

Ch.9 Bottom-up/top-down inference in And-Or-Graph

--- Event Recognition and Intent Prediction

Earley-Stolcke parsing algorithm

- $S \rightarrow A B C$
 - $S \rightarrow D E F$
 - $A \rightarrow a_1$ $A \rightarrow a_2$
 - $B \rightarrow b_1$ $B \rightarrow b_2$
 - $C \rightarrow c_1$ $C \rightarrow c_2$
 - $D \rightarrow a_1$ $D \rightarrow a_2$
 - $E \rightarrow b_1$ $E \rightarrow b_2$
 - $F \rightarrow c_1$ $F \rightarrow c_2$
- $a_1 b_1 c_1$
 - $a_1 b_1 c_2$
 - $a_1 b_2 c_1$
 - $a_1 b_2 c_2$
 - $a_2 b_1 c_1$
 - $a_2 b_1 c_2$
 - $a_2 b_2 c_1$
 - $a_2 b_2 c_2$

Earley-Stolcke parsing algorithm

The input is a1 b1 c1

State set 0	State set 1	State set 2	State set 3
$_0 \rightarrow .S$ Predicted $_0 S \rightarrow .A B C$ $_0 S \rightarrow .D E F$ $_0 A \rightarrow .a1$ $_0 A \rightarrow .a2$ $_0 D \rightarrow .a1$ $_0 D \rightarrow .a2$	Scanned $_0 A \rightarrow a1.$ $_0 D \rightarrow a1.$ Completed $_0 S \rightarrow A. B C$ $_0 S \rightarrow D. E F$ Predicted $_1 B \rightarrow .b1$ $_1 B \rightarrow .b2$ $_1 E \rightarrow .b1$ $_1 E \rightarrow .b2$	Scanned $_1 B \rightarrow b1.$ $_1 E \rightarrow b1.$ Completed $_0 S \rightarrow A B. C$ $_0 S \rightarrow D E. F$ Predicted $_2 C \rightarrow .c1$ $_2 C \rightarrow .c2$ $_2 F \rightarrow .c1$ $_2 F \rightarrow .c2$	Scanned $_2 C \rightarrow c1.$ $_2 F \rightarrow c1.$ Completed $_0 S \rightarrow A B C.$ $_0 S \rightarrow D E F.$ $_0 \rightarrow S.$

