NEURAL FOURIER MODELLING: A HIGHLY COMPACT APPROACH TO TIME-SERIES ANALYSIS

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

027

028

029

031 032

033

Paper under double-blind review

ABSTRACT

Neural time-series analysis has traditionally focused on modeling data in the time domain, often with some approaches incorporating equivalent Fourier domain representations as auxiliary spectral features. In this work, we shift the main focus to frequency representations, modeling time-series data fully and directly in the Fourier domain. We introduce Neural Fourier Modelling (NFM), a compact yet powerful solution for time-series analysis. NFM is grounded in two key properties of the Fourier transform (FT): (i) the ability to model finite-length time series as functions in the Fourier domain, treating them as continuous-time elements in function space, and (ii) the capacity for data manipulation (such as resampling and timespan extension) within the Fourier domain. We reinterpret Fourier-domain data manipulation as frequency extrapolation and interpolation, incorporating this as a core learning mechanism in NFM, applicable across various tasks. To support flexible frequency extension with spectral priors and effective modulation of frequency representations, we propose two learning modules: Learnable Frequency Tokens (LFT) and Implicit Neural Fourier Filters (INFF). These modules enable compact and expressive modeling in the Fourier domain. Extensive experiments demonstrate that NFM achieves state-of-the-art performance on a wide range of tasks (forecasting, anomaly detection, and classification), including challenging time-series scenarios with previously unseen sampling rates at test time. Moreover, NFM is highly compact, requiring fewer than **40K** parameters in each task, with time-series lengths ranging from 100 to 16K.

1 INTRODUCTION

034 Time series analysis is valuable in understanding the dynamics of systems and phenomena that evolve over time, and to address practical 037 problems in a range of domains. With rapidly increasing computational resources and data available for learning, neural-based modelling 040 approaches (Vaswani et al., 2017; Oord et al., 2016) have recently gained vast popularity in the 041 discipline. A number of sophisticated method-042 ologies and models have been developed, greatly 043 advancing performance on a variety of time-044 series tasks such as forecasting (Nie et al., 2022; 045 Wu et al., 2021), classification (Zhang et al., 046 2022; Raghu et al., 2023), anomaly detection 047 (Xu et al., 2021a; Chen et al., 2022a). Behind 048 this success, a remaining central question in time series modelling is how to capture meaningful information relevant to tasks from temporal pat-051 terns and generalize the dependencies ingrained



Figure 1: Illustration of Fourier-domain manipulations (left), including zero-padding/truncation (top) and zero-interleaving (bottom), and equivalent effects in the time domain (right).

within time-evolving data that are inherently diverse and intricate. To answer the question and thus
 progress the discipline further, in this work we study frequency representations, giving rise to a
 powerful time-series modelling scheme, Neural Fourier Modelling (NFM).

054 Given its well-known prevalence in the field of conventional signal processing, the adoption of frequency domain analysis to neural-based modelling is, unsurprisingly, not a unique idea in itself. 056 There have been a number of research efforts across multiple areas from time-series analysis to 057 computer vision (CV). Spectral representations are often harnessed as an alternative (Trabelsi et al., 058 2018; Choi et al., 2019) or as a complement (Yang & Hong, 2022; Woo et al., 2022) to time- or spatialdomain representations to capture information in a more compact and overarching form. Moreover, several recent studies have extended its applications to an efficient form of global convolution (Huang 060 et al., 2023; Lee-Thorp et al., 2021; Lin et al., 2023; Alaa et al., 2020; Yi et al., 2024a) and a 061 means of data augmentation (Xie et al., 2022; Xu et al., 2021b; Zhou et al., 2022b) for facilitating 062 invariances conducive to generalization in learning. While the existing works have successfully 063 utilized the frequency representation, it still remains under-explored with the lack of study on 064 frequency interpolation and extrapolation as a direct means of modelling data fully in the Fourier 065 domain. More specifically, our work is motivated by two aspects of frequency representations. (i) 066 Simply learning directly in the Fourier domain enables finding a function-to-function mapping - an 067 inductive bias towards learning resolution-invariance property. (ii) Fourier-domain manipulation and 068 a fundamental connection to its time-domain counterparts can provide a general means of learning in 069 the Fourier domain.

070

071 Fourier-domain Manipulation. As shown in Figure 1, there are two ways to manipulate time 072 series in the Fourier domain -1) padding/truncating and 2) interleaving the original Frequency 073 representation with zero coefficients, each of which resulting in resampling and extending the original 074 time-domain representation, respectively. Taking this into a view where a resultant time-domain representation caused by the frequency manipulation to a given input sequence is a desired target, 075 we can naturally reformulate the manipulation with zero frequency coefficients into a constructive 076 process of learning meaningful coefficients towards the target - i.e., frequency extrapolation and 077 interpolation. Notably, this view provides a comprehensive learning framework requiring no architectural modification to models. For example, time series can be readily modelled for forecasting 079 task from the frequency interpolation. Moreover, context learning (e.g., classification and regression) or representation learning can be made through the frequency extrapolation or directly imposing 081 reconstruction practice with an auxiliary modelling scheme such as Fourier-domain (Xie et al., 2022; 082 Zhang et al., 2022) masking. 083

Adopting the above insight as a core learning mechanism, we frame time series modelling into finding 084 an interpolation or extrapolation solution between input and target directly in the Fourier domain 085 and propose NFM. To achieve this, we introduce two main learning modules that operate directly in the Fourier domain and equip NFM with them. (i) A complex-valued learnable frequency token (LFT) is proposed to capture effective spectral priors and enable flexible frequency extension for 088 the frequency interpolation and extrapolation. (ii) Implicit neural Fourier filter (INFF) as a principal 089 processing operator is designed to realize an expressive continuous global convolution for learning the interpolation or extrapolation in the Fourier domain. We apply NFM to various datasets and 091 distinct tasks of scenarios with both normal (constant) and unseen discretization rate, and show that NFM achieves state-of-the-art performance in a remarkably compact form - forecasting with 27K, 092 anomaly detection with 6.6K, and classification with 37K parameters. 093

094

2 RELATED WORK

096

Frequency representations for time series modelling. There has been a growing attention for 098 processing time series and learning temporal dynamics of it with Fourier-domain information and/or through Fourier-domain operations in neural-based time series modelling methods. Autoformer 100 (Wu et al., 2021), FEDformer (Zhou et al., 2022b), and FourierGNN (Yi et al., 2024a) adopt a 101 frequency-based mixing mechanism as a main operator for learning temporal dependencies, where 102 the computation efficiency is ensured by operating in the Fourier domain with the Fast Fourier 103 Transform (FFT) algorithm. FiLM (Zhou et al., 2022a), BTSF (Yang & Hong, 2022), and TimesNet 104 (Wu et al., 2022) leverage frequency representations in conjunction with time-domain representation 105 to have a better realization of long-term dependency and global patterns (low frequency components). COST (Woo et al., 2021) and Autoformer exploit the Fourier domain as an explicit inductive bias, 106 utilizing it not only for efficient computation but for decomposing periodic patterns from complex time 107 series. TF-C (Zhang et al., 2022) enhances transferability of representations by integrating spectral



Figure 2: Overall workflow (forecasting scenario is exemplified) of the proposed NFM which deals with discrete signals as continuous-time elements in the compact function space through Fourier lens. NFM finds an interpolation/extrapolation from discrete input to target directly in the Fourier domain.

representations and introducing Fourier-domain augmentation in contrastive learning framework. The current utilization is capitalized on enabling efficient computation over long sequences and complementing the point-wise time-domain representations with the overarching representations. Differing from these works, Yi et al. (2024b) designs FreTS fully with Fourier operators and models time series fully in the Fourier domain for forecasting task. Our work follows in a similar vein to FreTS in the sense that we model time series directly in the Fourier domain and fully leverage on learning with the samples of functional representations of discrete signals, but with the distinct difference in the used core learning mechanism, frequency extrapolation/interpolation.

FITS. More recently, Frequency Interpolation Time Series analysis baseline (FITS) (Xu et al., 2023), a remarkably lightweight frequency-domain linear model, is introduced to address time-series problems. While NFM is designed leveraging the analogous principle (frequency interpolation and extrapolation) as FITS, there are several improvements essential for practicality - see details in **Appendix A**. In short, NFM is a generalization of FITS, that can 1) model both multivariate and univariate time series, 2) readily scale up, 3) be adaptive to variable-length inputs/outputs, and 4) be as compact as FITS and even surpass FITS's compactness while yielding better performance.

3 NEURAL FOURIER MODELLING

In this section, we provide an overview of NFM (Section 3.1) and introduce two main learning modules, Learnable Frequency Tokens (LFT) (Section 3.2) and Implicit Neural Fourier Filter (INFF) (Section 3.3). To begin, we first provide notations and necessary preliminaries below.

Notations. Considering a c-channel signal $x \in (\mathcal{D}, \mathbb{R}^c)$ and a target function $y \in (\mathcal{D}, \mathbb{R}^{d_y})$ defined on some temporal domain $\mathcal{D} \subset \mathbb{R}$, let $D_i = (x_i, y_i) \subset \mathcal{D}$ be a *i*th pair of input signal and target in domain \mathcal{D} . We denote $I_O = \{0, ..., O-1\}$ a set of indices for integer $O \ge 1$ and define $x[n \in I_N]|_{D_i} = x(n/f_x)|_{D_i}$ as a time series with N-point discretization at rate f_x over the timespan $[0, T_x)$, i.e., $N := T_x f_x$. The target y can be in any form, depending on tasks (see Appendix E), and its information spanning on its L-point latent representation $z[n \in I_L]|_{D_i}$ over an output timespan $[0, T_y(\geq T_x))$ at sampling rate f_y , where $L(\geq N) \coloneqq T_y f_y$. For notational simplicity, we drop $|_{D_i}$ and set T_x to a unit timespan and subsequently the input sampling rate $f_x = N$.

Preliminary: Relationship between input and output discretization. In our framework, it is 161 important to pay an attention to the relationship between the input and output discretization N and L as it removes ambiguities in formulating time series problems. For example, a task with $L \neq N$

173

178 179

181 182

183

195

196

199 200 201

would require modelling time series either across timespan or about the same timespan but at different discretization rate. To explicitly consider it, we denote an interpolation factor $m_{\tau} := T_y/T_x$ and an extrapolation factor $m_f := f_y/f_x$ and relate N and L with respect to these factors as follows:

$$\frac{L}{N} = \frac{T_y f_y}{T_x f_x} = m_\tau m_f \tag{1}$$

168 m_{τ} and m_f can be understood as the input-and-output ratio with respect to timespan and discretization 169 rate in modelling, respectively. NFM is designed to effortlessly switching its forward processing 170 between frequency interpolation when $m_{\tau} > 1$ and frequency extrapolation when $m_f > 1$ with 171 adoption of LFT in Section 3.2, which allows an easy specialization of NFM to various time series 172 tasks.

