions	Gov.UK	healthcare	D	28	8	1	1,731	0 0 0	160
Hospital Admissions	Gov.UK	healthcare	W	13	8	1	246	0 0 0	128
Redset	GitHub	cloud	5T	288	118	1	25,920	0 0 1	1,180
Redset	GitHub	cloud	15T	96	126	1	8,640	0 0 1	1,260
Redset	GitHub	cloud	H	24	138	1	2,160	0 0 1	1,380
UCI Air Quality	UCI	nature	H	168	1	4	9,357	0 3 0	80
UCI Air Quality	UCI	nature	D	28	1	4	389	0 3 0	44
UK COVID - Nation - Cumul.	Gov.UK	healthcare	D	28	4	3	729	5 0 0	240
UK COVID - Nation - Cumul.	Gov.UK	healthcare	W	8	4	3	105	5 0 0	48
UK COVID - Nation - New	Gov.UK	healthcare	D	28	4	3	729	5 0 0	240
UK COVID - Nation - New	Gov.UK	healthcare	W	8	4	3	105	5 0 0	48
UK COVID - UTLA - Cumul.	Gov.UK	healthcare	W	13	214	1	104	0 0 0	1,070
UK COVID - UTLA - New	Gov.UK	healthcare	D	28	214	1	721	0 0 0	2,140

1080 **F.2 TIME Benchmark**

1081 In this section, we evaluate the zero-shot forecasting capabilities of TS-ICL on the TIME benchmark
 1082 [33], covering univariate settings. This benchmark is particularly valuable as it provides a rigorous
 1083 framework for zero-shot evaluation, ensuring a total absence of data leakage across all compared
 1084 foundation models. Performance is assessed across a wide range of missingness patterns, sequence
 1085 lengths, and application domains (see Table 19 for details).

Table 19: Individual statistics of forecasting tasks across all datasets. Freq denotes the sampling frequency.

