

1 **Reviewer #2:** Thank you for the thorough review. Indeed, better theoretical guidance enables XNAS to naturally select
 2 several hyper parameters, achieving best performance and accelerate differential architecture search.
 3 Experiments: since the compared methods use different sets of training tricks, we adopted the exact ones reported
 4 in SHARPDARTS, which is the newest NAS compared. The additional tricks in SHARPDARTS vs DARTS are:
 5 1.AutoAugment 2.power cosine annealing learning rate, instead of cosine annealing.
 6 Since DARTS' code is published, a fair comparison is attainable:

CIFAR-10 error	DARTS repo	XNAS repo	FLOPS(M)	inference time (ms, 1080ti)
DARTS-V2 (36 channels)	2.76%	2.42%	550	4.2
XNAS-small (36 channels)	2.16%	1.81%	529	3.6

8 In Bold are the results of new experiments following the reviewer's request. It can be seen that while additional tricks
 9 do help, XNAS cells exhibit significantly better performance even without those. For clarity, we will add 'XNAS-plain'
 10 results in the camera-ready version. Plain-mode is now supported by XNAS repo, which will be published upon
 11 acceptance to enable further comparisons. FLOPs and inference time comparisons were out of this paper's scope, and
 12 were not reported by the majority of compared methods. We will report our results in the revised experiments section.
 13 EG/Wipeout ablation: Please find EG/Wipeout ablation below.
 14 Paper clarity: the revised version now includes clearer explanations and modified figures as remarked by the reviewer.

15 Specifically, $p_t \doteq p_t^{(j,k)}$ is the prediction of a forecaster, as "superscript indices are ignored for brevity" (line 90).
 16 For citation we simply used the command *cite* with NeurIPS style.

17 Fig.1 was replaced with an illustration of multiplicative vs additive updates for differential architecture search.

18 **Reviewer #3:** Thank you for your review. Apparently our contribution and novelty was not stated clearly.
 19 *main difference compared to DARTS method is that it has a wipe-out method* - While DARTS picks a generic optimizer
 20 for architecture optimization, we derive an optimizer tailor-made for NAS based on EG. We suggest regret-minimization
 21 as a novel criterion for NAS to derive XNAS, and prove it achieves optimality under this criterion. We explore the
 22 algorithm's appealing properties, like robustness to initialization and the usage of multiple theoretically-derived learning
 23 rates, and demonstrate these in different setups. XNAS is lighter and x3.3 faster than DARTS, thanks to simple
 24 multiplicative updates and less hyper-parameters involved. It outperforms DARTS in 7/7 datasets, including -42%
 25 error in CIFAR-10. The wipeout is merely a part of this synergy. Please find EG/Wipeout ablation below.

26 *wipeout ... has the risk of eliminating useful "expert"* - the wipeout's merit is the elimination of unuseful experts only.
 27 Lemma 1 states that clearly: "In XNAS, The optimal expert in hindsight cannot be wiped-out".
 28 In the light of the above, we ask the reviewer to reconsider his review for this work.

29 **Reviewer #4:** Thank you for the valuable remarks which help us make the revised version much clearer, and also for
 30 acknowledging the originality and the need for EG in the context of NAS.

31 Regret analysis assumptions: provided under conditions that are not obvious to hold - that is correct. In the analysis, we
 32 assume gradient-boundedness and convexity, as the former is enforced by applying gradient-clipping, the latter does not
 33 hold in general for DNNs. However, the convexity of the loss is a very common assumption for regret minimization and
 34 online learning in general. It is assumed when deriving similar theoretical guarantees for all previous NAS optimizers,
 35 e.g. gradient descent (ASAP) and ADAM (DARTS). In addition, we put to test the optimizer in several toy setups and
 36 simulations, deterministic and statistical (sections 8.1, 8.2, 8.3), investigate its properties and compare with common
 37 optimizers. Finally, extensive experimentation in NAS provides an empirical support for the algorithm derived under
 38 that assumption. Based on that, we believe that XNAS analysis has a high significance as well.

39 EG loss function: Can the authors provide the form of this loss function? - The loss function is the parent network's
 40 cross-entropy loss over the validation set. This loss is not defined explicitly for a predictor in an intermediate node.
 41 However, its back-propagated gradient is: $\nabla_{p_t} \ell_{\text{val}}(p_t)$. Therefore, unlike loss-based PEA methods like HEDGE, we
 42 derived a gradient-based NAS optimizer, via auxiliary-losses (Eq. 22). We will make this observation clearer.

43 Reward bounds, gradient clipping and optimal learning rates - the bounds are over the reward values (Alg.1, line 8),
 44 $|R_{t,i}| \leq \mathcal{L}$, as required by our theory and fulfilled by a simple gradient-clipping. We will make that clearer in the paper.
 45 The optimal learning rate derived from theorem 1 depends on that bound (line 163) and is used in our experiments
 46 (values-line 229). The successful usage of the theoretical learning rate in practice is a key property of XNAS.

47 EG/Wipeout ablation: The revised version will contain a section for the analysis of EG without wipeout. We include
 48 here only the summary and omit the full details and graphs due to the limited space.

49 Runtime: EG: 0.35 GPU-days. Faster than ADAM which is used in DARTS(1nd) and runs for 0.4 GPU-days with a
 50 significantly inferior performance (table 2). This is mostly due to fewer calculations and no momentum updates (section
 51 4.2.2). Wipeout contribution: -0.05 GPU-days. Wipeout takes effect around the last 30% of the run, and reduces this
 52 part's runtime by half. If one desires a further acceleration (e.g. a large dataset or many operations), while risking with
 53 the elimination of useful experts - the wipeout factor discussed in section 4.2.1 controls this trade-off.

54 Performance: EG is responsible for most of the improvement compared to previous methods. Based on a few runs, the
 55 mean error for CIFAR-10 is around 1.75% (50 channels). The wipeout contributes to the overall mean performance,
 56 and also reduces its standard deviation caused by hard selections of operations ('relaxation bias', line 128).