

1 We thank the reviewers for their comments and suggestions and for finding our theoretical considerations insightful and  
2 solid [R1, R3, R4], our experimental results valuable [R2, R3, R4] and our paper well-written [R1, R3, R4].

3 **About multi-step rollouts** We appreciate the suggestion of R1 and R3 on multi-step rollouts. We are curious as well  
4 but we opted to leave it to future work: the goal of the paper is to show that model-based methods can be used to  
5 directly learn the action-gradient of the value function, and we intentionally designed the simplest algorithm possible to  
6 show this point. Using multi-step methods for computing either the gradient of the TD-target or the value gradient itself  
7 is orthogonal to what we propose in MAGE, and would introduce additional complexity and hyperparameters (e.g.,  
8 rollout length); we believe this would distract the reader from the main contribution of the paper. Nonetheless, we agree  
9 that combinations of MAGE and these approaches can be particularly fruitful.

10 **R1:** \* The setup of MAGE follows common practices in (MB)RL: ensembles of models and the swish activation have  
11 been used since PETS [7], and, more recently, in (SotA) MBPO [20], whose training scheme is very similar to ours;  
12 the Huber loss has been shown to improve stability in DQN (2015). Those choices are known to be effective, but any  
13 reasonable setting is enough for MAGE to perform well (see, e.g., the robustness of MAGE to the choice of  $\lambda$  in Fig. 7  
14 in Appendix B.3 and the additional experiment with different model architectures in Fig. R1). Therefore, we did not try  
15 to cherry-pick the settings nor did any overly extensive hyperparameter search.

16 \* Concerning the experiments on the unknown reward function setting, notice that Fig. 6 (in Appendix B.2) contains  
17 results for all the tasks. We only included a single environment (Pusher-v2) in the main paper in order to save space.

18 \* We indeed show the asymptotic performance on HalfCheetah-v2 only since our focus is on sample-efficiency;  
19 nonetheless, we plan to include a larger number of steps (also for Swimmer-v2) in the final version of the paper.

20 \* We will include the suggested references into the paper. See also **About multi-step rollouts**.

21 **R2:** (1) The reviewer suggests that the paper should first "show that minimizing the TD-error is not  
22 a good choice". Notice, however, that despite being commonly used and thought of as "intuitive",  
23 the minimization of the TD-error in DPGs is just a heuristic and lacks proper theoretical justification ever since its  
24 first use in [42]. In contrast, our Proposition 3.1 inspires a reliable theoretically-grounded objective for learning a critic.

25 Furthermore, Fig. 3 shows indeed that minimizing the TD-error can lead to a critic being far away from the ideal one.  
26 (2) We find no contradiction in our line of reasoning: we state that, despite the objective suggested by Proposition 3.1  
27 is theoretically sound, there are practical issues that make the optimization problem hard. As written in the paper, we  
28 interpret the TD-error as a regularization term with almost no additional overhead, but we do not exclude the existence  
29 of other practical strategies for minimizing the bound in Proposition 3.1 without the use of the raw TD-error. We agree  
30 that in MAGE the minimization of the TD-error is also important but we clearly showed with our experiments that the  
31 overall objective is better than the TD-error alone.

32 (3) We did not write that "model-based RL has no advantage in terms of sample-efficiency than model-free RL". Instead,  
33 we said that this advantage is not *intrinsic* and it is instead "highly environment and algorithm-dependent". We find  
34 this statement uncontroversial in the light of other recent works in MBRL (see, e.g., [20,52]).

35 **R3:** \* We carefully positioned our approach in the literature and discussed similarities and differences w.r.t. most of  
36 the works listed by R3. In particular, we directly cite Balduzzi et al. (2015) for proposing a solution to the problem of  
37 value gradient learning. Notice, however, that MAGE provides a completely different and more direct way to address it.  
38 Moreover, in our Related Works section, we acknowledge Werbos who pioneered (in 1977!) value gradient learning,  
39 preceding Balduzzi et al's, Fairbank's and Weber et al's work by decades. Having said that, notice that these works lack  
40 our theoretical motivation and belong to significantly different algorithmic frameworks. Likewise, Sobolev training  
41 and related concepts in value gradients mainly share *technical tools* with MAGE, but they have different algorithmic  
42 implications. Therefore, we are truly downhearted by the assertion regarding the alleged limited novelty of MAGE.

43 \* Concerning the gradient stopping, we indeed use it: see Algorithm 1, in which target parameters  $\phi$  are employed  
44 during the computation of  $\hat{\delta}$ , a convenient notation for the gradient-stopping operation on the target  $r(s, a) + \gamma Q_\phi(\hat{s}, a')$   
45 and its gradient w.r.t. the action  $a$ . Nonetheless, as suggested, we will make the gradient stopping more explicit.

46 \* See **About multi-step rollouts**.

47 **R4:** We agree with the reviewer that the quality of the learned model is of paramount importance for most MBRL  
48 algorithms, whose performance, generally, deteriorates when the model is not enough expressive for a given task.

49 Thus, as suggested, we performed an additional experiment to investigate how the  
50 performance of MAGE (compared to the Dyna-TD3 baseline) is affected by the use of  
51 less powerful models. We evaluated two versions of the model with reduced capacity: (i)  
52 only 2 members in the ensemble, 4 hidden layers and 256 units per layer (**-small suffix**);  
53 (ii) no ensemble (a single model), 2 hidden layers, 256 units per layer (**-smaller suffix**).  
54 Recall that the original setting involves a more powerful model with 8 members in  
55 the ensemble, 4 hidden layers, 512 units per layer (no suffix). The results, reported in  
56 Fig. R1, show that MAGE is robust to the presence of a misspecified model: while a  
57 simpler but still quite capable model does no harm to MAGE, a significantly smaller  
58 model has a reasonable impact on the obtained average return. We will report results  
59 for all environments in the final version.

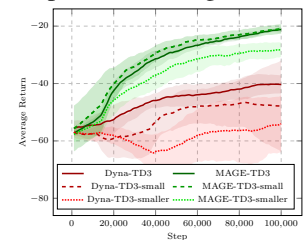


Figure R1: Pusher-v2 (5 runs, 95% c.i.).