

1 S4 Detailed model description

2 S4.1 Statistical model

3 The statistical model formally describes the framework and assumptions of the model for the kitchen
4 task scenario. First, an extension of the CSSM to correctly handle durative actions is described. Then
5 the overall DBN interpretation of the model is given. Afterward, the sensor model, duration model, and
6 action selection heuristics are discussed in detail.

7 S4.1.1 Durative actions

8 The model given in Appendix S2 assumes that actions take a single time step, where progression of time
9 is given by observation events. This requires that sensor data – usually some kind of real time signal –
10 is already partitioned into segments that correspond to individual actions, each segment representing a
11 single observation. In this case, the subscript t in y_t simply counts the number of actions that have been
12 observed. However, in general correct segmentation needs additional context information for disambigua-
13 tion [1], such as the possible sequence of actions. The assumption that the CSSM will be presented with
14 observations that represent complete action sightings, as for instance used in all but one of the previous
15 CSSM studies (Sec. 1.3.2), thus hides important aspects of the real-world inference problem.

16 Therefore, a model was chosen where multiple observations may correspond to a single action. This
17 model introduces real-valued random variables U and V that represent the starting time of an action A
18 and its duration. The action duration density $\tau(v | a, u) := p(V | A, U)$ defines the probability that action
19 A stops at time V if started at time U . Let $F_\tau(v | a, u)$ be the cumulative distribution function of τ . The
20 probability that an action’s duration lies in the interval $(v, v']$ then is $\zeta_{a,u}^{v,v'} := P(v < V \leq v' | a, u) =$
21 $F_\tau(v' | a, u) - F_\tau(v | a, u)$. If this duration is known to be greater than v , the corresponding conditional
22 probability is $P(v < V \leq v' | a, u, V > v) = \zeta_{a,u}^{v,v'} / \zeta_{a,u}^{v,\infty}$.

23 Note that durative actions will increase state space complexity.

24 S4.1.2 Overall CSSM model structure

25 For the probabilistic model of Sec. 1.2 (see Appendix S2 for details), a DBN with the structure given
26 in Fig. 1 was used. $Y_t = (W_t, Z_t)$ is the observation data for time step t . V_t is the associated time
27 stamp, required to be strictly increasing. As discussed in Appendix S2, state and action observations
28 were assumed to be conditionally independent, i. e., $p(y_t | a_t, s_t) = p(w_t | s_t) p(z_t | a_t)$.

29 $X_t = (A_t, D_t, G_t, S_t, U_t)$ defines the hidden state. For this study, G_t , the current goal, could be
30 assumed to be constant, so that $p(g_t | x_{t-1}) = [g_t = g_{t-1}]$. This allowed to efficiently precompute goal
31 distance values, which depend on G_t . The boolean flag D_t signals termination of A_{t-1} in the interval
32 $(v_{t-1}, v_t]$. It is a Bernoulli random variable defined by $p(D_t=1 | v_t, v_{t-1}, a_{t-1}, u_{t-1}) = \zeta_t / \zeta_{a_{t-1}, u_{t-1}}^{v_{t-1}, \infty}$ where
33 $\zeta_t := \zeta_{a_{t-1}, u_{t-1}}^{v_{t-1}, v_t}$. The variable D_t introduces a context-specific independence [2] into the DBN: if $d_t = 0$,
34 then A_t , S_t , and U_t carry over their values from the previous state and are independent of their other
35 parents. Otherwise, new values for A_t , S_t , and U_t are selected as follows:

36 A new action is selected according to γ by: $p(a_t | d_t = 1, s_{t-1}, a_{t-1}, g_t) = \gamma(a_t | g_t, s_{t-1}, a_{t-1})$. U_t
37 is the starting time of action A_t , given by $p(u_t | d_t = 1, u_{t-1}, v_t, v_{t-1}, a_{t-1}) = \bar{\tau}(u_t | u_{t-1}, v_t, v_{t-1}, a_{t-1})$,
38 where $\bar{\tau}(u_t | u_{t-1}, v_t, v_{t-1}, a_{t-1}) := [v_{t-1} < u_t \leq v_t] \tau(u_t | a_{t-1}, u_{t-1}) / \zeta_t$ is the truncated duration density,
39 constrained by $v_{t-1} < u_t \leq v_t$. S_t is the LTS state for time step t : either the result of applying the
40 new action to the previous state, or by carrying over the old state. For the purpose of this study,
41 actions could be assumed to be deterministic and with instantaneous effect, giving the simple model
42 $p(s_t | d_t = 1, s_{t-1}, a_t) = [s_t = a_t(s_{t-1})]$. As outlined for aLTS model development in Appendix S3,
43 non-deterministic effects for modeling complex interleaving were represented by multiple actions.

