

1 We thank the reviewers for their time and feedback. In this paper, we argue that intrinsic motivation can take many
2 forms. Inspired by humans, we propose a curiosity formulation based on multimodal association: searching for novel
3 associations to explore. We demonstrate that audio-visual curiosity shows promising results not only on standard Atari
4 environments but also on the realistic Habitat setting.

5 We highlight that the reviewers believe the paper is interesting (R1), novel (R4), “the first multimodal curiosity work”
6 (R2), and shares “intuitive inspiration and some promising results” (R3). There are concerns with regards to the
7 applicability of the approach and the failure cases, and R3 has concerns with respect to the formulation (specifically the
8 use of error as reward). We address all of these concerns below. However, the main contribution of our paper is the
9 introduction of multimodal curiosity. Independent of the specific formulation, this should be of great interest to the
10 NeurIPS community and pave the way for exploiting the richness of data for better performance.

11 **R1, R3: Formulation (R3: error as reward, R1: association is as hard, R1,3: couch-potato issue):** Error as
12 reward is not necessarily bad; the question is what error. Our approach is not to predict audio given visual features (or
13 vice-versa). Instead, given both audio and visual features, we classify whether they are aligned or not. We highlight that
14 this association classification error (i.e. association novelty) is fundamentally different from the prediction error based
15 formulation typically used. In particular, this has two implications:

16 (a) In the prior visual curiosity framework, the model predicts the future frame in high-dimensional space. We argue
17 that compared to higher-dimensional prediction, our classification formulation is less susceptible to the issues
18 mentioned. For example, if pressing a button produces three distinct sounds, we could learn to classify all of these
19 as associated, while an agent using future prediction error would always be curious.

20 (b) With error as reward, as Schmidhuber points out, “the problem is that in non-deterministic environments the
21 controller will focus on parts of the environmental dynamics which are inherently unpredictable.” On the other
22 hand, in our case, the discriminator focuses on deterministic aspects to solve the alignment classification problem.
23 This effectively helps ignore stochasticity! Therefore, while our approach would not overcome a purely random
24 environment, it would mitigate the couch-potato problem. In contrast, future prediction-based curiosity is attracted
25 to randomness. Our noise ablation, which adds noise to *both* the audio and the visual features, provides some
26 insight into this.

27 **R1, R4: Applicability/Generality:** This paper has taken a step forward in terms of real-scenario experiments compared
28 to prior work, which uses Atari as a standard benchmark. We extend to Habitat, which has many characteristics of the
29 real world: realistic audio and visual modalities, generated from physical processes, no clear separable sound effects,
30 and nontrivial associations (i.e. not one-to-one object-to-sound correspondence). We highlight that our method can
31 perform without direct visibility, in the presence of background noise, and with more than one audio source:

32 (a) **Visibility:** Direct visibility of an object is not required for association; context should be sufficient. For example,
33 hearing a microwave sound in a kitchen would be positively associated even if the microwave is out of view.

34 (b) **Background Noise:** We acknowledge that background noise could be an issue, but it would be so when there is
35 only background or random noise (as discussed in the above formulation response). If there are foreground audio
36 and visual signals, we can still learn associations in the presence of noise.

37 (c) **Multiple Audio Sources:** Multiple objects also render visual prediction hard. It requires object segmentation and
38 modeling relationships. Similarly, multiple audio sources would require segmentation.

39 **R1: Issue viewing supplementary material:** Thank you for bringing this to our attention. We cannot change the
40 format at this time, but you should be able to unzip it with `jar xvf supplementary material.zip`.

41 **R2: Which baseline from Burda et al.:** Our Burda et al. baseline uses random CNN features, which is stronger than
42 pixel prediction, as you mentioned. Our method uses the same random CNN features, as shown in Figure 2.

43 **R2: RND with audio baseline:** We agree that RND would strengthen the audio in baseline ablation. In our preliminary
44 experiments, RND with audio appears to perform similarly to the RND baseline without audio. We will include this
45 ablation in the final paper.

46 **R3: Definition of association:** We think of association as learning shared information between modalities, i.e. the
47 underlying physical processes that govern these different signals. We do implement this as learning alignment (are two
48 things cotemporal), but this is not the only way to learn associations. We will modify the text to clarify this.

49 **R3: Audio padding:** Yes, the padding to 1 second is done to make each sound equal length for feature computation.

50 **R3: Prior work on mitigating shortcomings:** Thanks for pointing this out and we will add this discussion and context.
51 Our approach, different from this body of prior work, looks at how multimodal data can mitigate these issues.

52 **R3: Terminology, prior work, typos, citations:** We are grateful to R3 for their detailed comments and will definitely
53 incorporate this feedback into the final paper.

54 **R4: Fixed random initialization:** We use random CNN features to be similar to prior work. This is the same feature
55 representation as used in our Burda et al. baseline.

56 **R4: Intrinsic or extrinsic rewards:** The agent only has access to intrinsic rewards, as described at L174-175. Extrinsic
57 rewards are used only for our evaluation of exploration.