

## 1 All reviewers

2 We would like to thank all the reviewers for their positive assessment of our work. We are also grateful for pointing out  
3 the typos and providing writing suggestions — we will incorporate them in the updated version.

## 4 Reviewer 1

- 5 • **Experimental results.**

6 We would like to politely disagree with your statement that “the proposed model seems to underperform the  
7 existing and simple baseline models on most datasets”. In the density estimation experiments (Tables 1 & 2),  
8 our method achieves the best or the second best score in 19/26 cases. RNN is the only model besides TriTPP  
9 that consistently achieves good results across all the datasets. In the NLL experiment TriTPP dominates the  
10 simpler baselines (IP, RP, MRP models) on Hawkes1, Hawkes2, SC, PUBG, Reddit-C, Reddit-S, Twitter  
11 (7/13), and achieves nearly identical scores on the remaining datasets (6/13). Similarly, we convincingly  
12 outperform the Hawkes model on Hawkes2, SC, IPP, MRP, RP, PUBG, Reddit-S, Yelp1, Yelp2 datasets (9/13),  
13 get slightly better scores on Taxi and Twitter (2/13) and get worse scores only on the Hawkes1 and Reddit-C  
14 datasets (2/13).

15 We would also like to reiterate that the parallelism of our model is beneficial not only for faster training on  
16 long sequences. Rather, fast parallel sampling opens new applications for TPPs that were completely off the  
17 charts for existing approaches (where sampling can be hundreds of times slower).

## 18 Reviewer 2

- 19 • **Choice of the block size.**

20 We found that increasing the block size beyond 16 for the TriTPP model did not improve the performance.  
21 In contrast, the RNN did benefit from larger hidden sizes (32 or 64). Therefore, we allowed the RNN to  
22 have more parameters in the density estimation experiment (Section 6.2), not to give an unfair advantage to  
23 TripTPP. In the scalability experiment (Section 6.1) we made sure that both models have approximately the  
24 same number of parameters.

- 25 • **Additional nonlinearities.**

26 We found that stacking the nonlinearities (splines) did not improve the performance on the validation set. One  
27 explanation for this can be that the splines do not benefit from depth, as their flexibility can be increased by  
28 simply increasing the number of knots.

- 29 • **Using SGD & tuning the optimization procedure.**

30 Thank you, we will run more experiments to see if this will allow us to improve the performance.

- 31 • **More examples & intuition for TPPs and triangular maps.**

32 Thank you for the excellent suggestion — we will add more examples as well as the following discussion to  
33 help the readers build intuition on this topic.

34 The function  $\Lambda$  captures the global trend by rescaling the time (e.g., if there are more events during the weekend  
35 than on a weekday). The pairwise difference matrix  $D$  computes the “rescaled” inter-event times after the  
36 global trend component  $\Lambda$  has been removed. The block-diagonal matrices  $B_1, \dots, B_L$  capture the interactions  
37 between events, which allows our model to capture self-exciting or self-correcting behavior. The function  
38  $\Phi$  models the distribution of the “raw” inter-event times (this may correspond to a recurring pattern, as in a  
39 renewal process). Finally, the cumulative sum matrix  $C$  sums up the transformed inter-event times to obtain  
40 the arrival times of the base unit Poisson process.

## 41 Reviewer 3

- 42 • **Sampling from a Poisson process.**

43 We would like to point out that we can sample from a unit-rate homogeneous Poisson process in parallel as  
44 following (using Pytorch as an example):

```
45 tau = torch.empty(N).exponential_(1.0) # draw inter-event times from Expo(1)
46 t = tau.cumsum(dim=-1) # obtain HPP arrival times
```

Both steps can be done in parallel in  $O(N)$  operations, including the cumulative sum (Blelloch, 1990. “Prefix sums and their applications”).