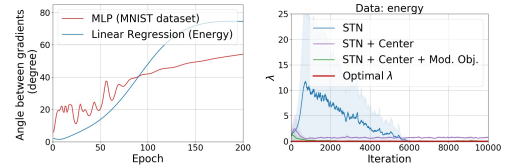1 We thank the reviewers for their detailed reviews and constructive feedbacks. We are encouraged that the reviewers
2 found our paper well-written & easy to follow [**R3 R4 R6**], theoretically motivated [**R3 R6**], empirically convincing
3 [**R3 R4 R6**], and novel [**R3**]. The major concerns and questions are addressed below; we will incorporate all suggestions
4 in the next revision and improve the notation (also include a table of notation) as suggested by **R5**.

5 [**R3 R5**] **Undesirable bilevel dynamics for uncentered parame-**
6 **terization.** In support of our claim that the training and validation
7 gradients are well aligned early in training, the left-hand figure shows
8 (for linear regression and an MLP) that the angle between the training
9 and validation gradients starts close to 0 and increases throughout
10 training. Based on the analysis in Section 3.1 (line 152), this align-



11 ment causes "spikes" in the hyperparameter values (Fig. 2 on page 7 and Fig. 6 on page 17), even with smaller learning
12 rates. The right-hand figure decouples the centered parameterization from all other changes to the training objective;
13 the centering by itself eliminates the spike, as our analysis predicts. This effect and the poor conditioning of the
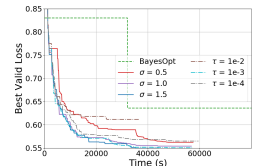14 Gauss-Newton Matrix, limit the optimization performance of uncentered representations.

15 [**R3**] **Smaller validation gradient $\Rightarrow$ sensitive to initial hyperparameters?** The uncentered parameterization indeed
16 results in larger hypernetwork gradients, but in a way that is unhelpful for the optimization as described above. Hence,
17 centered parameterizations are *less* sensitive to hyperparameter initializations.

18 [**R3 R6**] **Applicability to other bilevel problems and limitations of the STN approach.** STNs and $\Delta$-STNs consti-
19 tute a third general approach to gradient-based bilevel optimization, in addition to implicit differentiation (ID) and
20 unrolling. In relation to those two approaches, ($\Delta$-)STNs are much more efficient, because they amortize the inner
21 optimization and the response Jacobian. The range of applicability is slightly more restricted, as STNs require a single
22 (rather than per-example) inner objective and (like ID but unlike unrolling) require that the outer variables parameterize
23 the training objective. We note that NAS falls into this category and STN approaches can be applied. We will add this
24 discussion to the paper.

25 [**R5**] **Societal impact.** As almost any application of deep learning involves regularization hyperparameters, and
26 hyperparameter tuning is one stage of a much longer pipeline, any discussion of societal impacts would necessarily be
27 speculative. One predictable impact is to lessen the need for massive computing resources to tune hyperparameters.

28 [**R3 R4 R6**] **Learning rates and other hyperparameters.** Out of the three aforementioned gradient-based approaches
29 to bilevel optimization — implicit differentiation, unrolling, and ($\Delta$-)STNs — only unrolling is applicable to learning
30 rate tuning. However, even for unrolling, short-horizon objectives suffer from a severe bias [Wu et al, 2018] which
31 makes unrolling impractical for learning rate adaptation unless one uses huge amounts of computation (e.g. [Metz
32 et al., 2019]). Learning rate adaptation remains a hard problem in general. We note that it would be natural to wrap
33 Bayesian optimization (or random search) around the $\Delta$-STN to optimize the remaining hyperparameters that $\Delta$-STNs
34 are inapplicable to. BO suffers from high-dimensional search spaces and pushing the responsibility for regularization
35 hyperparameters onto the $\Delta$-STN can significantly reduce the dimensionality of the BO search space.

36 [**R3 R4**] **Impact of meta-parameters.** While it is always difficult to rigorously quantify
37 the ease of tuning, we found $\Delta$-STNs to perform well with a default set of meta-parameters
38 (e.g. fixed perturbation scale, choices of $\alpha_i$) across all tasks we investigated, unlike STNs.
39 $\Delta$-STNs are more forgiving of large perturbation scales than STNs, because they do not need
40 to model the entire response function over the likely range of hyperparameters. Furthermore,
41 $\Delta$-STNs are less sensitive than STNs to initial hyperparameters (e.g. figure 10), for the



42 reasons described in Section 3.1. The right-hand figure shows the effect of different entropy weights, where we grid
43 searched over $\{10^{-2}, 10^{-3}, 10^{-4}\}$ on all experiments. We will provide a more thorough sensitivity analysis in the next
44 revision.

45 [**R3 R5**] **Validation accuracy, and significance.** $\Delta$-STNs improve over STNs on validation accuracy as well as loss:
46 on AlexNet and VGG16, they improve accuracy by $0.5\%$ and $1.3\%$, respectively. In Table 1, we showed that $\Delta$-STNs
47 consistently outperform STNs with $4\% \sim 13\%$ improvements in the validation objective across all tasks. We believe
48 that these improvements are significant, since we are using the same set of regularization and data augmentation
49 hyperparameters, without modifying the architecture or introducing new techniques. We also showed that $\Delta$-STNs are
50 more stable and converge faster compared to STNs.

51 [**R3**] **Scalability.** We modified the statement made on line 249. On a ResNet-32 ImageNet classifier, the $\Delta$-STN
52 requires only a factor of 1.8 more time per epoch compared to training a single network (and this overhead can be made
53 arbitrarily small by updating the hyperparameters less frequently). We did not include ImageNet experiments because,
54 for proper baseline comparisons, it would cost thousands of dollars, but there is no computational obstacle to running
55 $\Delta$-STNs on ImageNet.