

1 **EK Benchmark May Not Be Correct (R1).** We do not derive the EK ground truth from noisy narrations. The EK  
2 dataset contains 1) (*noun, verb*) action tuple labels, 2) bounding box annotations of *active* objects, and 3) a list of  
3 *nouns* associated with an action (e.g., [’pan’, ’mushrooms’] for "put mushrooms in the pan"). We construct the ground  
4 truth by finding frames where the object category of a bounding box matches the *noun* of either 1) an action tuple or 2)  
5 one of the *nouns* from a *noun* list. We manually verified our annotations and confirmed that they were error-free.

6 **Claims of Not Training COBE on EK (R1, R2, R4).** We acknowledge that even though we do not directly train  
7 COBE on EK, we use a detector trained on EK to pseudo-label the frames of HowTo100M, which might ease the  
8 transfer to EK. We will revise our former claims to make this point clear. Also, note that we do *not* start from a detector  
9 pre-trained on EK when training on HowTo100M\_BB. To expand our evaluation beyond EK, we also manually labeled  
10 a test set of 9K HowTo100M frames (using a disjoint set from HowTo100M\_BB) with 171 unique (*noun, context*)  
11 tuples, where *context* can be one of: *noun, verb, adjective* or *adverb*. On this new set, COBE achieves 17.5 mAP,  
12 while Tuple Faster R-CNN yields 16.1 mAP. We will release all benchmark annotations upon publication of the paper.

13 **Additional Baselines (R1, R2, R3, R4).** As requested by R3, and R4, we use S3D trained on HowTo100M [Miech et  
14 al., CVPR20] with Faster R-CNN to construct a baseline for our task (as for the other baselines in Table 1 of our draft).  
15 While this baseline performs better than other action recognition baselines, its mAP accuracy is 7.9 and 10.6 lower than  
16 that achieved by COBE on the EK\_H evaluation set for (*noun, noun*), and (*noun, verb*), respectively. In response to  
17 R1’s question, we note that these results indicate that classification systems such as those of Miech et al. cannot be  
18 easily adapted to our problem, and that a specialized design such as COBE might be needed for better solving this task.  
19 Also, as suggested by R3, we replaced CTRL, with a word2vec embedding, which produced mAP accuracies of 15.8  
20 and 21.9 compared to the 16.9 and 24.7 of COBE w/ CTRL. We also tried the weighted averaging scheme suggested by  
21 R3, which yielded 14.9 and 22.1 in mAP. Lastly, per R2’s request, we implemented a baseline that “matches object  
22 detection outputs to noun tokens,” which achieved a 6.7 mAP on the EK\_H (compared to the 16.9 mAP of COBE).

23 **Object Detection Results (R2, R3).** As requested, we ran object detection experiments on 124K frames of EK (180  
24 object categories) by comparing COBE to a Faster R-CNN trained on HowTo100M\_BB for object detection. Both  
25 methods share the same architecture (except for the contextualized object branch). COBE outperforms this baseline by  
26 1.4 mAP. We also note that pre-training COBE on HowTo100M\_BB and then finetuning it on EK outperforms Faster  
27 R-CNN only trained on EK by 10.1 mAP. This highlights the benefit of pretraining on HowTo100M\_BB.

28 **Limitations of COBE (R3, R4).** R3 and R4 correctly point out that COBE relies on a predefined set of noun-centric  
29 object categories. However, because of the supervision from narration, our detector can recognize fine-grained  
30 contextual cues beyond the coarse categorical labels. For example, a standard object detector can predict that there  
31 is a “tomato” in an image, whereas COBE can predict that there is “a yellow tomato” or “a chopped tomato.” R4  
32 rightly points out that unseen tuple categories may still occur visually in the training frames even though they are not  
33 mentioned in the text. However, because COBE is supervised by textual captions, it would *not* be trained to predict  
34 the categories of those “unmentioned” tuples. Thus, in those cases, COBE is indeed inferring new tuple concepts. We  
35 fully agree with R3 and R4 that it would be great to eliminate the reliance on bounding box labels. However, this is  
36 challenging, particularly because HowTo100M captions are noisy. We intend to tackle this problem in our future work.

37 **Related Work (R4).** ViLBERT is better suited for tasks that require multi-modal inputs (e.g. VQA, VCR). Instead,  
38 COBE uses captions to supervise a visual detection model during training, but then operates on video frames alone  
39 during inference. We thank R4 for highlighting these related models, which we will cite and discuss in our final paper.

40 **EK Action Recognition Results (R2).** Unlike action recognition models, which are trained on manually annotated  
41 action labels, COBE is trained on noisy narrations. Thus, such a comparison would not be properly fair. However, we  
42 agree with R2 that such experiments would be useful and we will attempt to include them in our final draft.

43 **Technical Novelty (R1, R2, R3).** We acknowledge that our work offers limited contributions in terms of model design.  
44 Our main technical contribution is our large-scale training framework which leverages noisy narrations for learning  
45 object detectors that generalize to novel states. This is in contrast to prior work [43, 44] that focuses on noun-centric  
46 object detection, and that relies on smaller datasets with manually annotated text descriptions. As noted by R4, our idea  
47 is simple, and effective, and it demonstrates a novel object detection application of modern language models.

48 **Experimental Questions (R2, R4).** 1) Prior to our evaluation, we verified that our evaluation set does *not* contain  
49 frames with multiple tuple categories. However, we agree with R2 that for instances with multiple valid tuples, using a  
50 softmax would be suboptimal. To address this issue, we would use a sigmoid instead of a softmax in our formulation in  
51 L210. 2) The two tasks are indeed separated to include specialized baselines. 3) We considered a 3D backbone, but it  
52 was too costly to train it on HowTo100M\_BB. 4) Misalignment between captions and videos was not an issue because  
53 we used only the most confident detections to construct HowTo100M\_BB (see L118-L120). 5) We appreciate R4’s  
54 advice of using global features for capturing more context: we will try it. 6) We confirm that: a) Faster R-CNN is  
55 trained to predict the classes from EK, and b) COBE is trained end-to-end.