

1 We thank the reviewers and are happy that they found the solution simple (R1,R4) and the empirical results promising
2 (R1,R2,R3,R4) in a field of high relevance to the community(R3).

3 The questions raised by R3 are valid and we thank R3 for raising them since it would aid in making the text better. We
4 thank R3 for pointing us to the references of stiff ODEs in the field of numerical analysis. We would definitely add
5 these as references. The experimental setting for section 4.4 was similar to the one used in the examples/ode_demo.py
6 file of the official github repository of neural ode. Instead of using a 2-D ODE we used a 1-D ODE to be able to better
7 visualize and analyze the behavior. We used Eq 12 in the paper as mentioned in the numerical methods book by Chapra
8 et al on page 752. It was surprising that a Neural ODE was able to model a 2 dimensional model with ease while the
9 stiff equation was extremely hard for the Neural ODE to model. We agree that we have barely scratched the surface in
10 terms of stiff ODEs. We had mentioned in lines 287-291 explaining that it was an empirical observation regarding the
11 failure of neural ODE based models where STEER was able to somehow work better. The text wanted to highlight the
12 problems of stiff ODEs and demonstrate that it could become a real challenge moving forwards as the area of Neural
13 ODE matures. We would be open to removing the section until we understand more about the reason why STEER
14 performs better in those circumstances. As R3 points out the method still would have merits without the stiff ODE
15 section.

16 The concerns raised regarding the proof is also valid on closer inspection. The proof, though valid does not add much in
17 the discussion of the proposed STEER method since T is only drawn after a theta-parameter update has been conducted.
18 Our rationale was based on the fact that Picard's iteration enables us to prove existence of a unique solution for Neural
19 ODEs and the modified Picard's iteration enables us to show the existence of a unique solution for one with changes in
20 the end time. However, modifications to the Picard's iterates on closer inspection are not the exact same modifications
21 introduced by STEER due to the random sampling involved. We would be open to removing the discussion of the proof
22 from the paper.

23 As R4 pointed out, the eigenvectors can be complex or real. The ratio of the magnitudes of the real parts of the
24 eigenvalues defines the stiffness ratio. We will mention dataset details as pointed out for Tables 1 and 5. The units
25 of time of table 1 is hours. The sentence can be better phrased as RNNs have also been used for irregularly sampled
26 timeseries models with good results.

27 R4 and R3 have pointed out the comparison to fixed depth solvers such as resnets. Extensive comparisons between
28 resnets and ODEs with similar number of parameters were made in Neural ODE Chen et al. (2018) and Augmented
29 Neural ODE Dupont et al. (2019). STEER doesn't provide a huge performance improvement in terms of accuracy over
30 these 2 techniques. We would however include the results for a fair comparison. With respect to the optimum values of
31 b as R4 and R3 have pointed out. We have shown some experiments in the supplementary regarding the values of b for
32 optimal performance. The general observation has been that the greater the value of b the better the performance. In
33 case of feedforward models with fixed limits of integration $[0, 1]$ the best performance was obtained when $b = 0.99$. In
34 the case of generative models with continuous normalizing flows, however there were some numerical issues when we
35 tried values of b as high as feedforward models. The best results while also ensuring numerical stability was obtained
36 with $b = 0.5$ with fixed limits of integration $[0, 1]$. In case of irregularly sampled timeseries models, since the points are
37 irregularly sampled, hence the limits of integration change every time. In such a case rather than explicitly tuning b we
38 tune the ϵ parameter as mentioned in line 207 the main text of the paper. Similar to the observations in the preceding
39 cases, the best results are obtained using b as high as possible. In this case it implies ϵ as low as possible without
40 numerical instability. The best results were obtained with $\epsilon = 0.05$.

41 As R1 has rightly pointed out, the integration times are
42 arbitrary in the case of Neural ODEs. The dynamics of
43 these models rearrange to similar qualitative behaviors
44 even when the integration times are varied by instance.
45 We observe that STEER changes the behavior of these
46 models in qualitatively meaningful ways as shown in
47 Figure 1. The principled work by Finlay et al. (2020) is
48 based on concepts of optimal transport. We have shown
49 via experiments on generative modeling (Table 1) that
50 our technique can be used to augment their technique to
51 obtain even faster convergence times.

52 R2 points out the rationale behind the increase in the
53 performance in case of timeseries models and feedforward models. Previous studies Dupont et al. (2019) have
54 corroborated that simpler dynamics leads to better solutions. The temporal regularization is intended to improve the
55 dynamics for neural ODE flows as demonstrated via multiple experiments. It could include unintended improvements
56 in accuracy that we observe in Tables 2,4 and 5 and Figure 1. We intend to explore these aspects in future work.

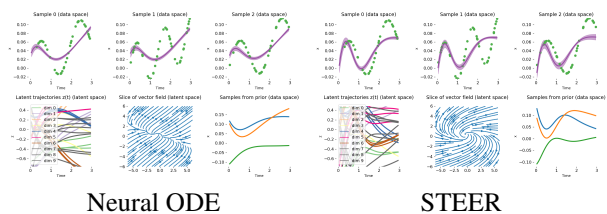


Figure 1: A simple addition of t to $\sin(t)$ changes the qualitative performance of Neural ODE vs STEER.