# Joint Discovery of Object States and Manipulating Actions

Jean-Baptiste Alayrac<sup>\*</sup>

Josef Sivic\*

Ivan Laptev\*

Simon Lacoste-Julien<sup>†</sup>

# Abstract

Many activities involve object manipulations aiming to modify object state. Examples of common state changes include full/empty bottle, open/closed door, and attached/detached car wheel. In this work, we seek to automatically discover the states of objects and the associated manipulating actions. Given a set of videos for a particular task, we propose a joint model that learns to identify object states and to localize state-modifying actions. Our model is formulated as a discriminative clustering cost. We assume a consistent temporal order for the changes in object states and manipulating actions, and learn the model without additional supervision. Our method is validated on a new dataset of videos depicting real-life object manipulations. We demonstrate the successful discovery of seven manipulating actions and corresponding object states. Moreover, we emphasize our joint formulation and show the improvement of object state discovery by action recognition and vice versa.

# 1. Introduction

Many of our activities involve changes in object states. We need to open a book to read it, to cut bread before eating it and to lighten candles before taking out a birthday cake. Transitions of object states are often coupled with particular manipulating actions (open, cut, lighten). Moreover, the success of an action is often signified by reaching the desired state of an object (whipped cream, ironed shirt) and avoiding other states (burned shirt). Recognizing object states and manipulating actions is, hence, expected to become a key component of future systems such as wearable automatic assistants or home robots helping people in their daily tasks.

Human visual system can easily distinguish different states of objects, such as open/closed bottle or full/empty coffee cup [6]. Automatic recognition of object states and state changes, however, presents challenges as it requires distinguishing subtle changes in object appearance such as



Figure 1: We automatically discover object states such as detached/attached car tire or empty/full coffee cup along with their corresponding manipulating actions by observing people interacting with the objects.

the presence of a cap on the bottle or screws on the car tire Despite much work on object recognition and localization, recognition of object states has received only limited attention in computer vision [17].

One solution to recognizing object states would be to manually annotate states for different objects, and treat the problem as a supervised fine-grained object classification task [10, 11]. This approach, however, presents two problems. First, we would have to decide *a priori* on the set of state labels for each object, which can be ambiguous and not suitable for future tasks. Second, for each label we would need to collect a large number of examples, which can be very costly.

In this paper we propose to automatically *discover* object states directly from videos with object manipulations. As state changes are often caused by specific actions, we attempt to jointly discover object states and corresponding manipulations. In our setup we assume that two distinct object states are temporally separated by a manipulating action. For example, the empty and full states of a coffee cup are separated by the "pouring coffee" action, see Figure 1. Equipped with this constraint, we develop an *unsupervised* clustering approach that jointly (i) groups objects with similar appearance and consistent temporal locations

<sup>\*</sup>WILLOW-SIERRA project-teams, Département d'Informatique de l'Ecole Normale Supérieure, ENS/INRIA/CNRS UMR 8548, Paris, France.

<sup>&</sup>lt;sup>†</sup>Department of CS & OR (DIRO), Université de Montréal, Montréal.

with respect to the action and (ii) finds similar manipulating actions separating different object states in the video. Our approach exploits the complementarity of both subproblems and finds a joint solution for states and actions. We formulate our problem by adopting the discriminative clustering loss [2]. To evaluate our method, we collect a new video dataset depicting real-life object manipulation in realistic video. Given this dataset for training, our method demonstrates successful unsupervised discovery of object states and manipulating actions. We also emphasize our joint formulation and show the improvement of object state discovery by action recognition and vice versa.

# 2. Related work

Below we review related work on person-object interactions, recognizing object states, action recognition and discriminative clustering that we employ in our model.

**Person-object interactions.** Many daily activities involve person-object interactions. Modeling co-occurrences of objects and actions have shown benefits for recognizing actions in [7, 15, 21, 28, 39]. Recent work has also focused on building realistic datasets with people manipulating objects, *e.g.* in instruction videos [1, 30] or while performing daily activities [31]. We build on this work but focus on joint modeling and recognition of actions and *object states*.

**States of objects.** Prior work has addressed recognition of object attributes [11, 26, 27], which can be seen as different object states in some cases. Differently from our approach, these works typically focus on classifying still images, do not consider human actions and assume an *a priori* known list of possible attributes. Closer to our setting, Isola *et al.* [17] discover object states and transformations between them analyzing large collections of still images downloaded from the Internet. In contrast, our method does not require annotations of object states. Instead, we use the dynamics of consistent manipulations to discover object states in the video in an unsupervised manner.

Action recognition. Most of the prior work on action recognition has focused on designing features to describe time intervals of a video using motion and appearance [25, 32, 35, 36]. This is effective for actions such as dancing or jumping, however, many of our daily activities are best distinguishable by their effect on the environment. For example, opening door and closing door can look very similar using only motion and appearance descriptors but their outcome is completely different. This observation has been used to design action models in [12, 13, 37]. In [37], for example, the authors propose to learn an embedding in which a given action acts as a transformation of features of the video. In our work we localize objects and recognize changes of their states using manipulating actions as a supervisory signal. Related to ours is also the work of Fathi et al. [12] who represent actions in egocentric videos by changes of appearance of objects (also called object states), however, their method requires manually annotated temporal localization of actions in training videos. In contrast, we focus on (non-egocentric) Internet videos depicting real-life object manipulations where actions are performed by different people in a variety of challenging indoor/outdoor environments. In addition, our model does not require supervision in terms of known temporal localization of actions at training time and learns both actions and object states jointly in an unsupervised manner.

Discriminative clustering. Our model builds on unsupervised discriminative clustering methods [2, 34, 38] that group data samples according to a simultaneously learned classifier. Such methods can incorporate (weak) supervision that helps to steer the clustering towards a preferred solution [3, 9, 19]. In particular, we build on the discriminative clustering approach of [2] that has been shown to perform well in a variety of computer vision problems [3]. It leads to a quadratic optimization problem where different forms of supervision can be incorporated in the form of (typically) linear constraints. Building on this formalism, we develop a model that jointly finds object states and temporal locations of actions in the video. Part of our object state model is related to [20], while our action model is related to [5], but we consider novel constraints that posed optimization challenges, and our joint model is novel.

**Contributions.** The contributions of this work are threefold. First, we develop a discriminative clustering model that jointly discovers states of objects and temporally localizes associated manipulating actions in the video. Second, we build a new challenging dataset depicting manipulating actions changing object states in unconstrained indoor and outdoor settings. Finally, we experimentally demonstrate that our model discovers key object states and manipulating actions from input videos in an unsupervised manner.

