

# Model-Free Incremental Learning of the Semantics of Manipulation Actions

Eren Erdal Aksoy<sup>a,\*</sup>, Minija Tamosiunaite<sup>a</sup>, Florentin Wörgötter<sup>a</sup>

<sup>a</sup>*Georg-August-Universität Göttingen, BCCN, Department for Computational Neuroscience  
Inst. Physics-3, Friedrich-Hund Platz 1, D-37077 Göttingen, Germany*

---

## Abstract

Understanding and learning the semantics of complex manipulation actions are intriguing and non-trivial issues for the development of autonomous robots. In this paper, we present a novel method for an on-line, incremental learning of the semantics of manipulation actions by observation. Recently, we had introduced the Semantic Event Chains (SECs) as a new generic representation for manipulations, which can be directly computed from a stream of images and is based on the changes in the relationships between objects involved in a manipulation. We here show that the SEC concept can be used to bootstrap the learning of the semantics of manipulation actions without using any prior knowledge about actions or objects. We create a new manipulation action benchmark with 8 different manipulation tasks including in total 120 samples to learn an archetypal SEC model for each manipulation action. We then evaluate the learned SEC models with 20 long and complex chained manipulation sequences including in total 103 manipulation samples. Thereby we put the event chains to a decisive test asking how powerful is action classification when using this framework? We find that we reach up to 100% and 87% average precision and recall values in the validation phase and 99% and 96% in the testing phase. This supports the notion that SECs are a useful tool for classifying manipulation actions in a fully automatic way.

---

\*Corresponding author

*Email address:* [eaksoye@physik3.gwdg.de](mailto:eaksoye@physik3.gwdg.de) (Eren Erdal Aksoy)

## 1. Introduction

One of the main problems in cognitive robotics is how to recognize and learn human demonstrations of new concepts, for example learning a relatively complex manipulation sequence like cutting a cucumber. Association-based or reinforcement learning methods are usually too slow to achieve this in an efficient way. They are therefore most often used in combination with supervised learning. Especially the Learning from Demonstration (LfD) paradigm seems promising for cognitive learning ([1, 2, 3, 4, 5]) because humans employ it very successfully. The problem that remains in all these approaches is how to represent complex actions or chains of actions in a generic and generalizable way allowing inferring the essential “meaning” (semantics) of an action irrespective of its individual instantiation.

In our earlier studies we introduced the “Semantic Event Chain” (SEC) as a possible descriptor for manipulation actions [6, 7]. The SEC framework analyzes the sequence of changes of the *spatial relations* between the objects that are being manipulated by a human or a robot. Consequently, SECs are invariant to the particular objects used, the precise object poses observed, the actual trajectories followed, or the resulting interaction forces between objects. All these aspects are allowed to change and still the same SEC is observed and captures the “essence of the action” as demonstrated in several action classification tests performed by us [6, 7, 8, 9].

In this paper, we address the problem of on-line, incremental learning of the semantics of manipulation actions observed from human demonstrations. We use the concept of SECs as the main processing tool to encode manipulations in a generic and compact way. Manipulations are continuous in the temporal domain but with event chains we discretize them by sampling only decisive key time points. Those time points represent topological changes between objects and the hand in the scene which are highly descriptive for a given manipulation. Our main intent here is to design a cognitive agent that can infer and learn frequently observed spatiotemporal changes embedded in SECs in an un-

31 supervised manner whenever a new manipulation instance occurs. The learning  
32 phase is bootstrapped only with the semantic similarities between SECs, i.e. the  
33 encoded spatiotemporal patterns, without using any prior knowledge about ac-  
34 tions or objects. Since we use computer vision methods to derive event chains,  
35 our approach for incremental learning of semantics is highly grounded in the  
36 signal domain. To the best of our knowledge, this is the first attempt to apply  
37 reasoning at the semantic level, while being fully grounded at the signal level,  
38 to learn manipulations with an unsupervised method. Note, here – on purpose  
39 – we do not include any object- or other information to show the power of  
40 our methods to fully automatically and in an unsupervised way extract action  
41 and object information. Clearly, in praxis, it will often make sense to include  
42 whatever additional knowledge is available to further ease action understanding.

43 The paper is organized as follows. We start with introducing the state  
44 of the art. We next provide a detailed description of each processing step;  
45 extraction of SEC representations and learning model-SECs for each observed  
46 manipulation. In the next section, we discuss experimental results from the  
47 proposed framework, which also includes validation and testing of the learned  
48 models. We finally conclude with a discussion.

## 49 **2. State of the Art**

50 Learning from Demonstration (LfD) has been successfully applied both at  
51 the control [1, 2, 10] as well as the symbolic level [3, 4, 5]. Although vari-  
52 ous types of actions can be encoded at the control level, e.g. trajectory-level,  
53 this is not general enough to imitate complicated actions under different cir-  
54 cumstances. On the other hand, at the symbolic level, sequences of predefined  
55 abstract action units are used to learn complex actions, but this might lead to  
56 problems for execution as many parameters are left out in a symbolic represen-  
57 tation. Although our approach with SECs is a symbolic-level representation,  
58 SECs can be enriched with additional decisive descriptors (e.g. trajectory, pose,  
59 etc.) and do not use any assumption or prior knowledge in the object or action

60 domain. Ideas to utilize relations to reach semantics of actions can be found as  
61 early as in 1975. For instance, [11] introduced the first approach about directed  
62 scene graphs in which each node identifies one object. Edges hold spatial in-  
63 formation (e.g., LEFT-OF, IN-FRONT-OF, etc.) between objects. Based on  
64 object movement (trajectory) information, events are defined to represent ac-  
65 tions. The main drawback of this approach is that the continuous perception  
66 of actions is ignored and is substituted instead by idealized hand-made image  
67 sequences. This approach, however, had not been pursued in the field any longer  
68 as only now powerful enough image processing methods became available from  
69 which object and action information can be extracted.

70 Still there are only a few approaches attempting to reach the semantics of  
71 manipulation actions in conjunction with the manipulated objects [12, 13, 14,  
72 15, 16]. The work in [12] is one of the first approaches in robotics that uses  
73 the configuration transition between objects to generate a high-level description  
74 of an assembly task from observation. Configuration transitions occur when  
75 a face-contact relation between manipulated and stationary environmental ob-  
76 jects changes. The work presented in [13] represents an entire manipulation  
77 sequence by an activity graph which holds spatiotemporal object interactions.  
78 The difficulty is, however, that very complex and large activity graphs need  
79 to be decomposed for further processing. In the work of [14], segmented hand  
80 poses and velocities are used to classify manipulations. A histogram of gradients  
81 approach with a support vector machine classifier is separately used to catego-  
82 rize manipulated objects. Factorial conditional random fields are then used to  
83 compute the correlation between objects and manipulations. Visual semantic  
84 graphs (inspired from our scene graphs) were introduced in [15] to recognize  
85 action consequences based on changes in the topological structure of the manip-  
86 ulated object. In [16] activity trees were presented to recognize actions using  
87 a minimal action grammar. Recent works such as [17] modeled human activi-  
88 ties employing the human skeleton information as well as roles of manipulated  
89 objects. Although all those works to a certain extent improve the classification  
90 of manipulations and/or objects, none of them extracts key events of individ-

91 ual manipulations or learns a descriptive semantic model to represent different  
92 manipulation tasks.

93 In this sense, to our best knowledge, our work is the first study to evaluate  
94 and learn the semantics of manipulations in an incremental and model free man-  
95 ner. The concept of semantic event chains has been successfully utilized and  
96 extended by others [18, 19, 20, 21, 22, 23] not only to represent manipulation  
97 actions but also to replicate them by robots. The work in [18] presented active  
98 learning of goal directed manipulation sequences, each was recognized using se-  
99 mantic similarities between event chains. Our scene graphs were represented  
100 with kernels in [19] to further apply different machine learning approaches. Ad-  
101 ditional trajectory information was used in [20] to reduce noisy events occur  
102 in SECs. Others [21, 22, 23] showed execution of various manipulations with  
103 different robots by using the key spatiotemporal points provided by SECs.

### 104 **3. Method**

105 In this method section we will present the core algorithmic components where  
106 are complex details will only be given in the Appendix. This should make  
107 reading easier, while still everything is present to implement this algorithm if  
108 desired.

#### 109 *3.1. Data Acquisition*

110 In this work, we investigate eight different manipulation actions: *Pushing*,  
111 *Hiding*, *Putting*, *Stirring*, *Cutting*, *Chopping*, *Taking*, and *Uncovering*. Fig. 1 (a)  
112 shows a sample frame for each manipulation action. All movies used in this  
113 study can also be found at [www.dpi.physik.uni-goettingen.de/~eaksoye/](http://www.dpi.physik.uni-goettingen.de/~eaksoye/MANIAC_DATASET)  
114 `MANIAC_DATASET`. The *Pushing* action shows how a hand can move objects  
115 around randomly. In the action of *Hiding*, some objects are made invisible by  
116 covering them with other objects. In the *Putting* action objects are taken from  
117 the supporting background and put on top of each other. The *Stirring* action  
118 represents a scenario in which a spoon is used to stir some liquid in a bucket. In

119 the *Cutting* action, a hand is cutting vegetables by moving a cutting tool back  
 120 and forth. In the *Chopping* action, a cutting tool follows a straight trajectory  
 121 to divide vegetables into parts. The *Taking* action represents a scenario where  
 122 some objects are taken down and put on the supporting background. In the  
 123 *Uncovering* action some objects are becoming visible after moving occluding  
 124 objects away.

125 We recorded 15 different versions for each of these manipulations by asking  
 126 5 different individuals to demonstrate each manipulation 3 times with different  
 127 objects in various scene contexts. Fig. 1 (b) depicts a sample frame from each

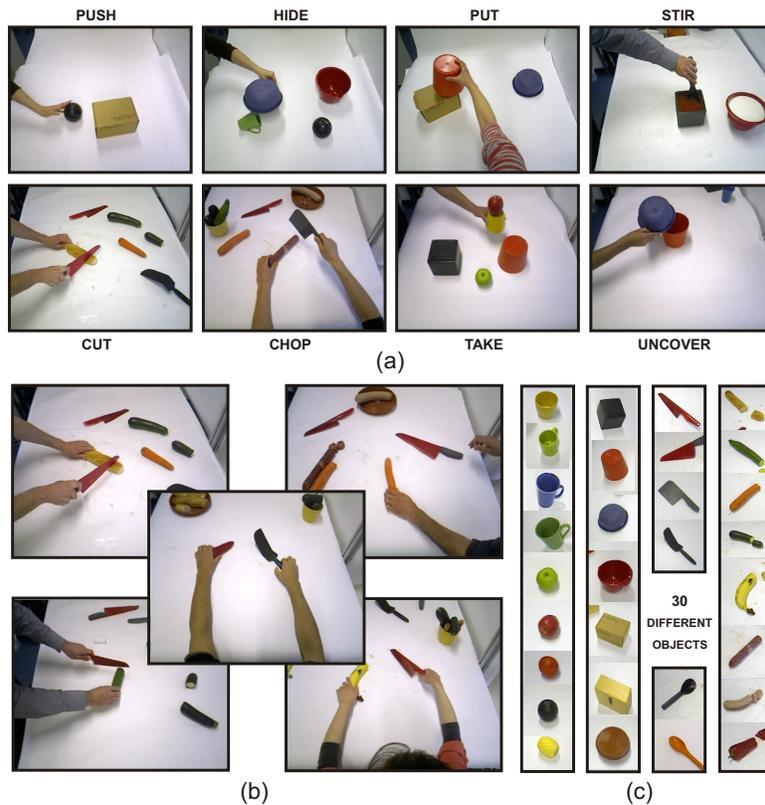


Figure 1: Eight different real action scenarios: *Pushing*, *Hiding*, *Putting*, *Stirring*, *Cutting*, *Chopping*, *Taking*, and *Uncovering*. (a) A sample original frame for each manipulation. (b) A sample frame from each demonstration of the *Cutting* action performed by 5 different individuals. (c) 30 different objects manipulated in all 120 manipulation demonstrations.

