SEAT: SPARSIFIED ENHANCEMENTS FOR ATTENTION MECHANISMS IN TIME SERIES TRANSFORMERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Transformer models excel in time series tasks due to their attention mechanisms. However, they often suffer from "block-like" attention patterns caused by high feature correlation, leading to feature confusion and reduced performance. In this study, we mathematically prove and quantify this limitation, demonstrating how it affects the sparsity of the attention matrix and hinders effective feature representation. To overcome this issue, we propose a novel, model-agnostic, and plugand-play method called SEAT (Sparsification-Enhanced Attention Transformer) that leverages frequency domain sparsification. By transforming time series data into the frequency domain, our method induces inherent sparsity, reduces feature similarity, and mitigates block-like attention, allowing the attention mechanism to focus more precisely on relevant features. Experiments on benchmark datasets demonstrate that our approach significantly enhances the accuracy and robustness of Transformer models while maintaining computational efficiency. This provides a mathematically grounded solution to inherent flaws in attention mechanisms, offering a versatile and effective approach for advancing time series analysis.

024 025 026

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

027 028

029 In long-term sequence forecasting (LTSF) tasks, attention mechanisms play a crucial role in capturing dependencies within time series data. However, existing approaches, particularly those employing "Point-wise Attention" exhibit significant limitations. This strategy maps various features 031 at each time step to embeddings, treating each as a token. Consequently, attention mechanisms allocate weights to each corresponding token. Experimental studies Zhang & Yan (2023) have revealed 033 that, in conventional Transformers for LTSF tasks, cross-dimensional dependencies are not explicitly 034 captured during the embedding process. This limitation adversely affects forecasting capabilities, as attention values tend to exhibit segmentation; adjacent data points often receive similar attention weights. The "Point-wise Attention" approach typically results in a block-like distribution of feature 037 map values within time series representations. This block-shaped pattern causes approximate values 038 of different features in the attention mechanism to become conflated, leading to model overfitting to noise. Consequently, the decision boundaries of the model become ambiguous, negatively impacting overall performance. 040

O41 Subsequent research has attempted to mitigate these issues by employing alternative attention strategies, such as using pairs of patches. For instance, PatchTST Nie et al. (2023) utilizes patching and channel independence within Transformer architectures to significantly enhance performance, demonstrating that Transformers retain considerable potential for improvement in time series forecasting when appropriately adapted. Pathformer Chen et al. (2024) adopts a patch-based method to divide the time series into various temporal resolutions. Through this multi-scale division, dual attention is applied to these patches, enabling the capture of global correlations and local temporal dependencies.

Despite these advancements, existing methods primarily focus on enhancing feature extraction with out fundamentally improving the quality of time series feature representation. Block-like attention
 may result in the attention weight matrix assigning significantly higher weights to certain time slices
 compared to others, leading to the aggregation of similar features within the attention output and
 increasing feature confusion. Such effects can adversely impact the model's learning capability and
 performance limits, as will be substantiated in the Methods section.

054 To address these challenges, we propose **SEAT**, a model-agnostic enhancement framework appli-055 cable to any Transformer architecture's input signals. SEAT employs frequency modelling across 056 the entire time series and introduces a finite energy representation in the frequency domain, cap-057 italizing on the sparse features present in time series data. This design enables the Transformer's 058 attention mechanism to focus on independent feature representations as "Channel-wise Attention." By explicitly modelling the Fourier transform to reconstruct features as incremental signals, SEAT enhances the model's ability to distinguish approximate features and reduce feature redundancy and 060 similarity. Consequently, this approach improves predictive performance in long-term forecasting 061 tasks across various Transformer-based models. In summary, our contributions are threefold: 062

- 1. We provide a theoretical analysis of the limitations in attention mechanisms within time series forecasting Transformers, investigating the causes of feature confusion and susceptibility to overfitting.
- 2. Based on rigorous mathematical proofs, we design SEAT, a sparse sensing enhancement framework tailored for time series attention, ensuring the independence and sparsity of input features within Transformer architectures.
 - 3. Our framework is decoupled from the underlying model architecture, offering plug-andplay functionality and compatibility with any existing Transformer-based architecture.

2 RELATED WORK

063

064

065

066

067

068

069

071 072

073

074 075 2.1 TIME-DOMAIN-BASED TRANSFORMERS

Time series forecasting has been significantly advanced by the integration of Transformer architectures Ashish (2017). Leveraging the self-attention mechanism, Transformer-based models have demonstrated exceptional performance in capturing long-range dependencies, a critical aspect for effective LTSF. Notable models in this domain include:

Autoformer Wu et al. (2021) introduced an auto-correlation mechanism specifically designed to leverage the inherent periodicity of time series data. This mechanism discovers dependencies and aggregates representations at the sub-series level, significantly improving the model's ability to utilize long-range information. By focusing on periodic patterns, Autoformer overcomes the information bottleneck that constrains the original Transformer architecture, thereby enhancing its forecasting performance for non-stationary time series.

Pyraformer Liu et al. (2022a) further addresses these limitations through a hierarchical time se-087 ries decomposition approach. Its pyramid structure exhibits multiresolution properties, allowing the 880 model to decompose time series into distinct temporal scales. At coarser resolutions, long-term 089 dependencies are captured, while finer resolutions discern intricate short-term variations. This hi-090 erarchical decomposition facilitates a nuanced understanding of temporal dynamics, enhancing the 091 model's ability to accurately forecast future trends and patterns. However, distinguishing long-term 092 and short-term features is often limited by window parameter selection and data resolution, necessi-093 tating model structure modifications for different datasets, which reduces robustness and increases susceptibility to overfitting. 094

 Nonstationary Transformer Liu et al. (2022b) introduces a dual approach to improve the modelling of non-stationary time series. It enhances data stationarity through techniques that mitigate non-stationarity within the time series and reformulates the Transformer's internal mechanisms to reintegrate non-stationary information. This twofold innovation significantly advances the Transformer's capability in handling non-stationary data, improving both predictability and forecasting performance.

iTransformer Liu et al. (2024) marks a significant advancement by surpassing traditional Transformer models in time series forecasting. It treats individual series as variate tokens, utilizing attention mechanisms to capture multivariate correlations and employing layer normalization and feed-forward networks to learn robust series representations. However, relying solely on time-domain tokens poses challenges in comprehensively portraying integral properties of time series data, such as overarching trends and periodic fluctuations. Additionally, the quadratic complexity and large parameter count inherent in Transformer models make them prone to overfitting, especially when applied to non-stationary time series.