#### 3. New dataset of manipulated objects

Due to the lack of open domain annotated datasets for the problem of discovering object states and manipulating actions we have built a new dataset that we detail below. Our dataset is composed of 630 video clips with average length of 30 seconds depicting seven different actions<sup>1</sup> manipulating five distinct objects<sup>2</sup>. The clips were obtained from different sources: the instruction video dataset introduced in [1], the Charades dataset from [31], and some additional videos downloaded from YouTube. The actions were chosen to cause visible state changes of objects. We focus on "third person" videos (rather than ego-centric) as such videos depict a variety of people in different settings

<sup>&</sup>lt;sup>1</sup>put the wheel on the car, withdraw the wheel from the car, place a plant inside a pot, open an oyster, open a refrigerator, close a refrigerator and pour coffee.

<sup>&</sup>lt;sup>2</sup>car wheel, flower pot, oyster, refrigerator and coffee cup.

and can be obtained at large-scale from YouTube. Detailed statistics of the dataset are given in Table 1.

Objects	Actions (#clips)	States	#Tracklets
wheel	{remove (47), put (46)}	{attached, detached}	5447
coffee cup	{ <b>fill</b> (57)}	{full, empty}	1819
flower pot	$\{$ <b>put plant</b> (27) $\}$	{full, empty}	2463
fridge	{ <b>open</b> (234), <b>close</b> (191)}	{open, closed}	7968
oyster	{ <b>open</b> (28)}	{open, closed}	1802

Table 1: Statistics of our new dataset of manipulated objects

Annotation. The goal of our dataset is to provide a benchmark for analyzing actions together with their associated changes in object states. The annotation described below is used to evaluate our method only, *i.e.* we do not use it at training time. We annotated the precise temporal extent of each action. We also defined a list of two states each object was going through during the action. Finally, we run automatic object detectors (pre-trained on ImageNet) for each involved object and track the detected object occurrences throughout the video. This use of an automatic object detector is motivated by two reasons: (i) annotating tracks is much faster than drawing object bounding boxes in each frame, and (ii) using real (and noisy) object detections provides a more realistic evaluation set-up. We have used a relatively low detector threshold to obtain a large number of detections with a good recall for each object state. This is to avoid situations where, for example, detecting full cup would be harder than empty cup that could bias the statistics of our data. Finally, every track is labeled with four possible states: state 1, state 2, ambiguous state or false positive detection. The ambiguous state covers the (not so common) in-between cases, such as cup half-full. In total, we have 19,499 fully annotated short object tracks, called tracklets, out of which: 35% cover state 1 or state 2, 25%are ambiguous, and 40% are false positives.

#### 4. Modeling manipulated objects

We are given a set of N clips that contain a common **manipulation** of the same *object* (such as "open an *oys*ter"). Note that manipulations in our dataset only occupy a small time interval within each clip. We also assume that we are given an *a priori* model of the corresponding object in the form of a pre-trained object detector [14]. Given these inputs, our goal is twofold: (i) precisely localize in time the extent of the action and (ii) spatially/temporally localize the manipulated object while identifying its states. This is achieved by jointly clustering the appearances of an object (such as a "oyster") appearing in all clips into two classes, corresponding to the two different states (such as "closed" and "open"), while at the same time temporally localizing a consistent "opening" action that separates the two states consistently in all clips. More formally, we formulate the problem as a minimization of a joint cost function that ties

together the action prediction in time, encoded in the assignment variable Z, with the object state discovery in space and time, defined by the assignment variable Y:

		Action localization		Object state labe	ling
		saliency of action		ordering + non ove	erlap
$Z {\in} \{0, 1\}^T$	s.t.	$\underbrace{Z\in\mathcal{Z}}$	and	$\underbrace{Y\in\mathcal{Y}}_{}$	
$\underset{Y \in \{0,1\}^{M \times 2}}{\text{minimize}}$		f(Z) + g	(Y) -	+ d(Z, Y)	(1)

where f(Z) is a discriminative clustering cost to temporally localize the action in each clip, q(Y) is a discriminative clustering cost to identify and localize the different object states and d(Z, Y) is a joint cost that relates object states and actions together. T denotes the total length of all video clips and M denotes the total number of tracked object candidate boxes (tracklets). In addition, we impose constraints  $\mathcal{Y}$  and  $\mathcal{Z}$  that encode additional structure of the problem: we localize the action with its most salient time interval per clip ("saliency"); we assume that the ordering of object states is consistent in all clips ("ordering") and that only one object is manipulated at a time ("non overlap"). In the following, we proceed with describing different parts of the model (1). We define and explain the cost function and model assumption first for the discovery of object states in Sec. 4.1, and then action localization in Sec. 4.2. We then describe and motivate the joint cost function d in Sec. 4.3. Finally, we describe the optimization procedure in Sec. 4.4.

### 4.1. Discovering object states

The goal here is to both (i) spatially localize the manipulated object and (ii) temporally identify its individual states. To address the first goal, we employ pre-trained object detectors. To address the second goal, we formulate the discovery of object states as a discriminative clustering task with constraints. We give details next. We obtain candidate object detections using standard object detectors pre-trained on large scale existing datasets such as ImageNet [8]. We assume that each clip n is accompanied with a set of  $M_n$ tracklets<sup>3</sup> of the object of interest.

We formalize the task of localizing the states of objects as a discriminative clustering problem where the goal is to find an assignment matrix  $Y_n \in \{0,1\}^{M_n \times 2}$ , where  $(Y_n)_{mk} = 1$  indicates that the *m*-th tracklet represents the object in state k. We also allow a complete row of  $Y_n$  to be zero to encode that no state was assigned to the corresponding tracklet. This is to model the possibility of false positive detections of an object, or that another object of the same class appears in the video, but is not manipulated and

<sup>&</sup>lt;sup>3</sup>In this work, we use short tracks of objects (less than one second) that we call tracklet. We want to avoid to have long tracks that continue across a state change of objects. By using the finer granularity of tracklets, our model has the ability to correct for detection mistakes within a track as well as identify more precisely the state change.



Figure 2: Given a set of clips that depict a given object being manipulated, we wish to automatically discover the main states that the object can take along with localizing the associated action. In this example, we show one video of someone filling up a coffee cup. The video starts with an empty cup (*state 1*), which is filled with coffee (**action**) to become full (*state 2*). Given imperfect object detectors, we wish to assign to the valid object candidates either the initial state or the final state (encoded in Y). Additionally, we want to localize the manipulating action (encoded in Z). Finally, we wish to enforce a consistency between states and actions through a distortion measure d(Z, Y).

thus is not undergoing any state change. In more details, our approach here is to minimize the following discriminative clustering  $\cos [2]$ :<sup>4</sup>

$$g(Y) = \min_{W_s \in \mathbb{R}^{d_s \times 2}} \underbrace{\frac{1}{2M} \|Y - X_s W_s\|_F^2}_{\text{Discriminative loss on data}} + \underbrace{\frac{\mu}{2} \|W_s\|_F^2}_{\text{Regularizer}}$$
(2)

where  $W_s$  is the object state classifier that we seek to learn,  $\mu$  is a regularization parameter and  $X_s$  is a  $M \times d_s$  matrix of features, where each row is a  $d_s$ -dimensional (state) feature vector storing features for one particular tracklet. The minimization in  $W_s$  actually leads to a convex quadratic cost function in Y (see [2]). The first term in (2) is the discriminative loss on the data that measures how easy the input data  $X_s$  is separable by the linear classifier  $W_s$  when the object state assignment is given by matrix Y. In other words, we wish to find a labeling Y for given object tracklets into two states (or no state) so that their appearance features X are easily separated by a linear classifier. To steer the cost towards the right solution, we employ the following constraints (encoded by  $Y \in \mathcal{Y}$  in (1)).