128 individual demonstration of the *Cutting* action to give an impression of the dif-  
129 ferences in demonstrations. There are in total 30 different objects manipulated  
130 in all 120 demonstrations. All manipulated objects are shown in Fig. 1 (c).

131 All manipulations were recorded with the Microsoft Kinect sensor which  
132 provides both color and depth image sequences. Colored objects are preferred  
133 to cope with the intrinsic limitations of the Kinect device. The central goal in  
134 these demonstrations is to learn a common archetypical SEC model for each  
135 manipulation including all possible variations in trajectory, pose, velocity, and  
136 object domains.

### 137 3.2. Segmentation and Tracking

138 The recorded image sequences are first pre-processed by a real-time image  
139 segmentation procedure to uniquely identify and track objects (including hands)  
140 in the scene. The segmentation algorithm is based on the color and depth in-  
141 formation fed from the Kinect device and uses phase-based optical flow [24] to  
142 track segments between consecutive frames. Data transmission between these  
143 different pre-processing sub-units is achieved with the modular system architec-  
144 ture described in [25]. Segmentation and tracking approaches are described in  
145 detail elsewhere [26, 27], therefore, details are omitted here.

### 146 3.3. Extracting Semantic Event Chains (SECs)

147 Each consistently segmented image is represented by a graph: nodes repre-  
148 sent segment centers and edges indicate whether two objects touch each other  
149 or not. By using the depth information we exclude the graph node for the back-  
150 ground segment, i.e. supporting surface, since it does not play any crucial role  
151 in the dynamics of the manipulation. By using an exact graph matching tech-  
152 nique, the framework discretizes the entire graph sequence into decisive main  
153 graphs. A new main graph is identified whenever a new node or edge is formed  
154 or an existing edge or node is deleted. Thus, each main graph represents a “key  
155 frame” in the manipulation sequence. All extracted main graphs form the core  
156 skeleton of the SEC, which is a matrix where rows are spatial relations (e.g.

157 touching) between object pairs and columns describe the scene configuration at  
 158 the time point when a new main graph has occurred.

159 Fig. 2 depicts the SEC representation with some sample *key frames* including  
 160 original images, respective segments (colored regions), and corresponding main  
 161 graphs for one of the *Cutting* action demonstrations. For instance, the first row  
 162 represents the spatial relations between graph nodes 9 and 6 which are hand  
 163 and knife, respectively. Note that, although the whole demonstration sample  
 164 has approximately 1000 frames, it is now represented by a  $3 \times 9$  matrix.

165 Possible spatial relations are *Not touching* ( $N$ ), *Touching* ( $T$ ), and *Absence*  
 166 ( $A$ ), where  $N$  means that there is no edge between two segments, i.e. graph  
 167 nodes corresponding to two spatially separated objects,  $T$  represents objects  
 168 that touch each other, and the absence of an object yields  $A$ . In the event  
 169 chain representation, all pairs of objects need to be considered once, however,  
 170 static rows which do not contain any change from  $N$  to  $T$  or vice versa are  
 171 deleted as being irrelevant. For instance, the relation between the left and right

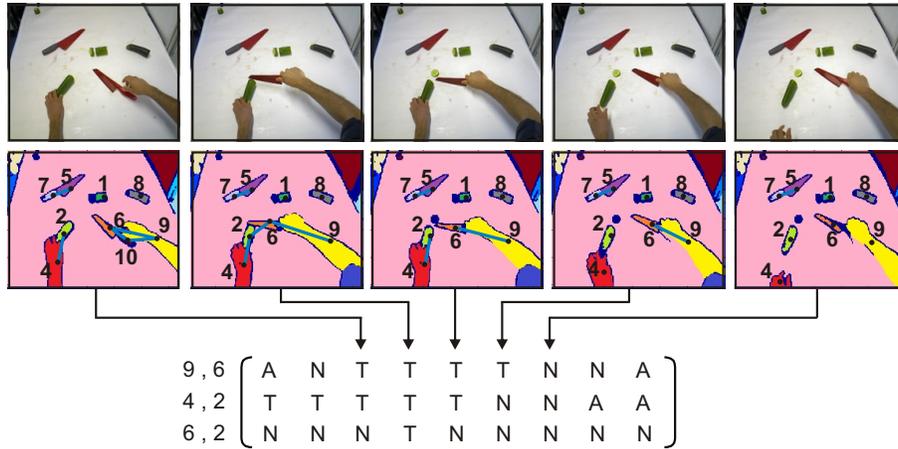


Figure 2: SEC representation for a sample *Cutting* action where a hand is cutting a cucumber with a knife. Each column corresponds to one *key frame* some of which are shown on the top with original images, respective segments (colored regions), and main graphs. Rows are spatial relations between object pairs, e.g. between the hand (9) and knife (6) in the first row. Possible spatial relations are  $N$ ,  $T$ , and  $A$  standing for *Not touching*, *Touching*, and *Absence*.

172 hand is always  $N$  and never switches to  $T$  to trigger an event, therefore, the  
173 respective row is ignored in the event chain. In Appendix A we introduce a de-  
174 noising process to cope with spurious spatial (rows) and/or temporal (columns)  
175 information propagated from noisy segmentation and tracking.

176 We note that there is no object recognition module included to identify  
177 graph nodes, i.e. segments, in the SEC framework. Event chains purely rely on  
178 spatial relational changes between segments in the temporal domain. The SEC  
179 extraction explained briefly in this section has been described in detail in [7].

### 180 3.4. Learning of Model SECs

181 The learning approach described next is an on-line unsupervised method to  
182 cluster observed SEC samples and to derive an archetypal SEC model for each  
183 cluster based on the semantic similarities between event chains. Each learned  
184 SEC model can then be used to describe a manipulation action.

185 Fig. 3 shows an overview of the proposed framework. The learning phase  
186 is triggered when a new manipulation experiment is observed; for example, a  
187 *Cutting* manipulation sample is introduced as the first experiment in Fig. 3.  
188 The new observed sample is represented by an event chain to be compared with  
189 the already learned SEC models. If there is no model existing, as in the case  
190 for this very first manipulation observation, the currently observed SEC sample  
191  $N$  is directly assumed as a new model  $M_1$ . Once a new manipulation example  
192 is acquired, e.g. a *Chopping* sample as the second experiment in Fig. 3, the  
193 framework measures semantic similarities between the new SEC sample  $N$  and  
194 the known model  $M_1$  in the spatiotemporal domain. We provide a detailed  
195 explanation of the similarity measure in Appendix B.

196 Semantic similarity values between the known models and the new sample  
197 are stored in a matrix, called the similarity matrix ( $\zeta_{semantic}$ ), which is then  
198 converted into a histogram ( $\mathcal{H}$ ) representing the distribution of similarities.  
199 We apply the conventional Otsu’s method introduced in [28] to the normalized  
200 version of the histogram to further compute a threshold  $\tau$ . See section 3.4.1  
201 for the details of the derivation of  $\mathcal{H}$  and  $\tau$  from  $\zeta_{semantic}$ . The gray box in

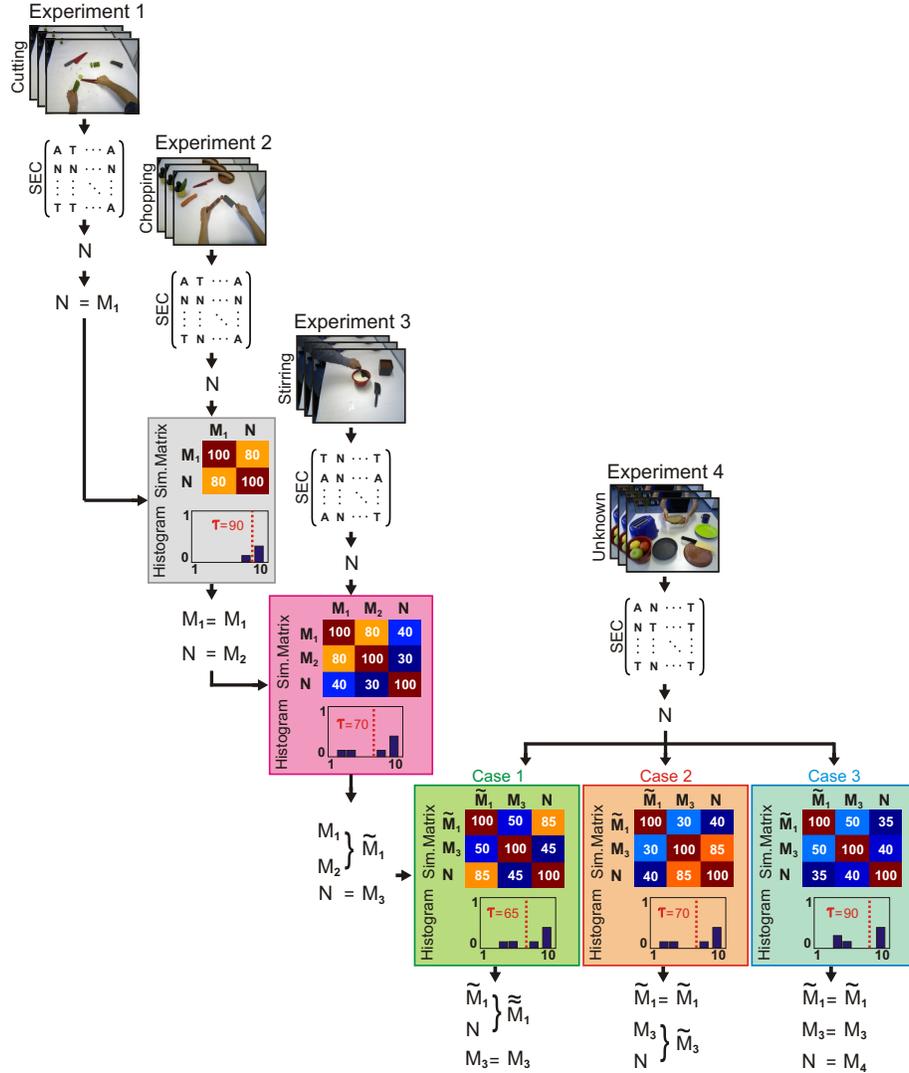


Figure 3: Overview of the proposed on-line learning framework.

202 Fig. 3 depicts extracted  $\zeta_{semantic}$  and  $\mathcal{H}$  in which the red dashed line indicates  
 203  $\tau$  computed between the first two experiments.

204 Threshold  $\tau$  is used for two purposes: First, we merge already learned SEC  
 205 models which have higher semantic similarities than  $\tau$ . Second, we compare the  
 206 currently observed SEC sample with the so far existing models. If the com-

207 parison yields a higher similarity than  $\tau$ , then the best fitting (highest similar)  
 208 model will be refined with the new SEC sample. Otherwise, a new model will  
 209 be created based on the SEC sample.

210 The comparison of the first two experiments  $N$  and  $M_1$  shown in the gray  
 211 box in Fig. 3 yields 80% semantic similarity which is less than  $\tau$  estimated as  
 212 90% (See Appendix B). Therefore, the *Chopping* sample  $N$  is considered as a  
 213 new SEC model  $M_2$ . We repeat the same procedure, i.e. computing  $\zeta_{semantic}$ ,  
 214  $\mathcal{H}$ , and  $\tau$ , once the next sample  $N$ , which is a *Stirring* experiment in this case,  
 215 is observed. As depicted in the purple box  $\tau$  drops below 80% which allows us  
 216 to update  $M_1$  with  $M_2$  yielding  $\widetilde{M}_1$ . As the *Stirring* demonstration still has less  
 217 similarities with any of the known models, a new model  $M_3$  is initialized with  
 218  $N$ .

