# From Stochastic Grammar to Bayes Network:

# Probabilistic Parsing of Complex Activity

**CVPR14 Submission #83**

**Supplementary Document**

In section 1, we provide a detail description of how the inference is implemented. We further discuss primitive action in section 2 and 3. Finally we show the full grammar of the toy assembly task used in our experiment.

**1. Inference by Message Passing**

**Input:**

- The constructed Bayes network.

- CPT $P(v_e \mid v_s)$ for every primitive v

- CPT $P(Z^v \mid v_s, v_e)$ for every primitive v, including special value $P(Z^v \mid !v)$

- Prior information $P(\exists M \mid \exists A)$ for every OR-rule A -> M | …

- Prior information $P(\exists S) = 1$

- Prior information $P(S_s \mid \exists S)$

- Prior information $P(Z^{end} \mid S_e, \exists S)$

**Step 0 :** For every composition A, recursively compute:

$$P(Z^A \mid !A) = \prod_{M \text{ in } A} P(Z^M \mid !M)$$

Note: since scaling the likelihood $P(Z^v \mid v_s, v_e)$ of primitive v does not change the inference result. In our implementation, we allow the value bigger than 1 and scale it so that $P(Z^v \mid !v) = 1$ for every primitive v. Then $P(Z^A \mid !A) = 1$ for every A. Then we can safely ignore them in following calculation.

**Step 1 – Forward Phase:**

26 • Forward Phase on primitive v: assume $P(v_s, Z^{pre(v)} | \exists v)$ is given, compute:

27 $$P(v_s, v_e, Z^{pre(v),v} | \exists v) = P(v_s, Z^{pre(v)} | \exists v) P(v_e | v_s) P(Z^v | v_s, v_e)$$

28 $$P(v_e, Z^{pre(v),v} | \exists v) = \sum_{t=1}^{T} P(v_s = t, v_e, Z^{pre(v),v} | \exists v)$$

29 • Forward Phase on composition A defined by A -> M N: given $P(A_s, Z^{pre(A)} | \exists A)$, compute:

30 $$P(M_s = t, Z^{pre(M)} | \exists M) = P(A_s = t, Z^{pre(A)} | \exists A)$$

31 For t is between 1 and T. This will apply for every t is used in all following formulas.

32 Recursively perform forward phase on M to get $P(M_e, Z^{pre(M),M} | \exists M)$

33 $$P(N_s = t, Z^{pre(N)} | \exists N) = P(M_e = t, Z^{pre(M)}, M | \exists M)$$

34 Recursively perform forward phase on N to get $P(N_e, Z^{pre(N),N} | \exists N)$

35 $$P(A_e = t, Z^{pre(A),A} | \exists A) = P(N_e = t, Z^{pre(N),N} | \exists N)$$

36 • Forward Phase on composition A defined by A -> M | N: given $P(A_s, Z^{pre(A)} | \exists A)$, compute:

37 $$P(M_s = t, Z^{pre(M)} | \exists M) = P(A_s = t, Z^{pre(A)} | \exists A)$$

38 $$P(N_s = t, Z^{pre(N)} | \exists N) = P(A_s = t, Z^{pre(A)} | \exists A)$$

39 Recursively perform forward phase on M and N to get $P(M_e, Z^{pre(M),M} | \exists M)$ and

40 $P(N_e, Z^{pre(N),N} | \exists N)$, then:

41 $$P(A_e = t, Z^{A,pre(A)} | \exists A) = P(\exists M | \exists A) P(Z^N |!N) P(M_e = t, Z^{M,pre(M)} | \exists M)$$

42 $$P(\exists N | \exists A) P(Z^M |!M) P(N_e = t, Z^{N,pre(N)} | \exists N)$$

43 • Start from $P(S_s | \exists S)$, perform forward phase on S, and recursively on other actions. The

44 output is $P(A_s, Z^{pre(A)} | \exists A)$ and $P(A_e, Z^{pre(A),A} | \exists A)$ for every A.