**Only one object is manipulated : non overlap constraint.** We assume that the manipulating modifies the state of only one object in the video. However, in practice, it is common to have multiple (spatially diverse) tracklets that occur at the same time, for example, due to a false positive detection in the same frame. To overcome this issue, we impose that at most one tracklet can be labeled as belonging to *state 1* or *state 2* at any given time. We refer to this constraint as "*non overlap*" in problem (1).

state  $1 \rightarrow Action \rightarrow state 2$ : ordering constraints. We assume that the manipulating action transforms the object from an initial state to a final state and that both states are present in each video. This naturally introduces two constraints. The first one is the ordering constraints on the labeling  $Y_n$ , *i.e.* the state 1 should occur before state 2 in each video. The second constraint imposes that we have at *least one* tracklet labeled as state 1 and at least one tracklet labeled as state 2. We call this last constraint the "**at least one**" constraint in contrast to forcing "exactly one" ordered prediction as previously proposed in a discriminative clustering approach on video for action localization [5]. While our novel constraint brings new optimization challenges (see Sec. 4.4), we show in the experimental section 5.2 that it is key to get good results for our problem.

### 4.2. Action localization

We now detail the cost and constraint for action localization. Given N input videos, our goal is to localize time intervals corresponding to *one common* action in all clips.<sup>5</sup> More formally, the n-th clip is composed of a video stream decomposed in  $T_n$  chunks of frames  $(x_t^n)_{t=1}^{T_n}$ . The goal is to find an assignment matrix  $Z_n \in \{0, 1\}^{T_n}$  for each clip n,

<sup>&</sup>lt;sup>4</sup>We concatenate all the decision variables  $Y_n$  into one  $M \times 2$  matrix Y.

<sup>&</sup>lt;sup>5</sup>Our action model is equivalent to the one of [5] applied to only *one action*. This model is incomplete as the clips in our dataset can contain other actions that do not manipulate the object of interest. Our key contribution is to propose a *joint formulation* that links this simple action model with the object state prediction model (see Sec. 4.3), thereby resolving the ambiguity on actions, as we will see in our experiments.

where  $Z_{nt} = 1$  encodes that the *t*-th time interval of video is assigned to an action and  $Z_{nt} = 0$  encodes that no action is detected in interval *t*. We further assume that every time interval *t* has an associated  $d_v$ -dimensional feature vector stored in the *t*-th row of the concatenated  $T \times d_v$  (video) feature matrix  $X_v$ . Similar to localizing object states, we formulate action localization as a discriminative clustering problem with the following cost

$$f(Z) = \min_{W_v \in \mathbb{R}^{d_v}} \underbrace{\frac{1}{2T} \|Z - X_v W_v\|_F^2}_{\text{Discriminative loss on data}} + \underbrace{\frac{\lambda}{2} \|W_v\|_F^2}_{\text{Regularizer}}, \quad (3)$$

where  $W_v$  is a classifier for the considered action that we seek to learn and  $\lambda$  is a regularization parameter. Similar to the labeling of object states, we wish to find Z with action labels for video frames so that the appearance features  $X_v$ of the action frames are easily separated by a linear classifier from the non-action frames in the video. To better steer the cost to the desired solution, we employ the following constraint (encoded by  $Z \in \mathcal{Z}$  in (1)).

Action saliency constraint. An action can often be well represented by a short and discriminative interval. We thus force our model to predict *exactly* one time interval for an action per clip. This can be interpreted as seeking the most salient interval that best describes the action, an approach for actions that was shown to be beneficial in a weakly supervised setting [5].

#### 4.3. Linking actions and object states

Actions in our model are directly related to changes in object states. We therefore want to enforce consistency between the two tasks. To do so, we design a joint cost function that operates on the action video labeling  $Z_n$  and the state tracklet assignment  $Y_n$  for each clip. We want to impose a constraint that the action occurs in between the presence of the two different object states. In other words, we want to penalize when state 1 is detected *after* the action happens (type I), or when state 2 is triggered *before* the action occurs (type II). In the following, we derive a symmetric penalty on both  $Y_n$  and  $Z_n$  that enforces this constraint in a soft manner.

**Joint cost definition.** We propose the following joint cost function for each clip:

$$d(Z_n, Y_n) = \sum_{y \in \mathcal{S}_1(Y_n)} [t_y - t_{Z_n}]_+ + \sum_{y \in \mathcal{S}_2(Y_n)} [t_{Z_n} - t_y]_+, \quad (4)$$

where  $t_{Z_n}$  and  $t_y$  are the times when the action  $Z_n$  and the tracklet y occurs in a clip n, respectively.  $S_1(Y_n)$  and  $S_2(Y_n)$  respectively contains the tracklets that have been assigned to state 1 and to state 2 in the n-th clip. Finally the  $[.]_+$  is the positive part function. In other words, the function simply penalizes every wrong decision on  $Y_n$  of (type I) (left sum in (4)) or (type II) (right sum in (4)) by the amount of time that separates the inconsistent tracklet decision with the single action decision for the clip.<sup>6</sup> The global joint cost is simply the sum over all clips weighted by a scaling hyperparameter  $\nu > 0$ :

$$d(Z,Y) = \nu \frac{1}{T} \sum_{n=1}^{N} d(Z_n, Y_n).$$
 (5)

# 4.4. Optimization.

Here we describe our solution to the joint objective (1).

**Convex hull relaxation.** Problem (1) is NP-hard due to the integer constraints. Inspired by the approach of [4] that was successful to approximate combinatorial optimization problems, we propose to use the tightest convex relaxation of the feasible subset of binary matrices by taking its convex hull. As our variables now can take values in [0, 1], we also have to propose a consistent extension for the different cost functions to handle fractional values as input. For the cost functions f and g, we can directly take their expression on the relaxed set as they are already expressed as (convex) quadratic functions. Similarly for the joint cost function d in (4), we use its natural bilinear relaxation:

$$d(Z_n, Y_n) = \sum_{i=1}^{M_n} \sum_{t=1}^{T_n} \left( (Y_n)_{i1} Z_{nt} [t_{ni} - t]_+ + (Y_n)_{i2} Z_{nt} [t - t_{ni}]_+ \right), \quad (6)$$

where  $t_{ni}$  gives the video time of tracklet *i* in clip *n*. This relaxation is equal to the function introduced in (4) on the integer points. However, it is not jointly convex in *Y* and *Z*, thus one has to design a good optimization technique to get good candidate solutions.