219 The threshold value is required to better assess the obtained semantic sim-  
 220 ilarities between models and the observed sample. Therefore, whenever a new  
 221 observation is available, the entire process of estimating a new  $\tau$  by determin-  
 222 ing  $\zeta_{semantic}$  and  $\mathcal{H}$  is repeated to decide on the fly whether the current SEC  
 223 sample belongs to one of the already learned manipulation models or whether it  
 224 represents a new manipulation. This is summarized with the fourth experiment  
 225 introduced as an *Unknown* demonstration in Fig. 3, the fate of which depends  
 226 on three possible cases. Case 1 and 2 are respectively standing for the processes  
 227 of refining the models  $\widetilde{M}_1$  and  $M_3$  with  $N$ , whereas Case 3 is representing the  
 228 initialization of a new model  $M_4$ .

229 In the following, we will describe how to compute the threshold and update  
 230 a learned model with a new SEC sample.

### 231 3.4.1. Computing the Threshold

232 Let  $\mathcal{M}$  be a set of learned SEC models at any observation time as

$$\mathcal{M} = \{m_1, m_2, \dots, m_n\} \quad , \quad (1)$$

233 where  $n$  is the total number of existing models. Semantic similarity values

234 between all learned models are stored in a matrix as

$$\zeta_{semantic} = \begin{bmatrix} \varphi_{1,1} & \varphi_{1,2} & \cdots & \varphi_{1,n} \\ \varphi_{2,1} & \varphi_{2,2} & \cdots & \varphi_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ \varphi_{n,1} & \varphi_{n,2} & \cdots & \varphi_{n,n} \end{bmatrix}, \quad 0 \leq \varphi_{i,j} \leq 100 \quad \text{and} \quad \varphi_{i,j} = \varphi_{j,i},$$

235 where  $\varphi_{i,j}$  holds the semantic similarity between models  $m_i$  and  $m_j$  and is  
 236 computed as described in Appendix B.

237 Semantic similarity matrix  $\zeta_{semantic}$  is then converted into a histogram  $\mathcal{H}$   
 238 representing the distribution of similarities as

$$\mathcal{H} = \{h_k : k \in [1, \dots, \lambda]\}, \quad (2)$$

$$h_k = \frac{1}{\eta} \sum_{i=1}^n \sum_{j=i}^n \delta_{i,j}, \quad (3)$$

$$\delta_{i,j} = \begin{cases} 1 & \text{if } \frac{\varphi_{i,j}}{\phi} \text{ is at bin } k \\ 0 & \text{else} \end{cases}, \quad (4)$$

239 where  $\lambda$  is the total number of bins each has a size of  $\phi$  which is chosen as  
 240 10 in our experiments and  $\eta$  is the normalization factor. Note that, since the  
 241 similarity matrix  $\zeta_{semantic}$  is symmetric, only half of the matrix is processed,  
 242 thus, the value of  $j$  changes from  $i$  to  $n$  in Eq. (3) and  $\eta$  is defined as  $n(n+1)/2$ .

243 The normalized histogram  $\mathcal{H}$  is now used to calculate the required threshold  
 244 using the conventional Otsu's method introduced in [28]. For this purpose, we  
 245 compute zero- and first-order cumulative moments of the normalized histogram  
 246 at each bin as

$$\omega(k) = \sum_{i=1}^k h_i \quad , \quad (5) \quad \mu(k) = \sum_{i=1}^k ih_i \quad . \quad (6)$$

248 The total mean value of the histogram is calculated as

$$\mu_T = \sum_{i=1}^{\lambda} ih_i \quad . \quad (7)$$

249 The variance of the histogram separability is then given by

$$\sigma_B^2(k) = \frac{[\mu_T \omega(k) - \mu(k)]^2}{\omega(k)[1 - \omega(k)]} \quad . \quad (8)$$

250 Otsu’s method yields a threshold value  $k^*$  for that bin at which the variance  
251  $\sigma_B^2$  is maximal; that is,

$$k^* = \arg \max_{1 \leq k < \lambda} (\sigma_B^2(k)) \quad . \quad (9)$$

252 The threshold  $k^*$  separates the histogram into two distinct regions. The  
253 left side of  $k^*$  indicates low semantic similarity between models in  $\mathcal{M}$ , and vice  
254 versa. As we are seeking for a threshold  $\tau$  to group similar manipulations, we  
255 take the average of the similarities falling into the right side of  $k^*$  as

$$\tau = \frac{1}{\eta_r} \sum_{i=k^*}^{\lambda} h_i \quad , \quad (10)$$

256 where  $\eta_r$  is the normalization term which is the total number of similarity  
257 values on the right side of  $k^*$ .

### 258 3.4.2. Updating Model SECs

259 Once the highest semantic similarity between a novel SEC sample and any  
260 of the known models is higher than the threshold  $\tau$ , the one model with *highest*  
261 similarity to the new SEC is now updated with this new SEC sample. To update  
262 a model, the learning procedure just needs to search for all common rows and  
263 columns observed in the new SEC sample.

264 Each model is initially created by assigning weight value of 1 for each row.  
265 Once a new SEC sample is observed, weights of each row in the model that  
266 match to a row in the new SEC are incremented. This way existing common  
267 rows between the matched model and the novel sample are receiving increasing  
268 weights. In the case of having additional rows in the new SEC sample, the  
269 model is extended by these rows, each of which is initiated again by giving them a  
270 weight of one. As the next step, we search for the common temporal information  
271 embedded in the columns of the event chains by employing a procedure very  
272 similar to that applied for extracting common rows. Finally, the model SEC  
273 consists of only those rows and columns observed frequently in the observed  
274 new SEC samples. A detailed explanation of the model updating procedure is  
275 given in Appendix C.

## 276 4. Results

277 In this section, we will first show experimental results from our proposed  
278 incremental learning framework. We will then continue with enrichment of  
279 each learned SEC model with object information. Next, validation and testing  
280 processes of the learned models will be given.

### 281 4.1. Learning

282 We apply the incremental learning and clustering framework described above  
283 to 8 different manipulation actions each of which has 15 versions, yielding in  
284 total 120 samples, as introduced in section 3.1. Manipulation tasks have vast  
285 variations in terms of manipulated objects, their poses, and followed trajectories  
286 as depicted in Fig. 1. The framework is first tracking each segment in the scene  
287 and extracting the corresponding SEC representation from a randomly observed  
288 manipulation sample. While observing more samples, different SEC models are  
289 learned or updated based on the threshold value.

290 When we let the framework run only once through 120 manipulation tasks  
291 by randomly choosing a sample at each time, it learns 22 model event chains.

292 Fig. 4 (a) shows the final computed semantic similarity matrix  $\zeta_{semantic}$  be-  
 293 tween each of the learned models. Low similarities between models indicate  
 294 how distinct those models are. The corresponding histogram representation  $\mathcal{H}$   
 295 with derived thresholds  $k^*$  and  $\tau$  is depicted in Fig. 4 (b). The threshold  $k^*$  sep-  
 296 arates the histogram into two distinct regions as depicted with the gray shade  
 297 and  $\tau$  is then calculated as 72 from Eq. (10). In Fig. 4 (c), we can see the  
 298 complete behavior of  $\tau$  during the learning cycle with 120 observation samples.  
 299 It is initiated with 100 and after updating with Eq. (10) at each observation it  
 300 starts to converge to approximately 72.

301 Fig. 5 (a) depicts all learned models with corresponding number of observa-  
 302 tion samples employed for updating each. The green dashed line indicates the  
 303 actual sample numbers as the ground truth. Although the framework learns in  
 304 total 22 models, only 7 of them, those in the red box, contain more than 10  
 305 samples and the rest hold at most 2 samples. Recalling the fact that the train-  
 306 ing set has 8 manipulations, we can state as one central result that 7 of them  
 307 are indeed found with high numbers of examples each. *Cutting* and *Chopping*  
 308 models are merged, though, but we will explain below that this actually “makes  
 309 more sense” than the naively (by us) assumed ground truth. Furthermore, we

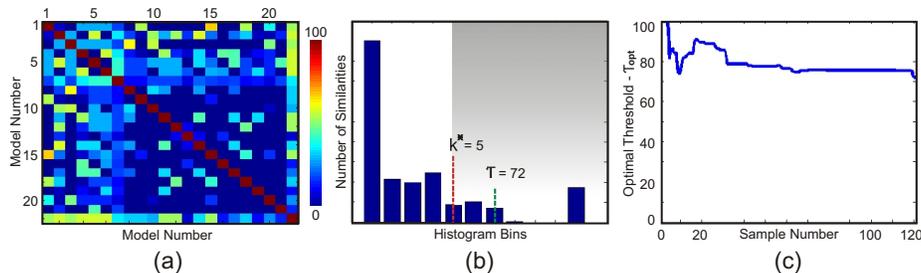


Figure 4: Thresholding. (a) Semantic similarity matrix  $\zeta_{semantic}$  computed between 22 learned SEC models. The scale bar on the right indicates the similarity values in percent. (b) Respective histogram representation  $\mathcal{H}$  with extracted  $k^*$  and  $\tau$  values. The threshold  $k^*$  separates the histogram into two distinct regions as depicted with the gray shade. (c) Development of  $\tau$  during the observation of 120 samples.  $\tau$  is initiated with 100 and after updating with Eq. (10) it starts to converge to 72.

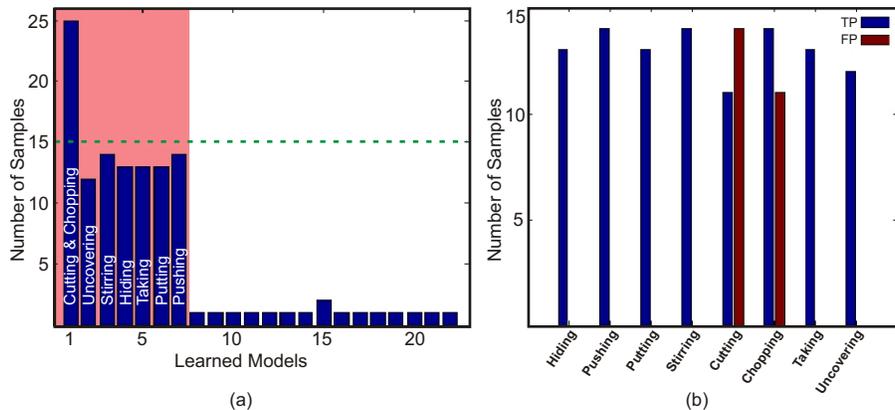


Figure 5: Number of learned models and clustering accuracy of observed samples. (a) Learned 22 SEC models with corresponding number of trained samples. The green dashed line indicates the actual sample numbers as the ground truth. (b) Number of true and false positive samples clustered in learned models with respect to the ground truth.

310 observe that only few demonstrated samples have either enormous variations or  
 311 noise, i.e. less semantic similarities than  $\tau$  with any other models, which leads  
 312 to the generation of the additional models outside the red box. As mentioned,  
 313 our framework produces a single model representing the *Cutting* and *Chopping*  
 314 manipulations together due to having high semantic similarities. It is because  
 315 both manipulations have the same fundamental action primitives, i.e. similar  
 316 columns in the event chains, and the only differences are mostly in the followed  
 317 trajectories and velocity of the movements which are not captured by SECs. See  
 318 Fig. B.14 in Appendix B as an example of high semantic similarities between  
 319 the *Cutting* and *Chopping* tasks. Thus, Fig. 5 (a) shows that without using any  
 320 human intervention the proposed learning framework can automatically retrieve  
 321 the demonstrated 8 manipulation types two of which are naturally merged.

322 As addressed in section 3.4, all manipulation samples used for updating the  
 323 same SEC model will have the same cluster label. In Fig. 5 (b), we show the  
 324 number of true and false positive samples falling into the same model with re-  
 325 spect to the ground truth. Except for the *Cutting* and *Chopping* manipulations,  
 326 none of the given manipulation samples is wrongly clustered. This means, for

327 instance, a given *Stirring* demonstration is used only for updating the *Stirring*  
 328 model, but not for the *Pushing* model, etc. However, since we have now only  
 329 one SEC model for the *Cutting* and *Chopping* manipulations, samples from both  
 330 manipulations will be used for the same model. As the ground truth expects  
 331 two different models, high false positives are observed for both.