108 2.2 FREQUENCY-DOMAIN-BASED TRANSFORMERS

110 Time-frequency transforms offer a promising avenue for transforming long-time series into sparse 111 representations. The development of frequency domain methods can be traced back to the introduction of Fourier transforms and wavelet transforms, which provided novel perspectives for analyzing 112 periodicity and global dependencies within time series data. Frequency-domain methods, tradi-113 tionally used for signal processing tasks, excel at decomposing time series into their constituent 114 frequencies, revealing periodic and seasonal components that are often crucial for accurate fore-115 casting. Early frequency domain approaches, such as FNet Lee-Thorp et al. (2021) and AFNO 116 Guibas et al. (2021), were primarily designed to enhance computational efficiency by leveraging 117 Fourier transforms to replace self-attention or token-mixing mechanisms, thereby achieving faster 118 computation. 119

There has been a growing trend towards integrating frequency-domain methods with sophisticated 120 techniques such as attention mechanisms to achieve superior predictive performance in time series 121 analysis. Traditional Transformer architectures applied to time series data rely on point-wise atten-122 tion mechanisms, where individual temporal points undergo attention computations and predictions 123 in isolation. As a result, these models often struggle to maintain and accurately model holistic, 124 global features that inherently encode crucial information for accurate forecasting. In contrast, FED-125 former Zhou et al. (2022) represents a paradigm shift by leveraging Transformer structures within 126 the spectral domain for feature extraction. This innovative approach enables FEDformer to more 127 effectively capture global characteristics vital for comprehending the intricate dynamics within time 128 series data. By harnessing the complementary strengths of both temporal and spectral representations, FEDformer fosters a deeper, more nuanced understanding of underlying patterns and trends, 129 thereby enhancing its predictive capabilities. 130

131 Similar to the objectives of numerous decomposition methods Wu et al. (2021); Zhou et al. (2021), 132 the adoption of spectral domain approaches aims to facilitate the decomposition of time series data 133 into distributions that are more conducive to learning. By transforming the time series into the 134 frequency domain, these methods enable a more effective decomposition of temporal dynamics, 135 transforming the data into representations that are easier for models to comprehend and utilize for prediction tasks. However, they share a common limitation: when dealing with complex time series, 136 most frequency models tend to prioritize learning low-frequency features while overlooking high-137 frequency features, exhibiting a frequency bias. For example, FITS Xu et al. (2024) incorporates 138 a low-pass filter in the frequency domain to capture essential time series information, but this in-139 evitably results in the loss of high-frequency components. This bias can hinder the model's ability 140 to fully capture the intricate dynamics present across all temporal scales, potentially impacting the 141 accuracy and robustness of forecasts. 142

Fredformer Piao et al. (2024) aims to alleviate frequency bias by equally learning features across
 different frequency bands. This approach helps prevent the model from overlooking lower amplitude
 features crucial for accurate predictions. In these previous studies, attention layers are designed to
 function directly in the frequency domain to enhance spatial or frequency representations.

147 Transformer-based models for long-term time series forecasting (LTSF) have primarily advanced through two approaches. Time-domain models focus on enhancing attention mechanisms and em-148 ploying patch-based training strategies. In contrast, frequency-domain models develop specialized 149 filters and address biases between high and low-frequency components to improve performance. 150 However, these methods mainly concentrate on updating and iterating Transformer architectures 151 without deeply considering the inherent input properties of time series data and their impact on at-152 tention mechanisms, making it difficult to effectively mitigate high feature confusion. In the follow-153 ing sections, we provide a mathematical definition of this issue and introduce our SEAT framework 154 to address it.

155 156 157

158

3 Method

- 159 3.1 PROBLEM STATEMENT
- In the realm of time series forecasting (TSF), the problem statement can be formally defined as follows: Given a set of data points $\mathbf{X} = \{x_{t1}, \dots, x_{tD}\}_{t=1}^{L} \in \mathbb{R}^{D \times L}$ within a lookback window of

183

184

186

a time series, where *L* represents the size of the window, $D \ge 1$ is the number of variables, and x_{tj} denotes the value of the *j*-th variable at the *t*-th time step. The objective of TSF is to predict the forecasting horizon $\hat{\mathbf{X}} = {\hat{x}_{(L+1)1}, \dots, \hat{x}_{(L+T)D}}_{t=L+1}^{L+T} \in \mathbb{R}^{D \times T}$.

166 3.2 SOLVING BLOCK-LIKE ATTENTION FROM THE INPUT SIDE

168 The block-like distribution of attention feature maps is observed in many Transformer-based time se-169 ries tasks, where the magnitudes of the feature values exhibit a block-like pattern. This phenomenon 170 arises from the model's confusion among different approximate feature values. The conventional 171 attention mechanism can be expressed by the following formula:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$
(1)

174 The phenomenon of "Block-like Attention" arises when contiguous elements within blocks of the 175 attention score matrix exhibit similar magnitudes of values calculated by $A_{ij} = Q_i \cdot K_i^T$. This 176 phenomenon signifies a clustering of high or low attention intensities within localized regions of the 177 matrix, leading to a distinct block-like pattern. The presence of such patterns can be indicative of 178 the model focusing its attention on specific subsets of input features or interactions, thereby pro-179 viding insights into the model's decision-making processes. We can quantify the degree of feature confusion using the following mathematical definition, where the default method for *Similarity* is 181 the cosine similarity function and f_i, f_j are features: 182

$$Sim(F) = \begin{cases} 0, & \text{if } N = 1\\ \frac{1}{N*(N-1)} \sum_{i \neq j} Similarity(f_i, f_j), & \text{if } N \ge 2 \end{cases}$$
(2)

¹⁸⁵ In the following, we define two primary attention mechanisms.

Point-wise Attention (Temporal Attention) Definition: Point-wise attention, also known as temporal attention, is a mechanism within deep learning architectures that assigns significance scores to individual data points or temporal instances within a sequence. This approach enables the model to dynamically adjust its focus on specific points in time, capturing nuances and salient features that may be crucial for downstream tasks. By selectively attending to these key points, point-wise attention enhances the model's ability to comprehend complex temporal patterns and dynamics within the data. Notably, the effectiveness of this approach has been demonstrated in seminal works (Ashish (2017); Hu et al. (2018); Nie et al. (2023)).

194 **Channel-wise Attention Definition:** Channel-wise attention, on the other hand, emphasizes the 195 significance of individual feature channels within a multidimensional tensor. In this framework, each 196 channel represents a unique feature map capturing distinct aspects of the input data. Channel-wise 197 attention aims to dynamically reweight these feature channels, allowing the network to concentrate 198 its representational power on the most informative and discriminative channels. Specifically, for 199 time series data, each variable token is embedded into a high-dimensional space, where channelwise attention operates to capture intricate multivariate correlations. By highlighting the channels 200 that are most relevant to the task at hand, this mechanism enhances the overall feature representation, 201 leading to improved performance in tasks such as classification, regression, or forecasting. Notable 202 papers have highlighted the effectiveness of channel-wise attention (Wang et al. (2017); Woo et al. 203 (2018); Liu et al. (2024)). 204

205 Compared to "Point-wise Attention" features, using multivariate data as tokens at the same time 206 step results in naturally similar tokens for adjacent time steps. Although the sampling is discrete, the time series varies continuously, leading to many similar features that yield high similarity scores. 207 Sparse feature representation facilitates the computation of a sparse attention matrix. This means 208 that the attention matrix of size N * N needs to meet a specified minimal error threshold $\epsilon > 0$, 209 such that $N_0 = \sum_{i,j} \mathbb{I}(A_{ij} > \epsilon)$ is significantly less than N^2 where $\mathbb{I}(\cdot)$ is an indicator function 210 that takes the value of 1 when the condition is true and 0 otherwise. Experiments show that sparse 211 and independent feature sets yield a smaller Sim value, indicating that sparse features exhibit lower 212 similarity compared to dense features. 213

Standard Attention mechanisms exhibit certain limitations, notably the susceptibility to noise inter ference during weight computation, potentially leading to distractions on irrelevant elements. Be yond Top-k Attention, other sparse Attention mechanisms have emerged, such as Fixed Factorized

Attention and Strided Attention (Child et al. (2019)). According to our proposed metric Sim, models with lower feature confusion exhibit enhanced feature-capturing capabilities, fostering stronger generalization abilities. While sparse attention mechanisms theoretically risk overlooking certain global contextual information, they can partially mimic the effects of global context through meticulous design and optimisation. For instance, the BigBird (Zaheer et al. (2020)) adeptly captures long-range dependencies by blending sparse attention, global attention, and random skip connections.

