

1 We would like to thank the reviewers for their feedback and their encouraging comments about the theoretical
 2 contribution [R1, R2, R3], novelty of the learning approach [R1, R2], convincing experiments [R1, R2, R3], and the
 3 paper being well-written [R1, R2]. Our responses are as follows.

4 **Response to Reviewer 1:**

5 **C1:** "The experimental results quite thoroughly demonstrate the value of the proposed method on a dataset, though
 6 perhaps more datasets could be used"

7 **Answer:** Thanks. To address this comment, we performed experiments on the Inria instructional video dataset (5 tasks
 8 and 30 videos per task). We used C3D features of video segments and ran all subset selection methods with the same
 9 settings discussed in the paper. The results in Table 1 (below) show significant improvement w.r.t. the state of the art.

10 **C2:** "If I have misjudged the novelty of the work, please explain why."

11 – **Answer:** We would like to mention that this is the first work on supervised subset selection that derives conditions for
 12 the exactness of a subset selection utility function (i.e., conditions under which subset selection recovers ground-truth
 13 representatives) and employs these conditions to design a loss function for representation learning, e.g., via DNNs. In
 14 fact, our work takes a major step towards a theoretically correct/motivated supervised subset selection framework. We
 15 also plan to release the code, if accepted.

16 We hope that our responses address the reviewer’s comments and kindly ask the reviewer to raise his/her rating.

17 **Response to Reviewer 2:**

18 **C1:** "It would be more convincing to include results for at least one other dataset"

19 **Answer:** Thanks. To address this, we performed experiments on the Inria instructional video dataset, which has 5 tasks
 20 and 30 videos per task. We used C3D features of video segments and ran all methods with the same settings discussed
 21 in the paper. Table 1 (below) shows that our method outperforms other algorithms with at least 5% on the entire dataset.

22 **C2:** "error bars (confidence estimates) should be provided for the scores in Tables 1 and 2."

23 – **Answer:** Thanks for the comment. We will update the tables to include score variances. We would like to mention
 24 that the variance of our method (SupUFL-L and SupUFL-L) for each task is less than 0.7% in Table 1 in the main paper.

Activity	Uniform	UFL	dppLSTM	SubmodMix	SupUFL-L	SupUFL-N
cpr	51.6	55.4	59.8	61.6	65.3	66.4
coffee	53.4	52.2	47.3	56.1	58.4	61.3
repot	50.8	56.5	67.2	60.8	74.8	74.3
jump car	52.7	55.9	57.4	57.5	62.8	66.7
changing tire	50.1	59.3	55.3	62.8	65.4	69.2
Average	51.7	55.9	57.4	59.8	65.3	67.6

25 **Response to Reviewer 3:** Table 1: F1 score (%) of different algorithms on the Inria dataset.

26 **C1:** "This tightness result is interesting but not novel"

27 **Answer:** While tightness of convex relaxation for some non-convex problems has been studied before, this is the first
 28 work studying the exactness for the uncapacitated facility location function. Moreover, we use our derived exactness
 29 conditions to design a loss function for representation learning in the supervised setting, hence, moving towards a
 30 theoretically correct/motivated supervised subset selection.

31 **C2:** "In experiments, is it possible to include a tweaked version of [59] as another baseline?"

32 **Answer:** Thanks for the question. The reason we cannot simply use [59] for supervised subset selection is that not
 33 only it requires knowing assignment of points to predefined categories (please see Remark 1), but also we need to
 34 know matching between ground-truth representatives of different videos, i.e., we need to know which ground-truth
 35 representative in each video belong to the same subactivity/category. Thus, to extend [59] to supervised subset selection
 36 we need to iteratively i) perform matching/clustering of ground-truth representatives across all videos, ii) assign
 37 segments to induced categories, iii) perform learning via [59]. In our initial experiments, this did not do well due to the
 38 difficulty of matching ground-truth representatives across videos (also not all videos have all subactivities). We will
 39 explain this further in the revised paper, if accepted.

40 **C3:** "In (7), there may be some trivial Θ , (say $\Theta = 0$ in some settings) which enforce all $f_{\Theta}(y)$ to be equal so that
 41 all the data points will be assigned to one cluster. In practice, a random initialization of Θ and an iterative gradient
 42 algorithm may result in a reasonably good Θ , but the problem in (7) is not what is being solved?"

43 **Answer:** Thanks for the question. In our case, each ground-truth representative of a dataset form its own group and
 44 our third loss function, $\mathcal{L}_3^{\ell}(\Theta)$ (equation (6)), enforces that different representatives must be separated by the adaptive
 45 margin. This prevents data to collapse to the same point (for which the total loss function is positive, whereas when all
 46 our three conditions are satisfied, the total loss will be zero).