Preliminary: Discrete Fourier Transform (DFT). DFT is a computational tool widely used to convert data in physical domain (e.g., time and spatial) to spectral representations. Given the finitelength sequence $x[n \in I_N]$, the DFT and its inverse (IDFT) for recovery to the original sequence, $\mathcal{F}(\cdot)$ and $\mathcal{F}^{-1}(\cdot)$, are defined as follows:

$$X[k] = \mathcal{F}(x) \coloneqq \sum_{n=0}^{N-1} x[n] e^{-j2\pi kn/N}$$
(2)

$$x[n] = \mathcal{F}^{-1}(X) \coloneqq \frac{1}{N} \sum_{k=0}^{N-1} X[k] e^{j2\pi kn/N}$$
(3)

184 where $k \in I_N$, and the capital letter X is a complex spectral representation of its time-domain variable 185 x. Importantly, each k frequency component represents the entirety of the sequence x summarized at different oscillation. The convolution theorem, in conjunction with the IDFT, leverages this 186 characteristic of the FT and provides an efficient way for performing convolution operations through 187 point-wise multiplication in Fourier domain (Rabiner & Gold, 1975; McGillem & Cooper, 1991). 188 We adopt this insight in NFM and introduce a new type of neural Fourier filters (NFFs) playing as a 189 continuous global token mixer that is both expressive and adaptive yet lightweight in Section 3.3. 190 Besides, for computation, we utilize an efficient algorithm of DFT, fast Fourier transform (FFT) 191 that offers $\mathcal{O}(N \log N)$ complexity at optimum, as well as its conjugate symmetry property (i.e., 192 $X[N-k] = X^*[k \in I_{K_N}]$. $K_N := |N/2| + 1$ is the number of the first half frequency components 193 of sequence length N. 194

3.1 OVERVIEW OF NFM

197 Linear frequency interpolation and extrapolation. We begin with rewriting IDFT in Eq.(3) for 198 the desired output of a model, z, as follows:

$$z[n \in I_L] = \frac{1}{L} \sum_{k=0}^{L-1} Z[k] e^{i2\pi kn/L}$$
(4)

202 One straightforward way to express Eq.(4) with respect to the given input sequence $x[n \in I_N]$ would 203 be by taking a linear system to the spectral representation such that $Z[k] = m_{\tau} m_f (W \mathcal{F}(x[n]) + b)$, 204 where $W \in \mathbb{C}^{K_L \times K_N}$ and $b \in \mathbb{C}^{K_L}$ are weights and bias term, respectively. Directly designing 205 $y = z = \mathcal{F}^{-1}(m_{\tau}m_f(\boldsymbol{W}\mathcal{F}(x)+b))$ with adoption of a heuristic low-pass filter to further reduce 206 the dimensionality of W and b yields exactly FITS. FITS greatly appreciates the simplicity and 207 lightweight-ness of the linear system and presents a solution to low-resources tasks like edge comput-208 ing. Nevertheless, it lacks in several aspects as a general solution to a range of time series analysis 209 (refer to **Appendix A**). 210

SIGNATE: NFM. We aim to build a general-purpose time series model as a composition $y = \mathcal{P} \circ \mathcal{M}(x)$, that learns mapping between infinite-dimensional spaces of input signal and target function directly in the Fourier domain, given discrete observations D. Especially, we introduce NFM for encoder $\mathcal{M}: \mathbb{R}^{N \times c} \to \mathbb{R}^{L \times d}$ that acts globally on the input signal x and seeks the mapping as a frequency interpolation/extrapolation to $z = \mathcal{M}(x)$. The target $y = \mathcal{P}(z)$ is evaluated from the latent representation z through a predictor (e.g., a local-to-local transformation $\mathcal{P}: \mathbb{R}^d \to \mathbb{R}^{d_y}$ and a global-to-local

216



Figure 3: Illustration of NFM architecture consisting of three main learning modules: 1) LFT block to 234 allow flexible frequency extension and provides effective spectral priors, 2) a plain MLP for channel 235 mixing, and 3) INFF module for effective token mixing with global convolution operation. The M in LFT block denotes frequency extension operation. 236

238 $\mathbb{R}^{L \times d} \to \mathbb{R}^{d_y}$). As depicted in **Figure 2** and **Figure 3**, the overall workflow of NFM is described 239 in two steps. Given input sequences projected to hidden dimension $\bar{x} \in \mathbb{R}^{N \times d}$ (see Appendix D), (i) it is first tailored to temporal embedding $z_0 \in \mathbb{R}^{L \times d} = \mathcal{F}^{-1}(Z_0) = \text{LFT}(\bar{x})$, in that its spectral 240 embedding $Z_0 \in \mathbb{C}^{K_L \times d}$ explicitly accounts for an extension of the original sequence with respect 241 to m_{τ} and m_{f} . (ii) Then, the low-level embedding tokens z_0 are polished iteratively through a 242 stack of l mixing blocks consisting of a channel-mixing module and global convolution, INFF, as in 243 $z_i = Mixer(z_{i-1})$ for $i = \{1, ..., l\}$. In the following sections, we detail out the two main modules, 244 LFT and INFF. 245

246 247

233

237

3.2 LEARNABLE FREQUENCY TOKENS

248 In our framework, we decouple the process of extending the spectral representations of original 249 sequence to that of target domain (i.e., $N \to L$ and $K_N \to K_L$), from learning the abstract 250 coefficients of Fourier interpolation/extrapolation. We denote such extended spectral representation with absence of the weighting coefficients $\bar{Z}_0 \in \mathbb{C}^{K_L \times d}$, and it is obtained by simply initializing the \bar{Z}_0 with zeros and rearranging $\bar{X} = \mathcal{F}(\bar{x})$ onto it with scaling for the extension such that 251 252 253 $\overline{Z}_0[|m_t k|] = m_\tau m_f \overline{X}[k]$ for all $k \in I_{K_N}$, where $|\cdot|$ is floor operation.

254 The result from the above extension is equivalent to the ones from applying zero-padding and/or 255 zero-interleaving to X. The operation itself is non-parametric and allows handling variable-length 256 input sequence. However, directly adopting Z_0 is not effective for learning since extending spectral 257 representations with zero coefficients itself does not bring in any information gain. One natural 258 solution to this is introducing extra embeddings that are learned to encapsulate certain abstraction 259 and priors in data during optimization. Indeed, it has become a canonical practice, especially in many 260 Transformer models (Devlin et al., 2018; Chen et al., 2022b; 2023; Wang & Chen, 2020), to enrich 261 the models' learning capability and performance.

262 Inspired by this, we introduce LFT that can be learned without a-priori and applied directly in 263 the complex Fourier domain as expressive spectral priors across sequences. Especially, we design 264 the LFT by representing the desired spectral priors as a composition of \mathcal{F} and an implicit neural 265 representation (INR) (Sitzmann et al., 2019b; Chen et al., 2021) $\phi \colon \mathbb{R} \to \mathbb{R}^d$ that maps a temporal 266 location $\tau \in [0, T_u)$ to abstract temporal priors v corresponding to that time location. Notably, the 267 LFT can be characterized as samples from continuous-time Fourier transform (CTFT) of the temporal priors, allowing sampling of the frequency tokens within bandwidth from any arbitrary temporal 268 locations without a need of re-training. We first obtain the learnable frequency tokens $V[k] \in \mathbb{C}^d$ by 269 sampling $v[n] = \phi(\tau_n)$ at temporal locations $\tau_n = \{n/f_u | n \in I_L\}$ and add it to \overline{Z}_0 to get Z_0 . This entire process of the LFT block is completed as follows:

$$V[k \in I_{K_L}] = \mathcal{F}(\text{InstanceNorm}(\phi(\tau_n))),$$

$$Z_0[k] = (\bar{Z}_0[k] + V[k]),$$

$$z_0[n] = \mathcal{F}^{-1}(Z_0[k])$$

274 275

282

283

289

290 291

272 273

276 We apply an instance normalization before applying DFT to v[n] to the LFT block as shown in **Figure** 277 **3**. This removes DC priors of each channel and a source of internal covariate shift, thereby helping 278 the LFT effectively learn the priors over spectrum - we find in the experiment that the energy of 279 the spectral priors without it is often largely concentrated in DC component, preventing effective 280 learning. The ϕ is expressed by MLP with a periodic activation function (Sitzmann et al., 2020) - see 281 **Appendix D**.

3.3 IMPLICIT NEURAL FOURIER FILTER

Neural Fourier filters. Upon convolution theorem and FT, a convolution kernel can be directly defined in the Fourier domain with arbitrary resolution and its size being as large as the length of inputs. This, then, gives rise to a way for instantiating an efficient global convolution operator $\mathcal{K}: (\mathcal{D}, \mathbb{R}^d) \to (\mathcal{D}, \mathbb{R}^d)$ as follows (Guibas et al., 2021; Li et al., 2020):

$$\mathcal{K}(z)(\tau) = \int_{\mathcal{D}} \kappa(\tau - s) z(s) \, \mathrm{d}s, \quad \forall \tau \in \mathcal{D}$$
(6)

(5)

$$\mathcal{K}(z)(\tau) = \mathcal{F}^{-1}(\mathcal{R} \cdot \mathcal{F}(z))(\tau), \qquad \forall \tau \in \mathcal{D}$$
(7)

292 where Eq.(6) and Eq.(7) express convolution operator on physical space of the signals and its 293 equivalent form with convolution theorem and the FT applied, respectively. $\mathcal{R} \equiv \mathcal{F}(\kappa)$ is a Fourier 294 filter defined directly in the Fourier domain and parameterized by a neural network. While a shallow 295 fully-connected network is generally sufficient to parameterize \mathcal{R} , much design concerns are put in 296 how to achieve properties of NFFs that are desirable for modelling. We summarize them and compare the existing NFFs from the designing perspective in Appendix B. In short, it is highly desirable 297 to have a NFF that is memory-efficient, length-independent (i.e., flexible), instance-adaptive, and 298 mode-adaptive (i.e., expressive and generalized across spectrum). We design INFF that satisfy these 299 properties below. 300

301 302

303 304

305

309

INFF. INFF is formulated for modulating the embedding tokens z globally in the Fourier domain to \hat{z} through a Fourier filter, $\mathcal{R}[k] \in \mathbb{C}^d$, as follows:

$$\hat{z} = \mathcal{F}^{-1}(\mathcal{R}(z_0) \odot \mathcal{F}(z)) \tag{8}$$

where \odot denotes Hadamard product. The computation of \mathcal{R} is conditioned on the initial spectral embedding Z_0 , for which a reason will be clarified later. Here, the use of an INR for encapsulating abstract spectral priors in **Section 3.2** is extended to define \mathcal{R} with implicit filter coefficients.

$$\mathcal{R}(z_0) \coloneqq \mathcal{W}(\mathcal{F}(\text{InstanceNorm}(\phi(\tau_n) + z_0))) \tag{9}$$

310 Recalling that $\mathcal{F}(\phi(\tau_n))$ implicitly models CTFT, in this formulation the Fourier filter has filter 311 coefficients defined uniquely for each period in spectrum, but with a single parameterization. That is, 312 the designed Fourier filter is in a compact form and can readily handle variable-length sequence with 313 unique frequency coefficients (i.e., mode-adaptive). Unfortunately, the INFF with the filter solely 314 defined by the $\mathcal{F}(\phi(\tau_n))$ would lack expressivity due to the reliance on depth-wise convolution (i.e., no channel mixing) and struggle to generalize across different instances. We further improve the filter 315 on both factors by aggregating its temporal coefficients with the initial temporal embedding tokens 316 z_0 and processing them through a complex-valued MLP, $\mathcal{W} \colon \mathbb{C}^d \to \mathbb{C}^d$. Note that we use ReLU 317 for non-linearity and have no bottleneck or expansion factor for the intermediate dimension. With 318 this, each feature of the updated kth filter coefficient now represents a mixture of all features of kth 319 implicit filter coefficient and all features of the kth spectral representation of the input embedding. 320

Finally, INFF is put together with a plain channel-mixing block and Layer Normalization to form a complete Mixer block. We configure the components as pre-channel-mixing and post-normalization with a skip connection in each Mixer block and stack multiple Mixer blocks followed by a final channel-mixing block to constitute a NFM backbone network $\mathcal{M}(\cdot)$ as shown in **Figure 3**.

Model (params)	NI (27	F M 7K)	FI (~0	TS .2M)	N-Li (∼0	inear .5M)	iTrans (~5	former .3M)	Patch (~8	nTST .7M)	Time (∼0	esNet .3B)
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTm1	0.345	0.375	0.357	0.380	0.366	0.383	0.371	0.401	0.353	0.382	0.400	0.406
ETTm2	0.250	0.311	0.250	0.313	0.258	0.315	0.276	0.337	0.256	0.317	0.291	0.333
ETTh1	0.407	0.420	0.407	0.420	0.413	0.422	0.503	0.491	0.413	0.434	0.458	0.450
ETTh2	0.356	0.400	0.334	0.382	0.343	0.389	0.405	0.430	0.331	0.381	0.414	0.427
Weather	0.227	0.269	0.241	0.280	0.254	0.288	0.255	0.289	0.227	0.264	0.259	0.287
Electricity	0.159	0.251	0.163	0.254	0.169	0.262	0.163	0.259	0.159	0.253	0.193	0.298
Traffic	0.391	0.260	0.411	0.280	0.433	0.290	0.376	0.270	<u>0.391</u>	0.264	0.620	0.336

Table 1: Long-term forecasting results averaged over 4 horizons. The best averaged results are in **bold** and the second best are <u>underlined</u>. Full result is available in **Appendix E**.

4 EXPERIMENTS

324

325

326 327 328

337 338

339

340

341

342

349

350 351

352

353

354

355

We conduct extensive experiments to demonstrate the effectiveness of NFM and its competence as a general solution to time series modelling. Note that we provide only a brief description about the setting for each experiment below, but one can find all details in the supplementary section (**Appendix** $C \sim E$). Our code is publicly available at: https://github.com/minkiml/NFM.

