

1 We thank all of the reviewers for their useful and insightful feedback. We are encouraged they found our work interesting
2 (R1, R2, R3, R4), novel (R2) and scientifically rigorous (R2, R3). We are happy all reviewers recognize the importance
3 of the problem we address. We are glad they found our work impactful in the design of new algorithms (R1), that our
4 empirical studies were convincing (R1, R2, R3) and that they proactively answer their questions (R3). We address
5 reviewer comments below and will incorporate all feedback.

6 @R1 "clear examples to explain the concepts": Our work stands on the intuition that an agent might benefit from using
7 its computation budget to spread information as fast as possible through the state space to all precursors of the current
8 state, rather than just the one state and action that brought it to the current state on this occasion. An **analogy to human**
9 **memory consolidation** may help: we are arguably good at retrospectively inferring (or inventing) the underlying
10 causes of our experience and ruminating over them so as to improve our future predictions; *we hypothesize agents*
11 *could also use retrospective knowledge about the world (backward models) to consolidate (hindsight planning) their*
12 *prospective knowledge (value functions)*. We will complement the **thought experiment** in Section 3 with this. Thank
13 you for the suggestion. @R1 "not clear what an abstract model is": we define a model as being abstract when we
14 remove inductive biases (e.g. structural constraints) in its construction and use, leaving it as a (learnable) black box. We
15 will expand in the paper.

16 @R1,R2 on related work: The closest work to ours is that of van Hasselt et. al., our results complement theirs: in
17 the control experiment we extend by showing backward planning is *more robust in dealing with rare events (e.g.*
18 *stochastic rewards) and different levels of stochasticity in the transitions;*@R2 van Hasselt et. al. compare against
19 replay – a non-parametric model-based approach, NOT against model-free learning. Goyal et. al. use imitation
20 learning on a generative model’s outputs so as to improve exploration by incentivizing the agent towards the high value
21 states on which the model was trained; in contrast, *we aim to formalize and tease apart the fundamental properties of*
22 *online hindsight planning*. Satija et. al. use backward value functions to encode constraints for solving constrained
23 MDPs (CMDPs) with safe policy improvements; though both our work and theirs employ some form of retrospective
24 knowledge, *the contents, purposes and uses differ*. @R4 we will add the missing citations and discuss how they relate
25 to our work. Thank you for spotting them.

26 @R3, R2 on explanatory details: **Details on experimental settings are in appendix D**. The prediction experiment is
27 run on a leveled state space s.t. transition dynamics between states generate bipartite graphs; we vary the no. of states
28 on each level and the no. of levels to generate different structural properties. We refer to **fan-in/fan-out** as *the no. of*
29 *predecessors/successors a state might have in the state space*. Thumbnails depict the phenomena of transitioning from
30 a larger no. of predecessors that "funnel" into a smaller no. of successors, and vice-versa. **Task** \equiv *value prediction*. We
31 will make the main paper more clear by adding more details from the appendix. Thank you for highlighting this.

32 @R2 *utility of forward planning when the future is predictable*: The claim rests on the prediction experiment; the
33 deterministic control setting is more conflated, performance difference too marginal, task too easy and backward
34 models are also in their best regime. @R2 *codebase*: Missing requirements can be installed from pip and the *import*
35 *bw_PAML_gt* safely deleted. @R2, R3 *backward planning being less harmful*: Great question! An erroneous forward
36 model that predicts an unreachable state will move the value of a real state towards the arbitrary value of the bootstrapped
37 state; an erroneous backward model will harmlessly distribute the value of a bootstrapped real state to an unreachable
38 arbitrary-value state (see concurrent work Jafferjee et. al. - "Hallucinating value..."). @R3 "tethered to a policy": Both
39 $\overleftarrow{p}_\pi(x_t, a_t | x_{t+1})$ and $\overleftarrow{p}_\pi(x_t | a_t, x_{t+1})$ are backward models. Eq (3) shows how the latter still **depends on π through**
40 **the stationary distributions η_π** . @R3 "using Bayes rule" **Yes** (see appendix A). @R1,R2,R3 we appreciate pointing
41 out the misplaced arrow in Fig.1’s legend, detailed feedback on typos and suggestions on how to improve clarity.

42 @R1,R2,R3,R4 "deep RL experiments": Deep RL comes with **confounding factors**, e.g. the environment-dependent
43 mixing time characterizing the correlated datastream forces the use of replay buffers to decorrelate experience and
44 target networks to stabilize learning for incremental-update algorithms; this in turn pushes learning in the off-policy
45 regime where convergence is only serendipitous; common testbeds are not informative of how different components of
46 algorithms influence learning. We have left as future work extensions to function approximation and more complex
47 problems, as we reckon (substantial) additional research is required in terms of testbeds and ablations to allow for
48 **scientifically relevant hypotheses**. Backward models came with their own peculiarities and available choices in terms
49 of estimation and use. Understanding and formalizing them was, in our view, **the first step in the right direction**,
50 i.e. transferring to more complex problems in a principled way. **Part of our scientific contribution** was this analysis,
51 revealing many options w.r.t model-learning objectives and planning-strategies (Section 3). The **goal** of our empirical
52 studies was to disentangle core properties of these approaches so as to inform on the design of complex testbeds, where
53 we can understand planning methods. We strongly believe that our contributions **set the stage for principled deep RL**
54 **investigations**.

55 We can clarify these points in the paper and strengthen the discussion of existing literature. We feel the technical
56 contribution is sound, surprising (to us at least) and potentially impactful to both scientific understanding and practice.