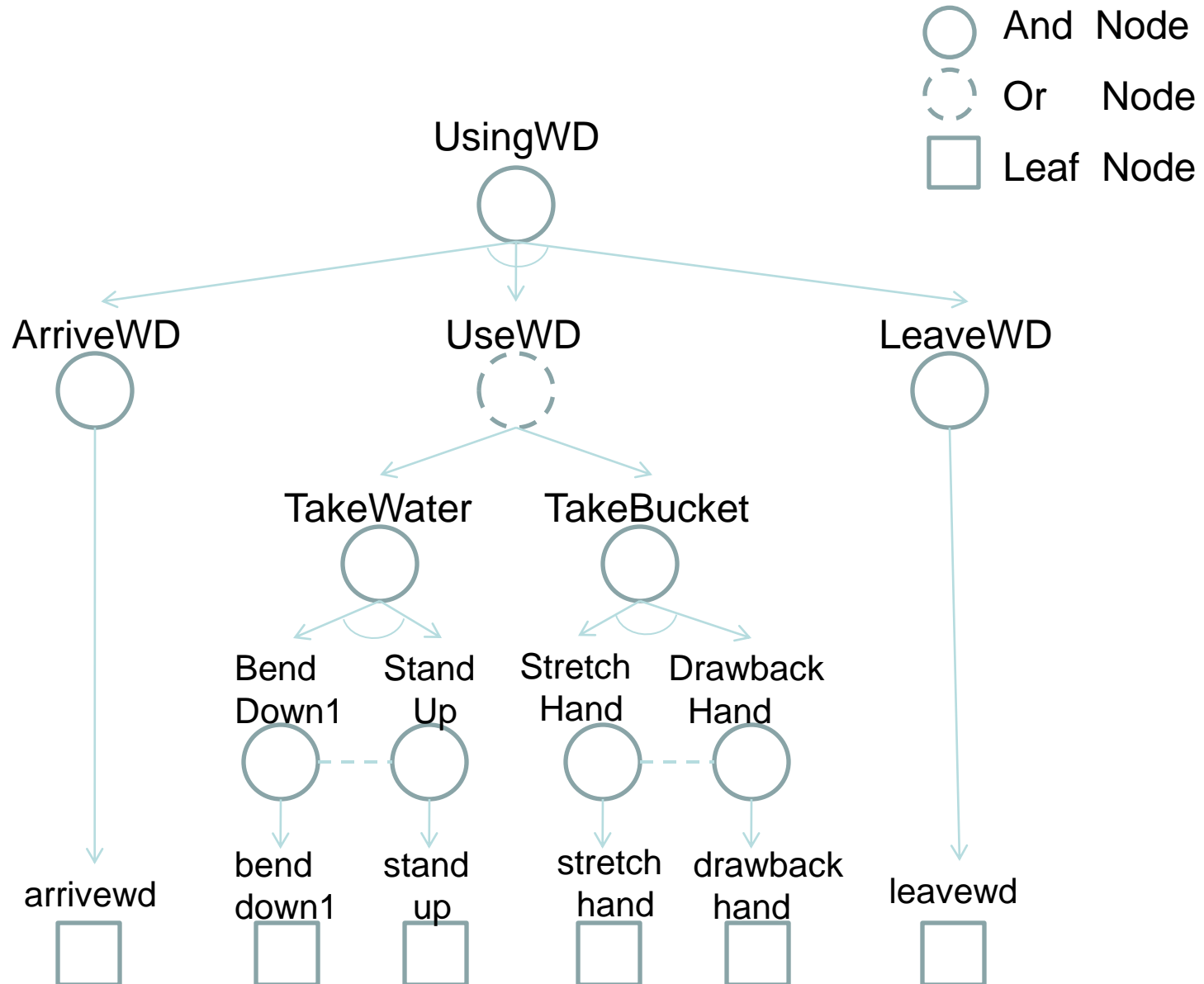
Grammars for Events

- UsingWD → ArriveWD UseWD LeaveWD
- ArriveWD → arrivewd
- UseWD → TakeWater
- UseWD → TakeBucket
- TakeWater → benddown1 standup
- TakeBucket → stretchhand drawbackhand
- LeaveWD → leavewd

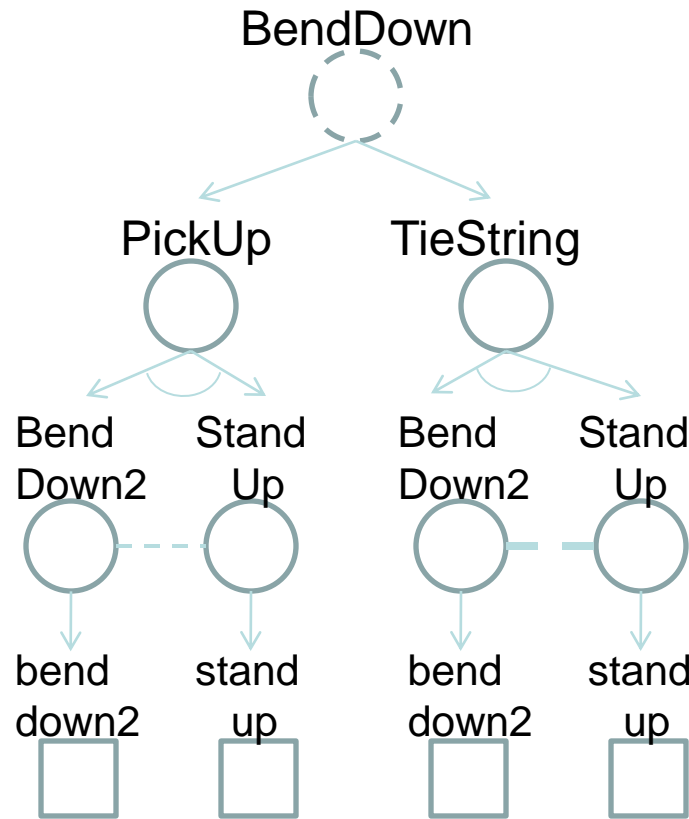
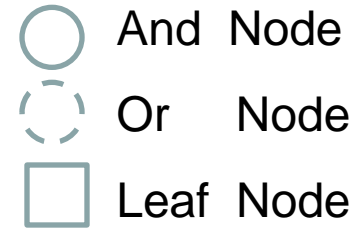
Grammars for Events

