

1 We'd like to thank the reviewers for their careful reading and valuable comments. First, we want to emphasize that the
 2 novelty of the proposed method, which addresses how to embed continuous time to differentiable functional domain
 3 that works with self-attention, comes with substantial theoretical derivations and superior practical performances. We
 4 believe this is a useful contribution for both functional representation learning and deep temporal sequence learning.

5 Second, we apologize for typos, grammar mistakes and unclear notations. They will be corrected in the final version.
 6 In Fig 1a (1c), the *SEARCH/ATC/VIEW/TRX* in legend stands for user actions of *search/add to cart/view/purchase*
 7 in online shopping. It shows that by combining time and event representations, the model captures useful time-event
 8 interactions, e.g. the products searched by a user gets higher attention weight when recommending the next product.

9 Third, we provide additional experiment results in Table 1. Previously, *Mercer* time embedding uses $k = 30$ as the
 10 degree for Fourier basis under each frequency in all experiments. We now treat k as a tuning parameter and report the
 11 best performances. For the *Bochner Inv CDF* method, we employ two SOTA flow-based inverse CDF learning methods,
 12 i.e. masked autoregressive flow (*MAF*) and non-volume preserving (*NVP*) transformation, in addition to *MLP*.

Dataset	Stack Overflow	Movielens		eCommerce	
Metric	Accuracy	Hit@10	NDCG@10	HIT@10	NDCG@10
Mercer	46.53(.20), [k=10]	82.92(.16)	60.88(.11), [k=5]	14.94(.31)	12.81(.22), [k=25]
Bochner	40.47(.65), [MLP]	81.60(.65)	60.60(.53), [MLP]	9.84(.86)	7.95(.94), [MLP]
Inv CDF	42.01(.39), [MAF]	82.52(.36)	60.80(.47), [MAF]	12.77(.65)	10.95(.74), [MAF]
	42.13(.38), [NVP]	82.37(.38)	60.55(.50), [NVP]	12.38(.68)	10.62(.73), [NVP]
PosEnc	41.09(.33)	82.45(.31)	59.05(.14)	12.49(.38)	10.84(.26)

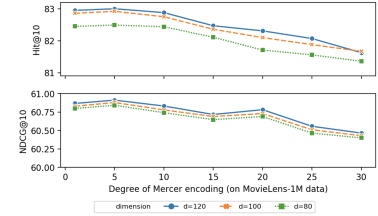


Table 1: Additional experiment results (converted to percentage by multiply-
 ing by 100). The model configurations are reported in the square brackets.

Figure 1: Sensitivity analysis on degree k under
 different time embedding dimensions.

13 **To reviewer #1. Q1: Why Mercer fails to outperform in Movielens dataset?** With the tuned Fourier basis degree k ,
 14 Mercer’s method consistently outperforms others across all tasks. While d , the dimension of time embedding, controls
 15 how well the bandwidth of $[\omega_{\min}, \omega_{\max}]$ is covered, k controls the degree of freedom for the Fourier basis under each
 16 frequency. When d is fixed, larger k may lead to overfitting issue for the time kernels under certain frequencies, which
 17 is confirmed by the sensitivity analysis on k provided in Figure 1 for the Movielens dataset.

18 **Q2: The insight/explanation/analysis of why other 3 methods do not perform well.** After employing the SOTA
 19 inverse CDF learning methods, *Bochner Inv CDF* achieves better performances than positional encoding and other
 20 baselines on Movielens and eCommerce dataset. This suggests the importance of having higher model complexity
 21 for learning $p(\omega)$ under Bochner’s Thm, which also explains why *Bochner Normal* fails in most cases since normal
 22 distribution has limited capacity in capturing complicated distributional signals. On the other hand, *Bochner Non-para*
 23 is actually the special case of *Mercer*’s method with $k = 1$ and no intercept. While Bochner’s methods originate from
 24 random feature sampling, Mercer’s method grounds in functional basis expansion. In practice, we may expect Mercer’s
 25 method to give more stable performances since it does not rely on distributional learning and sampling. However, with
 26 advancements in Bayesian deep learning and probabilistic computation, we may also expect *Bochner Inv CDF* to work
 27 well with proper distribution learning models (as we have shown above with the flow-based methods).

28 **To reviewer #2. Q1: Connections to other tasks that would make the proposed approach more generally appli-
 29 cable?** Besides recommender systems, temporal sequence learning has wide applications in clinical trial analysis,
 30 temporal network representation learning and reinforcement learning. The idea of functional representation learning, on
 31 the other hand, is not restricted to temporal sequence learning. The use of positional encoding with self-attention in NLP
 32 tasks requires a fixed input sequence length and a large number of parameters when sentences are long. The proposed
 33 approach can be easily adapted to functional position encoding, which does not suffer from the above drawbacks.

34 **To reviewer #3. Q1: Relations to related position/temporal encoding methods?** To the best of our knowledge, the
 35 other positional/temporal encoding methods either do not handle continuous time or are driven by heuristics (we will
 36 include the missing references). The original fixed positional encoding also takes the form of sinusoidal functions and
 37 can be thought of as a special version of *Bochner Non-para* method. While their frequencies are mostly chosen by
 38 insights, in our approach, the frequencies are either free model parameters or sampled from learnable distributions.

39 **Q2: Motivation for capturing the relationships in dot product.** The motivation mainly comes from the key-query
 40 inner product formulation of the self-attention model in (1).

41 **Q3-Q7:** In the additional experiments, the flow-based inverse CDF learning methods give much better performances
 42 than the three-layer *MLP*, which was originally chosen for illustration purposes. Given the limited space here, the
 43 detailed implementation for comparisons models will be described in full in the final version. An algorithm description
 44 (flow chart) will also be added to the appendix. We will address all remaining suggestions in the final version.