Joint optimization using Frank-Wolfe. When dealing with a constrained optimization problem for which it is easy to solve linear programs but difficult to project on the set, the Frank-Wolfe algorithm is well adapted and presents numerous advantages [18, 24]. It is exactly the case for our relaxed problem, where the linear program over the convex hull of the feasible integer matrices can be solved efficiently via dynamic programming. Moreover, [23] recently showed that the Frank-Wolfe algorithm with line-search converges to a stationary point for non-convex objectives at a rate of  $O(1/\sqrt{k})$ . We thus use this algorithm for the joint optimization of (1). As the objective is quadratic, we can perform exact line-search analytically, which speeds up convergence in practice. Finally, in order to get a good initialization for both variables Z and Y, we first optimize separately f(Z)and q(Y) (without the non-convex d(Z, Y)), which are both convex functions.

**Dynamic program for the tracklets.** In order to apply the Frank-Wolfe algorithm, we need to solve a linear

<sup>&</sup>lt;sup>6</sup>We see here that our saliency constraint modeling choice on actions (one action per clip) simplifies this joint formulation.

program (LP) over our set of constraints. The LP for localizing actions is simple (taking the max over all time steps in a clip). On the other hand, the previous dynamic programming solutions proposed in related prior work [5, 20] cannot handle our proposed constraints for the tracklet assignment problem. We thus derived a novel dynamic programming solution (described in Appendix B).

**Rounding method.** Once we get a candidate solution of the relaxed problem, we have to round it back to an integer solution in order to make predictions. [5] and [1] observed that using the learned  $W^*$  classifier to round gave better results than alternatives. We extend this approach to our novel joint setup by proposing the following new rounding possibility. We optimize problem (1) but fix the values of W in the discriminative clustering costs. Specifically, we minimize the following cost function over the integer points  $Z \in \mathcal{Z}$  and  $Y \in \mathcal{Y}$ :

$$\frac{1}{2T} \|Z - X_v W_v^*\|_F^2 + \frac{1}{2M} \|Y - X_s W_s^*\|_F^2 + d(Z, Y),$$
(7)

where  $W_v^*$  and  $W_s^*$  are the classifiers weights obtained at the end of the relaxed optimization. Because  $y^2 = y$  when y is binary, (7) is actually a *linear* objective over the binary matrix  $Y_n$  for  $Z_n$  fixed. Thus we can optimize (7) *exactly* by solving a dynamic program on  $Y_n$  for each of the  $T_n$ possibilities of  $Z_n$ , yielding  $O(M_nT_n)$  time complexity per clip (see Appendix C).

## 5. Experiments

We first describe the object tracking pipeline including the feature representation of tracklets and videos (Section 5.1). Then we present (Section 5.2) the experimental results that are divided into two main parts: (i) the evaluation of state discovery and (ii) the evaluation of temporal action localization.

### 5.1. Object representation and tracking

**Object detection and tracking.** In order to obtain detectors for the five objects, we finetune the FastRCNN network [14] with training data from ImageNet [8]. We use bounding box annotations from ImageNet when available (*e.g.* the "wheel" class). For the other classes, we manually labeled more than 500 instances per class. In our set-up with only moderate amount of training data, we observed that class-agnostic object proposals combined with FastRCNN performed better than FasterRCNN [29]. In detail, we use geodesic object proposals [22] and set a relatively low object detection threshold (0.4) to have good recall. We track objects using a generic KLT tracker from [3]. The tracks are then post-processed into shorter tracklets that last about one second. This is to avoid very long tracks that may contain several different object states.

**Object tracklet representation.** For each detected object, represented by a set of bounding boxes over the course of the tracklet, we compute a CNN feature from each (extended) bounding box that we then average over the length of the tracklet to get the final representation. We use ROI pooling [29] at block 5 of ResNet50 [16] with a spatial 2 by 2 subdivision. The ROI pooling allows to capture some context around the object which is important for some cases (*e.g. wheel* "on" or "off" the car). The resulting feature descriptor of each object tracklet is 8,192 dimensional.

**Representing video for recognizing actions.** Each video is divided into chunks of 10 frames that are represented by a motion and appearance descriptor averaged over 30 frames. For the motion we use a 2,000 dimensional bagof-word representation of histogram of local optical flow (HOF) obtained from Improved Dense Trajectories [36]. Following [1], we add an appearance vector that is obtained from a 1,000 dimensional bag-of-word vector of conv5 features from VGG16 [33]. This results in a 3,000 dimensional feature vector for each chunk of 10 frames.

#### 5.2. State discovery and localization of actions

**Experimental setup.** For each action of the dataset (see Table 1), we truncate the video 10 seconds before and after the annotated action of interest. This is to simplify the set-up to contain video clips that contain one action that changes the state of the object of interest. However, each clip may contain other actions that affect other objects or actions that do not affect any object at all (e.g. walking / jumping). This results in a dataset of 630 video clips that are about 20 seconds long and are annotated with temporal localization of 7 different actions together with about 20,000 object tracklets fully annotated with their states. We believe understanding this simplified set-up with coarsely pre-segmented clips is a necessary step towards the ambitious goal of relating actions and their effects on states of objects in continuous video streams. However, the proposed method is applicable outside of this set-up in a more realistic scenario where the clips are pre-selected automatically based on text metadata from, for example, natural language narration [1].

**Evaluation metric: average precision.** For all variants of our method, we use the rounded solution that reached the smallest objective during optimization. We evaluate these predictions with a precision score averaged over all the videos. A temporal action localization is said to be correct if it falls within the ground truth time interval. Similarly, a state prediction is correct if it matches the ground truth state.<sup>7</sup> Note that "precision" metric is reasonable in our set-up as our method is forced to predict in all videos, i.e. the recall level is fixed to all videos and the method cannot produce high precision with low recall. We provide

<sup>&</sup>lt;sup>7</sup>In particular, we count "*ambiguous*" labels as incorrect.