45

46 **Step 2 – Backward Phase: similar to Forward Phase**

47 • Backward Phase on primitive v: assume $P(Z^{post(v)} | v_e, \exists v)$ is given, compute:

48 $$P(v_e, Z^{v,post(v)} | v_s, \exists v) = P(Z^{post(v)} | v_e, \exists v) P(v_e | v_s) P(Z^v | v_s, v_e)$$

49     $$P(Z^{v,post(v)} \mid v_s, \exists v) = \sum_{t=1}^{T} P(v_e = t, Z^{v,post(v)} \mid v_s, \exists v)$$

50     •   Backward Phase on composition A defined by A -> M N: given $P(Z^{post(A)} \mid A_e, \exists A)$, compute:

51     $$P(Z^{post(N)} \mid N_e = t, \exists N) = P(Z^{post(A)} \mid A_e = t, \exists A)$$

52     Recursively perform backward phase on N to get $P(Z^{N,post(N)} \mid N_s, \exists N)$

53     $$P(Z^{post(M)} \mid M_e = t, \exists M) = P(Z^{N,post(N)} \mid N_s = t, \exists N)$$

54     Recursively perform backward phase on N to get $P(Z^{M,post(M)} \mid M_s, \exists M)$. Then:

55     $$P(A^{A,post(A)} \mid A_s = t, \exists A) = P(Z^{M,post(M)} \mid M_s = t, \exists M)$$

56     •   Backward Phase on composition A defined by A -> M | N: given $P(Z^{post(A)} \mid A_e, \exists A)$, compute:

57     $$P(Z^{post(M)} \mid M_e = t, \exists M) = P(Z^{post(A)} \mid A_e = t, \exists A)$$

58     $$P(Z^{post(N)} \mid N_e = t, \exists N) = P(Z^{post(A)} \mid A_e = t, \exists A)$$

59     Recursively perform backward phase on M and N to get $P(Z^{M,post(M)} \mid M_s, \exists M)$ and

60     $P(Z^{N,post(N)} \mid N_s, \exists N)$. Then:

61     $$P(Z^{A,post(A)} \mid A_s = t, \exists A) = P(\exists M \mid \exists A) P(Z^N \mid !N) P(Z^{M,post(M)} \mid M_s = t, \exists M)$$

62     $$+ P(\exists N \mid \exists A) P(Z^M \mid !M) P(Z^{N,post(N)} \mid N_s = t, \exists N)$$

63     •   Start from $P(Z^{end} \mid S_e, \exists S)$, perform backward phase on S and recursively on other action. The

64     output is $P(Z^{post(A)} \mid A_e, \exists A)$ and $P(Z^{A,post(A)} \mid A_s, \exists A)$.

65   **Step 3 – Compute Posterior Probability:** by multiplying forward and backward messages:

66     $$P(A_s, Z^{pre(A),A,post(A)} \mid \exists A) = P(A_s, Z^{pre(A)} \mid \exists A) P(Z^{A,post(A)} \mid A_s, \exists A)$$

67     $$P(A_e, Z^{pre(A),A,post(A)} \mid \exists A) = P(A_e, Z^{pre(A),A} \mid \exists A) P(Z^{post(A)} \mid A_e, \exists A)$$

68     $$P(A_s, Z \mid \exists A) = P(A_s, Z^{pre(A),A,post(A)} \mid \exists A) \prod_{M \text{ not in pre(A), A, post(A)}} P(Z^M \mid !M)$$

69 $$P(A_e, Z \mid \exists A) = P(A_e, Z^{pre(A), A, post(A)} \mid \exists A) \prod_{\text{M not in pre(A), A, post(A)}} P(Z^M \mid !M)$$

70 If v is a primitive we can have the joint distribution:

71 $$P(v_s, v_e, Z^{pre(v), v, post(v)} \mid \exists v) = P(v_s, v_e, Z^{pre(v), v} \mid \exists v) P(Z^{post(v)} \mid v_e, \exists v)$$

72 $$P(v_s, v_e, Z \mid \exists v) = P(v_s, v_e, Z^{pre(v), v, post(v)} \mid \exists v) \prod_{\text{M not in pre(v), v, post(v)}} P(Z^M \mid !M)$$

73 **Step 4 – Compute the happening probability:** Start from $P(\exists S \mid Z) = P(\exists S) = 1$

74 • For AND-rule A -> M N, assume $P(\exists A \mid Z)$ is given, compute:

75 $$P(\exists M \mid Z) = P(\exists N \mid Z) = P(\exists A \mid Z)$$

76 • For OR-rule A -> M | N, given $P(\exists A \mid Z)$, compute:

77 $$P(\exists M, Z \mid \exists A) = P(\exists M \mid \exists A) \sum_{t=1}^{T} P(M_e = t, Z \mid \exists M)$$

78 $$P(\exists N, Z \mid \exists A) = P(\exists N \mid \exists A) \sum_{t=1}^{T} P(N_e = t, Z \mid \exists N)$$

79 $$P(\exists M \mid Z) = P(\exists A \mid Z) \frac{P(\exists M, Z \mid \exists A)}{P(\exists M, Z \mid \exists A) + P(\exists N, Z \mid \exists A)}$$

80 $$P(\exists N \mid Z) = P(\exists A \mid Z) \frac{P(\exists N, Z \mid \exists A)}{P(\exists M, Z \mid \exists A) + P(\exists N, Z \mid \exists A)}$$

81 **Output:** For every action A in the grammar: $P(\exists A \mid Z)$, $P(A_s, Z \mid \exists A)$ and $P(A_e, Z \mid \exists A)$. If v is a

82 primitive we can have the joint: $P(v_s, v_e, Z \mid \exists v)$.

83 Optionally we can have 2 more steps:

84 **Step 5:** For every action A, we can compute $P(A_s \mid Z)$ and $P(A_e \mid Z)$. If v is a primitive then we also

85 have $P(v_s, v_e \mid Z)$.

86 Probability of label of a time step t being action A: $P(label_t = A \mid Z)$ can also be derived.

87 The calculation is shown in section 4.5 in the paper.

88 **Step 6 (optional & not in the paper) Compute the joint of the start and the end for every action:** This
89 can be done if needed. However the computational complexity will change.

90 • For every primitive v, compute $P(v_e, Z^v \mid v_s) = P(v_e \mid v_s)P(Z^v \mid v_s, v_e)$

91 • For every composition A -> M, N:

92 Recursively compute $P(M_e, Z^M \mid M_s)$, $P(N_e, Z^N \mid N_s)$. Then:

93 $$P(A_e = \beta, Z^A \mid A_s = \alpha) = \sum_{t=1}^{T} P(M_e = t, Z^M \mid M_s = \alpha)P(N_e = \beta, Z^N \mid N_s = t)$$

94 • For every composition A -> M | N:

95 Recursively compute $P(M_e, Z^M \mid M_s)$, $P(N_e, Z^N \mid N_s)$. Then:

96 $$P(A_e = \beta, Z^A \mid A_s = \alpha) = P(\exists M \mid \exists A)P(Z^N \mid !N)P(M_e = \beta, Z^M \mid M_s = \alpha)$$

97 $$+P(\exists N \mid \exists A)P(Z^M \mid !M)P(N_e = \beta, Z^N \mid N_s = \alpha)$$

98 For every value of $\alpha, \beta$ between 1 and T.

99 • Given $P(A_e, Z^A \mid A_s)$ for every A, we can compute the joint:

100 $$P(A_s, A_e, Z^{pre(A), A, post(A)} \mid \exists A) = P(A_s, Z^{pre(A)} \mid \exists A)P(A_e, Z^A \mid A_s)P(Z^{post(A)} \mid A_e)$$

101 $$P(A_s, A_e, Z \mid \exists A) = P(A_s, A_e, Z^{pre(A), A, post(A)} \mid \exists A) \prod_{M \text{ not in pre(A), A, post(A)}} P(Z^M \mid !M)$$

102 • Then for every A, we can compute $P(A_s, A_e \mid Z)$ similar to step 5.

103

104 **Computational Complexity:** The inference process starts from S and then performs on all symbols
105 recursively like a depth-first-search travel on the AND-OR tree representation of the grammar. (Note
106 that if the grammar does not have any OR-rule, it becomes a traditional message passing algorithm on a
107 linear chain). Each symbol is "visited" 4 times (4 above steps), there are calculations of vectors of size
108 Tx1 and matrices of size TxT in the step 1 and step 2. Overall the complexity is $O(KT^2)$ where K is the
109 number of symbol in the compiled grammar. With K=50 and T=1000, our Matlab implementation runs in
110 0.1 second on an average machine (CPU 2.5GHz, RAM 6GB).

