

1 We thank the reviewers for their constructive feedback.

2 **The use of LSTMs instead of 3D CNNs [R1, R2, R3].** We use LSTMs as an “agent” not only to make classification
3 predictions but more importantly to make *sequential gating decisions*—*dynamically* determine whether to compute
4 features at a finer scale conditioned on incoming video frames and historical information. In particular, the decision for
5 the t -th time step depends on previous observations and decisions—if the current frame is very similar to previous seen
6 frames (this is very common as there are a lot of redundant frames) the model is more likely to use coarse features as no
7 new information is provided. Therefore, we use LSTMs to learn the fine feature usage policy based on the history of
8 past interactions with the video, which is similar in spirit to a MDP process if we replace Gumbel-Softmax with RL. In
9 addition, the autoregressive nature of LSTMs allows LITEEVAL to save computation for online video recognition with
10 minimal modification, while it is not clear how to use uniform sampling for online recognition since it is hard to select
11 the optimal sampling rate for different videos as there is no information (*i.e.*, duration and fps) known beforehand about
12 incoming videos (sampling every 1 min might work for long videos but would be problematic for a 1-min-long video).

13 Existing 3D CNN models (I3D, S3D, *etc*) typically average prediction scores from N (*e.g.*, 10/25) uniformly sampled
14 snippets (stacked frames) as video-level predictions. There are a few disadvantages: (1) they operate on snippets,
15 requiring storing multiple frames for 3D convs, which is not feasible on low-power devices; (2) 3D CNNs are generally
16 computationally expensive (108 GFLOPs for a single snippet in I3D); (3) they produce one-size-fits-all models for all
17 videos regardless of their complexity; (4) the uniform sampling strategy for testing prevents them to be readily used in
18 online settings. In contrast, our model saves computation *for both online and offline settings* by using expensive features
19 as infrequently as possible. Note that we could use 3D CNNs as our feature extractor when computational budget is
20 sufficient. To study whether our framework is compatible with modern frameworks for video recognition, we adopted a
21 DPN model trained using the temporal segment network, which is a state-of-the-art framework for video recognition;
22 LITEEVAL offers 83.6% on ACTIVITYNET, confirming LITEEVAL supports features from different backbones.

23 **Comparisons with SlowFast [R1, R2, R3].** Thanks for pointing out this paper, which will be cited and discussed. But
24 we do like to stress that there are significant differences between our approach and SlowFast—(1) SlowFast produces
25 the same set of parameters for all videos whereas our approach allocates computational resources conditioned on input
26 videos; (2) SlowFast relies on the uniform sampling baseline, making it unsuitable for online recognition; (3) SlowFast
27 operates on video frames with the same spatial resolution (*i.e.*, 224×224) and uses lightweight CNNs for the Fast
28 pathway ($\sim 20\%$ computation) and heavy CNNs for the Slow pathway. In our model, we not only use a lightweight
29 CNN to extract coarse features but also reduce the input resolution, making the computation overhead of the coarse
30 features negligible (0.08 GLOPs). We are currently preparing a comparison.

31 **Combining coarse and fine features [R1].** As suggested by the reviewer, we also compare with the uniform sampling
32 and the LSTM baseline using both coarse and fine features on ACTIVITYNET. Although fusing two features does
33 slightly improve the performance of using fine features alone, we observe that LITEEVAL still achieves better results
34 compared to the uniform sampling (72.7% vs. 70.6%) and the LSTM baseline (72.7% vs. 71.5%).

35 **Second term in the loss/syncing cLSTM and fLSTM [R1].** The ablation study of synchronizing the LSTMs are
36 reported in Tab. 2. The 2nd term in the loss function controls the computational budget and results are shown in Tab. 3.

37 **Comparison with Skip-RNNs [R1].** Skip-RNNs achieved similar mAP as our LSTM baseline with slightly less
38 computation, since it processes *every frame* and saves computation by learning to skip updates of RNN models. Note
39 that the most expensive computation in a video recognition pipeline is feature extraction (7.82 GFLOPs for a single
40 frame with a ResNet101), and the computation incurred by LSTMs is negligible (0.005 GFLOPs per time step).

41 **Numbers in Tab 1 and Fig. 2 [R1].** Tab. 1 reports the performance of LITEEVAL in an offline setting while Fig. 2
42 summarizes online recognition results. Please refer to L221-L225 for more details.

43 **Gating [R3].** We vary the computational budget of LITEEVAL by adjusting γ and then compare with random selection
44 that uses similar computation budget during inference on ACTIVITYNET. The mAP (random vs. LITEEVAL) is 65.8%
45 vs. 72.7% (102 GFLOPs); 69.1% vs. 73.2% (183 GFLOPs). In addition to qualitatively visualized selected frames
46 in Fig. 4, we also visualized *quantitatively* fine feature usage statistics of selected classes in Fig. 3, we can see that
47 for simple classes like objects (*e.g.*, gorilla), LITEEVAL makes predictions with less fine feature usage, while more
48 computation is needed for more complicated classes (*e.g.*, marriage proposal). This confirms LITEEVAL is able to learn
49 useful gating decisions.

50 **Datasets [R3].** We chose FCVID and ACTIVITYNET for evaluation since videos in these two datasets are “untrimmed”
51 with an average duration over 100 seconds whereas videos in UCF-101 (~ 8 s), Kinetics (~ 10 s) and Something-
52 Something (~ 4 s) are all “trimmed”. Compared to extensive research efforts on action recognition on trimmed videos,
53 we believe long-term video understanding is also a very important and arguably a more challenging problem; resource-
54 efficient models for long videos are of great value for their deployment in real-world scenarios.