

1 We thank reviewers for their thoughtful feedback.

2 **Additional references (R1, R2, R3).** We reported in Table 2 the best existing self-supervised baselines at the date of
3 submission (including unpublished methods such as XDC). We outperform all the reviewers’ additional references that
4 target unsupervised learning (R2’s [1,2]), apart from results on a single benchmark (ESC-50) for R2’s [1] which we
5 indeed miss and will add (it appeared on arXiv one month before the deadline and is unpublished). Other references are
6 relevant but do not address the problem of unsupervised representation learning (e.g. “One model to learn them all” is
7 about supervised learning, “Multimodal transformer for unaligned multimodal language sequences” uses supervised
8 features...). We will add all mentioned references but we stress that they do not hurt our novelty claim nor our results.

9 **Fair comparison to previous work (R2, R4).** Comparison on equal grounds is a problem for all papers in this area,
10 and we try to be fair: (1) The **comparison with MIL-NCE (R4)** is performed in Table 1(a) since MIL-NCE is strictly
11 equivalent to our VT only architecture. With the exact same training setup and backbone, our FAC outperforms
12 MIL-NCE by 2% in both UCF and HMDB, while enabling a task (ESC-50) that would not be possible with MIL-NCE.
13 (2) The point of the paper is to demonstrate that **more modalities help (R2)**. Note that we outperform ELo which
14 uses an additional modality that we don’t (flow, which is really important for action recognition). (3) It is difficult
15 to compare our method on the **same data and architectures** as (i) FAC requires text and only HowTo has it (only
16 MIL-NCE uses this dataset, and we compare to this), (ii) IG65M is not public (**R1**), (iii) AVTS, XDC and ELo haven’t
17 released code, and (iv) other works use a variety of architectures. We beat XDC despite using less data (16 vs 21 years)
18 and fewer parameters (24 vs 33M), and ELo despite using a $2\times$ smaller model. (4) We did **another comparison to**
19 **XDC (R2, R4)** by running our VA model on the same data (AudioSet) and backbone (R(2+1)D-18) and outperform
20 them: 91.5 (vs. 91.2) on UCF and 70.1 (vs. 61.0) on HMDB (also matching R2’s [1] on UCF and beating on HMDB on
21 the same data). Note that R(2+1)D-18 actually outperforms S3D (used in our submission) so MMV does not simply
22 beat XDC due to a better backbone (**R4**). (5) **XDC beats their own fully supervised baseline** but we report a stronger
23 and more meaningful quantity – the best externally published performance for supervised transfer. Finally, as detailed
24 next, we would like to stress that our good performance (e.g. a significant boost of **5.5** point on HMDB) is not the only
25 contribution of the paper.

26 **Novelty, contributions and claims (R2).** We agree that our novelty does not lie in the loss which is indeed not novel.
27 The loss is not the only means to induce different structures on the embeddings, instead, we achieve that through
28 architecture design. For example, the FAC design allows us to navigate from the video-audio space (fine) to the
29 video-text space (coarse) (property (iii) of the MMV L28), which is not possible with the disjoint design. We validate
30 that claim in the supplemental video as highlighted in L313 through qualitative audio-to-text retrieval (we are not aware
31 of standard benchmarks for quantitative audio-text retrieval evaluation). This shows that we can influence the structure
32 of the embedding through architecture design and not only with losses. Furthermore, we show that FAC performs better
33 than the other considered designs in Table 1(b) on 3 out of 5 benchmarks. In addition, to the best of our knowledge
34 (acknowledged by R1 and R4), this is the first work to jointly learn from video, audio and text in a self-supervised
35 fashion. In the paper, we explore how to do that well at scale, propose various embedding strategies, and demonstrate
36 state-of-the-art performance on challenging downstream tasks which we deem to be an impactful contribution.

37 **Importance of the deflation contribution (R3).** We believe the deflation technique is an important contribution of
38 the paper as it allows video-trained models to do inference efficiently with image inputs. In particular, we show that the
39 image classification performance when using the deflated model is similar to using the original model with an inflated
40 input. In addition, we believe to be the first work to consider learning first from video to transfer to images and show
41 strong performance on two image benchmarks. We will clarify the impact of the deflation contribution in the paper.

42 **Baseline for deflation (R1).** R1’s proposed method is indeed a valid idea which we expect to perform on par or
43 better than our approach. However, this method effectively does partial finetuning, while we instead focused on linear
44 evaluation with frozen networks since it enables fast off-the-shelf evaluation on unseen image dataset and tasks.

45 **Positives and negatives for NCE (R1, R3).** For the text positive (L180), we strictly follow MIL-NCE [41] and refer
46 to this paper for details. In L173, we use all negatives coming from other elements in the batch. In total there are
47 $2 \times (N - 1)$ negative pairs ($N - 1$ audio negatives for the video and $N - 1$ video negatives for the audio).

48 **Model and code release (R2, R3).** We will release our pretrained models as well as the training scripts.

49 **Perfect audio-visual alignment (R3).** We agree there is no perfect alignment, but the point of L174 was to emphasise
50 that text is less aligned with the visual content than audio in general. Note that all competing audio-visual learning
51 approaches ignore the occasional misalignment and, like us, observe that learning is possible despite it. We will clarify.

52 **ASR for action recognition (R1).** Text from ASR provides semantic information about the visual content of the
53 videos for objects (e.g. tomato) or actions (e.g. cut). Interestingly even though there exists a domain gap between the