However, in the domain of time series analysis, Block-like attention may hinder the model's ability to perceive global characteristics, while multi-scale attention frequently Chen et al. (2024) grapples with overfitting issues. Skip-connected attention, while improving attention focus, may compromise temporal resolution and lead to information loss, thereby restricting model performance. Prior research (Wu et al. (2020); Zhou et al. (2021); Nie et al. (2023)) has primarily focused on designing attention architectures and patterns. In contrast, our work demonstrates that incorporating sparsity at the input end can potentially address the shortcomings of existing attention mechanisms.

In our study, we propose **SEAT** to leverage frequency domain transformation to induce sparsity in time series data, thereby optimizing the performance of attention mechanisms from the input end. This approach is agnostic to the underlying attention model architecture, functioning as a versatile, plug-and-play component tailored specifically to enhance attention functionality. Our method offers a novel perspective on enhancing attention mechanisms, emphasizing the potential benefits of sparsity induction from the frequency domain at the input stage.

235 236

3.3 FOURIER DOMAIN OF THE TIME SERIES

The key to effectively addressing the challenge of block-wise attention lies in the ability of timefrequency transformation to impart sparse representations to time-series data. In this context, we will embark on a rigorous mathematical derivation to prove that long time-series data exhibit sparse representations under frequency domain transformations. Furthermore, we posit that by integrating feature partitioning into patches with frequency domain transformation, we can achieve an even more sparse representation of time-series data than either method alone.

Theorem 1: Sparse Representation in the Frequency Domain Given a signal x(t) of time series data in L^2 space, if its Fourier transform X(f) is supported on a finite set of frequencies, then x(t)has a sparse representation in the frequency domain.

The Nyquist Sampling Theorem (Vaidyanathan (2001)), also known as the Shannon Sampling Theorem, fundamentally establishes the conditions under which a continuous-time signal can be accurately reconstructed from its discrete-time samples, ensuring that the Discrete Fourier Transform (DFT) yields equivalent results to the Continuous Fourier Transform (CFT) within the context of the sampled signal's representation. We assume that the sampled data of the time series satisfies the sampling theorem. This implies that the sampling frequency f_s used to acquire the discrete samples x[n] of the original continuous signal x(t) is sufficiently high, specifically $f_s \ge 2B$, where B is the maximum frequency content of the signal's spectrum X(f).

Furthermore, we assume that the signals we acquire adhere to the finite bandwidth theorem, which 255 states that the spectral content of the signal is confined within a limited frequency range. Given 256 that the signals we collect are inherently discrete, due to the nature of the sampling process, it is 257 crucial to consider the appropriate frequency resolution for the specific prediction task at hand. For 258 instance, in tasks involving daily analysis, frequency features at the minute level may be deemed 259 irrelevant or treated as noise. To ensure compliance with the finite bandwidth principle in practical 260 sampling scenarios, we can adopt suitable sampling strategies that capture the signal within a rea-261 sonable frequency bandwidth range, thereby facilitating effective and efficient data processing for 262 the intended prediction tasks.

Proof of Theorem 1: Given the assumptions mentioned above, we can deduce that in our long time-series analysis tasks, where the sampling theorem is satisfied, the use of the Discrete Fourier Transform (DFT) becomes equivalent to the Continuous Fourier Transform (CFT) for solving time-series problems. Thus, for these tasks, we can confidently employ the DFT as a valid and computationally efficient tool for analyzing and processing time-series data.

Assuming the signal x[n] is expressed as $x[n] = \sum_{m=1}^{M} A_m e^{j\omega_m n}$ where A_m represents the amplitude of the *m*-th frequency component, and ω_m denotes the angular frequency of the *m*-th frequency

component. It is further assumed that ω_m satisfies the condition $\omega_m = \frac{2\pi k_m}{N}$ where k_m is an integer and N is a constant that typically represents the total number of samples in the discrete domain.

273 The DFT of a discrete signal x[n] is given by:

274 275

276 277

278 279 280

$$X[k] = \sum_{n=0}^{N-1} x[n] e^{-j\frac{2\pi}{N}kn} = \sum_{m=1}^{M} A_m \sum_{n=0}^{N-1} e^{j(\omega_m - \frac{2\pi}{N}k)}, \quad k = 0, 1, \dots, N-1.$$
(3)

Under the frequency basis function we selected, the following equation can be calculated:

$$S_{m,k} = \sum_{n=0}^{N-1} e^{j(\omega_m - \frac{2\pi}{N}k)} = \begin{cases} N, & \text{if } k = k_m \\ \frac{1 - e^{j2\pi(k_m - k)}}{1 - e^{j\frac{2\pi}{N}(k_m - k)}} = 0, & \text{if } k \neq k_m \end{cases}$$
(4)

281 282 283

284

285

286 287

288

As a result, the frequency domain representation X[k] exhibits non-zero values only at $k = k_m$, where k_m corresponds to the indices of the frequency components present in the signal. Given that our time series can be arbitrarily long, the representation in the frequency domain of the "hidden space" of time series becomes sparse within frequency components.

3.4 FOURIER ATTENTION

Fourier attention operates by initially applying the Fourier Transform to the queries, keys, and values, subsequently conducting an attention mechanism within the frequency domain, and ultimately converting the outcomes back to the time domain via the inverse Fourier Transform. Let $F(\cdot)$ and $F^{-1}(\cdot)$ represent the Fourier Transform and inverse Fourier Transform, respectively. The Fourier attention mechanism can be formally expressed as follows:

296 297

298

299

300 301

302

 $Attention(Q, K, V) = F^{-1} \left(\sigma \left(\frac{F(Q)\overline{F(K)}^T}{\sqrt{d_q}} \right) F(V) \right), \tag{5}$

where $\sigma(\cdot)$ denotes the softmax function, and d_q is the dimensionality of the queries. Citing from previous research Zhang et al. (2022), calculating attention in the Fourier domain is equivalent to time-domain attention. Our method can be extended to attention in the frequency domain.

3.5 SEAT

Our proposed SEAT framework can be decomposed into several principal components, encompassing the Normalization Layer, SEAT Block, Feature Extraction Layer and Projection Layer.

Normalization Layer: Reversible instance normalization (RevIN)Kim et al. (2021), a generally applicable normalization-and-denormalization method with learnable affine transformation is well-known and widely used as the normalization layer. For better comparison of different models abilities, we use Revin normalization uniformly for all models.

SEAT Block: We drew a precise schematic diagram of the SEAT block structure in Fig 1. This
 novel framework is designed to induce sparsity in time series data through frequency domain trans formations. This model-agnostic approach serves as a versatile and easily integrated component,
 compatible with any underlying attention architecture. By focusing on sparsity induction via fre quency domain transformations, SEAT is specifically designed to enhance the overall performance
 of attention mechanisms from a new perspective in the input stage.