Implementation. A general-purpose time-series model allows one to address a range of tasks as well as various time series modalities without significant architectural modifications (i.e., no injection of task-specific inductive bias). To this end, we implement a single NFM backbone $\mathcal{M}(\cdot)$ and use it for all tasks (only differ by some hyper-parameters such as hidden size and the number of mixer blocks) by equipping it with a task-specific linear predictor $\mathcal{P}(\cdot)$ (simply a fully-connected layer). For all experiments, we use a single NVIDIA A100 GPU.

4.1 TIME SERIES MODELLING

We begin with showcasing the efficacy of NFM in three distinct time-series tasks, including long-term forecasting, anomaly detection, and classification under their conventional scenario ($f_x^{train} = f_x^{test}$). We follow the experimental setups: forecasting (Zhou et al., 2021), anomaly detection (Xu et al., 2021a), and classification (Romero et al., 2021). Refer to **Appendix D** for details.

356 **Long-term forecasting.** Long-term times-series forecasting task ($m_f = 1$ and $m_\tau > 1$) is con-357 ducted on 7 benchmark datasets over 4 horizons (96, 192, 336, and 720). In Table 1, The averaged 358 results of NFM on MSE and MAE metrics are compared with 5 SOTA forecasting models - FITS 359 (Xu et al., 2023), N-Linear (Zeng et al., 2023), iTransformer (Liu et al., 2023), PatchTST (Nie et al., 360 2022), TimesNet (Wu et al., 2022). While the performance of NFM is highly competitive with that of the dedicated forecasting models with hundreds or millions of trainable parameters, we put an extra 361 emphasis on its compactness. Especially, the number of parameters that need to be trained in NFM is 362 only around 27K for all forecasting cases, which sheds light on both low-resource on-device learning 363 and processing. This result stands out that of the simple linear models like N-Linear and FITS, and 364 other SOTA by large margin, for which the number of parameters is subject to the length of both 365 lookback and prediction horizon whereas the number of parameters in NFM is decoupled from the 366 length of input or target prediction. Please see Appendix E for more analysis about the results. 367

368 **Anomaly detection.** We frame the anomaly detection task as learning correct contexts (dominant 369 and normal contexts) and evaluating the validity of observations within the contexts. More specifically, 370 our approach is to learn the correct contexts by reconstructing complete sequences from their down-371 sampled counterparts (i.e., $m_f > 1$ and $m_{\tau} = 1$) - refer to Appendix D for detail. For quantitative 372 evaluation, we employ 4 popular anomaly detection benchmark datasets and compare the performance 373 with 6 baselines - vanilla Transformer (Vaswani et al., 2017), PatchTST (Nie et al., 2022), TimesNet 374 (Wu et al., 2022), ADformer (Xu et al., 2021a), N-linear (Zeng et al., 2023), and FITS (Xu et al., 375 2023). In **Table 2**, NFM shows effectiveness in anomaly detection task, presenting top-tier results in three datasets (SMD, MSL, and PSM) and third-tier result in SMAP. This result in NFM is achieved 376 with a compact model size of only around **6.6K** parameters, which follows that of FITS (1.3K) and 377 stands out that of the deep feature learning models (TimesNet and ADformer) by considerably large

Table 2: Anomaly detection results (F1-score) on 4 datasets, where higher F1-score indicates better performance. Full tabular result is available in Appendix E.

Model (params)	SMD	MSL	SMAP	PSM
Transformer (0.2M)	75.95	81.93	69.70	88.75
PatchTST (0.2M)	82.11	80.51	69.11	96.00
TimesNet (~28M)	83.27	81.70	73.23	97.30
ADformer*(4.8M)	76.38	81.78	71.18	83.14
N-Linear (10K)	81.81	81.18	67.50	95.77
FITS (1.3K)	81.67	80.77	64.07	96.60
NFM (6.6K)	84.32	82.46	70.88	97.51

The joint criterion in ADformer is replaced with the simple

reconstruction error to compute anomaly score for fair comparison.

Table 3: Classification accuracy (%) on Speech-Command. \sim denotes inapplicable (prohibitively slow) or computationally not possible on single GPU.

	S	peechComr	nand
Model (params)	MFCC	RAW (SR=1)	RAW (SR=1/2)
ODE-RNN (89K)	65.90	~	~
GRU- Δt (89K)	20.0	\sim	\sim
GRU-ODE (89K)	44.8	\sim	\sim
NCDE (89K)	88.5	\sim	\sim
NRDE (89K)	89.8	16.49	15.12
S4 (400K)	93.96	96.17	94.11
CKConv (100K)	95.27	71.66	65.96
Transformer (800K)	90.75	\sim	\sim
NFM (37K)	94.23	90.94	<u>90.30</u>

Table 4: Forecasting results (MSE) and performance drops (%) at different sampling rate. Three models, including PatchTST (Transformer-based), N-Linear (time-domain linear model), and FITS (frequency-domain linear model) are opted for comparison. The best performance is in **blue** and the least performance drop in red. Full tabular results over all horizons can be found in Appendix E.

Model	SR	ETT	m1	ETT	Γm2	Weather		
		96	720	96	720	96	720	
PatchTST N-Linear	1/4	0.394 (34.5 ↓)	0.433 (4.1 ↓)	0.232 (39.8 ↓)	0.382 (5.5 ↓)	0.212 (42.3 ↓)	0.329 (4.8 ↓)	
	1/6	0.434 (48.5 ↓)	$0.443~(6.5\downarrow)$	0.265 (59.6 ↓)	0.396 (9.4 ↓)	$0.236(58.4\downarrow)$	0.335 (6.7 ↓)	
N-Linear	1/4	0.436 (42.5 ↓)	0.469 (8.3 ↓)	0.288 (72.5 ↓)	0.409 (11.1 ↓)	0.268 (47.3 ↓)	0.351 (3.8 ↓)	
	1/6	0.482 (57.5 ↓)	0.471 (8.8 ↓)	0.281 (68.3 ↓)	$0.422~(14.7\downarrow)$	0.280 (53.8 ↓)	0.364 (7.7 ↓)	
EITC	1/4	0.348 (12.6 ↓)	0.426 (2.9 ↓)	0.198 (21.5 ↓)	0.356 (2 .0 ↓)	0.182 (7.7 ↓)	0.325 (1.2 ↓)	
FITS 1	1/6	0.368 (19.1 ↓)	0.437 (5.3 ↓)	$0.216(31.7\downarrow)$	$0.364~(4.3\downarrow)$	$0.195~(14.8\downarrow)$	0.329 (2.5 ↓)	
NEM	1/4	0.299 (4 . 5 ↓)	0.407 (0.2 ↓)	0.179 (11.9 ↓)	0.356 (2 . 0 ↓)	0.164 (6 . 5 ↓)	0.315 (1.0 ↓)	
NFIN	1/6	0.319 (11.5 ↓)	0.414 (2 .0 ↓)	0.189 (18.2 ↓)	0.358 (2.6 ↓)	0.173 (12.3 ↓)	0.318 (1.9 1)	

410 411

378

379

380

381

382

390

391 392 393

394

397

> margin. Note that the remarkably small model size of FITS is highly subject to cases (short length of input and target while their energy dominantly spans in low frequency regime).

412 **Classification.** We evaluate the effectiveness of NFM in time-series classification task ($m_f = 1$ 413 and $m_{\tau} = 1$) using SpeechCommand dataset (Warden, 2018) that provides both MFCC features (N = 161) and raw waveform (N = 16k). In **Table 3**, we report the classification accuracy (ACC) 414 and compare it with that of 8 different baselines (7 continuous-time models and 1 Transformer) - ODE-415 RNN (Rubanova et al., 2019), GRU- Δt (Kidger et al., 2020), GRU-ODE (De Brouwer et al., 2019), 416 NCDE (Kidger et al., 2020), NRDE (Morrill et al., 2021)), S4 (Gu et al., 2021), CKConv (Romero 417 et al., 2021), and vanilla Transformer (Vaswani et al., 2017). Overall, NFM yields competitive 418 performance in both cases of processing MFCC features - CKConv (95.27%) vs. NFM (94.23%) 419 vs. S4 (93.96%), and raw waveform - S4 (96.17%) vs. NFM (90.94%) vs. CKConv (71.66%) with 420 far smaller model size (**37K**) among all baselines. The neural ODE-based models are generally in 421 a compact form with a weight-tying architecture but struggle dealing with long series due to the 422 need of solving differential equations for a long step. Vanilla transformer suffers from huge memory 423 occupancy and thus is unable to process the raw waveform in a single GPU setting.

424 425

426

4.2 EVALUATION AT DIFFERENT DISCRETIZATION RATES

427 In practice, it is highly desirable for a model to have a resolution-invariance property that readily en-428 ables generalization of the learned solutions to unseen discretizations without significant performance lose. To this end, we now reveal this aspect of NFM, which learns function-to-function mappings, by 429 conducting the classification and forecasting tasks in a scenario where the observations are sampled 430 at different (unseen) sampling rate during testing time (SR= f_x^{test}/f_x^{train}). The classification result 431 in **Table 3** shows that NFM has the least performance drop, yielding only around $0.7\% \downarrow$ degradation

from the original performance on the input sequences sampled at unseen sampling rate of SR= 1/2, compared to NRDE (8.31% \downarrow), CKConv (7.95% \downarrow), and S4 (2.14% \downarrow). In the forecasting task, the similar result is drawn as shown in **Table 4**, where the performance degradation of NFM in all cases is considerably lower than that of the 3 baselines. We provide more details on the forecasting results in **Appedix E**.

438 4.3 EXPLORATION OF NFM

437

439

Effects of LFT and INFF. We examine the ef-440 fect of the proposed LFT in NFM (INFF + LFT) 441 by comparing a NFM with randomly initialized 442 complex-valued learnable weights as frequency 443 tokens (namely, *Naive*) and without frequency 444 tokens (namely, INFF-only). As shown in Fig-445 ure 4, the improvement on performance is re-446 markable, reporting 82.6 \rightarrow 90.9 (10.0% \uparrow) 447 and 79.4 \rightarrow 90.9 (14.5% \uparrow) compared to 448 *Naive* and *INFF-only*, respectively. In the case 449 of testing at different input resolution (SR=1/2), the improvement is more significant, yielding 450 $78.6 \rightarrow 90.4 \ (15.0\% \uparrow) \ \text{and} \ 74.1 \rightarrow 90.4$ 451 $(22.0\% \uparrow)$ respectively. These results from 452 INFF + LFT are achieved with ~ 8.6 times 453 smaller and ~ 1.16 times larger model size com-454 pared to the Naive and INFF-only, which shows 455



Figure 4: Comparison of NFM with different ablation cases on SpeechCommand dataset. PMU is peak memory usage during inference time.

that LFT is conducive to learning spectral priors in a compact form. Furthermore, we compare INFF 456 with 4 different SOTA NFFs, including FNO (Li et al., 2020), AFNO (Guibas et al., 2021), GFN (Rao 457 et al., 2021), AFF (Huang et al., 2023). The result in Figure 5 shows the overall performance of *INFF* 458 + LFT surpasses the others + LFT by notable margin. The instance-adaptive NFFs - AFNO (2.7% \downarrow) 459 and AFF (2.9% \downarrow), perform more robustly against different input resolution than the mode-adaptive 460 ones - FNO (11.6% \downarrow) and GFN (9.5% \downarrow). As analysed in Appendix B, INFF which is both instanceand mode-adaptive further improves the robustness ($0.7\% \downarrow$). This superiority of INFF comes only at 461 the cost of a minor complexity and memory usage increase. Besides, it is noteworthy that dealing with 462 different input resolution is originally not possible in the mode-adaptive ones (FNO and GFN) due to 463 the fixed-length operation without a heuristic low-pass filter, but becomes possible with the adoption 464 of LFT in NFM framework. The similar results are obtained on forecasting task (see Appendix E for 465 it and full tabular results). 466

467 Visualization on INFF. We study the be-468 haviour of INFF in terms of what representations 469 it potentially leads to be learned. Specifically, 470 one would expect to see INFF learning effec-471 tive filter coefficients such that INFF amplifies frequencies of relevant information while sup-472 pressing frequencies of irrelevant one. To con-473 firm this, we synthesize simple single channel 474 band-limited (up to Nyquist frequency $f_{nyquist}$) 475 signals by composing multiple frequency com-476 ponents. During generation, we assign class 477 labels to them according to certain combina-478 tions of the frequency components spanning in 479 a frequency range $[f_A, f_B(< f_{nyquist})]$. Please 480 refer to Appendix C for details about the gen-481 eration. For experiment, we sample 10 classes 482 of sequences of length N = 2000 at $f_x = N$ 483 with the class frequencies spanning in the range [320, 590]. Figure 5 provides a visualization of 484 the resulting INFF trained on the synthetic data, 485 where it is clearly shown that finding a correct



Figure 5: Visualization of INFF on synthetic data. The **top-left** figure shows the frequencies of input sequence, **bottom-left** the frequencies of filtered sequence, **bottom-right** the learned INFF's coefficients, and **top-right** the coefficients averaged over hidden dimension.