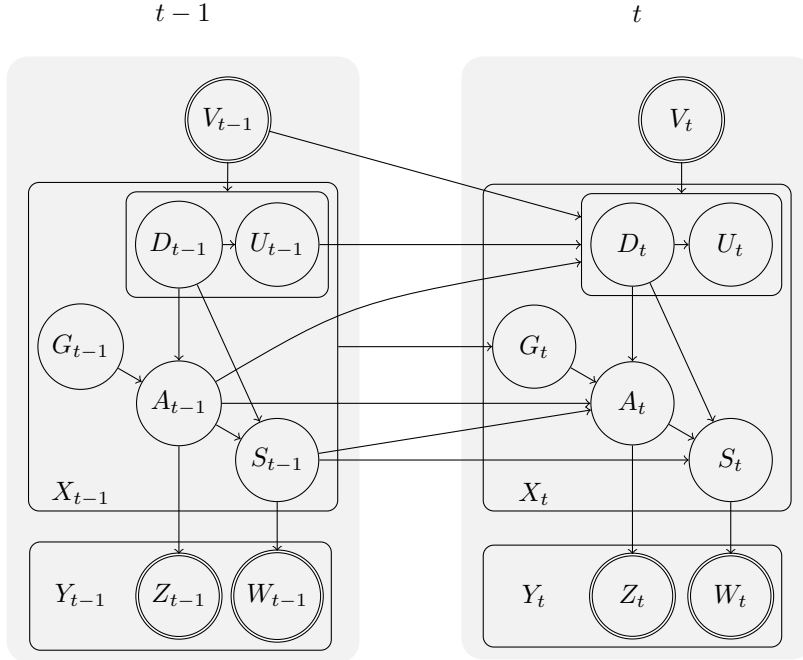


Figure 1. CSSM DBN structure. Boxes represent tuples of random variables. An arc starting / ending at a box (= a tuple) represents a set of arcs connected to the tuple’s components. Nodes with double outline signify observed random variables.

44 S4.1.3 Sensor models and scrambling

45 As construction of an elaborate observation model was not in the focus of this study, the following model
 46 was selected for its conceptual and computational simplicity: all actions a of a given class $c = class(a)$
 47 share the same observation distribution, such that $p(z|a) := p(z|c)$ and the distributions $p(z|c) :=$
 48 $N(z|\mu_c, \Sigma_c)$ were represented by multivariate normal distributions with unconstrained covariance Σ_c .
 49 The parameters μ_c, Σ_c were computed from the pool of all observations annotated with class c from all
 50 data sets. Although there is no reason to believe that the observation data is particularly well represented
 51 by this model, it was found to perform reasonably well in the baseline models, justifying its further use in
 52 this study. The 16 action classes can be seen in Fig. 6 or Fig. S3. For example, an action class is TAKE,
 53 while the actions belonging to it are take carrot, take bottle, take spoon etc.

54 A potential significant limitation of this model is its stationarity with respect to durative actions,
 55 which might be violated by the sensor data. A “run” for class c is a consecutive sequence of observations
 56 labeled with c . A run (often) represents an action that lasts several time steps. It might be the case that
 57 observations in the middle of a run are more “typical” for action class c than at the beginning (or end) of
 58 the run, where the transition to a different action class takes place and both kinds of motions are mixed.
 59 In this case, the expected probability of observations should increase from the run borders towards the
 60 run center. This effect should be specifically prominent for actions with long durations (which have long
 61 stable “center periods”). This could, for instance, impact the ability to correctly detect the start of such
 62 actions and eventually degrade recognition performance.

63 Conceptually, a better option to solve this problem is either to include temporal information into the
 64 observations [3] or – more general – to use a hierarchical approach, where the temporal structure of an
 65 action’s observations are represented by suitable sub-models [4]. A much simpler potential solution is

66 to scramble observations within a run. This will destroy order effects and other dependencies between
67 observations in a run. Clearly this approach is only possible in offline situations with data where the
68 runs are already delimited (either by annotations or by some suitable segmentation algorithm). However,
69 regarding the study objectives, the use of scrambling is considered a justified surrogate for defining
70 appropriate sub-models.