Feature Extraction Layer: This is a temporal feature extractor specifically based on transformer architecture since we designed SEAT to enhance the ac- curacy and robustness of attention mechanism.
Two primary types of attention mechanisms can be introduced: one represented by "Channel-wise
Attention" as in the case of iTransformer (Liu et al. (2024)), and the other by "Point-wise Attention" as exemplified by PatchTST (Nie et al. (2023)).

Regarding the whole processing pipeline of SEAT, as shown in Figure 1, we utilize the Revin normal ization technique Kim et al. (2021) to preprocess the input time series. Subsequently, the sequence
 undergoes transposition and is fed into the SEAT block. Within this block, each time point of the
 individual series is embedded into variable tokens, facilitating the application of the Fast Fourier

339 340



Figure 1: Overall Structure of SEAT.

341 Transform (FFT). This transformation converts the time series into the frequency domain. A linear 342 module is then employed to enable the model to learn robust and sparse representations of the sequence. Following this, an Inverse Fast Fourier Transform (IFFT) and skip connection are applied to 343 revert the sequence into the time domain. The resulting sparser representation of the sequence, once 344 obtained, undergoes transposition and is subsequently input into an attention-based feature extrac-345 tor. Ultimately, the sequence undergoes a de-normalization process and projection, yielding robust 346 predictions for Long-Term Sequence Forecasting (LTSF). 347

348 In summary, while Channel-wise Attention focuses on the interrelationships between different fea-349 ture channels, Point-wise Attention emphasizes dependencies between individual points or segments within the data. Both mechanisms are integral to enhancing model performance, depending on the 350 specific structure and demands of the task at hand. Our SEAT model utilizes the iTransformer as the 351 Temporal Feature Extractor, enabling it to serve as a versatile, plug-and-play component that can 352 seamlessly integrate with any state-of-the-art (SOTA) transformer to enhance the model's overall 353 performance. 354

355 356

357 358

4 EXPERIMENT

Datasets: we have undertaken extensive experiments on eight meticulously curated benchmarks, 359 notably including the ETT datasets, which are subdivided into four distinct subsets: ETTh1, ETTh2, 360 ETTm1, and ETTm2. Additionally, we have also employed the Weather, Exchange ECL, and Traffic 361 datasets, adhering to the precedents established in Zhou et al. (2021); Zeng et al. (2023); Hebrail & 362 Berard (2012); Zhao et al. (2019). These benchmarks, renowned for their rigour and comprehensive-363 ness, provide a robust framework for assessing the performance and effectiveness of our forecasting 364 models, particularly in the context of long-term horizon predictions.

Baselines Compared: Our proposed **SEAT** represents a model-agnostic approach that exhibits 366 broad applicability to any deep neural network architecture. In this study, we extensively com-367 pare the well-acknowledged and advanced Transformers and designed a plug-and-play experiment 368 to meticulously evaluate the efficacy of SEAT by integrating it into seven state-of-the-art Transform-369 ers designed specifically for time-series forecasting: iTransformer (Liu et al. (2024)), PatchTST (Nie 370 et al. (2023)), Crossformer (Zhang & Yan (2023)), Pyraformer (Liu et al. (2022a)), Autoformer (Wu 371 et al. (2021)), Informer (Zhou et al. (2021)), and Reformer (Kitaev et al. (2020)). This compre-372 hensive experiment serves to validate the generality and enhancement capabilities of SEAT when 373 applied to diverse yet sophisticated Transformer-based frameworks. The SEAT base model is typi-374 cally employed with iTransformer as backbone architecture without explicit indication in the Main 375 Result. We compute the MSE and MAE on Revin (Kim et al. (2021)) normalized data to measure different variables on the same scale. More details on experimental settings, including training de-376 tails and hyperparameters, are provided in the Appendix. Experiments are implemented in PyTorch 377 (Paszke et al. (2019)) and conducted on a single NVIDIA 4090 24G.

378 4.1 MAIN RESULT379

Table 1: Mean result of **SEAT** versus other SOTAs transformers. In eight benchmark datasets, our method achieved first place in six out of the mean squared error (MSE) metrics and seven out of the mean absolute error (MAE) metrics.

Models	SE Oi	AT 1rs	iTrans 20	former 24	Patch 20	nTST 23	Crossi 20	former 23	Pyraf 20	ormer 22	Autof 20	ormer 21	Info 20	rmer 21	Refo 20	ormer)20	
Metric	mse	mae	mse	mae	mse	mae	mse	mae	mse	mae	mse	mae	mse	mae	mse	mae	
ETTh1	0.436	0.433	0.454	0.447	0.469	0.454	0.529	0.522	0.865	0.731	0.518	0.500	1.078	0.813	0.961	0.757	
ETTh2	0.372	0.398	0.383	0.407	0.387	0.407	0.942	0.684	3.755	1.551	0.432	0.451	3.490	1.532	3.574	1.525	
ETTm1	0.396	0.400	0.407	0.410	0.387	0.400	0.513	0.496	0.750	0.615	0.583	0.513	0.948	0.717	0.928	0.688	
ETTm2	0.281	0.325	0.288	0.332	0.281	0.326	0.757	0.610	1.509	0.845	0.332	0.370	1.489	0.867	1.415	0.862	
Weather	0.249	0.277	0.258	0.278	0.259	0.281	0.259	0.315	0.278	0.342	0.317	0.359	0.723	0.605	0.485	0.500	
ECL	0.173	0.266	0.178	0.270	0.205	0.290	0.244	0.334	0.298	0.389	0.230	0.339	0.377	0.449	0.302	0.392	
Exchange	0.349	0.402	0.360	0.403	0.367	0.404	0.940	0.707	1.308	0.945	0.493	0.493	1.411	0.968	1.000	0.837	
Traffic	0.442	0.286	0.428	0.282	0.481	0.304	0.550	0.304	1.185	0.553	0.761	0.479	0.868	0.472	0.648	0.347	
1st count	6	7	1	1	1	0	0	0	0	0	0	0	0	0	0	0	

Main Results: Table 1 presents concise forecasting outcomes, with the best accurate predictions highlighted in **red** and the second-best <u>underlined</u>. A lower MSE/MAE signifies superior prediction accuracy. **SEAT** ensures the independence and sparsity of input features within Transformer architectures. SEAT optimizes the utilization of input data, facilitating more efficient and effective processing within the attention layers of the Transformer model. SEAT method has demonstrated significant advantages across eight benchmark datasets. Our proposed model aggregates sparser features in the input end of the attention module and achieves the best result in long temporal modelling, enhancing the model's ability to handle high-dimensional time series and mitigating overfitting.

4.2 PLUG AND PLAY EXPERIMENTS

Plug-and-play experiment: In our plug-and-play experimental results, we observe that SEAT sig-nificantly improves the performance of various Transformer models. This phenomenon is clearly visualized in Figure 2. These visualizations demonstrate that SEAT effectively enhances prediction accuracy, regardless of whether it is applied to iTransformer, PatchTST, or other models. Table 2 illustrates the improvements in prediction accuracy by contrasting the original models with those augmented by SEAT. By calculating the percentage improvements, we clearly demonstrate the ben-eficial effects brought by this powerful plugin. As a model-agnostic enhancement, SEAT delivered immediate performance gains across all the models we evaluated. SEAT ensures the independence and sparsity of input features within Transformer architectures, optimizing the utilization of input data and facilitating more efficient and effective processing within the attention layers of the Transformer model. The substantial improvements observed across diverse datasets underscore its broad applicability and adaptability, thereby reinforcing SEAT's value as an effective tool for LSTF task.