Figure 6: Scaling behavior (performance - the line plots) of NFM with respect to varying hidden dimension and depth on (a) ETTm1 (forecasting over the horizon of 96) and (b) SC-MFCC (classification) datasets. The bar plots represent the number of parameters computed at each set of hidden dimension and depth. Baseline performances are also included with the dashed horizontal lines for comparison.

solution comes with INFF learning its filter coefficients aligned with the frequency range [320, 590]
 of the class labels (relevant information).

Scaling behaviour and compactness of NFM. We present how scaling over hidden dimension (*d*) and depth affects NFM's performance in **Figure 6**. While it is clearly seen that increasing the NFM's scale consistently leads to improved performance, we observe that NFM rapidly reaches a regime of competitive performance despite having significantly fewer parameters compared to all baseline models (as indicated by the dashed horizontal lines). This highlights the compactness of NFM, which achieves high performance without the need for excessive parameter scaling.

513 514 515

507

508

509

510

511

512

5 CONCLUSION

516 517

517 In this work, we have introduced, NFM, modelling time series directly in Fourier domain by for-518 mulating the Fourier-domain data manipulation into Fourier interpolation and extrapolation. NFM 519 with *learnable frequency tokens* and *implicit neural Fourier filter* enjoys the intriguing properties of 520 the FT in learning, resulting in a remarkable compactness and continuous-time characteristics. Our 521 experiments demonstrate that NFM can be a powerful general solution to time series analysis across 522 a range of datasets, tasks, and scenarios and show that it is possible to achieve the state-of-the-art 523 performance with a remarkably compact model, compared to models with hundreds thousand and 524

Limitations and Future Work. Despite the fact that NFM models time series in function space through Fourier lens, the current implementation is not directly suitable for some other dynamic scenarios such as handling irregular time series. A reason for this is that the FFT for efficient transformation requires uniformly-sampled sequence thus not applicable while a naive algorithm with computing a DFT matrix and its pseudo inverse is not only too slow but also memory-intensive for every long and multivariate but irregular time series. While addressing this challenging scenario is highly valuable in time series analysis, we are currently improving NFM on this matter.

- 532 533 534
- References

Ahmed Abdulaal, Zhuanghua Liu, and Tomer Lancewicki. Practical approach to asynchronous
 multivariate time series anomaly detection and localization. In Proceedings of the 27th ACM
 SIGKDD conference on knowledge discovery & data mining, pp. 2485–2494, 2021.

- 538
- Ahmed Alaa, Alex James Chan, and Mihaela van der Schaar. Generative time-series modeling with fourier flows. In International Conference on Learning Representations, 2020.

540 541	Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. <u>Advances in neural information processing systems</u> , 31, 2018.
543 544 545	Wenchao Chen, Long Tian, Bo Chen, Liang Dai, Zhibin Duan, and Mingyuan Zhou. Deep varia- tional graph convolutional recurrent network for multivariate time series anomaly detection. In <u>International Conference on Machine Learning</u> , pp. 3621–3633. PMLR, 2022a.
546 547 548	Yinbo Chen, Sifei Liu, and Xiaolong Wang. Learning continuous image representation with local implicit image function. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 8628–8638, 2021.
549 550 551 552	Yinpeng Chen, Xiyang Dai, Dongdong Chen, Mengchen Liu, Xiaoyi Dong, Lu Yuan, and Zicheng Liu. Mobile-former: Bridging mobilenet and transformer. In <u>Proceedings of the IEEE/CVF</u> <u>Conference on Computer Vision and Pattern Recognition</u> , pp. 5270–5279, 2022b.
553 554 555 556	Yuxiao Chen, Jianbo Yuan, Yu Tian, Shijie Geng, Xinyu Li, Ding Zhou, Dimitris N Metaxas, and Hongxia Yang. Revisiting multimodal representation in contrastive learning: from patch and token embeddings to finite discrete tokens. In <u>Proceedings of the IEEE/CVF Conference on Computer</u> <u>Vision and Pattern Recognition</u> , pp. 15095–15104, 2023.
557 558 559	Hyeong-Seok Choi, Jang-Hyun Kim, Jaesung Huh, Adrian Kim, Jung-Woo Ha, and Kyogu Lee. Phase-aware speech enhancement with deep complex u-net. <u>arXiv preprint arXiv:1903.03107</u> , 2019.
560 561 562 563	Edward De Brouwer, Jaak Simm, Adam Arany, and Yves Moreau. Gru-ode-bayes: Continuous mod- eling of sporadically-observed time series. <u>Advances in neural information processing systems</u> , 32, 2019.
564 565 566	Ruizhi Deng, Bo Chang, Marcus A Brubaker, Greg Mori, and Andreas Lehrmann. Modeling continuous stochastic processes with dynamic normalizing flows. <u>Advances in Neural Information</u> <u>Processing Systems</u> , 33:7805–7815, 2020.
567 568 569	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. <u>arXiv preprint arXiv:1810.04805</u> , 2018.
570 571	Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured state spaces. In International Conference on Learning Representations, 2021.
572 573 574 575	John Guibas, Morteza Mardani, Zongyi Li, Andrew Tao, Anima Anandkumar, and Bryan Catan- zaro. Adaptive fourier neural operators: Efficient token mixers for transformers. <u>arXiv preprint</u> <u>arXiv:2111.13587</u> , 2021.
576 577	David Ha, Andrew M. Dai, and Quoc V. Le. Hypernetworks. In <u>International Conference on Learning</u> <u>Representations</u> , 2017.
578 579 580	Zhipeng Huang, Zhizheng Zhang, Cuiling Lan, Zheng-Jun Zha, Yan Lu, and Baining Guo. Adaptive frequency filters as efficient global token mixers. In <u>Proceedings of the IEEE/CVF International</u> <u>Conference on Computer Vision</u> , pp. 6049–6059, 2023.
582 583 584	Liming Jiang, Bo Dai, Wayne Wu, and Chen Change Loy. Focal frequency loss for image recon- struction and synthesis. In <u>Proceedings of the IEEE/CVF International Conference on Computer</u> <u>Vision</u> , pp. 13919–13929, 2021.
585 586 587	Patrick Kidger, James Morrill, James Foster, and Terry Lyons. Neural controlled differential equations for irregular time series. <u>Advances in Neural Information Processing Systems</u> , 33:6696–6707, 2020.
588 589 590	Chiheon Kim, Doyup Lee, Saehoon Kim, Minsu Cho, and Wook-Shin Han. Generalizable implicit neural representations via instance pattern composers. In <u>Proceedings of the IEEE/CVF</u> <u>Conference on Computer Vision and Pattern Recognition</u> , pp. 11808–11817, 2023.
592 593	Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Reversible instance normalization for accurate time-series forecasting against distribution shift. In International Conference on Learning Representations, 2021.

594 595	James Lee-Thorp, Joshua Ainslie, Ilya Eckstein, and Santiago Ontanon. Fnet: Mixing tokens with fourier transforms. <u>arXiv preprint arXiv:2105.03824</u> , 2021.
596 597 598	Zhe Li, Zhongwen Rao, Lujia Pan, Pengyun Wang, and Zenglin Xu. Ti-mae: Self-supervised masked time series autoencoders. <u>arXiv preprint arXiv:2301.08871</u> , 2023.
599 600 601	Zongyi Li, Nikola Borislavov Kovachki, Kamyar Azizzadenesheli, Kaushik Bhattacharya, Andrew Stuart, Anima Anandkumar, et al. Fourier neural operator for parametric partial differential equations. In <u>International Conference on Learning Representations</u> , 2020.
602 603 604 605 606	Shiqi Lin, Zhizheng Zhang, Zhipeng Huang, Yan Lu, Cuiling Lan, Peng Chu, Quanzeng You, Jiang Wang, Zicheng Liu, Amey Parulkar, et al. Deep frequency filtering for domain generalization. In <u>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</u> , pp. 11797–11807, 2023.
607 608 609	Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. itransformer: Inverted transformers are effective for time series forecasting. In <u>The Twelfth</u> <u>International Conference on Learning Representations</u> , 2023.
610	Clare D McGillem and George R Cooper. Continuous and discrete signal and system analysis. 1991.
611 612 613 614	Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. <u>Communications of the ACM</u> , 65(1):99–106, 2021.
615 616 617	James Morrill, Cristopher Salvi, Patrick Kidger, and James Foster. Neural rough differential equations for long time series. In <u>International Conference on Machine Learning</u> , pp. 7829–7838. PMLR, 2021.
618 619 620	Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In <u>The Eleventh International Conference on</u> <u>Learning Representations</u> , 2022.
622 623 624	Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. <u>arXiv preprint arXiv:1609.03499</u> , 2016.
625 626	Lawrence R Rabiner and Bernard Gold. Theory and application of digital signal processing. Englewood Cliffs: Prentice-Hall, 1975.
627 628 629 630	Aniruddh Raghu, Payal Chandak, Ridwan Alam, John Guttag, and Collin Stultz. Sequential multi- dimensional self-supervised learning for clinical time series. In <u>International Conference on</u> <u>Machine Learning</u> , pp. 28531–28548. PMLR, 2023.
631 632 633	Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In <u>International Conference</u> on <u>Machine Learning</u> , pp. 5301–5310. PMLR, 2019.
634 635 626	Yongming Rao, Wenliang Zhao, Zheng Zhu, Jiwen Lu, and Jie Zhou. Global filter networks for image classification. In <u>Advances in Neural Information Processing Systems</u> , 2021.
637 638 639	David W Romero, Anna Kuzina, Erik J Bekkers, Jakub Mikolaj Tomczak, and Mark Hoogendoorn. Ckconv: Continuous kernel convolution for sequential data. In <u>International Conference on</u> <u>Learning Representations</u> , 2021.
640 641	Yulia Rubanova, Ricky TQ Chen, and David K Duvenaud. Latent ordinary differential equations for irregularly-sampled time series. Advances in neural information processing systems, 32, 2019.
642 643 644 645	Lifeng Shen, Zhuocong Li, and James Kwok. Timeseries anomaly detection using temporal hierar- chical one-class network. <u>Advances in Neural Information Processing Systems</u> , 33:13016–13026, 2020.
646 647	Vincent Sitzmann, Michael Zollhoefer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In <u>Advances in Neural Information</u> Processing Systems, 2019a.

648 649 650	Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Contin- uous 3d-structure-aware neural scene representations. <u>Advances in Neural Information Processing</u> <u>Systems</u> , 32, 2019b.
652 653 654	Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. <u>Advances in neural information</u> processing systems, 33:7462–7473, 2020.
655 656 657 658 659	Ya Su, Youjian Zhao, Chenhao Niu, Rong Liu, Wei Sun, and Dan Pei. Robust anomaly detection for multivariate time series through stochastic recurrent neural network. In Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining, pp. 2828–2837, 2019.
660 661 662 663	Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. <u>Advances in Neural Information Processing Systems</u> , 33:7537–7547, 2020a.
664 665 666 667	Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. In <u>Proceedings of the 34th International Conference on Neural Information Processing Systems</u> , 2020b.
669 670 671	Chiheb Trabelsi, Olexa Bilaniuk, Ying Zhang, Dmitriy Serdyuk, Sandeep Subramanian, Joao Felipe Santos, Soroush Mehri, Negar Rostamzadeh, Yoshua Bengio, and Christopher J Pal. Deep complex networks. In International Conference on Learning Representations, 2018.
672 673 674 675	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. <u>Advances in neural information processing systems</u> , 30, 2017.
676 677 678	Yu-An Wang and Yun-Nung Chen. What do position embeddings learn? an empirical study of pre- trained language model positional encoding. In <u>Proceedings of the 2020 Conference on Empirical</u> <u>Methods in Natural Language Processing (EMNLP)</u> , pp. 6840–6849, 2020.
679 680 681	Pete Warden. Speech commands: A dataset for limited-vocabulary speech recognition. <u>arXiv preprint</u> <u>arXiv:1804.03209</u> , 2018.
682 683 684	Gerald Woo, Chenghao Liu, Doyen Sahoo, Akshat Kumar, and Steven Hoi. Cost: Contrastive learning of disentangled seasonal-trend representations for time series forecasting. In <u>International</u> <u>Conference on Learning Representations</u> , 2021.
685 686 687 688	Gerald Woo, Chenghao Liu, Doyen Sahoo, Akshat Kumar, and Steven Hoi. Cost: Contrastive learning of disentangled seasonal-trend representations for time series forecasting. <u>arXiv preprint</u> <u>arXiv:2202.01575</u> , 2022.
689 690 691	Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. <u>Advances in Neural Information Processing</u> <u>Systems</u> , 34:22419–22430, 2021.
693 694 695	Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In <u>The Eleventh International</u> <u>Conference on Learning Representations</u> , 2022.
696 697 698	Jiahao Xie, Wei Li, Xiaohang Zhan, Ziwei Liu, Yew-Soon Ong, and Chen Change Loy. Masked fre- quency modeling for self-supervised visual pre-training. In <u>The Eleventh International Conference</u> on Learning Representations, 2022.
700 701	Jiehui Xu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Anomaly transformer: Time series anomaly detection with association discrepancy. In <u>International Conference on Learning</u> Representations, 2021a.

