



Figure 1: Left to right: all mean, all median, ablation mean, ablation median.

1 **General.** Shortly after submission, we found an unintentional inconsistency in our evaluation protocol. Our M-Agents
 2 and their ablations were evaluated in a slightly different Atari setting from the others (episode terminated on life loss
 3 rather than game over), favoring them on some games. We reran experiments with the correct settings and (human-
 4 normalized) results are on Fig. 1. Despite some drop in absolute performance (eg, M-DQN outperforms C51 but is a bit
 5 less competitive with Rainbow), it does not change the ranking of algorithms (for all metrics) nor our conclusions.

6 **R1. Q1:** In Dopamine, IQN doesn't use double Q-learning (implicit_quantile_agent.py, ll. 192-196). **Q2:**
 7 Studying the homogeneity of the action gap is an interesting research direction that we didn't investigate, thanks for
 8 pointing it out. In the limit ($\epsilon = 1$), it is homogeneous (because infinite), but with numerical instabilities. We'll
 9 add a comparison with log-DQN in Fig. 2, to get at least some empirical insights. **Q3:** We do not see any obvious
 10 connections between Log-RL (related to the h -transform of Pohlen *et al.*) and M-RL (justified here through the lens of
 11 KL regularization), but that's also an interesting research direction. We also think that, even if different, both approaches
 12 could be combined to build an even stronger agent. **Q4:** Yes, exactly, we say why in footnote 6 (and this comparison is
 13 done in the ablation). We'll say it earlier. **Q5:** AL is indeed better than C51 on some metrics. Yet, this was not observed
 14 in [6], probably because the authors use RMS (to compare with the standard DQN) while we use Adam (to compare
 15 with M-DQN). Additionally, note that AL is a special case of M-DQN. We will rephrase appropriately.

16 **R2. Parameters:** To choose the parameters, we did a sensitivity analysis on a subset of games, that we'll add to the
 17 Appx. Yet, note that the parameters were not selected through a simple grid search. Values of β are on par with the ones
 18 provided by the analysis in [30], and γ is the same as the one found optimal in [6]. In terms of "easiness to tune", our
 19 empirical findings suggest that the most sensitive parameter is β , while it is rather easy to find working values for α
 20 and γ . **ALE:** our results are indeed not comparable with the original DQN ones, as discussed ll.217-231. Unfortunately, the
 21 authors used a proprietary version of ALE, so exact comparison to their results is not possible. **FQF:** thanks, we missed
 22 this paper. However, FQF does not use sticky actions and comparison is thus not straightforward. Our method readily
 23 applies to FQF (as for IQN), and we will try to add M-FQF results in the paper. If it is not possible, we will reformulate
 24 to soften the state-of-the-art claim. We believe this does not hinder the relevance of our approach, as M-IQN still
 25 outperforms Rainbow. **Action-gap:** by "quantifying", we mean "analytically quantifying": Bellemare et al. show that
 26 the action gap increases but not by how much, while we derive an actual value for the increase (Thm. 2).

27 **R4.** We trust that R4 has deeply misunderstood our contribution. Indeed, they state that "M-RL applies an entropy
 28 regularizer to the reward signal", and all the following (mostly negative) comments rely heavily on this statement. This
 29 is just wrong and we start the paper by stating otherwise (l.21-24, "We insist right away that this is different from
 30 maxent RL, that *subtracts* the scaled log-policy to *all* rewards [...]" while we *add* it to the *immediate* reward). **On
 31 novelty:** We strongly disagree, our contribution is not simply an instance of entropy-reg RL, as notably thoroughly
 32 discussed in Sec. 3 and related Appx. M-DQN *is different* from Soft-DQN, one just has to compare Eqs. (2) and (3)
 33 (paying attention to the *signs* of the different log terms). We're also absolutely certain that neither our algorithms nor
 34 our analysis are covered by Neu et al. **Most theoretical results are previously known.** Again, we strongly disagree.
 35 We're very clear about what our contributions are, and Thm. 1 and 2 are new (we could reevaluate this claim if a ref was
 36 provided). **Hard to believe** (about the theoretical result): we don't ask to believe our claim, as we provide the proofs,
 37 see Appx A.2 and Cor. 1 and 2 in Appx. A.3. Shortly, the implicit KL regularization avoids the error in the greedy step,
 38 that cannot be avoided when the KL regularization is explicit (the only case considered in [30]). **Clipping:** Obviously,
 39 we do not claim to be the first to clip a log term, but we provide this kind of details for the sake of reproducibility.

40 **R5. Linear case:** in this setting M-VI is strictly equivalent to MD-VI (see Thm. 1, noticing that in the linear case there
 41 is no error in the greedy step), so we refer to [30] for a study on tabular MDPs. Notice that this equivalence is lost with
 42 a non-linear parameterization (due to the necessary error in the greedy step with explicit KL regularization), hence the
 43 interest of the Munchausen principle for *deep* RL. **Comparison to Soft-DQN is missing:** no, it's not, it is provided in
 44 Fig. 3 (ablation). **Unsure about AL fairly compared:** we use $\beta = 0.9$ for both M-DQN and AL, and this choice is
 45 consistent with the AL paper, so we think the comparison to be fair. Also, notice that AL is a special case of M-DQN.
 46 **Montezuma:** this is a hard exploration game, and our method is not designed to favour exploration (there is no signal
 47 to reinforce, a discussion to a related issue is provided in Appx. B.2). **M-DQN** is defined l.31. **Eq. I.82:** this is exactly
 48 DQN, as β is defined to be greedy wrt q . **Bootstrap:** yes, we totally agree with this, it is actually the reason for the
 49 name "Munchausen" (l.18-26).