Figure 2: SEAT's Impact on Various Transformer Models' MSE Performance. The entire bars represent the original MSE scores, while the light green segments indicate the performance after applying the SEAT plugin, and the dark green segments represent the improvement specifically attributed to SEAT. These visualizations demonstrate the effectiveness of SEAT in enhancing the forecasting accuracy of diverse transformer models.

Table 2: Improvements of SEAT over different models with prediction lengths

	10010 2	• •••••	1010	meme	010		over a	merer	n mou	010 111	in proc	100101	Tionge	110	
	Models	iTrans	former	Patch	TST	Cross	Crossformer		former	Auto	former	Info	rmer	Refe	ormer
	Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
	Original	0.454	0.447	0.469	0.454	0.529	0.522	0.865	0.731	0.518	0.500	1.078	0.813	0.961	0.757
ETTh1	+SEAT	0.436	0.433	0.447	0.443	0.484	0.466	0.488	0.474	0.486	0.468	1.026	0.736	0.588	0.530
	Improvement	+4.0%	+3.3%	+4.7%	2.6%	+8.5%	+10.7%	+43.6%	+35.2%	+6.2%	+6.4%	+4.8%	+9.5%	+38.8%	+29.9%
	Original	0.383	0.407	0.387	0.407	0.389	0.416	3.755	1.551	0.432	0.451	3.490	1.532	3.574	1.525
ETTh2	+SEAT	0.372	0.398	0.374	0.402	0.394	0.416	0.433	0.434	0.459	0.449	0.677	0.571	0.459	0.451
	Improvement	+2.9%	+2.2%	+3.4%	1.2%	+58.2%	+39.1%	+88.5%	+72.0%	-6.3%	+0.4%	+80.6%	+62.7%	+87.2%	+70.4%
ETTm1	Original	0.407	0.410	0.387	0.400	0.513	0.496	0.750	0.615	0.583	0.513	0.948	0.717	0.928	0.688
	+SEAT	0.396	0.400	0.380	0.396	0.410	0.411	0.427	0.425	0.529	0.480	0.629	0.530	0.599	0.507
	Improvement	+2.7%	+2.4%	+1.8%	1.0%	+20.1%	+17.0%	+43.0%	+30.9%	+9.3%	+6.4%	+33.6%	+26.1%	+35.4%	+26.3%
ETTm2	Original	0.288	0.332	0.281	0.326	0.757	0.610	1.509	0.845	0.332	0.370	1.489	0.867	1.415	0.862
	+SEAT	0.281	0.325	0.279	0.326	0.292	0.332	0.302	0.336	0.304	0.342	0.389	0.406	0.318	0.348
	Improvement	+2.4%	+2.1%	+0.7%	+0.2%	+61.4%	+45.7%	+80.0%	+60.2%	+8.4%	+7.6%	+73.9%	+53.2%	+77.5%	+59.6%
	Original	0.258	0.278	0.259	0.281	0.259	0.315	0.278	0.342	0.317	0.359	0.723	0.605	0.485	0.500
weather	+SEAT	0.249	0.277	0.254	0.281	0.261	0.292	0.257	0.284	0.278	0.304	0.289	0.320	0.273	0.299
	Improvement	+3.5%	+0.4%	+1.9%	-0.4%	-0.8%	+7.3%	+7.6%	+17.0%	+12.3%	+15.3%	+60.0%	+47.1%	+43.7%	+40.2%
	Original	0.178	0.270	0.205	0.290	0.244	0.334	0.298	0.389	0.230	0.339	0.377	0.449	0.302	0.392
ECL	+SEAT	0.173	0.266	0.186	0.277	0.167	0.260	0.206	0.309	0.208	0.308	0.258	0.358	0.210	0.312
	Improvement	+2.8%	+1.5%	+9.3%	4.5%	+31.6%	+22.2%	+30.9%	+20.6%	+9.6%	+9.1%	+31.6%	+20.3%	+30.2%	+20.4%
	Original	0.360	0.403	0.367	0.404	0.940	0.707	1.308	0.945	0.493	0.493	1.411	0.968	1.000	0.837
Exchange	+SEAT	0.349	0.402	0.365	0.405	0.367	0.414	0.396	0.429	0.459	0.466	0.368	0.428	0.448	0.459
	Improvement	+3.1%	+0.2%	+0.3%	-0.2%	+61.0%	+41.4%	+69.7%	+54.6%	+6.9%	+5.5%	+73.9%	+55.8%	+55.2%	+45.2%
	Original	0.428	0.282	0.481	0.304	0.550	0.304	1.185	0.553	0.761	0.479	0.868	0.472	0.648	0.347
traffic	+SEAT	0.442	0.286	0.477	0.291	0.479	0.311	0.794	0.436	0.713	0.392	1.030	0.567	0.638	0.333
	Improvement	-3.3%	-1.4%	+0.8%	+4.3%	+12.9%	-2.3%	+33.0%	+21.16%	+6.2%	+18.2%	-18.7%	-20.1%	+1.4%	+4.0%

ATTENTION STUDY 4.3



Figure 3: Self-attention scores from iTransformer and SEAT trained on ETTh1. The left heatmap represents the iTransformer's performance, which is likely the most advanced model to our knowl-edge. We integrate SEAT as a plugin into the iTransformer, and the right heatmap clearly illustrates the significantly improved self-attention scores achieved after applying SEAT. The visualization re-sults demonstrate enhanced clarity and robustness in feature representation, thereby underscoring the significant efficacy of SEAT in augmenting the model's predictive capabilities.

Figure 3 presents a representative attention score map of the iTransformer for the Long-Term Short-Term Forecasting (LTSF) task. It is observable that the attention values in iTransformer tend to segment, with nearby data points sharing similar attention weights, resulting in ambiguous feature representations. Meanwhile, applying **SEAT** transforms the input signals into the frequency domain, leading to sparser and more distinct feature representations optimized for the attention mechanism. This transformation broadens the range of attention values and increases their distribution variance. The enhanced variance arises because frequency domain transformations decompose the time se-ries into orthogonal frequency components, effectively reducing feature correlation and improving feature distinguishability. With sparser and less redundant features, the attention mechanism can assign a wider range of weights, more accurately reflecting the true importance of diverse features. A broader range and higher variance in attention values enable the model to differentiate more distinct and relevant features, thereby reducing feature redundancy and mitigating overfitting to noise.

486 Consequently, the attention mechanism in **SEAT** more effectively captures both cross-time and 487 cross-dimension dependencies, resulting in clearer and more robust feature representations. These 488 improvements facilitate more precise pattern recognition and enhance the overall learning capacity 489 of the model, significantly boosting forecasting performance.

490 491 492

493

501 502

503 504

505

506

517

518

522

531

532

533

CONCLUSION 5

In this study, we introduce SEAT, a Sparse Sensing Enhancement Framework specifically designed 494 to optimize attention mechanisms for time series forecasting within Transformer architectures. 495 Through rigorous mathematical proofs and empirical analysis, we have demonstrated its feasibil-496 ity and effectiveness. By ensuring the independence and sparsity of input features, the framework 497 enhances both model performance and interpretability. Additionally, it overcomes the limitations of 498 existing attention mechanisms by offering seamless plug-and-play functionality and compatibility 499 with any Transformer-based architecture. This adaptability positions SEAT as a significant contri-500 bution to the field, paving the way for future advancements in attention-based models.