	Method	put wheel	remove wheel	fill pot	open oyster	fill coff.cup	open fridge	close fridge	Average
State discovery	(a) Constraints only	0.35	0.38	0.35	0.36	0.31	0.29	0.42	0.35
	( <b>b</b> ) Kmeans	0.25	0.12	0.11	0.23	0.14	0.19	0.22	0.18
	(c) Salient state only	0.33	0.60	0.19	0.25	0.14	0.43	0.39	0.33
	(d) At least one state only	0.51	0.65	0.33	0.48	0.28	0.45	0.35	0.44
	(e) Joint model	0.42	0.63	0.48	0.59	0.24	0.50	0.49	0.48
	( <b>f</b> ) Joint model + det. scores.	0.47	0.65	0.50	0.61	0.44	0.46	0.43	0.51
	(g) Joint + GT act. feat.	0.59	0.63	0.56	0.50	0.32	0.51	0.50	0.52
Action localization	(i) Chance	0.31	0.20	0.15	0.11	0.40	0.23	0.17	0.22
	(ii) [5]	0.22	0.13	0.15	0.07	0.33	0.35	0.21	0.21
	(iii) [5] + object cues	0.26	0.17	0.15	0.14	0.84	0.34	0.36	0.32
	(iv) Joint model	0.57	0.58	0.33	0.32	0.83	0.48	0.37	0.50
	(v) Joint + GT stat. feat.	0.72	0.66	0.41	0.43	0.84	0.50	0.45	0.57

Table 2: State discovery (top) and action localization results (bottom).

additional results using in alternative evaluation set-up using both precision and recall of the learnt object state and action classifiers in Appendix D.

**Hyperparameters.** In all methods that involve a discriminative clustering objective, we used  $\lambda = 10^{-2}$  (action localization) and  $\mu = 10^{-4}$  (state discovery) for all 7 actions. For joint methods that optimize (1), we set the weight  $\nu$  of the distortion measure (5) to 1.

State discovery results. Results are shown in the top part of Table 2. In the following, we refer to "State only" whenever we use our method without looking at the action cost or the distortion measure (1). We compare to two baselines for the state discovery task. Baseline (a) evaluates chance performance by running our "State only" method with random features for representing tracklets, using only constraints (at least one ordering + non overlap) to make predictions. Baseline (b) performs K-means clustering of the tracklets with K = 3 (2 clusters for the states and 1 for false positives). We report performance of the best assignment for the solution with the lowest objective after 10 different initializations. The next two methods are "State only" variants. Method (c) corresponds to a replacement of the "at least one constraint" by an "exactly one constraint" while the method (d) uses our new constraint. Finally, we report three joint methods that use our new joint rounding technique (7) for prediction. Method (e) corresponds to our joint method that optimizes (1). Method (f) is a simple improvement taking into account object detection score in the objective (details below). Method (g) is our joint method but using the action ground truth labels as video features in order to test the effect of having perfect action localization for the task of object state discovery. We first note that method (d) consistently outperforms (c), thus highlighting the importance of the "at least one" constraint for modeling object states. While the saliency approach (taking only the most confident detection per video) was useful for action modeling in [5], it is less suitable for our set-up where multiple tracklets can be in the same state. The joint approach with actions (e) outperforms the "State only" method (d)

on 4 actions out of 7 and obtains better average performance, confirming the benefits of joint modeling of actions and object states. Using ground truth action locations further improves results (cf. (g) with (e)). Our unsupervised approach (e) performs not much lower compared to using ground truth actions (g), except for the states of the coffee cup (empty/full). In this case we observe that a high number of false positive detections confuses our method (see a false positive example in the second row of Figure 3). A simple way to address this issue is to add the object detection score into the objective of our method, which then prefers to assign object states to higher scoring object candidates further reducing the effect of false positives. This can be done easily by adding a linear cost reflecting the object detection score to objective (1). We denote this modified method "(f)Joint model + det. scores". Note, in particular, the significantly improved performance on "fill coffee cup" and the best overall average performance.

Action localization results. We compare our method to three different baselines and give results in the bottom part of Table 2. Baseline (i) corresponds to chance performance, where the precision for each clip is simply the proportion of the entire clip taken by the ground truth time interval. Baseline (ii) is the method introduced in [5] used here with only one action. It also corresponds to a special case of our method where the object state part of the objective in equation (1) is turned off (salient action only). Interestingly, this baseline is actually worse than chance for several actions. This is because without additional information about objects this method localizes other common actions in the clip and not the action manipulating the object of interest. This also demonstrates the difficulty of our experimental set-up where the input video clips often contain multiple different actions. To address this issue, we also evalute baseline (iii), which complements [5] with the additional constraint that the action prediction has to be within the first and the last frame where the object of interest is detected, improving the overall performance above chance. Our joint approach (iv)consistently outperforms these baselines on all actions, thus



Figure 3: Typical failure cases for "*removing car wheel*" (top) and "*fill coffee cup*" (middle, bottom) actions. Yellow indicates correct predictions; red indicates mistakes. Top: the removed wheel is incorrectly localized (right). Middle: the "empty cup" is incorrectly localized (left). Bottom: In this case, both object tracklets are annotated as "ambiguous" in the ground truth as they occur *during* the pouring action and hence the predictions, while they appear reasonable, are deemed incorrect.

showing again the strong link between object states and actions. Finally, the approach (v) is the analog of method (g)for action localization where we use ground truth state labels as tracklet features in our joint formulation showing that the action localization can be further improved with better object state descriptors.

**Benefits of joint object-action modeling.** Interestingly, we observe that the joint modeling of object states and actions benefits both tasks. This effect is even more pronounced for actions. Intuitively, knowing perfectly the object states reduces a lot the search space for action localization. In addition, despite the recent major progress in object recognition using convolutional neural networks, action recognition still remains a hard problem with much room for improvement.

**Failure cases.** We observed two main types of failures, illustrated in Figure 3. The first one occurs when a false positive object detection consistently satisfies the hypothesis of our model in multiple videos (the top two rows in Figure 3). The second typical failure mode is due to ambiguous labels (bottom row in Figure 3). This highlights the difficulty in annotating ground truth for long actions such as "*pouring coffee*".



Figure 4: Qualitative results for joint action localization (middle) and state discovery (left and right) for 6 actions from our dataset (see Fig. 1 for *"fill coffee fup"* action).

# 6. Conclusion and future work

We have described a joint model that relates object states and manipulating actions. Given a set of input videos, our model both localizes the manipulating actions *and* discovers the corresponding object states without additional supervision. We have demonstrated that our joint approach improves performance of both object state recognition and action recognition. More generally, our work provides evidence that actions should be modeled in the larger context of goals and effects. While we have focused on modeling a single object state change per video, our work opens up the possibility of large-scale learning of manipulating actions from video sequences with multiple actions and objects. **Acknowledgments** This research was supported in part by a Google Research Award, ERC grants Activia (no. 307574) and LEAP (no. 336845), and the CIFAR Learning in Machines & Brains program.

