

1 Thank you for the thoughtful and careful reviews. We hope the AC nominates some of you for reviewer awards. We
2 start with common concerns and then respond to individual reviewer comments as space permits:

3 • **Common:** *There should be a baseline using MCTS and assuming access to simulator / common random numbers.*
4 There appears to be some imprecision in reviews about what this means. DirPG is not a planning algorithm, because
5 it does not use search (or a simulator) at test time, so standard MCTS / UCT at test time isn't applicable (L318-319).
6 Nevertheless, we designed an algorithm that uses MCTS-style tree search and a simulator to learn a policy similarly
7 to DirPG: A policy is initialized randomly. The algorithm explores k trajectories using UCT tree search with common
8 environment random numbers. After the k trajectories, one best trajectory is extracted by running without the
9 exploration bonus, and that trajectory is "distilled" into the policy by performing a gradient update to increase its
10 probability. Then environment stochasticity is re-sampled and the algorithm repeats.

11 This uses the same interface as DirPG, can be made to use the same number of environment interactions and gradient
12 updates via choosing k , and it makes the same assumptions about simulator and common random number availability.
13 We can make a key theoretical distinction, which is that this algorithm will behave asymptotically like DirPG as the ϵ
14 hyperparameter goes to ∞ , and thus it will be very risk-seeking. To see this, note that given sufficiently large k , UCT
15 will converge to the trajectory with maximal return under the fixed realization of noise shared across the trajectories,
16 which is the same trajectory that \mathbf{a}_{dir} will converge to in DirPG given sufficiently large search budget and $\epsilon \rightarrow \infty$.
17 Thus, DirPG offers a degree of control via ϵ that isn't available to this MCTS-style counterpart.

18 • **R2, R4:** *Surprised REINFORCE works well.* Ah, we should state more clearly that REINFORCE used the same
19 number of environment interactions-per-seed as DirPG. It uses the same common random numbers access as DirPG.

20 • **R1, R2:** *"If a simulator is available, one typically wouldn't use a model-free algorithm like REINFORCE. ... If
21 a heuristic and simulator are available (both typically assumed to be unavailable in RL), why wouldn't one use
22 heuristic search methods?"* There are state-of-the-art RL approaches to real problems (e.g., [1]) where a simulator
23 is available and researchers use REINFORCE. Heuristic search is not a silver bullet. The above work on solving
24 combinatorial optimization problems using RL is based on the premise that there is room for improvement over
25 traditional solvers. A potential future application of DirPG is to gain the benefits of learning in these domains without
26 having to throw out the benefits of heuristic search. We expect this to become increasingly relevant as RL is applied
27 to more non-traditional domains like combinatorial optimization, program synthesis, and sim-to-real.

28 • **R3** *"how much improvement can it bring compared with a perfect heuristic?"* With an uninformative heuristic, the
29 search for \mathbf{a}_{dir} reduces to sampling without replacement via Gumbel Top-K, which may bring some benefit over
30 independent sampling when policies are peaked [2], but generally the heuristic is important.

31 • **R1** *"First, if it is only applicable to tabular cases"* It's not.

32 • **R1,R2** *"Does the method extend to function approximation?"* Yes. See architecture details in Appendix D.3.

33 • **R1** *"to solve combinatorial bandit problems, there are other algorithms such as CUCB"* We agree there are lots of
34 algorithms for RL and combinatorial bandits. We say clearly in the paper that we are not claiming to be state-of-the-art
35 (e.g. L95-96). Please also note that the specific algorithm suggested is very similar to our "full bandit" baseline.
36 Compare Alg 1 in the CUCB paper to Appendix L614-617.

37 • **R3** *"Does DirPG also require the state to be finite?"* Continuous states are fine. Gumbels are applied to actions, not
38 states. Continuous actions may also be possible in future work (A^* sampling extends to Gumbel processes).

39 • **R3** *"gradient formula in l.66 and eq(5)"* Thank you! This is a typo. The LHS in l.66 should be expected reward of
40 the argmax like in the leftmost term in eq 14.

41 • **R4** *"how does this compare to simple finite differences?"* We tried implementing [3] using similar common random
42 numbers access as DirPG but couldn't get it to learn sensibly when given 5x the budget of interactions as DirPG.
43 However, we didn't have time to try all the tricks in the paper (e.g., rank transformations).

44 • **R4** *"this assumption in general seems to be a bit understated in the presentation"* That wasn't our intention. We felt
45 we said it clearly in L121-122 but understand your position and will increase its prominence.

46 • **R4** *"contrastive divergence"* Yes! Very interesting!

47 References

48 [1] Wouter Kool, Herke van Hoof, and Max Welling. Attention, learn to solve routing problems! In *International Conference on*
49 *Learning Representations*, 2019.

50 [2] Wouter Kool, Herke van Hoof, and Max Welling. Stochastic beams and where to find them: The gumbel-top-k trick for sampling
51 sequences without replacement. *arXiv preprint arXiv:1903.06059*, 2019.

52 [3] Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution strategies as a scalable alternative to
53 reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.