Dataset	Release Platform	Domain	Freq	Num. Series	Num. Variate	Avg Series Length	Short-term		Med-term		Long-term	
							Horizon	Num. Test Windows	Horizon	Num. Test Windows	Horizon	Num. Test Windows
Water Quality-Darwin	IMOS	Nature	15T	7	6	15,229	16 (4H)	3,780	96 (D)	630	288 (3D)	210
Current Velocity	IMOS	Nature	5T	1	6	26,486	36 (3H)	720	288 (D)	90	864 (3D)	30
Current Velocity	IMOS	Nature	10T	10	6	20,669	18 (3H)	7,200	144 (D)	900	432 (3D)	300
Current Velocity	IMOS	Nature	15T	5	6	8,503	12 (3H)	3,600	96 (D)	450	288 (3D)	150
Current Velocity	IMOS	Nature	20T	27	6	6,460	9 (3H)	19,440	72 (D)	2,430	216 (3D)	810
Current Velocity	IMOS	Nature	H	21	6	3,502	24 (D)	3,528	168 (W)	504	336 (2W)	252
CPHL	IMOS	Nature	15T	2	1	10,831	12 (3H)	240	96 (D)	30	288 (3D)	10
CPHL	IMOS	Nature	30T	2	1	14,687	12 (3H)	240	48 (D)	60	144 (3D)	20
CPHL	IMOS	Nature	H	4	1	4,971	24 (D)	112	168 (W)	16	336 (2W)	8
Coastal T-S	IMOS	Nature	5T	18	3	68,604	36 (3H)	6,480	288 (D)	810	864 (3D)	270
Coastal T-S	IMOS	Nature	15T	5	3	20,870	12 (3H)	1,800	96 (D)	225	288 (3D)	75
Coastal T-S	IMOS	Nature	20T	1	3	8,198	9 (3H)	360	72 (D)	45	216 (3D)	15
Coastal T-S	IMOS	Nature	H	24	3	5,489	24 (D)	2,016	168 (W)	288	336 (2W)	144
SG Weather	data.gov.sg	Nature	D	6	4	2,953	3 (3D)	2,928	7 (W)	1,272	14 (2W)	648
SG PM 2.5	data.gov.sg	Nature	H	1	5	38,688	24 (D)	460	72 (3D)	150	168 (W)	65
NE China Wind	GitHub	Nature	H	1	4	8,764	24 (D)	120	72 (3D)	40	168 (W)	16
Australia Solar	Pvoutput	Energy	H	1	3	35,064	24 (D)	315	72 (3D)	105	168 (W)	45
EPF Electricity	Academic	Energy	H	5	1	52,416	24 (D)	525	72 (3D)	175	168 (W)	75
OpenElectricity	OpenElec	Energy	5T	1	10	43,488	24 (2H)	1,680	96 (3H)	420	288 (D)	140
EWELD Load	Academic	Energy	15T	1	10	20,544	24 (6H)	560	96 (D)	140	672 (W)	20
SG Carpark	data.gov.sg	Transport	15T	354	1	14,332	16 (4H)	14,868	96 (D)	2,478	672 (W)	354
Finland Traffic	Digitraffic	Transport	15T	1	1	35,136	16 (4H)	186	96 (D)	31	672 (W)	4
Port Activity	Competition	Transport	D	99	2	2,127	30 (M)	2,376				
Port Activity	Competition	Transport	W	99	2	304	13 (Q)	792				
ECDC COVID	ECDC	Healthcare	D	9	1	1,117	30 (30D)	45				
ECDC COVID	ECDC	Healthcare	W	16	1	165	13 (Q)	64				
Global Influenza	WHO	Healthcare	W	15	4	205	13 (Q)	240				
Crypto	FRED	Finance	D	1	4	2,842	30 (M)	36				
US Term Structure	FRED	Finance	B	1	40	9,327	20 (4W)	1,400				
Oil Price	FRED	Finance	B	1	12	5,035	20 (4W)	420				
Job Claims	FRED	Finance	W	1	2	196	13 (Q)	8				
Uncertainty-1M	FRED	Economics	M	1	3	780	6 (6M)	21				
Housing Inventory	FRED	Economics	M	1	4	114	12 (A)	12				
JOLTS	FRED	Economics	M	1	6	297	12 (A)	30				
US Labor	FRED	Economics	M	1	14	380	12 (A)	70				
Vehicle Supply	FRED	Economics	M	1	6	391	12 (A)	30				
Auto Production-SF	FRED	Economics	M	1	1	367	12 (A)	5				
Commodity Prod.	FRED	Economics	M	32	1	325	12 (A)	160				
Commodity Import	FRED	Economics	M	8	1	697	12 (A)	40				
WUI-Global	FRED	Economics	Q	1	15	294	4 (A)	75				
Global Price	FRED	Economics	Q	1	60	142	4 (A)	300				
Vehicle Sales	FRED	Sales	M	1	10	596	12	50				
Online Retail II	Competition	Sales	D	1	1	739	30	6				
Supply Chain-Cust.	Competition	Sales	D	1	36	2,007	30	432				
Supply Chain-Loc.	Competition	Sales	D	1	51	2,007	30	612				
Azure2019-D	GitHub	CloudOPS	5T	989	3	8,627	288 (D)	8,901				
Azure2019-I	GitHub	CloudOPS	5T	492	3	8,630	288 (D)	4,428				
Azure2019-U	GitHub	CloudOPS	5T	78	3	1,406	48 (4H)	1,404				
Smart Mfg.	Competition	Industry	H	34	5	1,666	24 (D)	2,380	168 (W)	340	336 (2W)	170
MetroPT-3	Competition	Industry	5T	1	6	17,809	48 (4H)	216	288 (D)	36	576 (2D)	18

1086 **Setting.** The evaluation is conducted under a *primarily univariate forecasting* setting, where Time
 1087 considers short-, medium-, and long-context window sizes, resulting in 98 forecasting tasks and
 1088 approximately 110k windows to predict. Note, however, that some baselines (namely Chronos-2,
 1089 Toto, Moirai, and VisionTS) leverage look-back windows from others series at inference time,
 1090 operating in a multivariate forecasting regime [33].

