1 We greatly appreciate the feedback of the reviewers. We discuss the specific concerns of the reviewers below.

2 Reviewer 1:

3 • A theoretical upper-bound of the regret of Approx-Zooming-With-No-Arm-Similarity is stated in [7] as
4 $O(KT^{\frac{d+1}{d+2}})$, where $d$ is the dimension of the context space. Note that this regret bound is linear in the number
5 of arms. We will include this discussion into the paper.

6 • We will include empirical results of a gaussian process-based bandit in the final paper. As one needs to specify
7 the covariance/kernel matrix, we could either fit a Gaussian process to each arm separately when the metric is
8 unknown, or we can construct the covariance matrix given oracle access to a metric among the arms.

9 • Generalizing to higher dimension: Yes, indeed our algorithm and analysis does extend to the general $d$-
10 dimensional setting, and we will include that into the final paper. The only change required algorithmically is
11 in the subpartitioning/clustering step. Let us define $C_d(q)$ to be the number of balls of radius $r/q$ needed to
12 cover a ball of radius $r$, which scales exponentially in the dimension $d$, e.g. $q^d$. Since we are now estimating
13 the reward function $f$ over a $d$-dimensional context space, the number of sub-regions of the context space that
14 need to be clustered will be $C_d(2)$, and the number of samples needed to guarantee that the $k$-nearest neighbor
15 samples are within distance $\frac{1}{16}$ radius, will be equal to $\tilde{O}(kC_d(32))$. To compute $\hat{\mathcal{D}}$, we will instead have a
16 $d$-dimensional summation over the subset of the context space. Once $\hat{\mathcal{D}}$ is computed, then the clustering of
17 arms will have the same computational cost, i.e. linear in number of arms to be clustered. The analysis can be
18 modified to account for the $d$-dimensional setting, and the final regret bound will look like

$$O\Big(C_d(2)C_d(32)\sigma^2 L^{-2} K \ln(TK) + \min_{i_{\max}\in\mathbb{Z}_+} \big(LT2^{-i_{\max}} + \sum_{i=1}^{i_{max}-1} C_d(2)C_d(32)\sigma^2 L^{-1} M_i 2^i \ln(TK)\big)\Big),$$

19 where $M_i$ instead sums over an $\epsilon$-net of the context space for $\epsilon = 2^{-i}$, and thus we may expect $M_i$ to grow
20 exponentially in $i \times d$, although modified with respect to the distribution of the reward function and the finite
21 arms. The growth of $M_i$ will dominate the regret bound with respect to the dependence on the dimension $d$.

22 Reviewer 2:

23 • Line 160: We will include an example, in fact the setup we chose for the simulation illustrates this as it is a
24 periodic function that could be rearranged to be simply linear rather than the depicted zigzag.

25 • Adapting to the smoothness parameter of the reward function: Indeed our algorithm can be generalized to
26 Holder continuous reward functions. If the smoothness parameter is known, then modifying the algorithm is
27 straightforward. We will look into the techniques of Qian and Yang (2016) for adaptivity to the smoothness.

28 • Generalizing to higher dimension: see response to reviewer 1 above.

29 Reviewer 3:

30 • Regret bound scaling with Lipschitz constant: There are two terms in the final regret bound. The first term
31 comes from the cost of the initial clustering, and as you pointed out it scales inversely with $L^2$. However, the
32 first term scales as $\tilde{O}(K)$, logarithmic in $T$, whereas the second term has a polynomial dependence on $T$, e.g.
33 $\tilde{O}(\sqrt{KT})$ in the given examples. For large $T$ and $K = o(T)$, the second term will dominate, which does not
34 scale inversely with the Lipschitz constant. To give more inutition though, the inverse dependence on $L$ in the
35 first term (only logarithmic scaling wrt $T$) is due to the fact that the clusters are constructed to satisfy that the
36 bias in the reward due to different arms in the same cluster is controlled to be on the same order as the bias in
37 reward from different contexts in the same set. For small Lipschitz constant, the reward varies less across the
38 same size context width, and thus the algorithm requires that arm distances are measured more precisely to
39 guarantee that the bias of the new ball is bounded by $L$ times the context width. This initial clustering of the
40 algorithm could be modified to allow for looser clusters, where the bias due to different arms in a cluster is
41 larger than the context width. This would remove the inverse dependence on $L$ in the first term, but would add
42 more phases of subpartitioning that would contribute regret to the summation in the second term.

43 • The regret bound obtained in the contextual zooming algorithm when the metric among the arms is available is
44 $O(T^{\frac{d+1}{d+2}})$, where $d$ is the dimension of the joint context-arm space, i.e. $d = d_c + d_a$ where $d_c$ is the dimension
45 of the context space and $d_a$ is the dimension of the arm space. We will include this into the final paper.

46 We will also address the minor typos/comments in the revision as well. Thank you for your detailed feedback!