332 Fig. 6 shows how the clustering results for all 120 manipulation samples are  
 333 varying from observation to observation. Colors encode the cluster labels and  
 334 the ground truth for each cluster is given on the left. Note that time is pro-  
 335 gressing from left to right, thus the first observed sample is the one depicted  
 336 in cyan in the *Chopping* manipulation. As a consequence of merging models  
 337 with high semantic similarity, some clusters will merge once new observations  
 338 become available. Black ellipses depict when a sample switches from one cluster  
 339 to another. For instance, cyan clusters observed for the *Chopping* samples in the

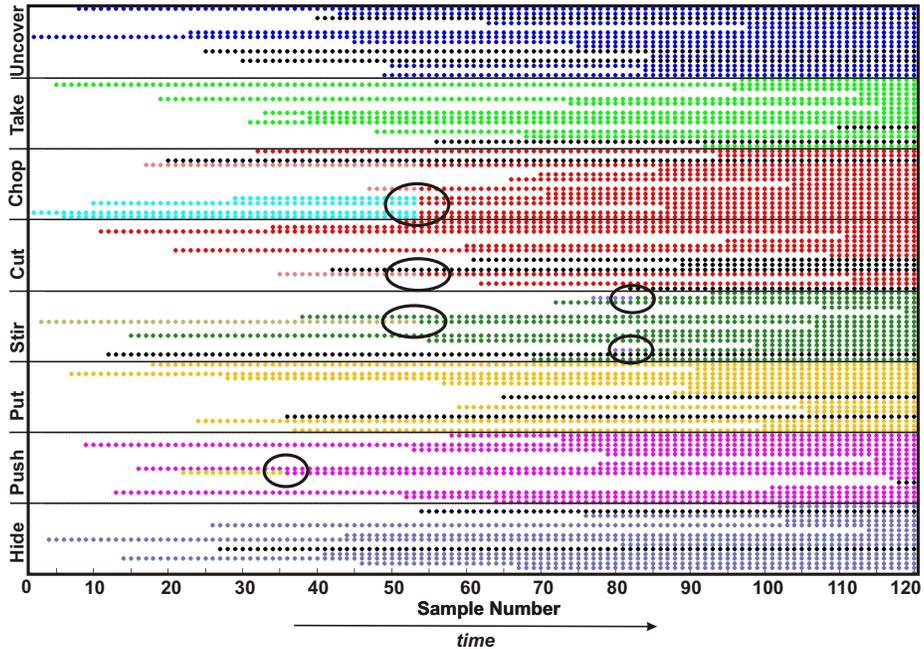


Figure 6: Clustering result of 120 manipulation samples. Colors encode the cluster labels and the ground truth for each cluster is given on the left. Noisy clusters are indicated in black. Black ellipses depict when a sample switches from one cluster to another.

340 beginning are turned into red clusters originally created for the *Cutting* task.  
 341 At the sample number 120 in the very right hand side we therefore observe 7  
 342 different colored clusters each from one learned model. This figure illustrates  
 343 that for some manipulations types the model is immediately converging to the  
 344 optimal solution, whereas for other models certain number of samples are re-  
 345 quired. Noisy clusters, which belong to the noisy models shown outside the red  
 346 box in Fig. 5 (a), are indicated by black dots.

347 To investigate the robustness of the framework, we repeat the same learn-  
 348 ing experiment explained above 100 times independently from each other and  
 349 compute differences between the learned models. In each trial, the framework  
 350 produces at least 21 and at most 23 various models. However, when we com-  
 351 pare all these models extracted in 100 trials, we see indeed 29 different ones, the  
 352 distribution of which is shown in Fig. 7 (a). Among those 29 models, it is again  
 353 the same 7 models introduced in Fig. 5 (a) which have high number of samples.  
 354 Furthermore, as indicated in Fig. 7 (b) we still do not obtain any false positives  
 355 among the clustered samples except for the *Cutting* and *Chopping* manipula-  
 356 tions due to the same reason as clarified above. Note that the red bars depict  
 357 the standard error of the mean for those which are not zero. Fig. 7 consequently

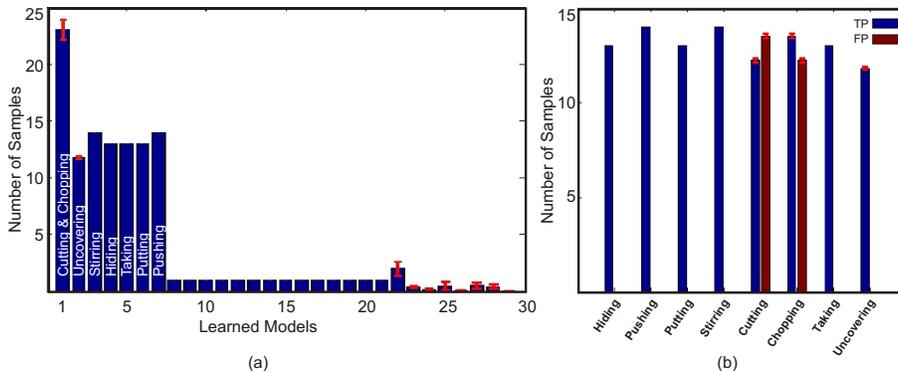


Figure 7: Total number of learned models and clustering accuracy after 100 independent trials. (a) Learned 29 SEC models with corresponding number of trained samples. (b) Number of true and false positive samples clustered in learned models with respect to the ground truth. Red bars depict standard error of the mean for those which are not zero.

358 proves that the learning approach is always converging to the same 7 models no  
 359 matter in which order the manipulation samples are provided.

360 We can now take a close look at some of those 7 SEC models explored  
 361 from demonstrated manipulation actions. Fig. 8 shows models for the *Cutting*  
 362 & *Chopping*, *Stirring*, and *Uncovering* manipulations with all derived states  
 363 introduced in Eq. (C.2) and the transition probabilities between each. States  
 364 and arrows given in red color correspond to the most commonly observed event  
 365 chain columns and their transitions with the highest probabilities as described  
 366 in Eqs. (C.3) and (C.4), respectively. On the left side of each model, we also  
 367 show weight values ( $\mathcal{W}$  from Eq. (C.1)) for each row in the states. It can be

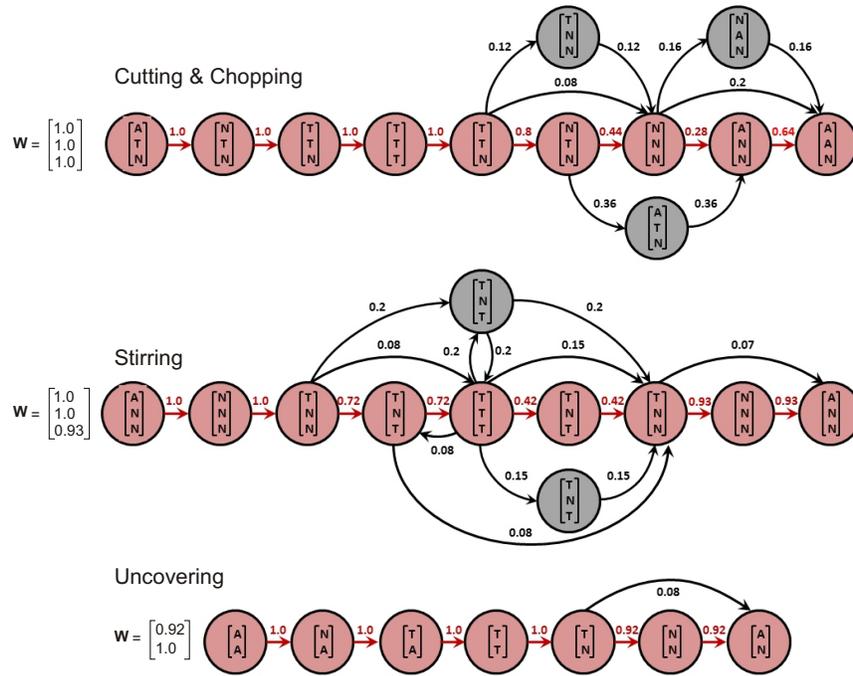


Figure 8: Complete learned SEC models for the *Cutting & Chopping*, *Stirring*, and *Uncovering* manipulations. Each state corresponds to one SEC column and arrows represent the transition probabilities from one state to the next. Those in red color correspond to the most commonly observed states, their transitions having the highest probabilities. Weight values  $\mathcal{W}$  on the left indicate how often each row in the states is obtained in the trained samples.

368 seen that in all 3 models some rows are quite commonly obtained in the trained  
 369 samples since their weights are close to 1, whereas this is not the same for  
 370 the state transitions. For instance, in the *Cutting & Chopping* model, there  
 371 exist three more states given in gray color which are particularly observed in  
 372 the second half of the action and cause drop of some state transitions to 0.28.  
 373 This is because even though each subject grasps a tool and cuts or chops an  
 374 object in the same temporal order, they leave the scene in different orders; for  
 375 example, one subject first removes the hand supporting the object to be cut  
 376 and then withdraws the hand holding the cutting tool whereas another subject  
 377 either does it the other way around or removes both hands at the same time.  
 378 Another reason of having extra states, thus smaller transition probabilities,  
 379 is the noise propagated from the segmentation and tracking components as  
 380 observed in the *Stirring* model. Nevertheless, we can now extract all these  
 381 variations that occurred due to the nature of manipulation or noise and pick  
 382 the most often observed states, i.e. states in red, as a representative model for  
 383 each manipulation action. Note that the learning process never ends and is open  
 384 to refine models incrementally whenever new samples are provided, just like the  
 385 assimilation process that happens in humans [29].

#### 386 4.2. Enriching Learned Model SECs

387 In this section, we will show how learned SEC models can be enriched with  
 388 additional object information.

389 During the updating process of model SECs, we determine correspondences  
 390 between rows of event chains as explained in section 3.4.2. Since each row in  
 391 an event chain holds relational changes between segments in the scene, the row  
 392 correspondences can also be used to calculate matchings between segments in  
 393 two event chains. We refer the interested reader to [7] for details of the segment  
 394 matching method. We now use this technique to extract segments, i.e. objects,  
 395 that play the same role in different versions of the same manipulations observed  
 396 during the learning phase.

397 Fig. 9 (a) shows the learned *Cutting & Chopping* model, columns of which

398 are the states indicated in red in Fig. 8. The framework now estimates which  
 399 segment is used as a main tool and which one as an object to be cut or chopped  
 400 in each observation. As explained in Appendix A, we refer to the hand as  
 401 the *manipulator* and to the object interacting with the hand as the *primary*  
 402 *object*, e.g. a knife or a cleaver. Other objects which are combined with the  
 403 primary object are called *secondary objects* like the cucumber to be cut. Note  
 404 that the second hand is almost always used to help the *manipulator*, hence it is  
 405 called the *supporter*. Fig. 9 (a) illustrates all matched *primary* and *secondary*  
 406 *objects* used for training of the *Cutting & Chopping* model. Fig. 9 (b) shows  
 407 the *primary* and *secondary objects* for the *Stirring* model. In this case, not  
 408 only a spoon but a knife and a spatula are also selected by subjects as the  
 409 *primary object* used for stirring. The *secondary object* is the stirred liquid and  
 410 the buckets are the *supporters*. As learned model SECs are refined with every  
 411 new observation, all these variations of the different objects will be attached to  
 412 the model, simultaneously. Note that segments representing the *manipulator*  
 413 and *supporter* are also matched, however, are not shown due to lack of space.

414 It is important to underline that the proposed framework is not utilizing any  
 415 object recognition method, hence, we are here strictly at the level of segments.  
 416 For the sake of simplicity, object images are shown instead of segments in Fig. 9.  
 417 It is evident that this unsupervised segment categorization process could be  
 418 coupled to object models, thus, providing access to object categorization, too.