111 If step 6 is performed, the complexity becomes $O(KT^3)$.

112    Note that even running in streaming mode, each inference is independent of each other. Hence the
113    inference rate does not need to be the same as video rate. In fact one can choose to only perform
114    inference when needed.

115

116    **2. Primitive action**

117    Calculating $D_v$ for primitive action v is only ingredient for the Bayes network that makes use of the test
118    sequence and can be the trickiest one to compute. We assume, for each primitive, there is a detector
119    that will output the TxT "heatmap" $D_v$ of the likelihood of the action for every possible interval. Ideally
120    if the action starts at $\alpha_0$ and ends at $\beta_0$ then $D_v[\alpha_0, \beta_0]$ would be high and $D_v[\alpha, \beta]$ would be low
121    for every other $\alpha, \beta$ value.

122    The detector is assumed to be black-box. It can be driven by explicitly detecting the start and the end of
123    the action (experiment in section 6.2 and 5). An alternative way is to perform sliding-window-detection
124    using statistics/features computed over the $[\alpha, \beta]$ segment (experiment in section 6.1). Note that the
125    calculation of $D_v[\alpha, \beta]$ can use the information of the entire input sequence if desired, not just the
126    $[\alpha, \beta]$ segment.

127    As the factor P(v.end | v.start) accounts for the duration of the action and the factors
128    $P(Z^v | v.start, v.end)$ accounts for the visual information of the action, the visual detector do not need
129    to concern about duration. Although one could derive a detector like that (or combine the 2 factors into
130    1 single factor $P(v.end, Z^v | v.start)$), we find that keeping these 2 factors separate is more flexible.
131    That way we could change them independently, and we can run in streaming mode, where visual
132    information is feed sequentially.

133    **Interpretation of the special value $D_v[-1, -1]$**

134    This value represents how likely the action v does not happen. Traditionally, non-maxima suppression
135    and thresholding are performed on the heatmap $D_v$ to obtain a set of detections. One can interpret the
136    $D_v[-1, -1]$ as the threshold: a high value means the action more likely does not happen. Informally
137    $D_v[\alpha, \beta] / D_v[-1, -1] > 1$ means segment $[\alpha, \beta]$ is a positive and the confidence is proportional with
138    that ratio value. Where $D_v[\alpha, \beta] / D_v[-1, -1] = 1$ basically means "nothing is known about $[\alpha, \beta]$" (we
139    made use of this in streaming mode). In our implementation, we choose $D_v[-1, -1]$ to be about the
140    expected detection score so that it has above properties (though if one has a way to check if the action
141    does not happen, it could also be incorporated).

142    Intuitively, $D_v[\alpha, \beta]$ and $D_v[-1, -1]$ put relative weights on the probability of sequence where the
143    action happens and the sequences where it does not, respectively, when we are considering the OR-

144    rule. For example if $D_v[\alpha, \beta] / D_v[-1, -1]$ is very big for some value $\alpha, \beta$ and these values are also

145    temporally consistent with overall activity's structure, this would contribute to increase the posterior

146    probability of the sequence where action v happens. Note that if the grammar does not have any OR-

147    rules, then the value $D_v[-1, -1]$ will not affect the inference result.

148

## 3. Special primitive action

150    One can include primitive actions with special duration factor or visual observation factor to serve

151    specific purpose.

152    **0-Duration Action:** This action always has duration to be 0; and $D_v[-1, -1] = 1$, and $D_v[\alpha, \beta] = f(Z^v)$.

153    It will not affect the inference result of action localization within a sequence of actions. However its

154    visual observation factor will affect the relative posterior probabilities between sequences where it

155    happens and the sequences where it doesn't.

156    **Dummy Action:** This action does not have a visual observation factor (or equivalently constant

157    likelihood value: $D_v[\alpha, \beta] = D_v[-1, -1] = 1$). It can serve as the gap between 2 actions in case we

158    assume the start time of the next action is not the same as the end time of the current one.

