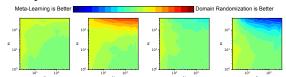1 To all reviewers, thank you very much for your thoughtful comments and suggestions.

2 **R#1 & R#2: Limitations of the linear regression.** Although linear regression has limited application, it removes the
3 optimization error (problem has an analytical solution), and enable us to study statistical and modeling trade-off. Our
4 extensive empirical study on meta-RL complements this and fully describe the trade-off in a more general setting.

5 **R#1:** *"...importance of similarity among the selected tasks..."* In Theorem 1&2, similarity in the tasks can be described
6 by $\text{Var}_\gamma(\mathbf{Q}_\gamma)$ and $\text{Var}_\gamma(\mathbf{Q}_\gamma\boldsymbol{\theta}_\gamma)$, implying that greater similarity leads to smaller statistical error.

7 **R#1:** *"...domain randomization, when enough samples are used, is a better alternative to meta-learning..."* In many
8 practical deep-RL scenarios (highly nonlinear, low sample-case), we find optimization error dominates and DRS
9 outperforms MAML. Although Theorem 4 shows that DRS needs some small amount of training data to be effective,
10 empirical results (Fig 1&2) suggest that the amount of data required is as small as a few rollouts.

11 **R#2:** *"...Theorems 1 and 2 are asymptotic..."*: Only the first sentence of each theorem is asymptotic, the rest (starting
12 with "specifically") holds for finite samples. Indeed, asymptotic statements are obtained via the limits of those finite
13 bounds as $M, N \to \infty$. Hence, the theorems are NOT asymptotic. We will remove the asymptotic parts for clarity.

14 **R#2:** *'Assumption 2 ... the per-task optimal models are centered around the corresponding optimal solutions."*: We
15 do NOT assume that the task optimal models are centered around the DRS/MAML optimal solutions, only that their
16 distance is bounded. This assumption can easily be dropped with the cost of including the distance as a term.

17 **R#3:** *'...trust region alone cannot justify why ... TRPO fares better..."* Thanks for this insightful comment. We agree
18 that the use of a trust region does not fully explain the behavior; we will investigate this and add further discussion.

19 **R#3: Other points.** We will add the links to the supplement. We will share the source code when the paper is public.

20 **R#4:** *"...as the learning rate goes to zero, MAML becomes DRS..."* We respectfully disagree that MAML with $\alpha = 0$
21 is the same as DRS. This statement is only true for the asymptotic objectives (Eq 9) and is not true in practice, when
22 there are finite samples. MAML uses some data for the inner optimization regardless if the parameters are updated
23 ($\alpha > 0$) or not ($\alpha = 0$). Our main contribution is analyzing MAML and DRS in the practical, finite-sample case.

24 **R#4:** *"...[1] reports consistent improvements over joint training..."* Our meta-RL experiments are more extensive. [1]
25 considers locomotion with varying reward, while we consider locomotion and manipulation with varying reward or
26 dynamics. On the HalfCheetahRandVel environment, the only one that overlaps, our result (Fig 6c) is consistent with
27 [1] as MAML slightly outperforms DRS. Evaluations on a wider range of environments results in a different conclusion.

28 **R#4:** *"...learning rates..."* For TRPO-MAML's inner learning rate & trust region size, we use values from [16] for
29 the locomotion environments (as we use their algorithm implementations) and [25] for the others; DRS+TRPO's trust
30 region size was set to be the same as TRPO-MAML. The comparison is fair and follows practice of previous work.

31 **R#4:** *"...Joint training plus fine-tuning ... never does as well as it seems to in this paper..."* Our contribution is to
32 introduce and discuss the trade-off between DRS and MAML. Existing literature compare them on a single point within
33 this trade-off, whereas we explore the entire spectrum. Hence, those works happen to experiment on one side of the
34 spectrum. Our paper does not conflict with existing works, it complements them by considering a larger picture.

35 **R#4:** *"... Figures 2, 5, and 6... Why do these figures not align with each other?..."* The one-sided Welch t-test obtains
36 an estimate of the distribution of the difference between the average rewards of DRS and MAML. The variance of that
37 distribution combines the variances of the two average rewards and is fairly large. Thus, even if the average reward for
38 MAML is visually above that of DRS (Figure 5), the probability that DRS is better (Figure 2) may not be small. In
39 other words, the difference in Figure 5 is not statistically significant considering the variances.

40 **R#4:** *"In Section 2, DRS does better with more training steps...Section 3, DRS does better with \*fewer\* training
41 tasks..."* In Section 3, the optimization error is 0 for both methods (problem has analytical solution), and only the
42 statistical and modeling error matter. MAML has smaller modeling error and so it is better with more data. In Section 2,
43 the optimization error is significantly larger for MAML because of its bilevel structure; it dominates the trade-off for a
44 wide range of training budgets. The difference between these behaviors is one of the main points of this paper. We
45 show that the trade-off between DRS and MAML is not straightforward and requires careful analysis.

46 **R#4:** *"...it's shown that DRS requires a large training budget to be successful..."* We showed that DRS needs some
47 data (can be very small). Results in Figures 5&6 suggest that this budget is very small (only a few rollouts).

48 **R#4:** *"...[2] has an analysis..."* For meta-linear regression, [2] derive the MAML and DRS estimates when the task-
49 specific losses are known. In contrast, we consider the practical scenario where the task-specific losses are not known,
50 analyzing the estimates, their finite-sample behaviors, and their meta-test performances. We will add a discussion.

51 **R#4:** *"... if $\alpha$ is reduced to 0.1 (or below) ..."* The adjacent figure
52 shows meta-linear regression results for $\alpha = 0.05, 0.1$. The same
53 conclusions hold, but the difference between MAML and DRS is
54 not as pronounced; MAML requires a larger data set to provide a
55 clear improvement over DRS after test-time optimization.



$\alpha = 0.05$, pre $\alpha = 0.1$, pre $\alpha = 0.05$, post $\alpha = 0.1$, post