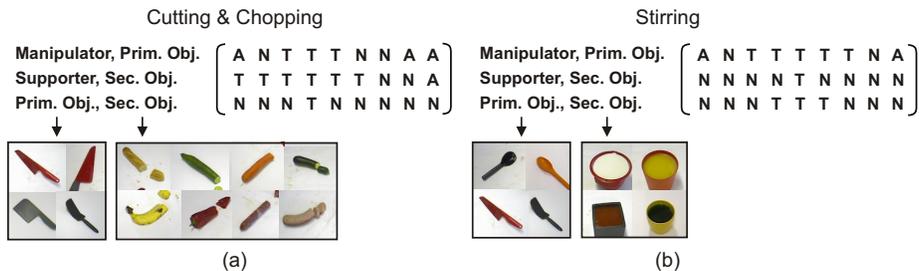


Figure 9: Learned *Cutting & Chopping* and *Stirring* models enriched with object information. Each column in the SEC model corresponds to one state indicated in red in Fig. 8. *Primary* and *secondary* objects are extracted from observed manipulations during the learning process.

419 *4.3. Validation and Testing*

420 A validation process of the learned 7 SEC models is performed with the  
 421 classification of all 120 training samples according to their semantic similarities  
 422 with the learned models. This step is required to show the clustering accuracy of  
 423 the training data but nothing unexpected will be observed here. We note that  
 424 the main and critical evaluation is then shown by the next following testing  
 425 experiment with a set of novel and complex manipulation *sequences*.

426 We label each SEC model as a different class and introduce a static threshold  
 427 chosen as 72 which is the converging value ( $\tau$ ) obtained during the learning  
 428 phase as depicted in Fig. 4 (c). Once the highest semantic similarity between  
 429 a training sample and any of the known models is higher than this threshold,  
 430 the sample is assigned to that class. The classification method has also a class  
 431 type called *Unknown* to detect samples that have low similarities with all known  
 432 models.

433 Fig. 10 (a) shows the confusion matrix depicting the classification accuracies  
 434 of the complete training data set with respect to the learned models. The  
 435 first impression that the figure conveys is that there is no misclassification of  
 436 any training data; for instance, 87% of the *Hiding* training manipulations are

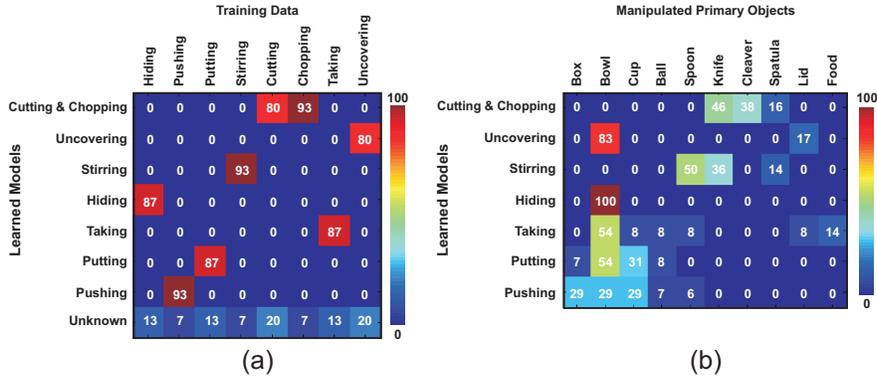


Figure 10: Confusion matrix showing (a) the classification accuracy for the complete training data set including in total 120 samples and (b) the usage rate of different objects primarily manipulated in the learned models.

437 correctly classified in the model *Hiding* and the rest is assigned as *Unknown*.  
438 As there is only one representative SEC model existing for both *Cutting* and  
439 *Chopping* manipulations, training samples from those are assigned within the  
440 same model *Cutting & Chopping*. The validation phase of the complete training  
441 set leads to 100% average precision and 87% average recall.

442 As addressed in section 4.2, we can also extract objects which are manipu-  
443 lated in a similar manner in different demonstrations of the same manipulation  
444 type. Fig. 10 (b) indicates the primary object types frequently manipulated in  
445 each classified training data. It is observed that objects like *Knife*, *Cleaver*, and  
446 *Spatula* are manipulated often in the *Cutting & Chopping* model class, whereas,  
447 due to its size, *Bowl* is the only preferred object in the *Hiding* manipulation to  
448 cover other objects. Fig. 10 consequently proves the high success rates of the  
449 discriminative and descriptive features of the learned 7 SEC models and their  
450 direct relations with manipulated objects.

451 To further evaluate the performance of the learned model SECs, we create  
452 a new testing set with 20 long chained actions which consist of in total 103  
453 different versions of the learned single manipulations such as *Cutting*, *Stirring*,  
454 and *Pushing*. We also introduce a new manipulation type called *Pouring* to  
455 measure the responses of the learned SEC models against a novel manipula-  
456 tion. In each chained action the subject has a certain task, e.g. “*making a*  
457 *sandwich*”, which involves execution of multiple single manipulations in various  
458 orders, either sequentially or parallelly. Fig. 11 depicts sample frames from two  
459 different chained action sequences in which subjects are performing the same  
460 task “*making a sandwich*” by using novel objects in various ways to increase the  
461 complexity of the scenes. We here apply an unsupervised, probabilistic method  
462 that measures the frequency of the changes in the spatial relations embedded  
463 in event chains to extract the main manipulator, e.g. hand, and to decompose  
464 the long chained actions into their primitive action components according to  
465 the spatiotemporal relations of the manipulator. Hence, also the decomposition  
466 process is model free and automatic. Since the decomposition issue is not in the  
467 core of the proposed framework, we omit the details here and refer the interested



Figure 11: Sample frames from two different long chained manipulation sequences which are used to test the learned SEC models. In these demonstrations subjects are performing the same task “*making a sandwich*” by using novel objects in various ways to increase the complexity of the scenes.

468 reader to [30].

469 Each single decomposed manipulation action is again analyzed as a classification task as described in the validation phase. Fig. 12 (a) indicates the highly  
 470 successful classification results of decomposed manipulations with respect to the  
 471 learned models. We receive minimum 83% accurate classification rate which is  
 472 for the *Stirring* manipulation and maximum 10% misclassification rate as ob-  
 473 served only for the *Pushing* manipulation. It is also significant to note that  
 474 the novel *Pouring* manipulation is never confused with any of the known SEC  
 475 models. In this testing phase, average precision and recall values are measured  
 476 as 99% and 96%, respectively.  
 477

478 Fig. 12 (b) shows the most often manipulated primary object types in each  
 479 classified test data. Compared to the results obtained in the training phase,  
 480 the major difference here is the high usage rates of the object type *Food* in the  
 481 *Hiding*, *Taking*, and *Putting* models. This is because making a sandwich by

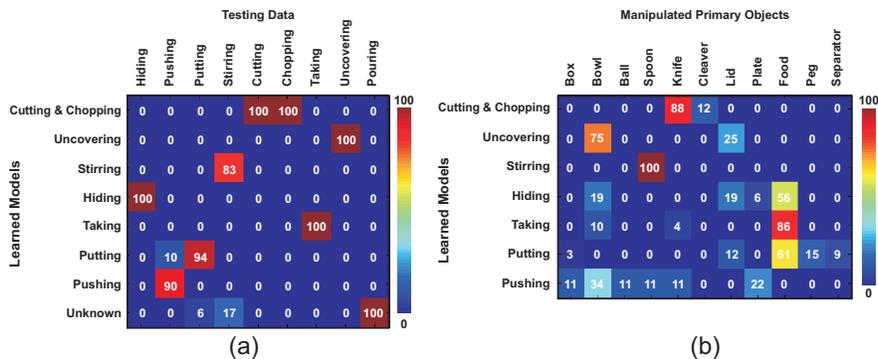


Figure 12: Confusion matrix showing (a) the classification accuracy for the complete testing data set including in total 103 samples and (b) the usage rate of different objects primarily manipulated in the learned models.

482 taking and putting cheese or bread slices on top of each other naturally results  
 483 in occlusions as expected by the *Hiding* model.

484 Note that all results shown in Figs. 10 and 12 are acquired in a fully auto-  
 485 mated, unsupervised manner and show that the learned SEC models are highly  
 486 accurate and discriminative to recognize manipulation actions which can even  
 487 be embedded in the long and complex chained demonstrations performed with  
 488 novel objects under different circumstances.

## 489 5. Discussion

490 The main contribution of our paper is a novel method for incrementally  
 491 learning the semantics of manipulation actions by observation. The proposed  
 492 learning framework is bootstrapped with the semantic relations (SECs) between  
 493 observed manipulations without using any prior knowledge about actions or  
 494 objects while being fully grounded at the sensory level (image segments). To  
 495 our best knowledge this is one of the first attempts in cognitive robotics to infer  
 496 descriptive semantic models of observed manipulations in a fully automated and  
 497 unsupervised manner.

498 One of the most fundamental advantages of the proposed framework is that

499 during the learning, when a new sample is observed, it is not compared with  
500 all previously acquired samples, which is an exhausting operation, but instead  
501 is compared only with the already learned models which are then updated ac-  
502 cordingly. This is of importance to allow the cognitive agent to use its memory  
503 in a more efficient way for lifelong learning, which is known as “Assimilation  
504 Process” in human cognition as originally defined by Piaget [29].

505 The proposed framework can be easily enriched with object information since  
506 event chains naturally group objects considering only their performed roles in a  
507 manipulation. As a strong contribution, we showed that objects, i.e. segments,  
508 can be categorized based on how an object is being manipulated, rather than  
509 by knowing what type of object it is. As shown in our previous works [9, 23],  
510 not only object but pose and the followed trajectory information can also be  
511 embedded into the SEC representations as further enrichment.

512 In this paper, we also introduced a new manipulation action data set with  
513 8 different manipulation tasks (e.g. *Cutting*, *Chopping*, *Stirring*, etc.), each of  
514 which consists of 15 different versions performed by 3 different human actors.  
515 This data set was used to learn an archetypal SEC model for each manipu-  
516 lation action. To further quantitatively evaluate the learned SEC models, we  
517 extended our data set with 20 long and complex chained manipulation sequences  
518 (e.g. “*making a sandwich*” or “*preparing a breakfast*”) which consist of in to-  
519 tal 103 different versions of these 8 manipulation tasks performed in different  
520 orders with novel objects under different circumstances. These data sets are  
521 publicly available and could be used for action/object benchmarking also of  
522 other methods.

523 In contrast to other well-known data sets, our new benchmark set captures  
524 manipulation activities from the subjects’s own point of view with a static RGB-  
525 D camera since we are interested in understanding the spatiotemporal interac-  
526 tions between the manipulated objects and hands. The conventional data sets,  
527 however, employ the entire human body configurations and movements as main  
528 features and therefore either do not involve hand-tool features [31, 17, 32] or  
529 are not rich to provide enough recordings required for the learning [15, 16].

530 The observed high accuracy of our method when classifying the unknown  
531 (long-sequence) test-data set support that the learned models are indeed dis-  
532 criminative and descriptive of these actions (and objects). The here shown  
533 experimental results also exhibit a similar behavior to that of the ontologies  
534 presented in [33, 15]. In these both studies manipulation actions were classified  
535 into six distinct structural categories (e.g. *Rearrange*, *Destroy*, *Break*, etc.) in  
536 which *Cutting* and *Chopping* manipulations were subsumed in the same category  
537 as the learned single Cutting & Chopping model in our framework.

538 As mentioned in the introduction, additional information, if available about a  
539 given action, will further improve action understanding. This notwithstanding,  
540 we believe that the current study strongly supports the power of the Semantic  
541 Event Chain framework, because here we have “pushed it to an extreme” by  
542 fully relying on model-free, unsupervised algorithms for clustering and classifi-  
543 cation. Therefore, we would hope that this study might stimulate the research  
544 community to adopt this framework in the future.