702 703 704	Qinwei Xu, Ruipeng Zhang, Ya Zhang, Yanfeng Wang, and Qi Tian. A fourier-based framework for domain generalization. In <u>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</u> , pp. 14383–14392, 2021b.
705 706 707	Zhijian Xu, Ailing Zeng, and Qiang Xu. Fits: Modeling time series with 10k parameters. <u>arXiv</u> preprint arXiv:2307.03756, 2023.
708 709 710	Ling Yang and Shenda Hong. Unsupervised time-series representation learning with iterative bilinear temporal-spectral fusion. In International Conference on Machine Learning, pp. 25038–25054. PMLR, 2022.
711 712 713 714	Kun Yi, Qi Zhang, Wei Fan, Hui He, Liang Hu, Pengyang Wang, Ning An, Longbing Cao, and Zhendong Niu. Fouriergnn: Rethinking multivariate time series forecasting from a pure graph perspective. <u>Advances in Neural Information Processing Systems</u> , 36, 2024a.
715 716 717	Kun Yi, Qi Zhang, Wei Fan, Shoujin Wang, Pengyang Wang, Hui He, Ning An, Defu Lian, Longbing Cao, and Zhendong Niu. Frequency-domain mlps are more effective learners in time series forecasting. <u>Advances in Neural Information Processing Systems</u> , 36, 2024b.
718 719 720 721	Gizem Yüce, Guillermo Ortiz-Jiménez, Beril Besbinar, and Pascal Frossard. A structured dictionary perspective on implicit neural representations. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 19228–19238, 2022.
722 723 724	Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In <u>Proceedings of the AAAI conference on artificial intelligence</u> , volume 37, pp. 11121–11128, 2023.
725 726 727 728	Xiang Zhang, Ziyuan Zhao, Theodoros Tsiligkaridis, and Marinka Zitnik. Self-supervised contrastive pre-training for time series via time-frequency consistency. <u>Advances in Neural Information</u> <u>Processing Systems</u> , 35:3988–4003, 2022.
729 730 731	Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In <u>Proceedings</u> of the AAAI conference on artificial intelligence, volume 35, pp. 11106–11115, 2021.
732 733 734	Tian Zhou, Ziqing Ma, Qingsong Wen, Liang Sun, Tao Yao, Wotao Yin, Rong Jin, et al. Film: Frequency improved legendre memory model for long-term time series forecasting. <u>Advances in</u> <u>Neural Information Processing Systems</u> , 35:12677–12690, 2022a.
735 736 737 738	Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In <u>International Conference</u> on Machine Learning, pp. 27268–27286. PMLR, 2022b.
739 740	Zachary Ziegler and Alexander Rush. Latent normalizing flows for discrete sequences. In International Conference on Machine Learning, pp. 7673–7682. PMLR, 2019.
741	
742	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
100	

A APPENDIX: COMPARISON WITH FITS



Figure 7: The number of parameters in NFM (dashed line: 27K) and FITS (colour bars) for different
length of the prediction horizons (lookback of 720 for all cases except for ETTh1 for which the
lookback window is 360). The number of parameters in FITS is to yield the results in the main
performance table.

We provide an additional discussion on a recent time series model, Frequency Interpolation Time
Series analysis baseline (FITS), that is introduced on the same principle of frequency-domain
manipulation as NFM.

FITS is designed with a single complex-valued linear layer to directly learn the linear interpolation or extrapolation from input frequency spectrum as analysed in Section 3.1. It is equipped with a heuristic low-pass filter which restricts the spectrum and further reduces the model complexity. While
FITS achieves an elevated degree of lightweight-ness with competing performance in some time series tasks, there present several limitations to be improved for broader utilization.

- FITS cannot model multivariate time series. With its simplified architecture, it relies on channel-independent modelling (Nie et al., 2022) which does not model correlation between channels, thus its application is limited to univariate scenarios only.
- FITS is a linear model operating directly in frequency domain, and thus its expressiveness and learning capacity are largely limited. This is especially disadvantageous when it comes to modelling large-scale dataset and more complex patterns. Moreover, the complexity of FITS itself without a low-pass filter increases exponentially with the length of inputs and target outputs (there is a trade-off between giving up some information and resolving the complexity with the use of a low-pass filter). This characteristic does not allow FITS to stay compact in many practical scenarios with more complex and/or long time series whose major frequencies spread in wide spectrum. This is well shown from the foreacsting results on electricity and traffic datasets (**Table 7**). In **Figure 6**, we also provide the overall number of parameters in NFM and FITS used in forecasting task. Taking any lower cut-off frequency for the low-pass filters in the datasets like ETTh2, electricity, and traffic to make FITS more compact leads to considerable performance drops.
- As discussed in the main body of this work, one natural advantage gained from modelling time series directly in the Fourier domain is that the learned function can be more robust to dealing with change in discretization of input data. These features, however, are difficult to

850

851 852 853

854

855

856

В

be fully exploited in FITS due to the static nature of the model (not able to deal with variable length of time series). Although the adoption of a low-pass filter mitigates this issue by allowing to handle any length of input sequence longer than the period of cut-off, defaulting a model with a heuristic low-pass filter is not always practical. It discards the information of all higher frequency components, and thus costs performance. Indeed, in our experiment, the performance with full spectrum is consistently better than that with restricted spectrum.

816 To address the aforementioned limitations in FITS, entirely renovating the model would be necessary. 817 In one sense, NFM is a complete renovation and generalization of FITS with a compact, lightweight, 818 and adaptive deep feature learning module. It is not only designed to model both multivariate and 819 univariate time series but scalable to learn more complex patterns in data. Notably, NFM can be 820 implemented in a very compact form (less than 40K for all tasks in the experiments) regardless of 821 the input and target output length and performs consistently well in various time series tasks. This 822 feature of NFM stands out the compactness of FITS and other linear models like (Zeng et al., 2023) -823 the compactness with respect to performance of FITS would only be better than NFM in some unique 824 cases where the processed input length is sufficiently short and/or an effective heuristic low-pass filter with low cut-off frequency can be chosen. 825



APPENDIX: ANALYSIS OF NEURAL FOURIER FILTER DESIGNS

Figure 8: Overview of different NFF designs.

NFFs can mainly be found in the context of operator learning to solve PDEs and global token mixing as an efficient alternative to attention mechanism in vision tasks. Here, we summarize the existing NFFs from design perspective as shown in **Figure 7** and compare them with NFM.

FNO (Li et al., 2020) is designed with per-mode matrix multiplication. Such per-mode parameterization allows the model to be expressive with large learning capacity. However, the model can also be easily over-parameterized and thus the advantage comes with a high risk of overfittig. Moreover, the inhibitive number of parameters can be incurred when the model needs to deal with long sequences. Besides, the per-mode parameterization essentially prevents the model from being adaptive to handling variable-size inputs unless a heuristic low-pass filter is integrated.

GFN (Rao et al., 2021) also adopts per-mode multiplication but of entry-wise operation (i.e., depthwise global convolution). While this features GFN with more efficient parameterization than FNO it raises another concern of landing no channel-mixing operation in learning and may limit the model's
expressivity. Additionally, due to the per-mode parameterization, GFN is limited to a scenario where
the input length is static as well, hence it is less suitable for time-series tasks. Note that FEDformer
(Zhou et al., 2022b) and FourierGNN (Yi et al., 2024a) are good examples of GFN adapted to time
series problems with random frequency masking and just a single set of filter coefficients, respectively,
which resolve the issue with the static input length but would still suffer from the lack of expressivity.

870 **AFNO** (Guibas et al., 2021) handles variable-size inputs and achieves channel mixing simply by 871 having a single parameterization (1×1 convolution + non-linearity + soft-shrinkage function) and 872 sharing its weights across all frequency modes in spectrum. AFNO can essentially be seen as a 873 generalization of FNO and GFN with a block-diagonal (i.e., multi-head) structure, Note that FreTS 874 (Yi et al., 2024b) adopts AFNO. However, it trades off expressivity (i.e., there is no principled way for AFNO with a single shared network to learn discriminative feature across frequency modes) 875 over efficiency and flexibility to handling variable-size inputs. Besides, a soft shrinkage function is 876 adopted in AFNO as a regularization to encourage sparsity, with which we, however, consistently 877 observe a performance drop on time series tasks in all NFFs including INFF as well. 878

The weights of the filter (i.e., filter coefficients) in the above works are all shared across different instances. Unfortunately, this poses a concern that the models can struggle to learn complex patterns across different instances that reside upon different underlying spectrum.

AFF (Huang et al., 2023) addresses this by altering the use of neural network from directly parameter ized filter coefficients to a hypernetwork (Ha et al., 2017) that yields the filter coefficients dynamically according to each input instance.

Based on the above analysis, one remaining gap in the current NFFs we found is how to design a filter that is aware of each mode separately in spectrum, namely "mode-adaptive", and modulate them locally in efficient way. As discussed, FNO and GFN achieve such mode-adaptivity by instantiating the filter coefficients separately for each mode. However, this way not only makes the operation static (due to fixed parameterization) but also come at the cost of significantly increasing the number of parameters (especially, when operating matrix multiplication as in FNO) with the input length. On the other hand, AFNO and AFF do not account for the mode-adaptivity while simply sharing filter coefficients for all modes.

Our proposed INFF essentially gives a solution to this matter without sacrificing other favourable NFFs' characteristics. We achieve this by leveraging an INR (Sitzmann et al., 2019a; Tancik et al., 2020b; Sitzmann et al., 2020), a neural-based technique to reformulate a direct representation of interests into a more compact form of implicit representation. Especially, in INFF, the per-mode parameterization of a Fourier filter is turned into a process of learning to encode an abstract and implicit representation of the whole filter coefficients in spectrum. Thus, it is possible to draw unique filter coefficients for any frequency modes within the bandwidth at the cost of only a single parameterization, which largely contributes to achieving high compactness in NFM.

901 902 903

904

905

906

907

908

909

910 911 912

C APPENDIX: SYNTHETIC DATA FOR INFF VISUALIZATION

We synthesize simple single-channel band-limited signals x of K classes by composing M = S + Rfrequency components within its Nyquist frequency, $f_{nyquist} = f_x/2$. For kth class, we first create a signal from a set of fixed S frequency components $\{f_1^k, ..., f_S^k\}$ of the class that are randomly chosen from a frequency range $[f_A, f_B(< f_{nyquist})]$. Then, we further combine the signal with different sets of R frequency components $\{f_1, ..., f_R\}$ drawn randomly from $U\{1, ..., f_{nyquist}\}$ to create different variants of the class signal. This generation of kth class signal is expressed as follows:

$$x^{k}(\tau) = \sum_{i=1}^{S} A_{i}^{k} sin(2\pi f_{i}^{k}\tau + \theta_{i}^{k}) + \sum_{j=1}^{R} A_{j} sin(2\pi f_{j}\tau + \theta_{j}) + \delta(\tau)$$
(10)

913 914

915 where $A \sim U(0, 1)$ and θ are amplitude and phase components, respectively, and $\delta \sim N(0, \sigma^2)$ is 916 gaussian noise. For the synthetic data used in experiment, we set θ to a constant value, K = 10, 917 S = 20, and R = 40 and for each class signal we generate 100 sequence samples of length N = 2000(with $T_x = 1$ and $f_x = N$).