1091 **Baselines.** The guaranteed absence of data leakage in the TIME benchmark allows us to expand
 1092 the comparison to a broader set of state-of-the-art foundation models that were previously excluded
 1093 in Section 5.2. Final aggregated results are reported in Figure 24 (and Table 20), using both proba-
 1094 bilistic (CRPS) and scale-normalized pointwise metrics (MASE; see Section D for a definition). We
 1095 compare TS-ICL against:

- Time Series Foundation Models (TSFMs): Chronos-2 [3], Timesfm2_5 [11], TiRex [4], Moirai2 [27], Toto [10], Chronos-bolt, Sundial [28], TimesFm [11], Vision.ts [7], and Moirai [43].

- Tabular Foundation Models (TFMs): TabPFNv2.5-TS [17] and TabICLv2-TS [34].

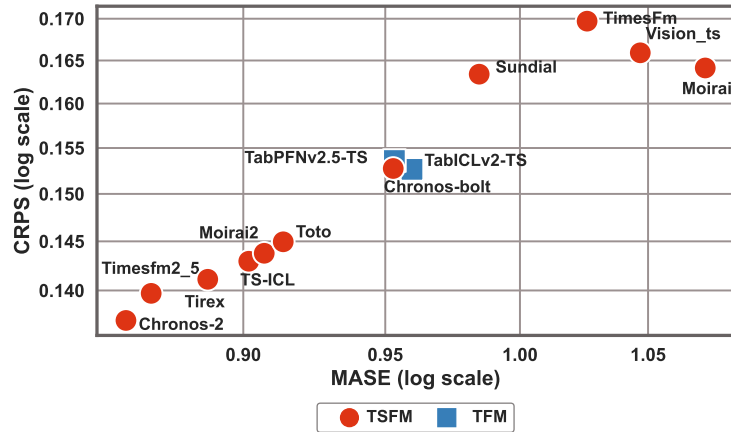


Figure 24: MASE-CRPS time trade-off (lower is better) on the TIME benchmark. The x-axis reports task-averaged MASE for each method, while the y-axis shows task-averaged CRPS for each method. Each point corresponds to a method, averaged across tasks.

Table 20: Detailed forecasting performance metrics aggregated (geometric mean \pm geometric std) across the 98 tasks of the *univariate setting* in the TIME benchmark. Best in **bold**.

	Time Series Foundation Models (TFSMs)						
	Chronos-2	TimesFM2.5	TiRex	TS-ICL	Toto	Moirai2	Bolt
MASE (\downarrow)	0.861\pm1.621	0.869 \pm 1.613	0.888 \pm 1.585	0.902 \pm 1.597	0.907 \pm 1.597	0.914 \pm 1.604	0.953 \pm 1.629
CRPS (\downarrow)	0.137\pm2.885	0.140 \pm 2.902	0.141 \pm 2.866	0.143 \pm 2.807	0.144 \pm 2.821	0.145 \pm 2.833	0.153 \pm 2.688
	Tabular FMs		Other TFSMs				
	TabPFNv2.5-TS	TabICLv2-TS	Sundial	TimesFM	Moirai	Vision_ts	
MASE (\downarrow)	0.953 \pm 1.601	0.960 \pm 1.590	0.985 \pm 1.675	1.026 \pm 1.593	1.073 \pm 1.598	1.047 \pm 1.603	
CRPS (\downarrow)	0.154 \pm 2.852	0.153 \pm 2.872	0.163 \pm 2.666	0.170 \pm 2.840	0.164 \pm 2.573	0.166 \pm 2.753	

1100 **Results.** TS-ICL demonstrates strong competitive performance on the Time benchmark, consistently ranking among the top-tier Time Series Foundation Models (TFSMs). As detailed in Table 20, 1101 while Chronos-2 maintains its position as the SOTA leader, TS-ICL achieves a highly comparable 1102 MASE of 0.902 and a CRPS of 0.143. Notably, the performance gap between TS-ICL model and 1103 Chronos-2 is minimal, with TS-ICL trailing by only 4.7% in MASE and 4.3% in CRPS in relative 1104 terms. This narrow margin is further evidenced by the pairwise win rate analysis (Figures 25 1105 and 26), where TS-ICL successfully outperforms Chronos-2 on 29.6% of tasks for CRPS and 1106 22.4% for MASE. 1107