## 545 **Appendices**

546 We here provide three appendices each describes details of individual algo-  
547 rithmic steps in details. The first appendix introduces the de-noising process  
548 to filter out noisy spatiotemporal relations in the event chains. In the next ap-  
549 pendix, the detailed description of the similarity measure between event chains  
550 is given. The last appendix highlights the updating process of a learned SEC  
551 model with a novel SEC sample.

### 552 **Appendix A. De-noising of SECs**

553 Due to some early vision problems such as illumination variations or occlu-  
554 sions observed in the segmentation and tracking phases, extracted event chains  
555 can contain noisy spatial (rows) and/or temporal (columns) information. To  
556 prevent noisy event chain elements to propagate further to the next learning  
557 stage, we apply a de-noising process to the extracted raw SECs. The de-noising

558 process is based on reasonable action descriptive assumptions (rules) introduced  
559 in [33], which are as follows:

- 560 1. only *single hand* manipulations are considered;
- 561 2. the hand can manipulate, i.e. *touch*, only one object at a time;
- 562 3. the manipulation can take place at the touched object itself (the one mentioned  
563 in rule 2) or only one other object can be a target, with which the first one  
564 interacts, i.e. *touches*;
- 565 4. before and after the manipulation the hand is *free* and *not-touching* anything;
- 566 5. before and after the manipulation the hand is *not in the scene*.

567 The first two rules guarantee that there is only one hand and at most one  
568 object interacting with the hand, which we call *manipulator* and *primary object*,  
569 respectively. Other objects, which are combined with the primary object, are  
570 called *secondary objects*. The third rule assures that *manipulator*, *primary* and  
571 *secondary objects* are the only ones having direct interaction with each other  
572 affecting the dynamics of the manipulation. The last two rules define the natural  
573 start and end points of the manipulation.

574 The de-noising process checks whether all those rules are satisfied in the  
575 SEC representation. For instance, the first two rules require that the event  
576 chain must have a row holding spatial relations between the *manipulator* and  
577 *primary object* and last three rules define these relations as:

$$\textit{manipulator, primary object} \quad \left[ A \quad N \quad T \quad \dots \quad T \quad N \quad A \right], \quad (\text{A.1})$$

578 where the *manipulator* is first absent (A) in the scene (*rule 5*), then appears  
579 but does not touch (N) the *primary object* (*rule 4*). Next, the *manipulator*  
580 touches (T) *primary object* to apply a certain task on it (*rule 3*). Depending  
581 on the manipulation, the temporal length of the touching (T) relation can vary.  
582 Finally, the *manipulator* releases (N) the *primary object* (*rule 4*) and leaves (A)  
583 the scene (*rule 5*).

584 Since segments, i.e. graph nodes, are not identified as objects in event chains,  
585 we do not know which segment corresponds to the *manipulator* or *primary*

586 *object*. Therefore, we apply a probabilistic reasoning to estimate segment roles  
 587 in the manipulation. Probability values for each segment are assigned based on  
 588 similarities of their relations with Eq. (A.1) and the frequency of their touching  
 589 relations. See Appendix B for similarity calculation between SEC rows. In this  
 590 regard, all rows in the event chain are compared with Eq. (A.1) and the most  
 591 similar one is taken as the best candidate for the *manipulator* and the *primary*  
 592 *object*.

593 Fig. A.13 (a-b) shows a noisy raw event chain with corresponding *key frames*  
 594 extracted from a *Putting* manipulation sample where a hand is putting a cup on  
 595 a box. For instance, the first and second rows of the SEC given in Fig. A.13 (b)  
 596 are similar to Eq. (A.1), however, the second row has a higher probability to be  
 597 a better candidate due to having more touching relations. Therefore, segments  
 598 4 and 1 in the second row have the highest likelihood to be the *manipulator*

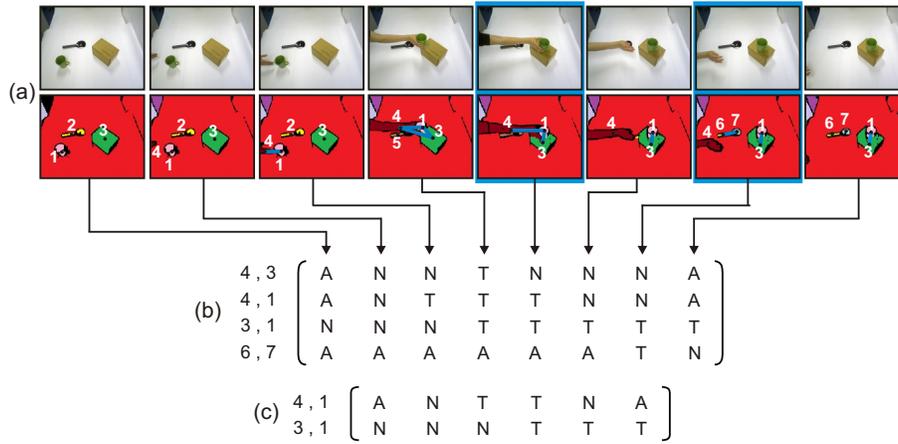


Figure A.13: SEC representation for a sample *Putting* action where a hand is putting a cup on a box. (a) Extracted 8 *key frames* with original images, corresponding segments (colored regions), and main graphs. (b) Respective SEC where each *key frame* corresponds to one column. Rows are spatial relations between object pairs, e.g. between the hand (4) and box (3) in the first row. Possible spatial relations are *N*, *T*, and *A* standing for *Not touching*, *Touching*, and *Absence*, respectively. (c) De-noised SEC after applying action descriptive rules. First and last rows as well as repetitive key frames, shown in blue frames, are removed from the raw SEC in (b).

599 and the *primary object*. Since *rule 5* constrains the *manipulator* to appear in  
600 the scene later, we choose segment 4 as the *manipulator* and segment 1 as the  
601 *primary object*.

602 Once the *manipulator* and *primary object* are estimated, the de-noising pro-  
603 cess is concluded by examining the second and third rules once more. Since  
604 the second rule does not allow the *manipulator* to interact with any other ob-  
605 ject other than the *primary object*, such rows can be considered as noise to be  
606 omitted. In this manner, the first row in the SEC given in Fig. A.13 (b) can  
607 be ignored as the *manipulator* (segment number 4) is also touching the box  
608 (segment number 3) which is not the *primary object*. Note that the hand is  
609 here accidentally touching the box while putting the cup. Recalling the third  
610 rule, we can ignore any segment which does not have any interaction with the  
611 *manipulator* or *primary object*. In this sense, the forth row of the SEC in  
612 Fig. A.13 (b) is omitted because segment 6 and 7 represent the spoon which is  
613 occluded by the *manipulator* and *primary object* and not playing any role in the  
614 manipulation. Fig. A.13 (c) shows the final de-noised SEC representation for  
615 the *Putting* action in Fig. A.13 (a). Note that de-noised event chain includes less  
616 columns since redundant duplicate (repetitive) columns observed after deleting  
617 noisy rows (indicated in blue frames in Fig. A.13 (a)) are also removed.

618 It is important to underline that the de-noising process considers temporal  
619 interactions between entire segments in the manipulation to solve illumination or  
620 occlusion based early vision problems which can not be solved without reasoning  
621 at a higher level.

## 622 **Appendix B. Measuring Semantic Similarity**

623 Once event chains are extracted in the observation phase, their semantic sim-  
624 ilarities need to be compared to further explore whether they describe the same  
625 type of manipulation. In [7], we introduced a method to measure semantic sim-  
626 ilarities and here we describe an updated version which is more robust against  
627 noisy spatiotemporal information coming from the early vision stage. To better

628 explain the semantic comparison we will use sample demonstrations from the  
629 *Cutting* and *Chopping* manipulations which are shown in Fig. B.14 (a-b) with  
630 extracted de-noised SECs including some sample key frames with respective  
631 segments and graphs. Note that even though those two samples have different  
632 perspectives and contain different number and types of objects, the dimensions  
633 of the event chains are accidentally the same. This is of no importance as our  
634 proposed method does not rely on dimensions, allowing to compare arbitrarily  
635 long manipulations.

636 To calculate the semantic similarity between two manipulations, spatial and  
637 temporal aspects are being analyzed in two separate steps. In the first step,  
638 we compare spatial information, i.e. relational changes in each row, and in the  
639 following second step the temporal information, i.e. the order of columns, is  
640 considered. In both steps we apply a standard sub-string search algorithm. To  
641 achieve this, we first perform a data-compression on the original chain ( $\xi_o$ ) by  
642 simply scanning each row of  $\xi_o$  from left to right and substitute “changes” by  
643 combining their values into a two-digit format. For example a change from  
644 *Not touching* to *Touching*, hence from  $N$  to  $T$ , is now encoded by  $NT$ . When  
645 nothing has changed, a double digit like  $TT$ , is removed. This compressed event  
646 chain, represented by  $\xi_c$ , lost all temporal information and is used only for the  
647 spatial-relational analysis in the first step. The original chain ( $\xi_o$ ) will then be  
648 used for the temporal analysis in the second step.  $\xi_o$  and  $\xi_c$  of the *Cutting* and  
649 *Chopping* actions are given in Fig. B.14 (a-d).

650 Let  $\xi_c^1$  and  $\xi_c^2$  be the sets of rows for the two manipulations, written as a  
651 matrix (e.g. Fig. B.14 (c) and B.14 (d)):

$$\xi_c^1 = \begin{bmatrix} r_{1,1}^1 & r_{1,2}^1 & \cdots & \cdots & r_{1,\gamma_1^1}^1 \\ r_{2,1}^1 & r_{2,2}^1 & \cdots & r_{2,\gamma_2^1}^1 & \\ \vdots & \vdots & \ddots & \vdots & \\ r_{m,1}^1 & r_{m,2}^1 & \cdots & \cdots & \cdots & r_{m,\gamma_m^1}^1 \end{bmatrix},$$

652 and

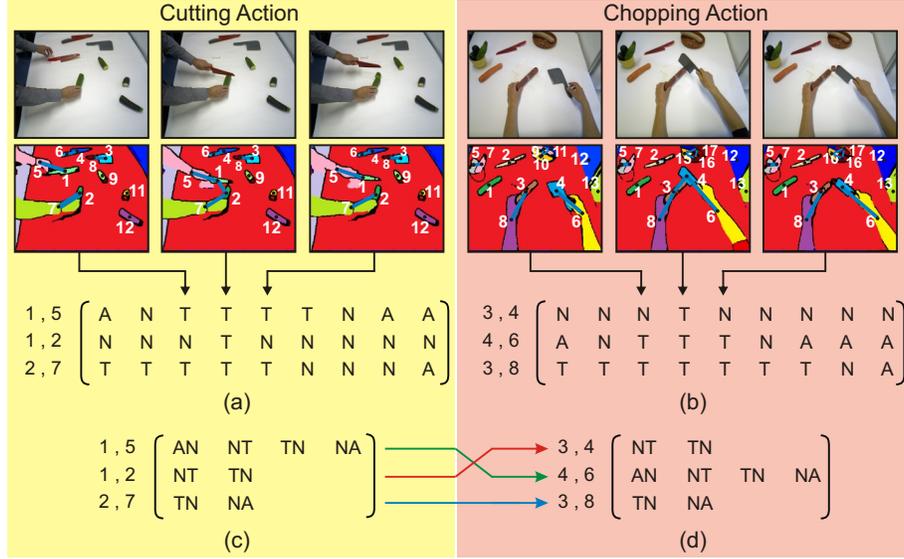


Figure B.14: Two sample manipulation action scenarios: “Cutting a cucumber with a knife” (on the left) and “Chopping a sausage with a cleaver” (on the right). (a-b) Extracted de-noised SECs ( $\xi_o$ ) with some sample original key frames including respective segments and main graphs. (c-d) Corresponding compressed SECs ( $\xi_c$ ). Colored arrows show row matchings.