71 As alternative to the IMU sensors, which give continuous-domain observations of the A state com-
72 ponent, a location-based model was set up, giving categorical observations (place names) of the S state
73 component. Motivation for introducing this alternative observation model was to demonstrate that
74 CSSMs allow an exchange of observed X component without requiring adaptation of the system model.
75 In addition, it was of interest to:

- 76 (i) Test in how far a switch from continuous (IMU) to discrete (place names) observation modalities
77 introduces additional challenges for inference.
- 78 (ii) Provide a non-peaked observation model that assigns a non-negligible probability to every state, in
79 order to test susceptibility of inference method to this observation model property.

80 With respect to these objectives, a very simple model was considered sufficient, where the locations of
81 the protagonist (3 places) and the food (6 places) were used (see Tbl. S5). For the conditional probability
82 $p(z_t | s_t)$ the value 10^{-6} was used in case the locations observations z_t did not match the locations found
83 in s_t and essentially 1 otherwise. The observations themselves were computed from the aLTS model by
84 stepwise execution of the aLTS ground actions recorded in the annotations and using the location slot
85 values of the observed objects in the resulting state. While being realistic with respect to the temporal and
86 causal structure of the underlying human activity, these synthetic observations are unrealistic concerning
87 the precision of temporal alignment and error model; they can be expected to exaggerate the achievable
88 precision. Nevertheless, as this was a convenient mechanism to obtain data for new observation models
89 for simple comparison purposes, its use was considered legitimate.

90 S4.1.4 Duration model

91 For simplification it was assumed that all actions of a given class share the same action duration distri-
92 bution and that the duration of an action does not depend on the time the action has started. Therefore,
93 $\tau(v | a, u) = \tau^*(v - u | class(a))$ where $\tau^*(d | c)$ is the class-specific duration distribution. As there is yet
94 no prior knowledge on action durations, it was necessary to use the durations found in the training data
95 as proxy.

96 Note that duration distributions $\tau^*(d | c)$ with large or even infinite support increase the branching
97 factor. In order to determine this effect on inference performance, an instance based model (with finite
98 support) and a parametric duration model (infinite support) was built. Let $(d_{c,1}, \dots, d_{c,n_c})$ be the set
99 of durations observed for actions of class c . The instance based model was given by the corresponding
100 empirical distribution function, so that $F_{\tau^*}(d | c) := n_c^{-1} \sum_{i=1}^{n_c} [d_{c,i} \leq d]$, where F_{τ^*} is the distribution
101 function of τ^* . For the parametric models, an approach was chosen that would provide a trade-off between
102 technical simplicity and flexibility. First, a lognormal model for the pooled observation was built, see
103 Fig. 2 for a comparison with the actual durations and a kernel density estimate. Then those classes were
104 determined whose means were significantly (at the .05 α level) different from the pooled mean by fitting
105 an ANOVA model (analysis of variance). For these special classes as well as for the remaining pooled
106 classes, the distribution giving the maximum likelihood was selected from a set of candidate distributions
107 whose parameters were fitted to the observed class durations. The set of candidate distributions was
108 Cauchy, exponential, gamma, geometric, lognormal, negative binomial, normal, Poisson, and Weibull.

109 The parametric models were expected to have an impact on performance for two reasons: (i) they
110 might not provide a good fit to the empirical data, (ii) since they are continuous, they provide a greater
111 number of possible durations, thereby increasing the branching factor and thus inference complexity.

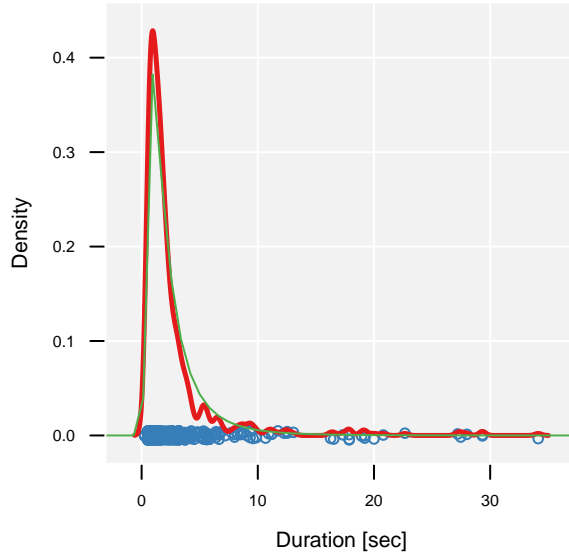


Figure 2. Kernel density estimate of pooled action durations vs. fitted lognormal density. Blue circles are the observed action durations. The red line is the kernel density estimate, the green line the fitted lognormal density.