- References
- Vaswani Ashish. Attention is all you need. Advances in neural information processing systems, 30: I, 2017.
- Peng Chen, Yingying ZHANG, Yunyao Cheng, Yang Shu, Yihang Wang, Oingsong Wen, Bin Yang, 507 and Chenjuan Guo. Pathformer: Multi-scale transformers with adaptive pathways for time series 508 forecasting. In The Twelfth International Conference on Learning Representations, 2024. 509
- 510 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse 511 transformers. arXiv preprint arXiv:1904.10509, 2019. 512
- 513 Dazhao Du, Bing Su, and Zhewei Wei. Preformer: predictive transformer with multi-scale segmentwise correlations for long-term time series forecasting. In ICASSP 2023-2023 IEEE International 514 Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1–5. IEEE, 2023. 515
- 516 John Guibas, Morteza Mardani, Zongyi Li, Andrew Tao, Anima Anandkumar, and Bryan Catanzaro. Efficient token mixing for transformers via adaptive fourier neural operators. In International Conference on Learning Representations, 2021. 519
- 520 Georges Hebrail and Alice Berard. Individual Household Electric Power Consumption. UCI Ma-521 chine Learning Repository, 2012.
- Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE* 523 conference on computer vision and pattern recognition, pp. 7132–7141, 2018. 524
- 525 Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Re-526 versible instance normalization for accurate time-series forecasting against distribution shift. In 527 International Conference on Learning Representations, 2021. 528
- 529 Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. arXiv preprint arXiv:2001.04451, 2020. 530
 - James Lee-Thorp, Joshua Ainslie, Ilya Eckstein, and Santiago Ontanon. Fnet: Mixing tokens with fourier transforms. arXiv preprint arXiv:2105.03824, 2021.
- 534 Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X. Liu, and Schahram Dustdar. 535 Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and fore-536 casting. In International Conference on Learning Representations, 2022a. 537
- Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers: Exploring 538 the stationarity in time series forecasting. Advances in Neural Information Processing Systems, 35:9881-9893, 2022b.

- 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 *The Twelfth International Conference on Learning Representations*, 2024.
- Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64
 words: Long-term forecasting with transformers. In *The Eleventh International Conference on Learning Representations*, 2023.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Xihao Piao, Zheng Chen, Taichi Murayama, Yasuko Matsubara, and Yasushi Sakurai. Fredformer:
 Frequency debiased transformer for time series forecasting. *arXiv preprint arXiv:2406.09009*, 2024.
- PP Vaidyanathan. Generalizations of the sampling theorem: Seven decades after nyquist. *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*, 48(9):1094–1109, 2001.
- Fei Wang, Mengqing Jiang, Chen Qian, Shuo Yang, Cheng Li, Honggang Zhang, Xiaogang Wang, and Xiaoou Tang. Residual attention network for image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3156–3164, 2017.
- Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 3–19, 2018.
- Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition trans formers with auto-correlation for long-term series forecasting. *Advances in neural information processing systems*, 34:22419–22430, 2021.
- Sifan Wu, Xi Xiao, Qianggang Ding, Peilin Zhao, Ying Wei, and Junzhou Huang. Adversarial sparse transformer for time series forecasting. *Advances in neural information processing systems*, 33: 17105–17115, 2020.
- Zhijian Xu, Ailing Zeng, and Qiang Xu. FITS: Modeling time series with \$10k\$ parameters. In *The Twelfth International Conference on Learning Representations*, 2024.
- Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago
 Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for
 longer sequences. *Advances in neural information processing systems*, 33:17283–17297, 2020.
- Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 11121–11128, 2023.
- Xiyuan Zhang, Xiaoyong Jin, Karthick Gopalswamy, Gaurav Gupta, Youngsuk Park, Xingjian Shi, Hao Wang, Danielle C Maddix, and Yuyang Wang. First de-trend then attend: Rethinking attention for time-series forecasting. *arXiv preprint arXiv:2212.08151*, 2022.
- Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency
 for multivariate time series forecasting. In *The Eleventh International Conference on Learning Representations*, 2023.

- Liang Zhao, Olga Gkountouna, and Dieter Pfoser. Spatial auto-regressive dependency interpretable
 learning based on spatial topological constraints. *ACM Trans. Spatial Algorithms Syst.*, 5(3), aug 2019. ISSN 2374-0353.
- Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.
 Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings* of the AAAI conference on artificial intelligence, volume 35, pp. 11106–11115, 2021.
- Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency
 enhanced decomposed transformer for long-term series forecasting. In *International conference* on machine learning, pp. 27268–27286. PMLR, 2022.

APPENDIX А

SAMPLING THEOREM A.1

598 If the signal x(t) possesses a limited bandwidth, with its spectrum X(f) = 0 for all |f| > B, then it can be perfectly reconstructed from samples taken at a sampling frequency f_s that satisfies the condition $f_s \ge 2B$. The reconstruction is achieved using the following formula:

 $x(t) = \sum_{n=-\infty}^{\infty} x[n] \operatorname{sinc}\left(\frac{t-nT}{T}\right), \quad \text{where} \quad T = \frac{1}{f_s}.$

594

595 596

597

600

606

607

608

609

610

605

In this equation, x[n] represents the discrete samples of the signal taken at intervals of T, which is the reciprocal of the sampling frequency f_s . The sinc function, defined as $\operatorname{sinc}(x) = \frac{\sin(\pi x)}{\pi x}$, plays a crucial role in interpolating between the samples to reconstruct the continuous signal $\ddot{x}(t)$. This process is a direct consequence of the Nyquist-Shannon sampling theorem, which ensures that a bandlimited signal can be uniquely determined from its samples taken at a sufficiently high rate.

Therefore, we assume that the sampled data of the time series satisfies the sampling theorem. This 611 implies that the sampling frequency f_s used to acquire the discrete samples x[n] of the original 612 continuous signal x(t) is sufficiently high, specifically $f_s \ge 2B$, where B is the maximum frequency 613 content of the signal's spectrum X(f). Given this assumption, the discrete samples contain enough 614 information to perfectly reconstruct the original continuous signal x(t) using the reconstruction 615 formula provided earlier.

616 617 618

A.2 PATCHING AND SEAT MAKES SPARSER FEATURE REPRESENTATION FOR TIME SERIES

619 Consider a feature set F, where F_s denotes a set of features $\{f_1, f_2, \ldots, f_N\}$. Each feature $f_i =$ 620 $(a_{i,1}, a_{i,2}, \ldots, a_{i,N})$ in F is characterized by N dimensions. To investigate the Patching algorithm through an inductive approach, we will initiate our discussion by focusing on the merging of a pair 621 of features. We hypothesize that the N^{th} and $(N-1)^{th}$ dimensions exhibit the highest degree 622 of similarity. A single step in the patching process involves removing f_N , resulting in f_p , which 623 represents one step of patching by eliminating highly correlated features. 624

625 Our objective is to evaluate whether this merging process, colloquially referred to as "Patching", 626 leads to a sparser feature matrix by comparing the mathematical expressions before and after the 627 operation. The new feature after patching can be defined as $f_i = (a_{i,1}, a_{i,2}, \dots, a_{i,N})$ in F_p . 628

$$Sim(F) = \frac{\sum_{i \neq j} \sum_{k=1}^{N} a_{i,k} a_{j,k}}{N(N-1)}$$
(7)