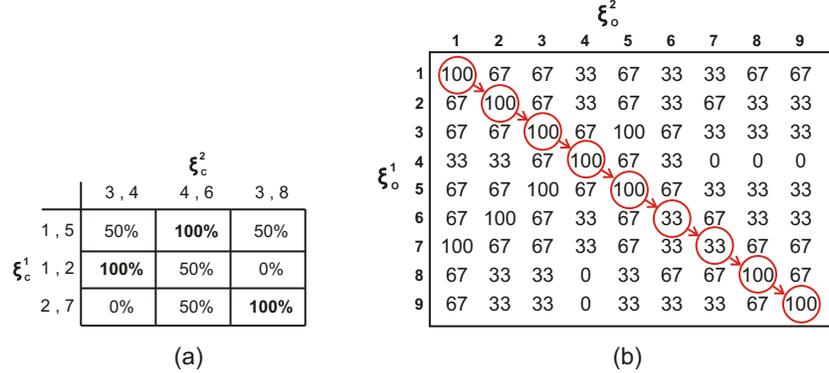


Figure B.15: Similarity matrices between the *Cutting* and *Chopping* samples given in Fig. B.14. (a) Spatial similarity matrix  $\zeta_{spatial}$  indicates possible correspondences between rows (see colored arrows in Fig. B.14). (b) Temporal similarity matrix  $\zeta_{temporal}$  with LCS matchings indicated in red circles shows correspondences between columns.

$$\xi_c^2 = \begin{bmatrix} r_{1,1}^2 & r_{1,2}^2 & \cdots & \cdots & \cdots & \cdots & r_{1,\gamma_1}^2 \\ r_{2,1}^2 & r_{2,2}^2 & \cdots & \cdots & r_{2,\gamma_2}^2 & & \\ \vdots & \vdots & \ddots & \vdots & & & \\ r_{k,1}^2 & r_{k,2}^2 & \cdots & \cdots & \cdots & r_{k,\gamma_k}^2 & \end{bmatrix},$$

653 where  $r_{i,j}$  represents a relational *change* between a segment pair

$$r_{i,j} \in \{AN, AT, NA, NT, TA, TN\},$$

654 where  $A$ ,  $N$ , and  $T$  stand for *Not touching*, *Touching*, and *Absence*, respectively. The lengths of the rows are usually different and given by indices  $\gamma$ .

655 The first step is comparing the rows of the compressed event chains ( $\xi_c^1$  and  
656  $\xi_c^2$ ) accounting for a possibly shuffling of rows in different versions of the same  
657 manipulations. Therefore, each row of  $\xi_c^1$  is compared with each row of  $\xi_c^2$  in  
658 order to find the highest similarity. The comparison process searches for equal  
659 entries of one row against the other using a standard sub-string search, briefly  
660 described next. Assume that we compare the  $a^{th}$  row of  $\xi_c^1$  with the  $b^{th}$  row of  
661  $\xi_c^2$ . If row  $a$  is shorter or of equal length than row  $b$  ( $\gamma_a^1 \leq \gamma_b^2$ ), the  $a^{th}$  row of  $\xi_c^1$   
662 is shifted  $\gamma_b^2 - \gamma_a^1 + 1$  times to the right. At each shift its entries are compared  
663 with the one of the  $b^{th}$  row of  $\xi_c^2$  and we get as a result set  $F_{a,b}$  defined as:  
664

$$F_{a,b} = \{f_t : t \in [1, \gamma_b^2 - \gamma_a^1 + 1]\},$$

$$f_t = \frac{100}{\gamma_b^2} \sum_{i=1}^{\gamma_a^1} \delta_i, \quad (\text{B.1})$$

665 where  $\gamma_b^2$  is the normalization factor and  $i$  is the row index and with

$$\delta_i = \begin{cases} 1 & \text{if } r_{a,i}^1 = r_{b,i+t-1}^2 \\ 0 & \text{else} \end{cases}, \quad (\text{B.2})$$

666 where the set  $F_{a,b}$  represents all possible similarities for every shift  $t$ , given by  
 667  $f_t$ , which holds the normalized percentage of the similarity calculated between  
 668 the shifted rows.

669 As usual for sub-string searches, we are only interested in the maximum  
 670 similarity of every comparison hence we define:

$$M_{a,b} = \max(F_{a,b}),$$

671 For the case  $\gamma_a^1 > \gamma_b^2$ , a symmetrical procedure is performed by interchanging  
 672 all indices of Eqs. (B.1), (B.2) above.

673 Spatial similarity values between all rows of  $\xi_c^1$  and  $\xi_c^2$  are stored in a matrix  
 674  $\zeta_{spatial}$  with size  $m \times k$  as

$$\zeta_{spatial} = \begin{bmatrix} M_{1,1} & M_{1,2} & \cdots & M_{1,k} \\ M_{2,1} & M_{2,2} & \cdots & M_{2,k} \\ \vdots & \vdots & \ddots & \vdots \\ M_{m,1} & M_{m,2} & \cdots & M_{m,k} \end{bmatrix}.$$

675 The final similarity value ( $\psi_{spatial}$ ) between the rows of two compressed  
 676 event chains is calculated by taking the mean value of the highest similarities  
 677 across both rows and columns of  $\zeta_{spatial}$  as

$$\psi_{spatial} = \frac{1}{m} \sum_{i=1}^m \max_j(M_{i,j}), \quad j \in [1, \dots, k], \quad (\text{B.3})$$

678 if

$$\max_j(M_{i,j}) = \max_t(M_{t,j}), \quad t \in [1, \dots, m] \quad . \quad (\text{B.4})$$

679 The spatial similarity matrix  $\zeta_{spatial}$  indicates possible correspondences be-  
 680 tween rows of  $\xi_c^1$  and  $\xi_c^2$  used to compute temporal similarity in the second step.  
 681 Note that there can be more than one correspondences between each row and  
 682 all existing permutations need to be considered in the second step, separately.  
 683 If there is a size differences between event chains, extra rows with no correspon-  
 684 dences will be omitted here, but penalty values will then be applied at the end  
 685 of the second step.

686 The complete similarity matrix ( $\zeta_{spatial}$ ) between the *Cutting* and *Chopping*  
687 samples ( $\xi_c^1$  and  $\xi_c^2$ ) is given in Fig. B.15 (a) which shows that first row of  $\xi_c^1$ ,  
688 i.e. 1, 5, corresponds to the second row of  $\xi_c^2$ , i.e. 4, 6. The same reverse relation  
689 exists between the second row of  $\xi_c^1$  and the first row of  $\xi_c^2$ . Therefore, rows of  
690 the second event chain will be resorted by simply interchanging first and second  
691 rows to initiate the second step, i.e. temporal analysis of the method.

692 In the following second step, we use the time sequence, encoded in the order  
693 of columns in the original event chains, to find the best matching permutation  
694 and thereby arrive at the final semantic similarity. To this end we will now  
695 compare columns of resorted  $\xi_o^2$  with that of  $\xi_o^1$ . Note that by contrast to rows,  
696 columns of event chains are never shuffled unless they represent different types  
697 of actions. Therefore, the column orders of type-similar event chains have to  
698 be the same. The comparison procedure of columns is very similar to the one  
699 for the rows. Since the lengths of the columns are the same, no shift-operator  
700 is required and columns are directly compared index-wise. Similarity values  
701 between all columns of  $\xi_o^1$  and  $\xi_o^2$  are stored in a matrix  $\zeta_{temporal}$  with the size  
702 of  $u \times v$  as

$$\zeta_{temporal} = \begin{bmatrix} \theta_{1,1} & \theta_{1,2} & \cdots & \theta_{1,v} \\ \theta_{2,1} & \theta_{2,2} & \cdots & \theta_{2,v} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{u,1} & \theta_{u,2} & \cdots & \theta_{u,v} \end{bmatrix},$$

703 where  $u$  and  $v$  are the lengths of columns in  $\xi_o^1$  and  $\xi_o^2$ .

704 Once similarities between columns are calculated, we use “Longest Common  
705 Subsequence, (LCS)” in order to guarantee that the order of columns is the  
706 same. LCS is generally used to explore the longest sequence existing in both  
707 input samples sequences. Columns of event chains are used as sequences for this  
708 task and LCS matching is computed based on similarities in  $\zeta_{temporal}$ . Since  
709 the number of sequences is constant, the problem is solvable in polynomial time  
710 by dynamic programming. Fig. B.15 (b) shows  $\zeta_{temporal}$  with LCS matchings  
711 indicated in red circles for the *Cutting* and *Chopping* samples  $\xi_o^1$  and  $\xi_o^2$ .

712 The temporal similarity value  $\psi_{temporal}$  between the columns of two event  
713 chains is then calculated by taking the mean value of the similarities given by  
714 LCS matching  $L_i$  as

$$\psi_{temporal} = \frac{1}{u} \sum_{i=1}^u L_i \quad , \quad (B.5)$$

$$L_i = \begin{cases} 100 & \text{if } \theta_{i,j} = 100 \\ 0 & \text{else} \end{cases} \quad , \quad (B.6)$$

715 where  $i$  and  $j$  are the matching column indices between  $\xi_o^1$  and  $\xi_o^2$ .

716 Note that due to noisy segmentation and tracking, size of  $\xi_o^1$  and  $\xi_o^2$  can be  
717 different. Therefore, size differences between event chains are used as a penalty  
718 to prevent false similarities. The final semantic similarity is then computed as

$$\psi_{final} = \frac{r_1 c_1 \psi_{temporal}}{r_1 c_1 + \frac{r_2 c_2 - r_1 c_1}{\rho}} \quad , \quad r_1 < r_2 \quad \text{and} \quad c_1 < c_2 \quad , \quad (B.7)$$

719 where  $\rho$  is the penalty value and  $r_1$ ,  $c_1$ ,  $r_2$ , and  $c_2$  are the number of rows and  
720 columns of  $\xi_o^1$  and  $\xi_o^2$ , respectively. The final  $\psi_{final}$  value between the *Cutting*  
721 and *Chopping* samples in Fig. B.15 is calculated as 78% by using Eqs. (B.7),  
722 (B.6), and (B.7) with  $\rho = 1$ . The best matching permutation is further used for  
723 categorizing objects as described in [7].

## 724 Appendix C. Model Updating

725 Let  $\xi_m$  and  $\xi_n$  be two matrices representing a SEC model and a new SEC  
726 sample with sizes of  $p \times q$  and  $k \times l$ , respectively. The two matrices can be  
727 written as

$$\xi_m = \begin{bmatrix} r_{1,1}^m & r_{1,2}^m & \cdots & r_{1,q}^m \\ r_{2,1}^m & r_{2,2}^m & \cdots & r_{2,q}^m \\ \vdots & \vdots & \ddots & \vdots \\ r_{p,1}^m & r_{p,2}^m & \cdots & r_{p,q}^m \end{bmatrix} \quad \text{and} \quad \xi_n = \begin{bmatrix} r_{1,1}^n & r_{1,2}^n & \cdots & r_{1,l}^n \\ r_{2,1}^n & r_{2,2}^n & \cdots & r_{2,l}^n \\ \vdots & \vdots & \ddots & \vdots \\ r_{k,1}^n & r_{k,2}^n & \cdots & r_{k,l}^n \end{bmatrix} ,$$

728 where  $r_{i,j} \in \{A, N, T\}$  is representing the spatial relations between each  
 729 segment pair as described in section 3.3.

730 Each model  $\xi_m$  is initially assigned with a set of weights  $\mathcal{W}$  as

$$\mathcal{W} = [w_1, w_2, \dots, w_p]^T , \quad (\text{C.1})$$

731 for representing the appearance frequency of each row, which leads to ex-  
 732 traction of all common rows observed in most of SEC samples. Each weight  
 733 value  $w_i$  is initialized to 1. We first compare each row of  $\xi_n$  with each row of  
 734  $\xi_m$  to find identical matches and to further increment the corresponding weight  
 735 values of the matched rows again by 1. This step is required since rows can be  
 736 shuffled in the new observation sample  $\xi_n$ . While comparing rows, we search  
 737 for only equal relational changes rather than temporal lengths of relations as  
 738 explained in Appendix B. In the case of  $k > p$ , all novel rows observed in  $\xi_n$   
 739 will then be appended to  $\xi_m$  with weight values  $\{w_{p+1}, \dots, w_k\}$  initialized to 1.  
 740 Common rows are then those with weights higher than  $\frac{|\mathcal{W}|}{2}$ . Next, the order of  
 741 rows in  $\xi_n$  is resorted considering the order of their best matches with common  
 742 rows in  $\xi_m$ . The sorting process yields the same row numbers in  $\xi_n$  and  $\xi_m$ ,  
 743 which is required for analyzing columns as described next.

