

1 We would like to sincerely thank all the reviewers for reading the paper carefully and their very valuable feedback.

2 **General comments:**

3 **Visualization of vanishing/exploding gradients:** Similarly, as in Fig. 4 from [2], we plotted  $L_2$ -norms of the gradients  
 4  $\frac{\partial \mathcal{L}}{\mathbf{x}_t}$  as a function of time  $t$  at the beginning of the optimization and after 100 iterations (first two plots in Fig. 1). We  
 5 see that the norm of the gradient for ODEtoODE barely changes, while for the NeuralODE it converges to 0 as we  
 6 backpropagate through time (we observed convergence to 0 also for other methods, yet for the clarify of the picture we  
 7 did not present additional curves). The plots were created for Humanoid training from the paper, where ODEtoODE  
 8 was **the only method** that successfully trained the agent.

9 **Presentation:** in the final version we will simplify Sections 1-4 and provide definitions of ISO/Gated-ODEtoODE  
 10 in the Introduction. As suggested, we will improve citations' style. In Section 5.2.3, we will provide more detailed  
 11 discussion, and incorporate new results (see: Table from Fig. 1). Training ODEtoODEs does not require any additional  
 12 modules since the parameters governing the evolution of weight matrices are unconstrained and therefore can be handled  
 13 by standard backpropagation in the supervised setting and standard GD-approach with ES-gradients for the RL one. We  
 14 do agree that this requires clarification, thus in the final version we will incorporate an algorithmic box on ODEtoODEs  
 15 in the experimental section for clarity.

16 **Reviewer 1:**

17 **ANODEV2:** Thank you for pointing this out. The HyperNets networks we compare against in Sec. 5.2 (see: Table 1,  
 18 2) are precisely ANODEV2 architectures (as indicated by our citation [56] to ANODEV2 in 1.265). This is also the case  
 19 for the ES experiments, and thus citation [26] (which is a typo) in 1.238 of Sec. 5.1 should be replaced by [56].

20 **Pros & cons of orthogonality:**

21 Even though in principle deep architectures leveraging orthogonal matrices can hurt accuracy by restricting model  
 22 capacity, a rich line of work ([2,30,40]) on orthogonal RNNs and related methods shows that they do not, if designed  
 23 correctly. In fact, the main motivation behind these architectures are accuracy improvements due to training not suffering  
 24 from vanishing/exploding gradients - one of the critical problems in training deep machine learning systems. We also  
 25 demonstrate it exhaustively in the empirical section, covering both: the supervised and RL setup (the latter usually  
 26 not touched in the papers on the subject). To summarize, ODEtoODEs outperform accuracy-wise other methods (in  
 27 particular highly competitive ANODEV2) in **11 out of 15** considered supervised tasks, as we clearly state in the paper  
 28 (1.272). In the RL setting ODEtoODE is **the only method** that manages to train the most difficult Humanoid task and,  
 29 as we demonstrate in "General comments" section above, it is precisely due to its ability to stabilize gradients' lengths.

30 **Theoretical results:** We are not aware of **any** results in the literature on neural networks (even in the simpler setting with  
 31 shallow discrete-architectures & supervised learning) on the convergence to global minimum under weak assumptions  
 32 from our Theorem 1 and to the best of our knowledge we are the first to show convergence to local minima with  
 33 depth-independent bounds in the more challenging ES scenario. Given that, we interpret our theoretical results as the  
 34 strength of the paper. Lemma 4.1 is implicitly used in the proof of our Theorem 1 (see: Eq. 21 in the Appendix) and  
 35 techniques used to prove it (Sec. 8) play key role in establishing depth-independent bounds. In the final version of  
 36 the paper we will explicitly refer to it. Lemma 4.1. is also a natural continuous extension of the analogous results  
 37 for the orthogonal RNNs (see for instance: [2]), so can be used in various analogous continuous variants, not only  
 38 ODEtoODEs. We will clarify it in the final version.

39 **ResNets in supervised setting:** We run these additional experiments and present results in Table from Fig. 1.

40 **Reviewer 3:**

41 Thank you very much for the review.

42 **Mathematical notation:** we will incorporate suggested improvements and clarify in the final version (in particular, by  
 43 "strong convergence results/guarantees" we mean "compelling"). **L.102:** This is a generic notation with parameters  
 44 encapsulated in  $\psi$ . We will clarify this in the final version.

45 **Visualization of gradients' norms:** See, our paragraph on visualization in "General comments" section.

46 **Reviewer 4:**

47 Thank you very much for the review.

48 **Visualization of gradients' norms:** See, our paragraph on visualization in "General comments" section.

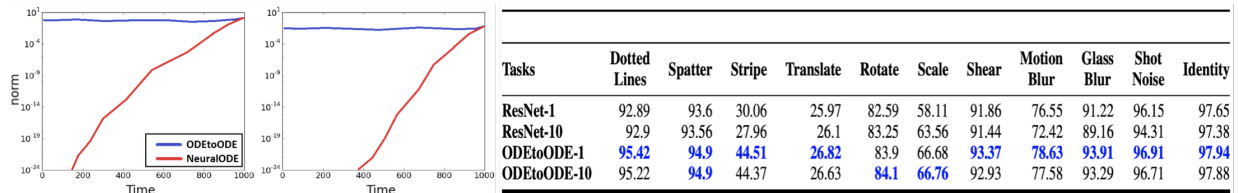


Figure 1: Showcasing gradient vanishing & extra experiments with ResNets. First two plots present norm of the gradient  $\frac{\partial \mathcal{L}}{\mathbf{x}_t}$  as a function of  $t$  for Humanoid training after 0 (first plot) and 100 iterations (second plot). Time = 1000 corresponds to  $T = 1.0$  from the paper. Table presents additional results on ResNets and comparison with ODEtoODEs. **Bolded blue** correspond to best results and **they are all** for variants of ODEtoODEs.