

1 We thank all reviewers for their thorough assessment of our paper and suggestions for improvement.

2 **Response to reviewer #1**

3 • **Experiments, clarity and additional feedback.** Thank you for your insightful feedback. We agree with you and in  
4 the revised paper, we will add an ablation study of RNNs and GRV, implement the combination of benchmarks you  
5 suggested, and add more simulation studies to provide a more complete evaluation of our method. We note though that  
6 RNNs are not necessarily superior, because the models in other methods can also capture the data generating functions.  
7 Finite sample regularization improves the estimation because it adapts all the initial models to solve the nonparametric  
8 estimating equation (EE), rather than after the initial models are trained and potentially overfits to the training data.  
9 Note that validating causal models is tricky because we do not observe the counterfactuals. We will definitely improve  
10 the paper clarity in the revision. As you suggested, we will expand the Appendix to provide a notation table, add several  
11 sections introducing the needed background from the various fields such that the paper becomes more accessible and  
12 our derivations are easier to follow. Thank you for the excellent suggestions.

13 • **Relation to prior work.** Both GRV and Targeted Regularization (TR) [Shi, et al.] are motivated by solving the EE.  
14 However, note that the EIC of DTR in GRV is much more complicated than the EIC of ATE in TR. The EIC of ATE  
15 only contains two outcome models and one propensity score model. The EIC of a DTR has many terms as shown  
16 in Equation (6-8). This is because the value function estimation involves the DTR, the outcomes, and the treatment  
17 assignment mechanism over time. We need to solve the EE by adapting a sequence of outcome and propensity score  
18 models over time. There is no direct implementation of TR which can solve the EE across time steps jointly in the DTR  
19 learning setting. The GRV regularizer is designed to simultaneously solve the complex EE of DTR learning, which  
20 involves models at different time steps; this requires a more technical derivation than TR. In addition, re-estimating  
21 nuisance models in DTR learning is more challenging than in ATE estimation. The current related work section gives  
22 an overview of the estimator classes and relegates the discussion of the longitudinal TMLE for OPE to the Appendix.  
23 In the revision, we will make the connection with TR and TMLE clear in the main body of the paper. Thank you.

24 **Response to reviewer #2**

25 Thank you for the useful suggestions. We agree and we will expand the Appendix to include a discussion of advantages  
26 and disadvantages as compared to the other methods, including the running time, stability etc. Thank you.

27 **Response to reviewer #3**

28 • **Motivations and clarity.** Thank you for your comments. We agree with you and we will improve the presentation  
29 of our method and the adopted architecture in Section 4.3. Due to space limitations, we already relegated some of  
30 the discussion about the network architecture to the Appendix. The objective function in Equation (9) has two terms,  
31  $\mathcal{L}_t$  and  $\mathcal{R}_t$ . The term  $\mathcal{L}_t$  includes the standard loss function for the outcome regression models and propensity score  
32 models in causal inference. In DTR learning, we also need sequential regression to estimate the long-term outcomes;  
33 this is why we have the last term in  $\mathcal{L}_t$ . The regularizer  $\mathcal{R}_t$  is motivated by Theorem 1. Minimizing  $\mathcal{R}_t$  encourages the  
34 fluctuated outcome models to solve the optimal score equation i.e. nonparametric estimating equation, on the R.H.S of  
35 Equation (10), so that the resulting estimator is efficient for estimating the value function. In the revised manuscript, we  
36 will improve the presentation in Section 4.3 and explain clearly the motivation of each term in our objective function.

37 • **Doubly robust (DR) and TMLE.** The DR estimator achieves efficiency by applying the DR estimating functional  
38 on some pre-trained initial outcome and propensity score models. An example of DR estimator is given in Equation  
39 (4), where the inverse propensity score (IPS) is multiplied by the cumulated outcome and  $Q$  functions. In GRV, the  
40 initial estimators are optimized to solve the nonparametric estimating equation jointly with an additional parameter  
41 epsilon; this provides additional flexibility to solve the estimating equation. Our estimator in Equation (11) is stable  
42 because the IPS is multiplied by a small epsilon. Similarly, it is well known that solving the estimating equation using  
43 an additional epsilon in TLME can give a more robust estimator than DR in ATE estimation. The parameter epsilon  
44 enables trading-off between finite sample stability and asymptotic efficiency. However, TMLE is not robust when the  
45 initial models are not well specified or estimated. This drawback exists even in the simplest setting of ATE estimation  
46 and becomes even more severe in DTR learning which involves many outcome and propensity score models. Our  
47 method adapts all the initial estimators and the epsilon to solve the estimating equation jointly by regularizing the  
48 models in training. By reducing the dependence on the quality of the initial models, our method is thus more robust  
49 than TMLE for DTR learning. The connection between our method and TMLE is discussed above Section 4.3 and in  
50 the Appendix. We will further improve the discussion and make these points clear in the revision. Thank you.

51 **Response to reviewer #4**

52 Thank you for your useful feedback. In the submitted paper, we have included two simulation studies to evaluate the  
53 method performance (in terms of both DTR evaluation and optimization) across various sample sizes. We note that,  
54 unfortunately, it is impossible to use offline datasets to benchmark the methods; existing DTR or policy learning papers  
55 are also based on simulation studies. We use the MIMIC III dataset for illustration only, to show that the learnt DTR  
56 will lead to gradually decreasing the treatment, which corresponds to the clinical literature. The real DTR application  
57 can only be demonstrated well through real-world experiments and simulation. In a real-world test-bed, the clinician  
58 can check if his/her treatment decision agrees with the DTR before prescribing it. In the revision, we will also explain  
59 how such test-beds can be built to evaluate DTR. Thank you.