## References

- J.-B. Alayrac, P. Bojanowski, N. Agrawal, I. Laptev, J. Sivic, and S. Lacoste Julien. Unsupervised learning from narrated instruction videos. In *CVPR*, 2016. 2, 6, 10, 11
- [2] F. Bach and Z. Harchaoui. DIFFRAC: A discriminative and flexible framework for clustering. In *NIPS*, 2007. 2, 4, 12
- [3] P. Bojanowski, F. Bach, I. Laptev, J. Ponce, C. Schmid, and J. Sivic. Finding actors and actions in movies. In *ICCV*, 2013. 2, 6
- [4] P. Bojanowski, R. Lajugie, F. Bach, I. Laptev, J. Ponce, C. Schmid, and J. Sivic. Weakly supervised action labeling in videos under ordering constraints. In *ECCV*, 2014. 5, 10, 11
- [5] P. Bojanowski, R. Lajugie, E. Grave, F. Bach, I. Laptev, J. Ponce, and C. Schmid. Weakly-supervised alignment of video with text. In *ICCV*, 2015. 2, 4, 5, 6, 7, 11, 12
- [6] T. F. Brady, T. Konkle, A. Oliva, and G. A. Alvarez. Detecting changes in real-world objects: The relationship between visual long-term memory and change blindness. *Communicative and Integrative Biology*, 2006. 1
- [7] V. Delaitre, J. Sivic, and I. Laptev. Learning person-object interactions for action recognition in still images. In *NIPS*, 2011. 2
- [8] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image database. In CVPR, 2009. 3, 6
- [9] C. Doersch, S. Singh, A. Gupta, J. Sivic, and A. A. Efros. What makes Paris look like Paris? *SIGGRAPH*, 2012. 2
- [10] K. Duan, D. Parikh, D. Crandall, and K. Grauman. Discovering localized attributes for fine-grained recognition. In *CVPR*, pages 3474–3481, 2012. 1
- [11] A. Farhadi, I. E. Lim, D. Hoiem, and D. Forsyth. Describing objects by their attributes. In *CVPR*, 2009. 1, 2
- [12] A. Fathi and J. M. Rehg. Modeling actions through state changes. In CVPR, 2013. 2
- [13] B. Fernando, E. Gavves, M. J. Oramas, A. Ghodrati, and T. Tuytelaars. Modeling video evolution for action recognition. In *CVPR*, 2015. 2
- [14] R. Girshick. Fast R-CNN. In ICCV, 2015. 3, 6
- [15] A. Gupta, A. Kembhavi, and L. S. Davis. Observing humanobject interactions: Using spatial and functional compatibility for recognition. *PAMI*, 31(10):1775–1789, 2009. 2
- [16] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6
- [17] P. Isola, J. J. Lim, and E. H. Adelson. Discovering states and transformations in image collections. In CVPR, 2015. 1, 2
- [18] M. Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In *ICML*, 2013. 5
- [19] A. Jain, A. Gupta, M. Rodriguez, and L. Davis. Representing videos using mid-level discriminative patches. In *CVPR*, pages 2571–2578, 2013. 2

- [20] A. Joulin, K. Tang, and L. Fei-Fei. Efficient image and video co-localization with Frank-Wolfe algorithm. In *ECCV*, 2014. 2, 6, 10
- [21] H. Kjellström, J. Romero, and D. Kragić. Visual objectaction recognition: Inferring object affordances from human demonstration. *CVIU*, 115(1):81–90, 2011. 2
- [22] P. Krhenbhl and V. Koltun. Geodesic object proposals. In ECCV, 2014. 6
- [23] S. Lacoste-Julien. Convergence rate of Frank-Wolfe for nonconvex objectives. arXiv preprint, 2016. 5
- [24] S. Lacoste-Julien and M. Jaggi. On the global linear convergence of Frank-Wolfe optimization variants. In *NIPS*, 2015.
  5
- [25] I. Laptev, M. Marszalek, C. Schmid, and B. Rozenfeld. Learning realistic human actions from movies. In *CVPR*, 2008. 2
- [26] D. Parikh and K. Grauman. Relative attributes. In ICCV, 2011. 2
- [27] G. Patterson, C. Xu, H. Su, and J. Hays. The SUN attribute database: Beyond categories for deeper scene understanding. *IJCV*, 2014. 2
- [28] H. Pirsiavash and D. Ramanan. Detecting activities of daily living in first-person camera views. In CVPR, 2012. 2
- [29] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 6
- [30] O. Sener, A. Zamir, S. Savarese, and A. Saxena. Unsupervised semantic parsing of video collections. In *ICCV*, 2015. 2
- [31] G. A. Sigurdsson, G. Varol, X. Wang, A. Farhadi, I. Laptev, and A. Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *ECCV*, 2016. 2
- [32] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In NIPS, 2014. 2
- [33] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, 2014. 6
- [34] S. Singh, A. Gupta, and A. A. Efros. Unsupervised discovery of mid-level discriminative patches. In ECCV, 2012. 2
- [35] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3D convolutional networks. In *ICCV*, 2015. 2
- [36] H. Wang and C. Schmid. Action recognition with improved trajectories. In *ICCV*, 2013. 2, 6
- [37] X. Wang, A. Farhadi, and A. Gupta. Actions ~ transformations. In CVPR, 2016. 2
- [38] L. Xu, J. Neufeld, B. Larson, and D. Schuurmans. Maximum margin clustering. In *NIPS*, 2004. 2
- [39] B. Yao, X. Jiang, A. Khosla, A. L. Lin, L. Guibas, and L. Fei-Fei. Human action recognition by learning bases of action attributes and parts. In *ICCV*, 2011. 2

# A. Outline of Supplementary Material

This supplementary material provides additional details and quantitative results. The organization is as follows. In Appendix B, we describe the dynamic program that we use to solve the linear program over the track constraints defined in Section 4.1 as needed for the Frank-Wolfe optimization algorithm. In Appendix C, we detail how we implement the *new* joint rounding method (7) that was briefly introduced in Section 4.3. In Appendix D, we complement the results from Section 5.2 with an additional set of results that use an alternative evaluation setup, where we evaluate the quality of the classifiers that are implicitly learned by our method.

# **B.** Dynamic program for the tracklets

The track constraints defined in Section 4.1 introduce new challenges compared to the previous related work [1, 4, 20]. Recall that there are three main components in the constraints. First, we assume that only one object is manipulated at a given time. Thus *at most one* tracklet can be assigned to a state at a given time. This constraint is referred to as the **non-overlap constraint**. Second, we have the **ordering constraint** that imposes that *state 1* always happens *before state 2*. The last constraint imposes that we have **at least one** tracklet labeled as *state 1* and *at least one* tracklet labeled as *state 2*. We need to be able to minimize linear functions over this set of constraints in order to use the Frank-Wolfe algorithm. More precisely, as the constraints decompose over the different clips, we can solve independently for each clip *n* the following linear problem:

where  $C_n \in \mathbb{R}^{M_n \times 2}$  is a cost matrix that typically comes from the computation of the gradient of the cost function at the current iterate. In order to solve this problem, we use a dynamic program approach that we explain next. Recall that we are given  $M_n$  tracklets  $(y_i)_{i=1}^{M_n}$  and our goal is to output the  $Y_n$  matrix that assigns to each of these tracklets either state 1, state 2 or no state at all while respecting the constraints. The whole method is illustrated in Figure 5 with a toy example.