(6)

631 632 633

634 635 636

645

629 630

$$Sim(F_p) = \frac{\left(\sum_{i \neq j, i, j \leq N-2} \sum_{k=1}^{N} a_{i,k} a_{j,k}\right) + 2\sum_{k=1}^{N} \sum_{i=1}^{N-2} a_{i,k} a_{c,k}}{(N-1)(N-2)}$$
(8)

The assumption can be rewritten as follows:

$$\sum_{k=1}^{N} a_{N,k} a_{N-1,k} > \sum_{k=1}^{N} a_{i,k} a_{j,k}, \quad \text{for all } i \neq j$$
(9)

$$Sim(F) = \frac{\sum_{i \neq j, i, j \leq N-1} \sum_{k=1}^{N} a_{i,k} a_{j,k} + 2 \sum_{i < N} \sum_{k=1}^{N} a_{i,k} a_{N,k} + 2 \sum_{k=1}^{N} a_{N-1,k} a_{N,k}}{N * (N-1)}$$
(10)

Since the similarity between f_N and f_{N-1} is the largest, this leads to a significant contribution that 646 is explicitly included in Sim(F) but diminished in $Sim(F_p)$, therefore $Sim(F) \ge Sim(F_p)$. The 647 Patching operation results in a sparser feature matrix.

The integration of feature segmentation into patches and frequency domain transformations has yielded a sparser representation of time series data compared to utilizing either method in isolation.
Previous research (Nie et al. (2023); Du et al. (2023)) has employed patching techniques to divide univariate time series into patches, which can be either overlapped or non-overlapped. The key feature in PatchTST, characterized by channel independence and patching, can also be interpreted as feature independence and a sparser feature representation in our context.

A.3 EXPERIMENT SETTING

we have standardized the parameters across all models to ensure a fair comparison on a uniform platform (time-series-library). Specifically, we have fixed the input dimension to 96 and varied the prediction horizon for time series forecasting, encompassing lengths of [96, 192, 336, 720]. The batch_size is set to 32, learning_rate is set to 1e-3,d_model is set to 512 and dropout is set to 0.1.

Table 3: Performance comparison in terms of forecasting error metrics. This table illustrates the performance comparison among different models in terms of forecasting error metrics, adhering to a unified setting to ensure fairness. The mse and mae values highlight the accuracy of the predictions. The best-performing results are highlighted in **red**, while the second-best results are marked in <u>blue</u> with underlining. Lower MSE/MAE values signify higher predictive accuracy. The incorporation of SEAT into various benchmark attention models demonstrates significant performance enhance-ments, showcasing SEAT's effectiveness in improving the forecasting capabilities of Transformer-based models as a model-agnostic plugin.