D APPENDIX: IMPLEMENTATION AND EXPERIMENTAL DETAILS

D.1 IMPLICIT NEURAL REPRESENTATIONS IN NFM

INR is parameterized by a neural network and trained to represent a target instance as a continuous function that maps grid-based representations (e.g., spatial coordinates or temporal locations), $\tau \in \mathbb{R}^{d_{in}}$, to the corresponding feature representations. A general formulation of INRs, $\phi \colon \mathbb{R}^{d_{in}} \to \mathbb{R}^{d}$, is based on a *L* layers MLP and can be expressed as follows (Yüce et al., 2022):

925 926 927

918

919 920

921 922

923

924

931

932

933

934

935

936

 $z_{INR}^{(0)} = \gamma(\tau),$ $z_{INR}^{(l)} = \alpha^{(l)} (\boldsymbol{W}^{(l)} z_{INR}^{(l-1)} + b^{(l)}), (l = 1, ..., A - 1)$ $\phi(\tau) = \boldsymbol{W}^{(L)} z_{INR}^{(A-1)} + b^{(A)}$ (11)

where $\gamma : \mathbb{R}^{d_{in}} \to \mathbb{R}^{h_0}$ is an initial feature encoding function, and $\mathbf{W}^{(l)} \in \mathbb{R}^{h_l \times h_{l-1}}$, $b^{(l)} \in \mathbb{R}^{h_l}$, and $\alpha^{(l)}$ are weights, bias, and an element-wise non-linear activation, respectively. For the INR layer in LFT and INFF modules of NFM, we opt for SIREN Sitzmann et al. (2020) - a MLP with $\alpha = sin(\cdot)$ and $z_{INR}^{(0)} = sin(w_0(\mathbf{W}^{(0)}r + b^{(0)}))$, where w_0 is a constant that controls the frequency region of the activation.

Implementing the $\phi(\cdot)$ solely based on SIREN requires putting a care on finding a good w_0 to deal with spectral bias - a tendency of MLPs towards prioritizing learning low-frequency components of the features (Rahaman et al., 2019). In order to alleviate this, as in (Kim et al., 2023) we further incorporate Fourier features (Tancik et al., 2020a; Mildenhall et al., 2021) into the $\phi(\cdot)$ by replacing the $\gamma(\tau) = sin(w_0(\mathbf{W}^{(0)}\tau + b^{(0)}))$ with $[sin(2\pi a_1\tau), cos(2\pi a_1\tau), ..., sin(2\pi a_{h_0/2}\tau), cos(2\pi a_{h_0/2}\tau)]$ where we sample $\{a_i\}_{i=1}^{h_0/2} \sim \mathcal{N}(0, 128)$.

In the experiments, we find that small ϕ works sufficiently well and do not see any noticeable improvements with larger ϕ in performance. For all tasks, we implement ϕ in LFT and INFF with the same settings - the temporal locations τ sampled equidistantly from the range [-1, 1], dimension of the input temporal location $d_{in} = 1$, the number of layers A = 3, and hidden unit dimension $h_0 \rightarrow h_1(32) \rightarrow h_2(32) \rightarrow h_3(d)$, where h_0 varies with datasets.

949 D.2 INPUT PROJECTION

951Instead of the projecting input temporal features x to initial hidden embeddings \bar{x} solely through a952matrix multiplication and passing down to the main processing modules, we combine it with non-953linear projection of a periodic activation using the same formulation as SIREN shown in Appendix954D.1.

955

948

950

956 957

$$\bar{x} = \mathbf{W}_l x + \mathbf{W}_n^{(2)} (\sin(w(\mathbf{W}_n^{(1)} x + b_n^{(1)})))$$
(12)

958 where $W_l \in \mathbb{R}^{d \times c}$, $W_n^{(1)} \in \mathbb{R}^{* \times c}$, $W_n^{(2)} \in \mathbb{R}^{d \times *}$, $b_n^{(1)} \in \mathbb{R}^*$, and w is a frequency scaling factor. 959 From the experiments, we find that relying solely on the projection through W_l was not effective for 960 both channel-independent and multi-channel cases - this could be due to the inherent information 961 sparsity of time series when working out of the features of each temporal location independently 962 (Li et al., 2023; Nie et al., 2022). A natural alleviation as a very common practice to this, especially 963 with Transformer-based models but not limited to, is to employ patchification (e.g., Nie et al. (2022)) before the projection. However, this brings in a concern. While patchification reliefs the information 964 sparsity in the inputs by forcing to reform them with a strong inductive bias of some "locality", the 965 process is essentially sensitive to the change in temporal structure (e.g., resolution and arrangement) 966 of the input time series. Thus, this factor hinders learning continuous-time characteristics in models. 967

968 Regarding learning, the projection in Eq.12 can be seen to be enriching each temporal feature 969 (point-wise input token) independently through initial channel mixing at different periods. In the 970 experiments, we see that the performance of NFM improves by a noticeable margin across all 971 datasets and tasks. We use both sine and cosine activation and set w simply to 1 (with higher w the 976 performance tends to degrade in our experiments).

972 D.3 GENERAL IMPLEMENTATION AND HYPERPARAMETER CONFIGURATIONS 973

Table 5: Hyperparameter settings for different experiments. ADs denotes all anomaly detection datasets. Note that the batch size of the forecasting datasets is set large as channel-independence is applied.

Params	ETTm1&2	ETTh1&2	Weather	Electricity	Traffic	SpeechC	ommand	ADs
						Raw	MFCC	
Epochs	40	40	40	40	40	300	300	150
Batch	1792	896	1680	1648	1648	160	240	128
Optimizer	Adam	Adam	Adam	Adam	Adam	Adam	Adam	Adam
Weight Decay	-	-	-	-	-	-	-	-
LR scheduler	-	-	-	cosine	cosine	cosine	cosine	-
Learning rate ($\times e^{-4}$)	2.0	1.5	2.5	3.5 ightarrow 1.5	3.5 ightarrow 1.5	7.5 ightarrow 3.5	5.0 ightarrow 2.5	1.0
Dropout	$0.05 \sim 0.35$	$0.05 \sim 0.35$	0.15	0.15	0.15	0.05	0.25	-
Patience	6	6	6	3	3	30	30	10
Mixer blocks	1	1	1	1	1	2	2	1
Hidden size (d)	36	36	36	36	36	32	32	8
h_0	32	32	32	32	32	32	32	16
Predictor $\mathcal{P}(\cdot)$	$d \rightarrow 1$	$d \rightarrow 1$	$d \rightarrow 1$	$d \rightarrow 1$	$d \rightarrow 1$	$d \rightarrow 10$	$d \rightarrow 10$	$d \rightarrow 1$

For all experiments, we use the same NFM backbone without a single architectural modification and equip it with a task-specific feature-to-feature linear projection. The hyper-parameter settings (empirically chosen) are specialized for each task as shown in **Table 5**.

995 D.4 EXPERIMENTAL SETUP: FORECASTING

⁹⁹⁶ Long-term times-series forecasting task ($m_f = 1$ and $m_\tau > 1$) is conducted on 7 benchmark ⁹⁹⁷ datasets (Zhou et al., 2021), including 4 *ETTs* (7 channels), *Weather* (21 channels), *Electricity* (321 ⁹⁹⁸ channels), and *Traffic* (862 channels) datasets. We use a lookback window of 720 (N = 720) for ⁹⁹⁹ all datasets except for ETTh1 for which N = 360 in NFM to make prediction over 4 horizons ¹⁰⁰⁰ {96, 192, 360, 720}. For data setup, see **Appendix D.8**.

We adopt the canonical modelling strategy (Xu et al., 2023; Zeng et al., 2023; Nie et al., 2022) of 1) channel-independence modelling and 2) a simple normalization trick to inputs and the final outputs for dealing with distribution shift caused by non-stationarity in complex time series data. For the latter, we adopt a recent popular choice of it, Reversible Instance Normalization (RevIN) (Kim et al., 2021). Meanwhile, in NFM, we also found that the normalization trick with only mean statistics works far better than with RevIN for some datasets, leading to more stable training and no early saturation to low performance regime. We use only mean statistics as the normalization trick in forecasting on ETTm2 and ETTh2 datasets.

1009

974

975

976

991

992

993 994

> Forecasting baseline results. We compare the result of NFM with that of 5 SOTA forecasting 1010 models, including FITS (Xu et al., 2023), N-Linear (Zeng et al., 2023), iTransformer (Liu et al., 1011 2023), PatchTST (Nie et al., 2022), TimesNet (Wu et al., 2022). For fair comparison, we collect 1012 the best performed results (of single prediction head but not channel-wise prediction heads) and 1013 adopted them after conducting a confirmation experiment using their official implementation (FITS¹, 1014 PatchTST², and N-Linear³). We replace them with our result if the difference was over $\pm 5\%$. Note 1015 that, as different models will have different regime for optimal lookback windows (although in 1016 general, sufficiently longer lookback would yield better results with larger the receptive field), we do not necessarily equalize the length of lookback window, and accordingly we do not compare 1017 the performance of the models with respect to different lookback window. For iTransformer⁴, we 1018 report new forecasting results made on the lookback window of 720 as the originally reported results 1019 were made with the lookback window of 96 and the performance of iTransformer is better in this 1020 setting. For TimesNet⁵, we use the lookback window of 96 (same as the original work) as we observe 1021

1025 ⁴https://github.com/thuml/iTransformer

⁵https://github.com/thuml/timesnet

¹https://anonymous.4open.science/r/FITS

²https://github.com/yuqinie98/PatchTST

³https://github.com/cure-lab/LTSF-Linear

consistently better performance with 96 over 720. Besides, we account only for the main results of the forecasting models made with "a single predictor" instead of channel-wise predictors for a fair comparison and due to its impracticality. For FITS, we use the hyperparameters that the official work used for their final results but **not 10K parameter setup** (cut-off frequency and output frequency of around 100), because with which the performance of FITS is no longer comparable with SOTAs.

1032 **Forecasting baseline results at different sampling rate.** Regarding the experiments of forecasting with different sampling rate in Section 4.2, surprisingly, none of the forecasting baseline models can 1033 1034 deal with varying-length time series inputs. This is due to their prediction head whose parameters are subject to the initial input length, thus not allowing the estimation directly from downsampled input 1035 sequences. While this can be overcome in FITS - by having downsampled factors up to the training 1036 cut-off frequency or by posing a lower cut-off frequency during testing time, this task is not directly 1037 applicable to PatchTST and N-Linear that working in time domain. A common practice to this is 1038 to resample the downsampled input sequences to match the resolution back (i.e., by upsampling). 1039 We simply do this by interpolating the downsampled input sequences through zero-padding them in 1040 frequency domain (equivalent to applying a sinc kernel in time domain). Note that the resolution of 1041 the target prediction remains unchanged (e.g., the case of forecasting on ETTm1 with SR=1/4 can be 1042 seen as predicting ETTm1 from ETTh1).

1043

1048

1031

1044 D.5 EXPERIMENTAL SETUP: CLASSIFICATION

For classification, we follow the same experimental setup used in CKconv (Romero et al., 2021) and S4 (Gu et al., 2021). For data setup, see **Appendix D.8**.

