

1 We thank the reviewers for their highly useful feedback, which we will incorporate in the revised version. To summarize,  
2 the main points are:

- 3 • Highlight relationship with existing representation formalisms
- 4 • Clarify how the ideas can be extended to practical settings

## 5 1 Reviewer 1

6 **Put “Related work” in part 2.** We agree we can improve the connection with other formalisms, we will emphasize  
7 this.

## 8 2 Reviewer 2

9 **Isn’t the first half of Lemma 1 a well known result?** You are right in that the proof makes direct use of existing  
10 linear programming results. There is a slight difference, however, since Lemma 1 is specifically about the value  
11 polytope: the domain of the functional is nonconvex, and the proof requires first an extension to its convex hull. The  
12 statement regarding deterministic policies is also polytope-specific. We will make sure to emphasize what is novel in  
13 the statement.

14 **Consider saying why Lemma 1 is useful.** The usefulness is the interesting property of the value polytope (which we  
15 then make use of). We will emphasize this point also.

16 **Connection to Wasserstein distance and model-based RL.** In that space, we know of Farahmand, Barreto, Nikovski  
17 (2017) and follow-up work. This is an interesting connection, although to the best of our knowledge there is no model  
18 learning equivalent of the value polytope or adversarial value functions. E.g. we would need a small set of models w.r.t.  
19 which we want low modelling error. Please let us know in your revised review if there are specific papers that might  
20 have a closer connection.

21 **Can you clarify why you used model-based algorithms in experiments?** Model-free experiments require dealing  
22 with stochasticity in the results, error bars, etc. and would give a murkier picture of the role of the representation. We  
23 could learn AVFs using Monte-Carlo samples from the different adversarial policies, and learn these policies through  
24 sample-based policy gradient, but we feel our setup gets to the point more clearly. We will highlight this.

25 **The experiments are preliminary / challenges when going to function approximation.** To clarify, our experiments  
26 use tabular information (e.g.,  $V^\pi$ ) to produce function approximation, so we read this comment as “how to learn a  
27 representation when tabular information is not available”.

28 The path to a “deep” method is clear to us, but not short. There are a few challenges to overcome: 1) how should we  
29 represent adversarial policies? 2) what is the effect of non-uniform distributions on the learned representation? this  
30 effectively changes the norm in Equation 1; 3) what is the effect of changing the sampling distribution for  $\delta$ ? and 4) we  
31 would like to make use of Bellman updates to learn AVFs, but these requires a good off-policy learning method. 1–3  
32 might be answered by considering the auxiliary tasks perspective (Section 3.2). We will add a section discussing these  
33 points.

## 34 3 Reviewer 3

35  **$\phi$  becomes part of the computation through its gradient.** This is a fair point. At first glance the difference in our  
36 perspectives seems to be about how the optimization process discovers the representation, rather than the space of  
37 possible representations, but if you are making a more specific distinction please let us know in the revised review.  
38 Either way, we will make sure to discuss this point.

39 **Figure 1 (Right)** is a cartoon version of a value polytope for a 2-state MDP. The axes are a good suggestion.

40 Thank you for catching bugs in the math, we will fix.