Mod	lels ^{i T}	rans 20	former 124	r -	+SEAT	F	Patch 20	nTST 23	+SI	EAT	Cross 20	forme)23	r +	SEA	т	Pyrafe 202	ormer 23	+SI	EAT	Auto 20	former 022	+8	EAT	Info 20	rmer 21	+SI	EAT	Refo 20	rmer 20	+SF	EAT	Ы	m
Met	ric 1	nse	mae	n	nse ma	ie r	mse	mae	mse	mae	mse	mae	ms	se r	nae	mse	mae	mse	mae	mse	mae	mse	mae	mse	mae	mse	mae	mse	mae	mse	mae	mse	
	96 <mark>0</mark> .	386	0.405	0.3	377 0.39	96 0.	.414	0.419	0.389	0.404	0.423	0.448	0.4	16 0.	.426	0.790	0.696	0.424	0.432	0.485	5 0.468	0.41	3 0.433	1.006	0.780	0.848	0.641	0.826	0.687	0.554	0.505	16.97%	16
E 1	92 0.	441	0.436	0.4	425 0.42	23 0.	.460	0.445	0.439	0.435	0.471	0.474	0.4	78 0.	.461	0.738	0.662	0.478	0.464	0.554	0.517	0.52	5 0.478	0.927	0.721	0.977	0.718	0.906	0.732	0.555	0.517	11.48%	4
E 3	36 0.	487	0.458	0.4	466 0.44	43 0.	.501	0.466	0.474	0.453	0.570	0.546	0.49	99 0.	.469	0.922	0.766	0.509	0.485	0.503	3 0.491	0.49	1 0.469	1.197	0.878	1.058	0.765	1.022	0.786	0.600	0.530	17.46%	ł
7	20 0.	503	0.491	0.4	475 0.43	71 0.	.500	0.488	<u>0.484</u>	0.480	0.653	0.621	0.54	43 0.	.510	1.011	0.799	0.539	0.514	0.531	0.523	0.50	9 0.491	1.182	0.873	1.220	0.819	1.091	0.821	0.642	0.570	16.349	ł
2	96 0.	297	0.349	0.2	2 <u>89</u> 0.3.	39 0.	.302	0.348	0.288	0.341	0.745	0.584	0.32	20 0.	.361	1.324	0.870	0.328	0.372	0.350	0.397	0.36	5 0.389	2.535	1.272	0.440	0.458	2.410	1.234	0.379	0.398	43.189	5
Ē	92 0.	380	0.400	0.3	<u>372 0.39</u>	<u>90</u> 0.	.388	0.400	0.357	0.388	0.877	0.656	0.3	77 0.	.402	4.264	1.608	0.422	0.420	0.430	0.440	0.44	7 0.440	3.557	1.522	0.787	0.608	4.812	1.930	0.472	0.449	45.909	5
Ξ 3	36 0.	428	0.432	0.4	413 0.42	25 0.	.426	0.433	0.425	0.434	1.043	0.731	0.4.	38 0.	.440	5.218	1.944	0.486	0.462	0.463	3 0.471	0.49	7 0.475	4.388	1.725	0.734	0.604	4.245	1.638	0.484	0.467	45.289	ž
	20 0.	427	0.445	0.4	415 0.43	39 0.	.431	0.446	0.426	0.446	1.104	0.763	0.4	<u>39 0.</u>	.461	4.213	1.780	0.498	0.483	0.484	0.497	0.52	3 0.492	3.478	1.609	0.747	0.614	2.828	1.296	0.499	0.488	43.459	į
- Ta 👌	<i>9</i> 6 0.	334	0.368	0.:	<u>532</u> 0.30	65 0.	.329	0.367	0.329	0.366	0.404	0.426	0.3.	34 O.	.373	0.609	0.540	0.371	0.391	0.519	0.490	0.47	0.444	0.765	0.655	0.427	0.427	0.749	0.575	0.479	0.450	20.829	ļ
Ē	92 0.	377	0.391	0.2	573 <u>0.3</u>	<u>83 0.</u>	367	0.385	0.356	0.380	0.450	0.451	0.3	/6 0.	.390	0.606	0.537	0.422	0.418	0.567	0.503	0.54	0.492	0.696	0.586	0.554	0.493	0.879	0.668	0.551	0.481	16.079	2
물건	30 0.	420	0.420	0.4	100 0.40	450	.399	0.410	0.390	0.403	0.552	0.513	0.44	42 U.	424	0.864	0.000	0.428	0.450	0.558	5 0.507	0.55	7 0.490	1.098	0.799	0.620	0.525	1.017	0.736	0.640	0.528	22.62	
	20 0.	191	0.439	0.4	174 0.44	+50.	175	0.459	0.444	0.433	0.000	0.385	0.4	37 0.	261	0.919	0.718	0.480	0.405	0.085	0.331	0.55	0.488	1.234	0.827	0.914	0.075	1.005	0.774	0.728	0.309	22.34	2
- E .	0 0.	250	0.204	0.1	1/4 0.2:	0/ <u>U.</u>	241	0.239	0.170	0.200	0.287	0.300	0.1	/0 U.	215	0.415	0.480	0.194	0.2/2	0.234	+ 0.515	0.19	0.215	0.487	0.551	0.204	0.295	0.555	0.504	0.199	0.270	33.17	7
Ea	36 0	311	0.309	0.3	307 0 3	130	305	0.302	0.248	0.307	0.414	0.492	0.20	06.0	345	1 207	0.009	0.200	0.312	0.292	0.349	0.20	3 0 360	1 567	0.392	0.314	0.370	1 221	0.833	0.284	0.327	30 34	ł
Ξ ₇	20 0	412	0.0407	0.7	105 0 30	000	402	0.400	0.304	0.345	1 730	1.042	0.1	24.0	106	3 676	1 4 1 1	0.315	0.113	0.357	10.073	0.42	3 0 4 1 2	3 333	1 353	0.400	0.527	2 670	1 206	0.354	0.303	48 58	ć
(26 0	174	0.214	0.1	160 0 2	100	177	0.218	0.185	0.231	0 158	0.230	0.1	72.0	221	0.210	0.204	0.172	0.222	0.235	5 0 304	0.10	5 0 243	0.333	0.303	0.000	0.258	0.354	0.424	0.101	0.125	16.25	ī
- Ē i	92 0	221	0.254	0.2	214 0.2	560	225	0.259	0.213	0.255	0.206	0.277	0.2	310	275	0.240	0.314	0.223	0.263	0.301	0.354	0.25	3 0 293	0.858	0.652	0.267	0.306	0.314	0.384	0.242	0.280	15.89	ê
r sat	36 0	278	0.296	0.2	269 0.29	93 0	278	0.297	0.269	0.295	0.272	0.335	0.2	83.0	312	0.295	0.356	0.278	0.301	0.329	0.367	0.29	9 0.321	0.664	0.583	0.317	0.340	0.590	0.566	0.293	0.314	17.13	o
₹ 7	20 0.	358	0.347	0.3	353 0.34	490.	354	0.348	0.349	0.345	0.398	0.418	0.3	58 0.	359	0.367	0.404	0.355	0.351	0.403	3 0.410	0.36	1 0.359	1.035	0.793	0.370	0.374	0.681	0.627	0.364	0.358	19.52	•
-	96 0.	148	0.240	0.1	38 0.2	34 0.	.181	0.270	0.157	0.248	0.219	0.314	0.13	36 0.	.230	0.285	0.377	0.178	0.285	0.196	5 0.313	0.17	3 0.279	0.340	0.420	0.219	0.325	0.274	0.372	0.190	0.294	24.78	7
H 1	92 0.	162	0.253	0.1	159 0.25	53 0.	188	0.274	0.169	0.262	0.231	0.322	0.1	56 0.	.248	0.289	0.384	0.197	0.301	0.220	0.330	0.21	1 0.308	0.378	0.452	0.241	0.346	0.293	0.389	0.200	0.303	21.19	9
Ĕ 3	36 0.	178	0.269	0.1	174 0.20	<mark>68</mark> 0.	.204	0.293	0.188	0.281	0.246	0.337	0.1	72 0.	.266	0.308	0.400	0.212	0.315	0.246	6 0.350	0.21	0.312	0.387	0.460	0.271	0.372	0.321	0.408	0.217	0.318	21.19	9
7	20 0.	225	0.317	0.2	222 0.30	<mark>09</mark> 0.	.246	0.324	0.232	0.316	0.280	0.363	0.20	03 0.	.295	0.309	0.396	0.239	0.334	0.259	0.362	0.23	8 0.333	0.403	0.464	0.301	0.390	0.318	0.400	0.235	0.333	16.67	3
8 9) 6 0.	086	0.206	0.0)87 0.20	07 0.	.088	0.205	0.086	0.206	0.256	0.367	0.0	88 0.	.211	0.898	0.798	0.106	0.234	0.158	3 0.288	0.12	$2\ 0.249$	0.974	0.794	0.122	0.250	0.820	0.758	0.121	0.247	50.06	ē
- E 1	92 <u>0.</u>	177	<u>0.299</u>	0.1	179 0.30	04 <mark>0.</mark>	.176	0.299	0.176	0.298	0.470	0.509	0.18	84 0.	.307	1.065	0.870	0.205	0.326	0.268	3 0.379	0.23	4 0.353	1.314	0.942	0.215	0.339	0.900	0.782	0.237	0.349	44.35	2
523	36 0.	331	0.417	0.3	<u>324 0.4</u>	<u>16</u> 0.	.301	0.397	0.331	0.417	1.268	0.883	0.30	67 0.	.443	1.287	0.957	0.353	0.431	0.434	1 0.487	0.41	5 0.473	1.403	0.988	0.368	0.450	1.080	0.877	0.403	0.464	39.48	2
<u><u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u></u></u>	20 0.	847	0.691	0.8	<u>306 0.68</u>	<u>82</u> 0.	.901	0.714	0.868	0.701	1.767	1.068	0.82	29 0.	.694	1.980	1.155	0.922	0.726	1.111	0.818	1.06	3 0.787	1.954	1.146	0.768	0.675	1.199	0.929	1.029	0.775	27.75	2
ు ⁶) 6 0 .	395	0.268	0.3	<u>396</u> 0.20	67 0.	.462	0.295	0.457	0.277	0.522	0.290	0.4	88 0.	.337	0.839	0.466	0.844	0.471	0.613	3 0.385	0.68	0.372	0.877	0.475	0.973	0.517	0.641	0.344	0.617	0.321	-1.63	7
E 1	92 0.	417	0.276	0.4	144 0.28	890.	.466	0.296	0.466	0.278	0.530	0.293	0.44	<u>45 0.</u>	.278	1.086	0.602	0.701	0.384	0.953	3 0.612	0.70	0.396	0.791	0.431	0.924	0.516	0.638	0.344	0.633	0.340	7.939	ł
_ <u></u> 3	36 0.	433	0.283	0.4	145 0.28	<u>86</u> 0.	.482	0.304	0.482	0.293	0.558	0.305	0.4	69 O.	.305	0.850	0.470	0.843	0.460	0.691	0.432	0.69	9 0.379	0.921	0.495	1.112	0.622	0.636	0.340	0.636	0.326	-1.139	7
7	20 <mark>0</mark> .	407	<u>0.302</u>	0.4	<u>+81</u> 0.3	DU 0.	.514	0.522	0.501	0.314	0.589	0.328	0.5	15 0.	.325	1.244	0.675	0./89	0.430	0.785	0.488	0.77	0.422	0.884	0.487	1.112	0.612	0.675	0.359	0.665	0.343	3.70%	e

A.4 CODE OF ETHICS

We have read and understood the ICLR Code of Ethics, as outlined on the conference website. We fully acknowledge the importance of adhering to these ethical guidelines throughout all aspects of my participation in ICLR, including paper submission, reviewing, and discussions.