112 S4.1.5 Action selection model

113 As primary goal-driven action selection feature, the goal distance feature f_δ as discussed in Appendix S2
 114 was chosen. Goal distances were computed by an exhaustive process, computing the set of LTS states
 115 reachable from the initial state and then for all reachable states the shortest path to a goal state. For
 116 simplicity, unit costs were used for actions. This approach will become intractable with models of increas-
 117 ing state space complexity. In addition, it is not yet established that goal distance indeed best describes
 118 human action selection preferences. Approximate distance values therefore might achieve the same result,
 119 for a much lower computational cost. Two approximations were considered in this study:

- 120 • A goal distance heuristic f_h , that assigns heuristic distances to LTS states based on values of LTS
 121 state variables. The trial task’s script consisting of 14 serial task steps was used to define a map
 122 from LTS state s to remaining script steps $h(s)$ based on its state variable values (see Tbl. S1 for
 123 this map). $h(s)$ then was used as distance value. (Most script tasks required several LTS actions for
 124 their realization.) The motivation of this heuristic is the idea that prior knowledge on the typical
 125 coarse-grained sequential structure of everyday activities should be easy to obtain [5,6]. Analyzing
 126 f_h thus should provide insight into the usefulness of such knowledge.
- 127 • A restricted goal distance feature $f_{\bar{\delta}}$. Here, only those LTS states were considered that are visited
 128 when using the annotations as exact observations. All other states discovered during inference
 129 received the nominal goal distance 100. Restricted goal distance $f_{\bar{\delta}}$ should give an upper limit to
 130 the gain achievable by a goal distance measure.

131 To gain insight on the effect of weight factors, each of these features was tested with the weight values
 132 $\lambda_i = -(2^k), k \in \{0, \dots, 4\}$, using exponential probing. To limit experimental complexity, we refrained
 133 from evaluating interactions between features by only assigning one of $f_\delta, f_h, f_{\bar{\delta}}$ a non-zero weight.

134 In the special case $\lambda_\delta = \lambda_h = \lambda_{\bar{\delta}} = 0$ all applicable actions receive the same selection probability,
 135 resulting in uniform selection. (This *locally uniform* selection strategy is not the same as giving each

136 possible action *sequence* the same probability – such a *globally uniform* action selection model requires a
137 considerable more complex feature computation.)

138 S4.2 iLTS model development

139 Building the inference LTS (iLTS, cf. Appendix S3) requires two steps: choosing an appropriate modeling
140 language and implementing the model. The following two sections justify the choice and present the model
141 development process used.

142 S4.2.1 Modeling language

143 As pure execution speed was not the primary concern of this study, Common Lisp was chosen as im-
144 plementation language for the inference procedures. This allowed a very convenient way to represent
145 actions in a PDDL / STRIPS like fashion as S-expressions, providing considerable latitude regarding the
146 expressive power available for defining precondition terms and effect expressions, as essentially any valid
147 Lisp expression could be used. For instance, the action of taking the food from the cutting board was
148 defined by the following expression:

```
149 (:action (take food cutting_board)
150          :precondition (and (eq ($ location self) ($ location cutting_board))
151                            (eq ($ location food) 'cutting_board)
152                            (> ($ available hands) 0))
153          :effect ((($ location food) 'hands
154                  ($ available hands) (- ($ available hands) 1)))
```

155 State variables are referred to by symbol lists, such as (`$ location self`). The precondition states
156 that the protagonist (`self`) and the cutting board must be at the same location, the food must be on
157 the cutting board, and the protagonist needs to have at least one free hand. The effect states that the
158 location of the food will now be one of the protagonist’s hand and that he has one less free hand available.