Classification baseline results. Learning a function-to-function mapping given discrete signals 1049 is a desirable property for models to generalize models across different challenging scenarios of 1050 time series analysis. Several works have addressed it, modelling continuously evolving dynamics in 1051 hidden states (Chen et al., 2018; Rubanova et al., 2019; Morrill et al., 2021; De Brouwer et al., 2019) 1052 and a stochastic mapping (Ziegler & Rush, 2019; Deng et al., 2020). On the other hand, this property 1053 can also be achieved by simply modelling data directly in Fourier domain with the fact that the DFT 1054 provides a function of frequency as samples of a new "functional" representation of the time-domain 1055 discrete signals. Intuitively, working explicitly with the frequency samples of discrete signals can be 1056 equivalently seen as working with their continuous-time elements in function space. To show this 1057 (especially, in the scenario of testing at different sampling rate), we compare the results of NFM with 8 baseline models (7 continuous-time models and 1 Transformer), including ODE-RNN (Rubanova 1058 et al., 2019), GRU- Δt (Kidger et al., 2020), GRU-ODE (De Brouwer et al., 2019), NCDE (Kidger 1059 et al., 2020), NRDE (Morrill et al., 2021)), S4, CKConv, and vanilla Transformer (Vaswani et al., 2017). The results were collected and verified using the implementations in $CKconv^6$ and $S4^7$. 1061

1062

D.6 EXPERIMENTAL SETUP: ANOMALY DETECTION

1064 We employ 4 popular anomaly detection benchmark datasets, including SMD (Server Machine Dataset, (Su et al., 2019)), PSM (Pooled Server Metrics, (Abdulaal et al., 2021)), MSL (Mars Science Laboratory rover, (Su et al., 2019)) SMAP (Soil Moisture Active Passive satellite, (Su et al., 2019)), 1067 and follow the well-established experimental protocol in (Shen et al., 2020) and the evaluation 1068 methodology in (Xu et al., 2021a). We use the same setup used in the work (Xu et al., 2021a) - a 1069 window of length 100 for all datasets and baselines, and anomaly ratio (%) of 1.0 (MSL), 1.0 (PSM), 1070 1.0 (SMAP), and 0.5 (SMD). Besides, since we workin in physical space of the time series, we apply channel-independent modelling and input normalization trick just like in forecasting task. For data 1071 setup, see Appendix D.8. 1072

- 1073 1074
- 1074
- 1076
- 1077

¹⁰⁷⁸ 1079

⁶https://github.com/dwromero/ckconv

⁷https://github.com/state-spaces/s4/tree/main

Anomaly detection baseline results. We compare the results of NFM with 6 SOTA baselines, including vanilla Transformer (Vaswani et al., 2017), PatchTST (Nie et al., 2022), TimesNet (Wu et al., 2022), ADformer (Xu et al., 2021a), N-linear (Zeng et al., 2023), and FITS (Xu et al., 2023).
For the baseline results, we follow the same evaluation methodology used in (Wu et al., 2022; Xu et al., 2021a) and produced the results, using the implementations in TimesNet, ADformer⁸, and FITS. Meanwhile, it is important to point out that there were some flaws (in training/validation/training data assignment, computing anomaly threshold, and estimation) in their (all three) official experimental codes that affect the final results. The fixed code samples are available in our repository, and all results in anomaly detection were made again with the fixed codes.



D.7 Optimizations



Figure 9: Training and testing framework for anomaly detection task in NFM.

1119 **Forecasting.** Given an input lookback sequence of length N, $x[n \in I_N]$ ($f_x = N$ and $T_x = 1$), the 1120 aim in forecasting task is to predict a future horizon of a desired length. While most of existing models 1121 (Zeng et al., 2023; Liu et al., 2023; Nie et al., 2022; Wu et al., 2022) directly outputs only the target 1122 horizon given the input lookback sequence, NFM yields a whole extrapolated sequence (frequency-1123 domain interpolation with $m_t = L/N$ of length L = N + horizon, $\hat{y}[n \in I_L]$ $(f_y = f_x$ and $T_y = f_y$ 1124 $T_x + \frac{horizon}{f_y}$). During training, we supervise NFM over the whole extrapolated sequence in both time 1125 domain and frequency domain against its ground truth $y[n \in I_L] = \{x[N \in I_N], x[N], \dots, x[L-1]\}$. 1126 Note that we observe that NFM in forecasting task performs consistently better when optimized in 1127 both time domain and frequency domain. We opt for the standard time-domain loss function \mathcal{L}_{TD} , 1128 **MSE**, and for the frequency domain loss function \mathcal{L}_{FD} we modify the **focal frequency loss** (Xie 1129 et al., 2022; Jiang et al., 2021) used for the recovery of image spectrum and adapt it for time series. 1130

1131

1113

1132

⁸https://github.com/thuml/Anomaly-Transformer

1137

1144

where Y_{Real} and Y_{Imag} are the real and imaginary part of the frequency representation of the 1145 sequence y, and λ (we set $\lambda = 0.5$) controls the contribution of the frequency-domain loss. As 1146 seen, the distance metric in \mathcal{L}_{FD} considers both amplitude (the contribution of each frequency 1147 component to the time-domain signal) and phase (temporal delay introduced by each frequency 1148 component) information by operating on both real and imaginary parts of the complex frequency 1149 representation. Importantly, one can see that the time-domain objective is local as applied point 1150 wise and the frequency-domain objective is global with the fact that each frequency component is 1151 a summary about the entirety of the sequence at different periods. Hence, this encourages more 1152 faithful construction for the extrapolated sequence as a whole. We argue that this is a reason for NFM 1153 working the best when both domains' objective is incorporated. 1154

 $\mathcal{L}_{Forecasting} = \underbrace{\lambda \frac{1}{L} \sum_{n=0}^{L-1} ||\hat{y}[n] - y[n]||_2}_{\mathcal{L}_{TD}} + (1 - \lambda) \frac{1}{K_L} \sum_{k=0}^{K_L - 1} ((\hat{Y}_{Real}[k] - Y_{Real}[k])^2 + (\hat{Y}_{Imag}[k] - Y_{Imag}[k])^2)^{1/2}}_{C}$

(13)

1155 **Classification.** In classification setup, given pairs of the sequence x of length N and the class label $y \in \{1, \ldots, K\}$, the NFM backbone yields the latent features $z[n \in I_L] = y(n/f_y)$, that resides on 1156 the same timespan $(T_y = T_x)$ as x, with respect to the class information. We train the NFM backbone 1157 with a linear classifier (global average pooling + a fully-connected layer) and optimize them using 1158 the standard classification loss function, Cross Entropy. 1159

1160 **Anomaly detection.** In anomaly detection task, the aim of NFM is to learn the dominant contexts 1161 (i.e., normal contexts) of the input sequences in unsupervised manner (i.e., no anomaly label is 1162 available). Then, during testing time, the learned sequence-wise normal contexts are used as a 1163 standard to establish a decision boundary for anomaly, and the elements of the sequences are 1164 evaluated within the corresponding contexts. To achieve this, we frame the objective of NFM as 1165 a context learning (see Figure 8) and train NFM to learn as faithful contexts as possible. More 1166 specifically, given the sequence $x[n \in I_N]$, we downsample (at equidistant sampling rate against the original discretization) it by a downsampling factor dr. Denoting the downsampled sequence 1167 $x_d[n \in I_{N_d}]$ where $N_d = N/dr$ (i.e., $f_{x_d} = f_x/dr$ and $m_f = dr$), the NFM takes in x_d as input and 1168 is trained to reconstruct the full original sequence $\hat{x}[n \in I_N]$ on both MSE loss (\mathcal{L}_{TD}) and the focal 1169 **frequency loss** (\mathcal{L}_{FD}) as follows: 1170

1171 1172

$$\mathcal{L}_{AD} = \lambda \frac{1}{N} \sum_{n=1}^{N-1}$$

$$\mathcal{L}_{AD} = \underbrace{\lambda \frac{1}{N} \sum_{n=0}^{N-1} ||\hat{x}[n] - x[n]||_{2}}_{\mathcal{L}_{TD}} + \underbrace{(1-\lambda) \frac{1}{K_{N}} \sum_{k=0}^{K_{N}-1} ((\hat{X}_{Real}[k] - X_{Real}[k])^{2} + (\hat{X}_{Imag}[k] - X_{Imag}[k])^{2})^{1/2}}_{\mathcal{L}_{FD}}$$
(14)

1179 1180 1181

We set $\lambda = 0.5$, and the \mathcal{L}_{FD} is only applied during training time and not used in any steps of 1182 anomaly detection. Besides, it is noteworthy that during the testing time, inputs can be any of original 1183 or downsampled version of candidate sequence with resolution-invariance property of NFM, and 1184 the evaluation of the sequence points for normality is made on the full length between the original 1185 sequence and the restored sequence of the input candidate. Additionally, we note that NFM requires 1186 no single architectural modification to adopt the above formulation or changing the above formulation 1187 to full reconstruction.

1188 D.8 SUMMARY OF DATASETS

1189

Table 6: Summary of data settings. SC: SpeechCommand, AD: Anomaly detection, and CLS: Classification.

1193 Tasks Dataset channels Train / Val / Test N/LNum.Class Domain 1194 1195 (720, 720, 720, 720) / 7 ETTm1 60%/20%/20% Temperature (816, 912, 1056, 1480) 1196 (720, 720, 720, 720) / 1197 ETTm2 7 60%/20%/20% Temperature (816, 912, 1056, 1480) 1198 (360, 360, 360, 360) / ETTh1 7 60%/20%/20% Temperature 1199 Forecasting (452, 552, 696, 1080) 1200 (720, 720, 720, 720) / ETTh2 7 60%/20%/20% Temperature (816, 912, 1056, 1480) 1201 (720, 720, 720, 720) 1202 Weather 21 70%/20%/10% Weather (816, 912, 1056, 1480) 1203 (720, 720, 720, 720) / Electricity 321 70%/20%/10% Electricity (816, 912, 1056, 1480) 1205 (720, 720, 720, 720) / Traffic 862 70%/20%/10% Transportation (816, 912, 1056, 1480) 1206 1207 50/100 2 SMD 38 80%/20%/ -Server Machine 2 MSL 55 80%/20%/ -50/100 Spacecraft 1208 AD SMAP 25 80%/20%/ -50/100 2 Spacecraft 1209 2 25 PSM 80%/20%/ -50/100 Server Machine 1210 CLS 10 SC-raw 1 70%/15%/15% 16000/-Speech 1211 SC-MFCC 20 70%/15%/15% 161/-10 Speech 1212

E APPENDIX: ADDITIONAL EXPERIMENTS AND ANALYSIS

Here, we provide extra experimental results and analysis omitted in the main work due to the limited work space.

1218 1219 1220

1213 1214 1215

1216 1217

E.1 FULL FORECASTING RESULTS

1222 **Discussion on the number of parameters.** With the prevalence of chunk-to-chunk prediction in the 1223 forecasting community, one practice in the deep forecasting baselines is to **adopt a prediction head** that acts on temporal dimension. Note that the linear models (FITS and N-Linear) naturally fall in 1224 this as they by themselves are the prediction head operating on the input time series. Due to this, they 1225 scale poorly with the length of horizons as well as the length of lookback window. For example, more 1226 than 95% of the learnable weights in PatchTST for L = 720 surprisingly belongs to the single "wide" 1227 prediction head of 8.3M parameters against 0.4M parameters in its Transformer-based backbone. 1228 In contrast, the prediction head used in NFM is feature-to-feature projection and the number of 1229 parameters in NFM is *completely independent* from the length of input sequence and prediction 1230 horizon. Importantly, we highlight that this aspect completely decouples the contribution of the linear 1231 head in modelling sequence and further validates the effectiveness of NFM to modelling temporal 1232 dependency unlike the other deep forecasting baselines that are equipped with a wide linear predictor. 1233

1234

- 1235
- 1236

1237

1238

1239

1240

Table 7: Full long-term forecasting results, where the best results are in **bold** and the second best are <u>underlined</u>. The number of parameters of the baselines are computed based on their original hyper-parameter setting.

		NI (27	F M 7K)	FI (~0	TS .2M)	N-L (∼0	inear .5M)	iTrans (∼5	former .3M)	Patcl (~8	nTST .7M)	Time (~0	esNet .3B)
		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
-	96	0.286	0.338	0.309	0.352	0.306	0.348	0.319	0.367	<u>0.293</u>	<u>0.346</u>	0.338	0.375
^m L	192	0.326	0.364	0.338	0.369	0.349	0.375	0.347	0.389	0.333	0.370	0.374	0.387
ET	336 720	0.362 0.406	0.384 0.414	$\frac{0.366}{0.415}$	0.385 0.412	0.375	0.388 0.422	0.382 0.437	0.409 0.440	0.369 0.416	0.392	0.410 0.478	0.411 0.450
	06	0 160	0.250	0.163	0.254	0.167	0.255	0.180	0.274	0.166	0.256	0.187	0.267
_m2	192	0.221	0.291	0.103	0.291	0.221	0.293	0.243	0.316	0.223	0.296	0.249	0.207
ET	336	0.271	0.326	0.268	0.326	0.274	0.327	0.299	0.352	0.274	0.329	0.321	0.351
	720	0.349	0.378	0.349	<u>0.379</u>	0.368	0.384	0.382	0.404	<u>0.362</u>	0.385	0.408	0.403
-	96	0.363	0.389	0.372	0.395	0.374	<u>0.394</u>	0.392	0.423	0.370	0.400	0.384	0.402
ΨĽ	192	0.404	0.413	0.404	0.413	0.408	0.415	0.428	0.447	0.413	0.429	0.436	0.429
표	336 720	0.420	0.422	0.427	0.427	0.429	$\frac{0.427}{0.453}$	0.494	0.488	$\frac{0.422}{0.447}$	0.440	0.491	0.469
	720	0.442	0.437	0.424	0.440	0.440	0.435	0.077	0.000	0.447	0.400	0.521	0.500
h2	96 102	0.281	0.341	0.271	0.337	0.277	0.338	0.304	0.364	$\frac{0.274}{0.341}$	0.337	0.340	0.374
TT	336	0.377	0.387	0.354	0.395	0.344	$\frac{0.381}{0.400}$	0.432	0.451	$\frac{0.341}{0.329}$	0.382	0.452	0.414
щ	720	0.414	0.455	0.378	0.423	0.394	0.436	0.441	0.469	0.379	0.422	0.462	0.468
5	96	0.154	0.203	0.169	0.224	0.182	0.232	0.173	0.227	0.149	0.198	0.172	0.220
athe	192	0.198	0.246	0.213	0.261	0.225	0.269	0.219	0.262	0.194	0.241	0.219	0.261
We	336	0.245	0.281	0.259	0.296	0.271	0.301	0.283	0.310	0.245	$\frac{0.282}{0.224}$	0.280	0.306
	720	0.312	0.331	0.321	0.340	0.338	0.348	0.344	0.355	0.314	0.334	0.365	0.359
city	96	0.131	0.222	0.135	0.231	0.141	0.237	0.132	$\frac{0.227}{0.252}$	0.129	0.222	0.168	0.272
ctri	192 336	0.147	0.240	$\frac{0.149}{0.165}$	$\frac{0.244}{0.261}$	0.154	0.248	0.155	0.252	0.147	0.240	0.184	0.289
Ele	720	0.194	0.284	$\frac{0.103}{0.203}$	0.293	0.210	0.297	0.195	0.289	0.197	0.290	0.220	0.320
	96	0.367	0 249	0.386	0.269	0.410	0.279	0 344	0.254	0.360	0 249	0 593	0.321
ffic	192	0.377	0.252	0.398	0.274	0.423	0.279	0.366	$\frac{0.234}{0.265}$	0.379	0.256	0.617	0.321
Trat	336	0.392	0.259	0.410	0.278	0.435	0.290	0.381	0.273	0.392	0.264	0.629	0.336
	720	<u>0.427</u>	0.279	0.448	0.296	0.464	0.307	0.413	0.287	0.432	0.286	0.640	0.350