159    **Negative-Duration Dummy Action:** similar to above dummy action, except its duration is allowed to be

160    negative. Including this between 2 actions allows them to overlap each other. This will be useful for

161    approaches that recognize activity by recognizing overlapping segments of that activity.

162    **Waiting Action:** in the Human-Robot collaboration application that we applied our method, the robot

163    delivers the bins to the human operator. In case the human needs a specific bin that the robot has not

164    yet delivered, he will have to wait. Therefore we designed a special "waiting" action that can starts at

165    any moment in time but only ends when that bin is delivered (if the bin is already available, the action, if

166    starts, will end immediately and have duration of 0).

167

## 4. Toy assembly grammar

| 169 | S | → | Body, (Wheel \| null), NWT, sticker |
| 170 | NWT | → | NWT_AB ~ 60%  \| NWT_C ~ 40% |
| 171 | NWT_C | → | Nose_C, ((Wing_C, Tail_C) \| (Tail_C, Wing_C)) |
| 172 | NWT_AB | → | Nose_AB, (WT_A \| WT_B) |
| 173 | WT_A | → | (Wing_A, Tail_A) \| (Tail_A, Wing_A ) |
| 174 | WT_B | → | (Wing_B, Tail_B) \| (Tail_B, Wing_B ) |
| 175 | Body | → | body1, body2, body3, body4 |
| 176 | Wheel | → | wheel1, wheel2 |
| 177 | Nose_AB | → | nose_ab1, nose_ab2, nose_ab3, nose_ab4 |

| 178 | Nose_C | → | nose_c1, nose_c2, nose_c3 |
| 179 | Wing_A | → | wing_a1, wing_a2, wing_a3 |
| 180 | Wing_B | → | wing_b1, wing_b2, wing_b3, wing_b4 |
| 181 | Wing_C | → | wing_c1, wing_c2, wing_c3, wing_c4, wing_c5, wing_c6 |
| 182 | Tail_A | → | tail_a1, tail_a2, tail_a3 |
| 183 | Tail_B | → | tail_b1, tail_b2, tail_b3, tail_b4 |
| 184 | Tail_C | → | tail_c1, tail_c2, tail_c3, tail_c4, tail_c5, tail_c6 |
| 185 | | | |

186    Map between primitive actions and corresponding bins:

| 187 | body1 | 5 | |
| 188 | body2 | 5 | |
| 189 | body3 | 3 | |
| 190 | body4 | 4 | |
| 191 | wheel1 | 3 | |
| 192 | wheel2 | 3 | |
| 193 | nose_ab1 | 3 | |
| 194 | nose_ab2 | 4 | |
| 195 | nose_ab3 | 3 | |
| 196 | nose_ab4 | 3 | |
| 197 | nose_c1 | | 3 |
| 198 | nose_c2 | | 4 |
| 199 | nose_c3 | | 3 |
| 200 | | | |
| 201 | wing_a1 | | 3 |
| 202 | wing_a2 | | 1 |
| 203 | wing_a3 | | 4 |
| 204 | wing_b1 | | 3 |
| 205 | wing_b2 | | 1 |
| 206 | wing_b3 | | 1 |
| 207 | wing_b4 | | 4 |
| 208 | wing_c1 | | 3 |
| 209 | wing_c2 | | 2 |
| 210 | wing_c3 | | 2 |
| 211 | wing_c4 | | 1 |
| 212 | wing_c5 | | 1 |
| 213 | wing_c6 | | 2 |
| 214 | tail_a1 | 3 | |
| 215 | tail_a2 | 5 | |
| 216 | tail_a3 | 4 | |
| 217 | tail_b1 | 3 | |

| 218 | tail_b2 | 5 |
| --- | --- | --- |
| 219 | tail_b3 | 5 |
| 220 | tail_b4 | 4 |
| 221 | tail_c1 | 3 |
| 222 | tail_c2 | 2 |
| 223 | tail_c3 | 2 |
| 224 | tail_c4 | 5 |
| 225 | tail_c5 | 5 |
| 226 | tail_c6 | 2 |
| 227 | sticker | 2 |
| 228 | | |
| 229 | | |
| 230 | | |
| 231 | | |
| 232 | | |