159 The state model was created by traversing all action definitions, collecting the state variable references,
160 and compiling this to a Lisp `defstruct` containing one structure slot for each state variable. At the same
161 time, action definitions were compiled to two pairs of Lisp functions, the first computing the precondition
162 value, the second one for applying the effect to a state. State variable references were replaced by
163 the corresponding slot access functions for the state `defstruct`; the value assignment in the effect was
164 performed via `setf`. For efficiency considerations, only atomic values (symbols, characters, integers, and
165 floating point numbers) were allowed as values for state variables. Therefore, states could be considered
166 as words of constant length, allowing the use of tries as dictionaries for handling state sets, providing
167 significantly faster access than hash tables in the Lisp environment [7] used in this study.

168 S4.2.2 Model development process

169 A two-stage process was used for iLTS model development. First a *feasible solution* \mathcal{A}^* for the domain
170 model (the set of action definitions) was constructed iteratively, and then this solution was refined to
171 the final model \mathcal{A} using local model modifications. For constructing the feasible solution, the domain
172 expert used an iterative procedure for successively building action sets $\mathcal{A}_0, \mathcal{A}_1, \dots, \mathcal{A}_S =: \mathcal{A}^*$ of increasing
173 complexity.

174 Constructing the feasible solution begins by setting $\mathcal{A}_0 := \emptyset$. For the current action set \mathcal{A}_s , the
175 model developer applies all action sequences obtained from the annotations in parallel and identifies the
176 smallest t where the action $a_{i,t}$ of some dataset i fails at time t . For instance, the first action `wash hands`
177 from the annotations of subject 1 (cf. Tbl. S3) is executed in the iLTS of \mathcal{A}_0 , failing because this action
178 is not included in $\mathcal{A}_0 = \emptyset$. If no action sequence fails, the feasible solution $\mathcal{A}^* := \mathcal{A}_s$ has been found.

179 Otherwise either (a) $a_{i,t}$ refers to an action not contained in \mathcal{A}_s (as in the example), or (b) $a_{i,t}$ is in
180 \mathcal{A}_s , but its precondition is not met in the current iLTS state after applying the preceding actions (for
181 instance because `take` is executed twice, but was previously modeled to allow only one object in hand,
182 cf. Tbl. S3, $t = 4-8$). In case (a), the model developer adds a new action schema with precondition and
183 effect matching the informal semantics of the domain. In case (b), the model developer either relaxes the
184 precondition or extends the effect of some other action such that the precondition of $a_{i,t}$ is met. It is the
185 task of the domain expert to judge which of the refinement strategies to apply for reconciling a failure.
186 The developer repeats the process with $\mathcal{A}_{s+1} := \mathcal{A}_s$ until a feasible solution has been found.

187 In the second step, the domain expert would refine \mathcal{A}^* to the final action set \mathcal{A} by adding preconditions
188 in order to limit the number of possible plans. (For instance, in the case of the given trial setting, the
189 protagonist would sit down at the table only after the food has been cooked). Note that this two stage
190 process – consisting of the iterative construction of a feasible solution and then applying local search to
191 refine the feasible solution towards a local optimum \mathcal{A} – is similar to heuristic optimization strategies for
192 solving integer programming problems.

193 References

- 194 1. Duda RO, Hart PE, Stork DG (2003) Pattern Classification. Wiley, second edition.
- 195 2. Boutilier C, Friedman N, Goldszmidt M, Koller D (1996) Context-specific independence in bayesian
196 networks. In: Proceedings of the Twelfth International Conference on Uncertainty in Artificial
197 Intelligence (UAI). San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., pp. 115–123.
- 198 3. Milner B (1996) Inclusion of temporal information into features for speech recognition. In: Proceed-
199 ings of the Fourth International Conference on Spoken Language. IEEE, volume 1, pp. 256–259.
- 200 4. Ostendorf M, Digalakis V, Kimball OA (1995) From hmms to segment models: A unified view of
201 stochastic modeling for speech recognition. IEEE Transactions on Speech and Audio Processing 4:
202 360–378.
- 203 5. Rashidi P, Cook DJ (2013) COM: A method for mining and monitoring human activity patterns in
204 home-based health monitoring systems. ACM Transactions of Intelligent Systems and Technology
205 4: 64:1–64:20.
- 206 6. Wyatt D, Philipose M, Choudhury T (2005) Unsupervised activity recognition using automatically
207 mined common sense. In: Proceedings of the 20th national conference on Artificial Intelligence
208 (AAAI). AAAI Press, pp. 21–27.
- 209 7. Clozure Associates. Clozure common lisp. URL <http://cc1.clozure.com/>. Accessed: 2014-03-26.