**Non-overlap data structure for the tracklets.** We first preprocess the tracklets to build an auxiliary data-structure that is used to enforce the non-overlap constraint between the tracklets, as illustrated in Figure 5a. First, we sort and index each tracklet by their beginning time, and add two fictitious tracklets:  $y_0$  as the starting tracklet and  $y_f$  as the ending tracklet. These two tracklets are used to start and terminate the dynamic program. If all the tracklets were sequentially ordered without any overlap in time, then we could simply make a decision for each of them sequentially as was done in previous work on action localization for example (one decision per time step) [4]. To enforce the non-overlap constraint, we force the decision process to choose *only one* possible successor among the group of overlapping valid immediate successors of a tracklet. For each tracklet  $y_i$ , we thus define its (smallest) set of "valid successors" as the earliest tracklet  $y_j^8$  after  $y_i$  that is also non-overlapping with  $y_i$ , as well as any other tracklet  $y_l$  for l > j that is overlapping with  $y_j$  (thus giving the earliest valid group of overlapping tracklets). The valid successors are illustrated by red dotted arrows in Figure 5a. For example, the valid successors of  $y_1$  are  $y_3$  (the earliest one that is non-overlapping) as well as  $y_4$  (which overlaps with  $y_3$  thus forming an overlapping group). Skipping a tracklet in this decision process means that we assign it to zero (which trivially always satisfies the non-overlapping constraint); whereas once we choose a tracklet to potentially assign it to state 1 or 2, we cannot visit any overlapping tracklet by construction of the valid successors, thus maintaining the non-overlap constraint.

**Dynamic program.** The dynamic programming approach is used when we can solve a large problem by solving a sequence of inclusive subproblems that are linked by a simple recursive formula and that use overlapping solutions (which can be stored in a table for efficiency). In terms of implementation, [4] encoded their dynamic program as finding an optimal path inside a cost matrix. This approach is particularly suited when the update cost rule depends *only* on the arrival entry in the cost matrix as opposed to be transition dependent. As we will show below, we can encode the solution to our problem in a way that satisfies this property. We therefore use the framework of [4] by casting our problem as a search for an optimal path inside a cost matrix  $\tilde{C}_n$  illustrated in Figure 5b, and where the valid transitions encode the possible constraints.

One main difference with [4] is that we have to deal with the challenging at least one constraint in the context of ordered labels. To do so, we can filter further the set of valid decisions by using "memory states" that encode in which of the following three situations we are: (i) that state 1 has not yet been visited, (ii) that state 1 has already been visited, but state 2 has not yet been visited (and thus that we can either come back to state 1 or go to state 2) and (iii) that both states have been visited. These memory states can be encoded by interleaving complete rows of 0s in between columns of  $C_n$  stored as rows, to obtain the  $5 \times M_n$  matrix  $\tilde{C}_n$ . These new rows encode the three different memory states previously described when making a prediction of 0 for a specific tracklet, and we enforce the correct memory semantic by only allowing a path to move to the same row or the row immediately below, except for state 1 which can also move directly to state 2 (two rows below), and the middle "between state 1/2" row, where one can go up one row additionally to state 1. Finally, the valid transitions between columns (tracklets) are given by the valid successors data structure as given in Figure 5a to encode the non-overlap constraints. Combining these two constraints (at least one ordering and nonoverlap), we illustrate with grey arrows in Figure 5b the possible transitions from the states along the path in red. To describe the dynamic program recursion below, we need to go the opposite direction from the successors, and thus we say that  $y_i$  is a *predeces*sor of  $y_i$  if and only if  $y_i$  is a successor of  $y_j$ .

To perform the dynamic program, we maintain a matrix  $D_n$  of the same size as  $\tilde{C}_n$  where  $D_n(k, i)$  contains the minimal valid path cost of going from  $(1, y_0)$  to  $(k, y_i)$  inside the cost matrix  $\tilde{C}_n$ . To define the cost update recursion to compute  $D_n(k, i)$ , let

<sup>&</sup>lt;sup>8</sup>Earliest means the smallest j.



(a) Non-overlap data structure for tracklets

(b) Cost matrix  $\tilde{C}_n$  for the dynamic program

Figure 5: In (a), we provide an illustration of a possible situation for the tracklets.  $y_0$  and  $y_f$  are two fictitious tracklets that encode the beginning and end of the video. Each tracklet is indexed based on its beginning time. The time overlap between tracklets is shown by the grey color. We specify for each tracklet its possible successors by the dotted red arrows (see main text). Finally an admissible labeling is illustrated by yellow tags where  $y_1$  and  $y_5$  have both been assigned to state 1 and  $y_6$  to state 2. In (b), we give an illustration of our approach to solve (8) with a dynamic program. We display the modified cost matrix  $\tilde{C}_n$  (see main text). A valid path has to go from the green dot  $(y_0)$  to the red dot  $(y_f)$ . The light yellow entries show part of the  $C_n$  matrix that are inserted in  $\tilde{C}_n$ , whereas white entries encode the rows of 0s that are inserted to impose the **at least one** ordering constraint. The red arrows specify an example optimal path inside the matrix. The red entries display the tracklets that have been assigned to state 1  $(y_1 \text{ and } y_5)$  or state 2  $(y_6)$  (equivalent to putting ones in the appropriate corresponding entries in  $Y_n$ ). Finally, the grey arrows display the possible valid transitions that can be made for the entries along the red path, for clarity. We see for example that from  $(2, y_1)$ , there are 6 possible transitions: two column choices from the two red arrows from  $y_1$  in (a) encoding the non-overlap constraint; and three row choices encoding the valid transition from "state 1" (corresponding to the choice "state 1", 0 or "state 2" for the next tracklet) encoding the "at least one" ordering constraint.

P(k, i) be the set of tuples (l, j) for which it is possible to go from (l, j) to  $(k, y_i)$  according to the rules described above. The update rule is then as follows:

$$D_n(k,i) = \min_{(l,j) \in P(k,i)} D_n(l,j) + \tilde{C}_n(k,y_i).$$
(9)

As we see here, the added cost depends *only* on the arrival entry  $\tilde{C}_n(k, y_i)$ . We can therefore use the approach of [4] and only consider entry costs rather than edge costs. Thanks to our indexing property (tracklets are sorted by the beginning time), we can update the dynamic program matrix by filling each column of  $D_n$  one after the other. Once this update is finished, we back-track to get the best path by starting from the ending track (predecessors of  $y_f$ ) at the last row (to be sure that both states have been visited) that has the lowest score in the  $D_n$  matrix. The total complexity of this algorithm is of order  $\mathcal{O}(M_n)$ .