1108 Beyond this head-to-head comparison, TS-ICL shows a clear superiority over the entire category 1109 of Tabular Foundation Models (TFMs), significantly outperforming both TabPFN and TabICL. It 1110 also surpasses several established TFSMs, including TimesFM, Moirai, and Sundial. By matching 1111 the performance of much larger, specialized models within such a tight margin while offering a 1112 more efficient architecture, TS-ICL establishes itself as a robust and scalable alternative for high- 1113 performance zero-shot forecasting.

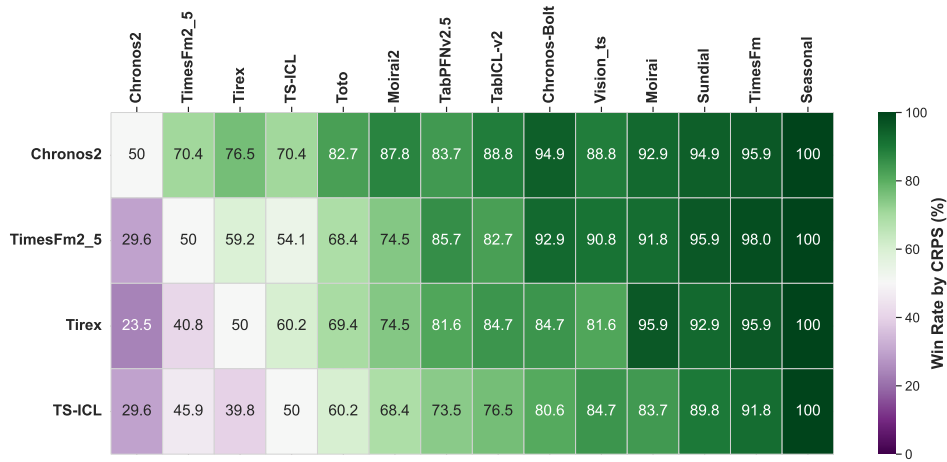


Figure 25: Pairwise win rates for the top-4 models against all other forecasters on the TIME benchmark. Each entry indicates the fraction of tasks where a method outperforms another according to the CRPS.

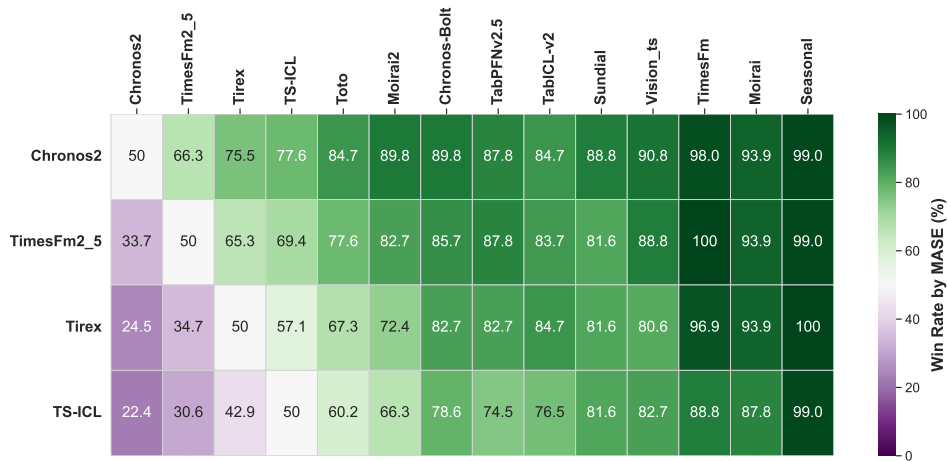


Figure 26: Pairwise win rates for the top-4 models against all other forecasters on the TIME benchmark. Each entry indicates the fraction of tasks where a method outperforms another according to the MASE.