E.2 FULL ANOMALY RESULTS

Table 8: Time series anomaly detection results on 4 datasets. The higher the three metrics, including the precision (P), recall (R), and F1-score (F1) in percentage, are, better the performance.

Model (params)		SMD			MSL			SMAP			PSM		
	Р	R	F1										
Transformer (0.2M)	68.40	85.37	75.95	88.11	76.55	81.93	89.37	57.12	69.70	99.96	79.80	88.75	
PatchTST (0.2M)	80.33	83.96	82.11	84.33	77.01	80.51	92.22	55.27	69.11	98.78	93.38	96.00	
TimesNet (~28M)	82.67	83.88	83.27	87.45	76.65	81.70	89.07	62.17	73.23	98.42	96.20	97.30	
ADformer*(4.8M)	68.79	85.86	76.38	88.68	75.87	81.78	91.85	58.11	71.18	99.98	71.15	83.14	
N-Linear (10K)	78.94	84.89	81.81	86.55	76.43	81.18	89.85	54.05	67.50	98.47	93.22	95.77	
FITS (1.3K)	79.90	83.52	81.67	86.85	75.49	80.77	88.47	50.22	64.07	98.74	94.55	96.60	
NFM (6.6K)	86.82	81.96	84.32	88.72	77.11	82.46	90.12	58.87	70.88	98.92	96.91	97.51	

* The joint criterion in ADformer is replaced with the simple reconstruction error to compute anomaly score for fair comparison.

E.3 FULL RESULTS OF FORECASTING AT DIFFERENT INPUT RESOLUTION

In practice, it is not rare to encounter a scenario where a system undergoes or necessitates a change in sampling rate for signals being monitored by a model. Such change does not affect the underlying temporal dynamic of the signals (i.e., the same solution) but brings in a positional alternation in the sequences of our observations, greatly affecting the performance of the model. To this end, we conduct forecasting on the input time series sampled at "unseen" discretization rate (this experiment is made for the first time in our work). Overall, the full results in Table 9 demonstrates the resolution-invariance property of NFM that can be highly valuable in practical applications.

1299										
1300	Dataset	Horizon	N	FM	FI	TS	N-L	inear	Patch	ıTST
1301			1/4	1/6	1/4	1/6	1/4	1/6	1/4	1/6
1302		06	0.299	0.319	0.348	0.368	0.436	0.482	0.394	0.434
1303		90	(4 . 5 ↓)	$(11.5\downarrow)$	$(12.6\downarrow)$	$(19.1\downarrow)$	$(42.5\downarrow)$	$(57.5\downarrow)$	$(34.5\downarrow)$	$(48.5\downarrow)$
120/		192	0.339	0.349	0.363	0.377	0.444	0.481	0.386	0.412
1304	ETTm1		$(4.0 \downarrow)$	$(7.1 \downarrow)$	$(7.4 \downarrow)$	$(11.5 \downarrow)$	$(27.2 \downarrow)$	$(37.8 \downarrow)$	$(15.9 \downarrow)$	$(23.7\downarrow)$
1305		336	(0.305)	(2, 1, 1)	(4.9)	$(7.1 \downarrow)$	$(21.9 \bot)$	$(32.5 \bot)$	$(6.8 \downarrow)$	$(11.7 \bot)$
1306			0.407	0.414	0.426	0.437	0.469	0.471	0.433	0.443
1307		720	(0 . 2 ↓)	(2 . 0 ↓)	$(2.9\downarrow)$	$(5.3\downarrow)$	(8.3 ↓)	(8.8 ↓)	$(4.1\downarrow)$	$(6.5\downarrow)$
1308		96	0.179	0.189	0.198	0.216	0.251	0.281	0.232	0.265
1309		50	(11.9 ↓)	(18.2 ↓)	$(21.5\downarrow)$	$(31.7\downarrow)$	$(50.3\downarrow)$	$(68.3\downarrow)$	(39.8 ↓)	$(59.6\downarrow)$
1000		192		0.239	0.240	0.258	0.267	0.284	0.264	0.287
1310	ETTm2		$(4.1 \downarrow)$ 0 281	$(0.1 \downarrow)$	$(10.0 \downarrow)$ 0.288	$(10.2 \downarrow)$ 0.300	$(20.8 \downarrow)$ 0.307	$(20.3 \downarrow)$ 0.331	$(18.4 \downarrow)$ 0.297	$(20.7 \downarrow)$ 0.311
1311		336	(3.7 .1)	(6.0.1.)	(7.5.1.)	(11.9.1.)	(12.0 l.)	$(20.8 \downarrow)$	(8.4.1.)	(13.5.1)
1312		720	0.356	0.358	0.356	0.364	0.409	0.422	0.382	0.396
1313		120	(2 . 0 ↓)	(2 . 6 ↓)	(2 . 0 ↓)	$(4.3\downarrow)$	$(11.1\downarrow)$	$(14.7\downarrow)$	$(5.5\downarrow)$	(9.4 ↓)
1010		0.0	0.164	0.173	0.182	0.195	0.268	0.280	0.212	0.236
1314		96	(6.5 ↓)	(12.3 ↓)	$(7.7\downarrow)$	(14.8 ↓)	$(47.3\downarrow)$	$(53.8\downarrow)$	$(42.3\downarrow)$	$(58.4\downarrow)$
1315		192	0.209	0.219	0.222	0.233	0.271	0.288	0.261	0.282
1316	Weather	102	$(5.6\downarrow)$	$(10.6\downarrow)$	(4 . 2 ↓)	(9.4 ↓)	$(20.4\downarrow)$	$(28.0\downarrow)$	$(34.5\downarrow)$	$(45.4\downarrow)$
1017		336	0.253	0.258	0.265	0.272	0.296	(12.71)	0.278	0.292
1317			(3.3 ↓) 0 315	(0.3↓) 0318	$(4.3 \downarrow)$	(∂.0 ↓) 0.329	$(9.2 \downarrow)$ 0.351	$(13.7 \downarrow)$ 0.364	$(13.5 \downarrow)$ 0.320	$(19.2 \downarrow)$ 0.335
1318		720	$(1.0 \downarrow)$	(1.9.1)	(1.2.1)	(2.5.1)	(3.8.1)	(7.7.1)	(4.8.1.)	(6.7.1)
			(=:• • •)	(=.• • •)	(=:= _)	(=:3 4)	(0.0 4)	(···· ψ)	(=··) ()	(

Table 9: Full forecasting results (MSE) and performance drops (%) at different testing-time sampling rate, where $SR = f_x^{test} / f_x^{train}$. The best performance is in **blue** and the least performance drop in **red**.

Analysis. Especially, we observe that the models operating fully on frequency domain (NFM and FITS) are much more robust to the change in sampling rate as learning a function-to-function mapping, than those operating on time domain (PatchTST and N-Linear). Interestingly, the performance degradation with unseen sampling rate tends to be more significant in forecasting over short horizons (with the same lookback length) and relatively minor over long horizons in all models. This tendency could be a strong indication that the forecasting models, including NFM, become more reliant on the global context (or similarly, low frequency regime) of the lookback window and less focusing on the local details (or similarly, high frequency regime) in the window as the predicting horizon gets longer. In this sense, the results imply that NFM is not the one that leverages local features in optimal way but is less prone to the local variations than the others. In the future, integrating a mechanism that encourages learning more of local features (high-frequency information) into models could potentially improve the forecasting ability in the long horizon cases as well as the short horizon cases.

1350 E.4 FULL TABULAR RESULTS ON COMPARISON OF NFM WITH DIFFERENT ABLATION CASES

Table 10: Tabular results of the ablation study on ETTm1. We use the same set up used in Table 5 for all cases and the number of heads = 2 for AFNO and AFF.

Horizon	Metric	INFF			LFT			
110/12011		×	Naive	LFT	FNO	AFNO	GFN	AFF
96	MSE	0.296	0.292	0.286	0.289	0.291	0.303	0.292
	MAE	0.348	0.343	0.338	0.347	0.347	0.349	0.346
	FLOP (G)	0.076	0.076	0.082	0.069	0.090	0.066	0.090
	PMU (GB)	0.029	0.030	0.030	0.030	0.032	0.025	0.033
	Params (M)	0.021	0.036	0.027	0.435	0.020	0.029	0.020
192	MSE	0.335	0.330	0.326	0.330	0.338	0.338	0.340
	MAE	0.370	0.364	0.364	0.367	0.374	0.372	0.372
	FLOP (G)	0.085	0.085	0.089	0.080	0.099	0.072	0.099
	PMU (GB)	0.031	0.033	0.033	0.033	0.035	0.027	0.036
	Params (M)	0.021	0.039	0.027	0.484	0.020	0.030	0.020
336	MSE	0.372	0.366	0.362	0.365	0.371	0.383	0.369
	MAE	0.396	0.387	0.384	0.390	0.391	0.394	0.393
	FLOP (G)	0.096	0.096	0.101	0.089	0.112	0.081	0.113
	PMU (GB)	0.036	0.036	0.037	0.038	0.039	0.030	0.040
	Params (M)	0.021	0.039	0.027	0.557	0.020	0.033	0.020
720	MSE	0.421	0.421	0.406	0.431	0.413	0.427	0.407
	MAE	0.422	0.419	0.414	0.424	0.414	0.420	0.421
	FLOP (G)	0.119	0.119	0.128	0.115	0.146	0.105	0.147
	PMU (GB)	0.047	0.048	0.048	0.051	0.038	0.051	0.053
	Params (M)	0.021	0.043	0.027	0.557	0.020	0.039	0.020

Table 11: Tabular results of the ablation study on SpeechCommand. We use the same set up used in **Table 5** for all cases and the number of heads = 2 for AFNO and AFF.

SR	Metric	INFF			LFT			
511		×	Naive	LFT	FNO	AFNO	GFN	AFF
1.0 0.5	ACC (%) ACC (%) FLOP (G) PMU (GB)	79.4 74.1 0.480 0.183	82.6 78.6 0.480 0.183	90.9 90.4 0.487 0.188	84.8 75.0 0.383 0.395	84.7 82.4 0.478 0.175	86.2 78.0 0.330 0.135	88.4 85.8 0.478 0.178
	Params (M)	0.031	0.185	0.188	16.4	0.031	0.133	0.031



E.5 STATISTICAL SIGNIFICANCE ON THE MAIN RESULTS WITH DIFFERENT RANDOM SEEDS



Figure 10: Statistical significance on the main experimental results computed by repeating the main experiments with different random seeds. The number in the title of forecasting results (a) \sim (g) indicates prediction horizon, and the legend in each box plot $mean \pm std$.