# C. Joint cost rounding method

Recall that we propose to use a convex relaxation approach in order to obtain a candidate solution of main problem (1). Thus, we need to *round* the relaxed solution afterward in order to get a valid integer solution. We propose here a new rounding that is adapted to our joint problem. We referred to this rounding as the **joint cost rounding** (see Section 4.4 of main paper). This rounding is inspired by [5, 1]. They observe that using the learned  $W^*$ 

classifier to round gives them better solutions, both in terms of objective value and performance. We propose to use its natural extension for our joint model. We first fix the classifiers for actions  $W_a$  and for states  $W_s$  to their relaxed solution  $(W_a^*, W_s^*)$ and find, for each clip n, the couple  $(Z_n, Y_n)$  that minimizes the joint cost (7). To do so, we observe that we can enumerate all  $T_n$ possibilities for  $Z_n$ , and solve for each of them the minimization of the joint cost with respect to  $Y_n$ . The minimization with respect to  $Y_n$  can be addressed as follows. First, we observe that the distortion function (6) is bilinear in  $(Z_n, Y_n)$ . Let  $Z_n$  be a  $T_n \times 1$  vector, and let  $\mathbf{1}_2$  be a vector of ones of length 2. We can actually write:  $d(Z_n, Y_n) = \text{Tr}((B_n Z_n \mathbf{1}_2^{\top})^{\top} Y_n)$  for some matrix  $B_n \in \mathbb{R}^{M_n \times T_n}$ . Thus, when  $Z_n$  is fixed, the joint term  $d(Z_n, Y_n)$  is actually a simple *linear* function of  $Y_n$ . In addition, the quadratic term in  $Y_n$  coming from (2) is also *linear* over the integer points (using the fact that  $y^2 = y$  for  $y \in \{0, 1\}$ ). Thus, when  $Z_n$ ,  $W_a^*$  and  $W_s^*$  are fixed, the minimization over  $Y_n$  is a linear program (8) that we solve using our dynamic program from the previous section. The final algorithm is given in Algorithm 1. Its complexity is of order  $\mathcal{O}(T_n M_n)$ .

# D. Additional evaluation using another metric

In the main paper, we focus on the task of action time localization and state discovery predictions for a given video. Thus, we

	Method	put wheel	remove wheel	fill pot	open oyster	fill coff.cup	open fridge	close fridge	mAP
	(a) Chance	0.31	0.34	0.29	0.20	0.19	0.31	0.30	0.28
State discovery	(c) Salient state only	0.27	0.53	0.26	0.35	0.16	0.35	0.36	0.33
	(d) At least one state only	0.34	0.58	0.28	0.38	0.29	0.32	0.36	0.36
	(e) Joint model	0.32	0.62	0.32	0.57	0.25	0.43	0.47	0.42
	$(\mathbf{f})$ Joint model + det. scores.	0.37	0.64	0.47	0.67	0.28	0.49	0.52	0.52
	(g) Joint + GT act. feat.	0.33	0.62	0.36	0.59	0.26	0.46	0.48	0.44
Action localization	(i) Chance	0.34	0.21	0.16	0.12	0.40	0.25	0.18	0.23
	(ii) [5]	0.25	0.16	0.13	0.09	0.36	0.26	0.19	0.21
	(iii) [5] + object cues	0.30	0.21	0.16	0.16	0.71	0.33	0.27	0.31
	(iv) Joint model	0.50	0.40	0.22	0.23	0.75	0.34	0.28	0.39
	$(\mathbf{v})$ Joint + GT stat. feat.	0.59	0.50	0.27	0.27	0.77	0.38	0.31	0.44

Table 3: Evaluation of object state and action classifiers (average precision).

action localization.

### Algorithm 1 Joint cost rounding for video n

Get  $W_s^*$  and  $W_a^*$  from the relaxed problem. Initialize  $Z^*$ ,  $Y^*$  and val<sup>\*</sup> =  $+\infty$ . # Loop over all possibilities for  $Z_n$  (saliency) for  $t \text{ in } 1 : T_n$  do  $Z \leftarrow \operatorname{zeros}(T_n, 1)$ # Set the t-th entry of Z to 1  $Z_t \leftarrow 1$ # Definition of the cost matrix  $C_n \leftarrow \frac{1}{2M} (\text{ones}(M_n, 2) - 2X_s W_s) + \frac{\nu}{T} B_n Z \mathbf{1}_2^\top$ # Dynamic program for the tracks  $Y_{\min} \leftarrow \arg\min_{Y \in \mathcal{Y}_n} \operatorname{Tr}(C_n^T Y)$ # Cost computation  $\begin{aligned} & \cos t_{Z} \leftarrow \frac{1}{2T} \| Z - X_a W_a \|_F^2 \\ & \cos t_Y \leftarrow \frac{1}{2M} \| Y_{\min} - X_s W_s \|_F^2 \\ & \cos t_{ZY} \leftarrow \frac{\nu}{T} d(Z, Y_{\min}) \end{aligned}$ # Update solution if better val  $\leftarrow \cos t_Z + \cos t_Y + \cos t_{ZY}$ if  $val < val^*$  then  $Z^* \leftarrow Z$  $Y^* \leftarrow Y_{\min}$  $val^* \leftarrow val$ end if end for return  $Z^*, Y^*$ 

and evaluate them with an AP (average precision) score. For actions, this score is simply the area under the precision-recall curve for the given action, where positives are all time intervals that were labeled as part of the action and negatives are all other intervals. For object states, we average the AP scores over the two states to get the final AP score. We compare the different methods that uses discriminative clustering objective together with the chance baselines in Table 3. Results are in line with what is shown in the result section (Section 5.2) of the main paper. This confirms the superiority of the joint approach both for state discovery *and* 

have reported performance in Table 2 with a metric that was evaluating the quality of our prediction obtained from rounding our relaxed solutions. We provide here a complementary study that evaluates the quality of the classifier weights that are implicitly learned by all methods based on discriminative clustering a.k.a. DIFFRAC [2]. This study is important to demonstrate that it is not only the constraints that give us better prediction performance but that we are also learning better models with the joint approach. In addition, it also evaluates the quality of the score obtained from the classifier weights. This is important, notably if we want to use these classifiers to predict on unseen videos. Therefore for each method that uses discriminative clustering cost, we keep the weights obtained at the last iteration of the optimization algorithm