744 The following step covers the temporal information embedded in the columns  
 745 of  $\xi_n$  and  $\xi_m$ , and is similar to the previous approach explained for rows. We  
 746 here assume that each column in an event chain is a *state* defining one action  
 747 primitive. Hence, we seek for all primitives derived from new observations and  
 748 compute the transition between them. Let  $\mathcal{S}_m$  be a set of existing states in the  
 749 current model  $\xi_m$  :

$$\mathcal{S}_m = \{s_1, s_2, \dots, s_q\} \quad , \quad (\text{C.2})$$

750 where each  $s_i = \{r_{j,i}^m : j \in [1, \dots, p]\}$ . We now compare each state, i.e.  
 751 column, in the sorted version of  $\xi_n$  with those in  $\xi_m$  by employing the same  
 752 approach as defined for the temporal analysis in Appendix B. In the case of  
 753 having more states in  $\xi_n$ , i.e.  $l > q$ , all novel states are appended to  $\mathcal{S}_m$ , and  
 754 then transitions between each state are calculated. We assign a probability  
 755 value  $P_{i,j}$  defining the transition from  $s_i$  to  $s_j$ , which is incremented when two  
 756 states are consecutive, i.e.  $s_j = s_{i+1}$  in  $\xi_n$ .

757 Following state transition calculation, the learned model  $\xi_m$  is refined with  
 758 the new states  $\hat{\mathcal{S}}_m$  having the maximum transitions between each; that is,

$$\hat{\mathcal{S}}_m = \{s_{\alpha_1}, s_{\alpha_2}, \dots, s_{\alpha_l}\} \quad , \quad (\text{C.3})$$

$$\alpha_{t+1} = \arg \max_j (P_{\alpha_t, j}) \quad , \quad (\text{C.4})$$

759 where  $\alpha_0 = 1$  is for the initial state and  $P_{i,j} = 0$  is the termination condition  
 760 of the state sequence.

761 Note that in the process of creating a new model,  $\hat{\mathcal{S}}_m$  will directly be equal  
 762 to the states of  $\xi_n$ . In the case of merging similar models, i.e. those with high  
 763 semantic similarity, one of the models will be assumed as  $\xi_n$  to employ the same  
 764 refinement procedure explained above. It is also important to note that all SEC  
 765 samples used for updating the same model  $\xi_m$  will be assigned with the same  
 766 cluster label which yields self-clustering of observed SEC samples.

## 767 Acknowledgements

768 The research leading to these results has received funding from the Euro-  
 769 pean Community’s Seventh Framework Programme FP7/2007-2013 (Specific  
 770 Programme Cooperation, Theme 3, Information and Communication Technolo-  
 771 gies) under grant agreement no. 270273, Xperience.

772 **References**

- 773 [1] S. Schaal, Is imitation learning the route to humanoid robots?, Trends in  
774 Cognitive Sciences 3 (1999) 233–242.
- 775 [2] A. Billard, S. Calinon, F. Guenter, Discriminative and adaptive imitation  
776 in uni-manual and bi-manual tasks, Robot. Auton. Syst. 54 (2006) 370–384.
- 777 [3] M. Pardowitz, S. Knoop, R. Dillmann, R. D. Zöllner, Incremental Learning  
778 of Tasks From User Demonstrations, Past Experiences, and Vocal Com-  
779 ments, IEEE Transactions on Systems, Man and Cybernetics – Part B:  
780 Cybernetics 37 (2) (2007) 322–332.
- 781 [4] S. Ekvall, D. Kragic, Robot learning from demonstration: a task-level plan-  
782 ning approach, International Journal of Advanced Robotic Systems 5 (3)  
783 (2008) 223–234.
- 784 [5] R. Cubek, W. Ertel, Learning and Execution of High-Level Concepts with  
785 Conceptual Spaces and PDDL, in: 3rd Workshop on Learning and Plan-  
786 ning, ICAPS (21st International Conference on Automated Planning and  
787 Scheduling), 2011.
- 788 [6] E. E. Aksoy, A. Abramov, F. Wörgötter, B. Dellen, Categorizing object-  
789 action relations from semantic scene graphs, in: IEEE International Con-  
790 ference on Robotics and Automation (ICRA), 2010, pp. 398–405.
- 791 [7] E. E. Aksoy, A. Abramov, J. Dörr, K. Ning, B. Dellen, F. Wörgötter,  
792 Learning the semantics of object-action relations by observation, The In-  
793 ternational Journal of Robotics Research 30 (10) (2011) 1229–1249.
- 794 [8] E. E. Aksoy, B. Dellen, M. Tamosiunaite, F. Wörgötter, Execution of a  
795 dual-object (pushing) action with semantic event chains, in: Proceedings  
796 of 11th IEEE-RAS International Conference on Humanoid Robots, 2011,  
797 pp. 576–583.

- 798 [9] E. E. Aksoy, M. Tamosiunaite, R. Vuga, A. Ude, C. Geib, M. Steedman,  
799 F. Wörgötter, Structural bootstrapping at the sensorimotor level for the  
800 fast acquisition of action knowledge for cognitive robots, in: IEEE Interna-  
801 tional Conference on Development and Learning and Epigenetic Robotics  
802 (ICDL-EPIROB), 2013.
- 803 [10] S. Niekum, S. Chitta, A. Barto, B. Marthi, S. Osentoski, Incremental se-  
804 mantically grounded learning from demonstration, Robotics Science and  
805 Systems (RSS) 2013.
- 806 [11] N. Badler, Temporal scene analysis: Conceptual descriptions of object  
807 movements, Ph.D. thesis, University of Toronto, Canada (1975).
- 808 [12] K. Ikeuchi, T. Suehiro, Toward an assembly plan from observation, part  
809 I: Task recognition with polyhedral objects, IEEE Trans. Robotics and  
810 Automation 10 (3) (1994) 368–385.
- 811 [13] M. Sridhar, G. A. Cohn, D. Hogg, Learning functional object-categories  
812 from a relational spatio-temporal representation, in: Proc. 18th European  
813 Conference on Artificial Intelligence, 2008, pp. 606–610.
- 814 [14] H. Kjellström, J. Romero, D. Kragić, Visual object-action recognition: In-  
815 ferring object affordances from human demonstration, Comput. Vis. Image  
816 Underst. 115 (1) (2011) 81–90.
- 817 [15] Y. Yang, C. Fermüller, Y. Aloimonos, Detection of manipulation action  
818 consequences (mac), in: International Conference on Computer Vision and  
819 Pattern Recognition (CVPR), 2013, pp. 2563–2570.
- 820 [16] D. Summers-Stay, C. Teo, Y. Yang, C. Fermüller, Y. Aloimonos, Using  
821 a minimal action grammar for activity understanding in the real world,  
822 in: Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International  
823 Conference on, 2012, pp. 4104–4111.

- 824 [17] H. S. Koppula, R. Gupta, A. Saxena, Learning human activities and ob-  
825 ject affordances from rgb-d videos, *The International Journal of Robotics*  
826 *Research* 32 (8) (2013) 951–970.
- 827 [18] D. Martinez, G. Alenya, P. Jimenez, C. Torras, J. Rossmann, N. Wantia,  
828 E. E. Aksoy, S. Haller, J. Piater, Active learning of manipulation sequences,  
829 (in press), in: *IEEE International Conference on Robotics and Automation*  
830 *(ICRA)*, 2014.
- 831 [19] G. Luo, N. Bergstrom, C. Ek, D. Kragic, Representing actions with kernels,  
832 in: *2011 IEEE/RSJ International Conference on Intelligent Robots and*  
833 *Systems (IROS)*, 2011, pp. 2028–2035.
- 834 [20] R. Vuga, E. E. Aksoy, F. Wörgötter, A. Ude, Augmenting semantic event  
835 chains with trajectory information for learning and recognition of manipu-  
836 lation tasks, in: *22nd International Workshop on Robotics in Alpe-Adria-*  
837 *Danube Region (RAAD)*, 2013.
- 838 [21] M. Wächter, S. Schulz, T. Asfour, E. E. Aksoy, F. Wörgötter, R. Dillmann,  
839 Action sequence reproduction based on automatic segmentation and object-  
840 action complexes, in: *IEEE/RAS International Conference on Humanoid*  
841 *Robots (Humanoids)*, 2013.
- 842 [22] J. Papon, T. Kulvicius, E. E. Aksoy, F. Wörgötter, Point cloud video object  
843 segmentation using a persistent supervoxel world-model, in: *IEEE/RSJ*  
844 *International Conference on Intelligent Robots and Systems*, 2013.
- 845 [23] M. J. Aein, E. E. Aksoy, M. Tamosiunaite, J. Papon, A. Ude, F. Wörgötter,  
846 Toward a library of manipulation actions based on semantic object-action  
847 relations, in: *IEEE/RSJ International Conference on Intelligent Robots*  
848 *and Systems*, 2013.
- 849 [24] K. Pauwels, N. Krüger, M. Lappe, F. Wörgötter, M. M. Van Hulle, A  
850 cortical architecture on parallel hardware for motion processing in real time,  
851 *Journal of Vision* 10.

- 852 [25] J. Papon, A. Abramov, E. E. Aksoy, F. Wörgötter, A modular system  
853 architecture for online parallel vision pipelines, in: IEEE Workshop on  
854 Applications of Computer Vision (WACV), 2012, pp. 361–368.
- 855 [26] A. Abramov, K. Pauwels, J. Papon, F. Wörgötter, B. Dellen, Depth-  
856 supported real-time video segmentation with the kinect, in: Applications  
857 of Computer Vision (WACV), 2012 IEEE Workshop on, 2012, pp. 457–464.
- 858 [27] A. Abramov, E. E. Aksoy, J. Dörr, K. Pauwels, F. Wörgötter, B. Dellen, 3d  
859 semantic representation of actions from efficient stereo-image-sequence seg-  
860 mentation on GPUs, in: 5th International Symposium 3D Data Processing,  
861 Visualization and Transmission, 2010, pp. 1–8.
- 862 [28] N. Otsu, A Threshold Selection Method from Gray-level Histograms, IEEE  
863 Transactions on Systems, Man and Cybernetics 9 (1) (1979) 62–66.
- 864 [29] J. Piaget, The Origins of Intelligence in the Child, Routledge, London, New  
865 York, 1953.
- 866 [30] E. E. Aksoy, M. Tamosiunaite, F. Wörgötter, Decomposition of long manip-  
867 ulation actions (under review), Computer Vision and Image Understanding.
- 868 [31] C. Schuldt, I. Laptev, B. Caputo, Recognizing human actions: a local svm  
869 approach, in: Pattern Recognition, 2004. ICPR 2004. Proceedings of the  
870 17th International Conference on, Vol. 3, 2004, pp. 32–36 Vol.3.
- 871 [32] A. Gupta, L. Davis, Objects in action: An approach for combining action  
872 understanding and object perception, in: Computer Vision and Pattern  
873 Recognition, 2007. CVPR '07. IEEE Conference on, 2007, pp. 1–8.
- 874 [33] F. Wörgötter, E. E. Aksoy, N. Krüger, J. Piater, A. Ude, M. Tamosiunaite,  
875 A simple ontology of manipulation actions based on hand-object relations,  
876 IEEE Transactions on Autonomous Mental Development 5 (2) (2013) 117–  
877 134.