- BendDown → Pickup
- BendDown → TieString
- Pickup → benddown2 standup
- TieString → benddown2 standup

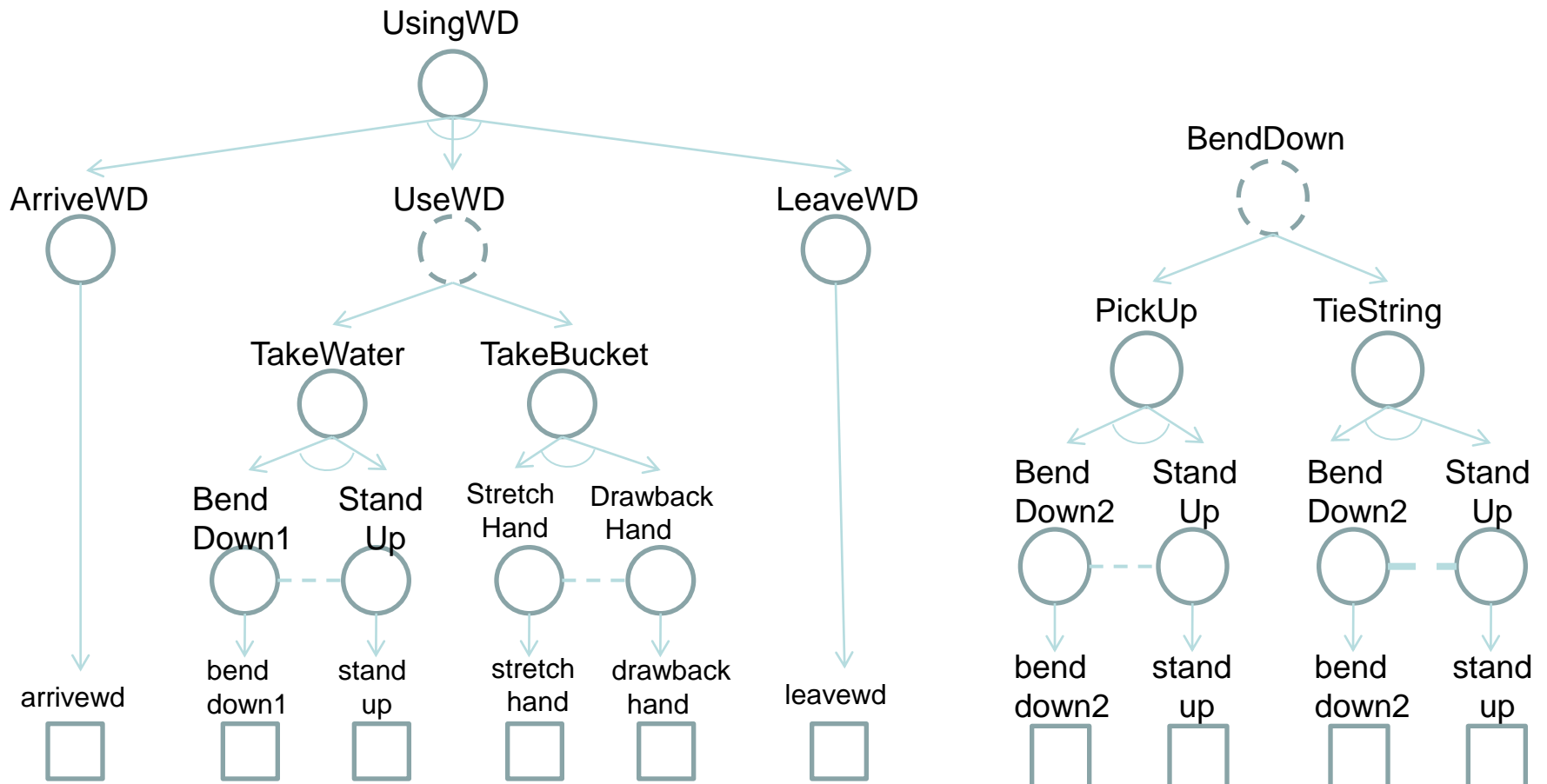
AOGs for Events



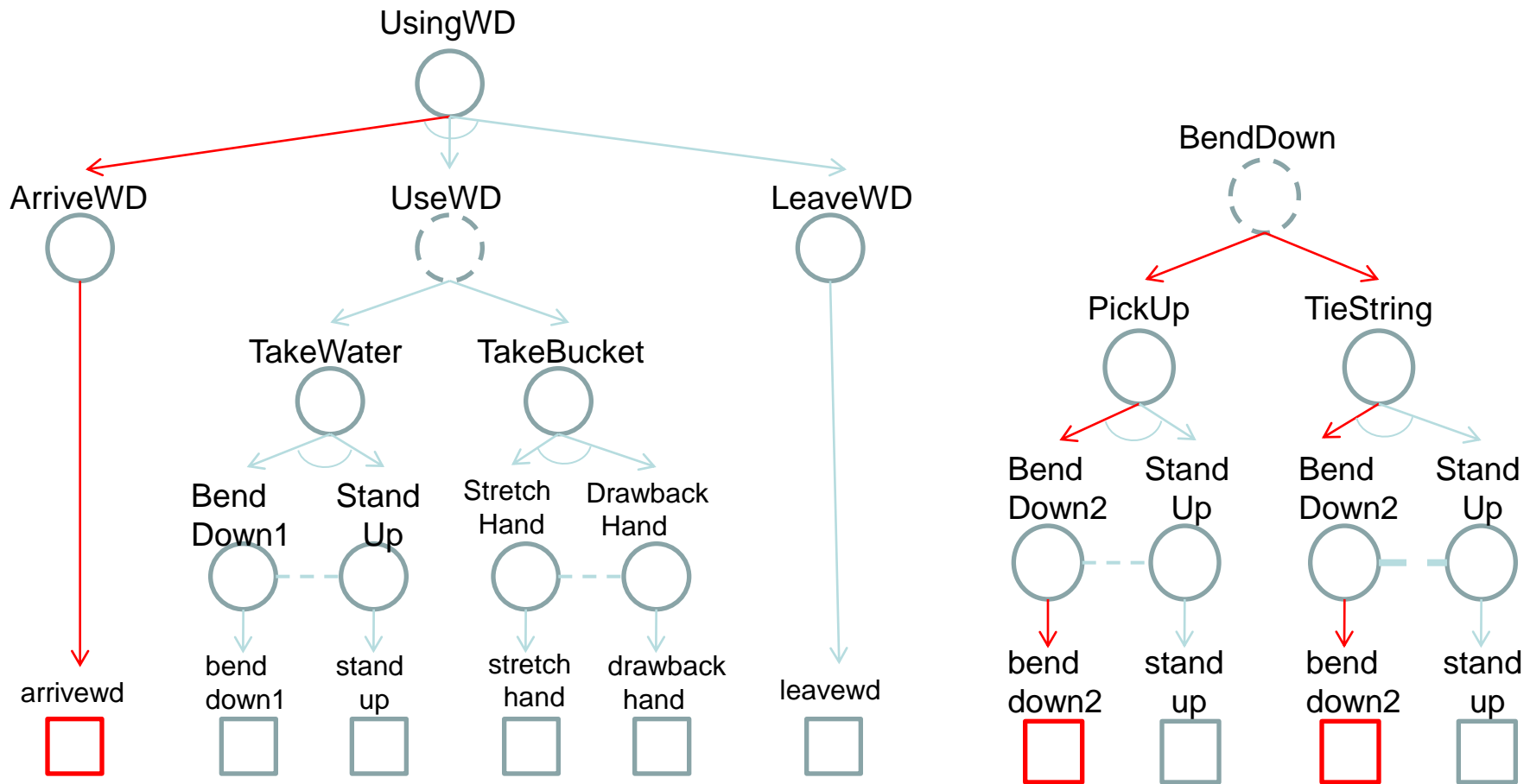
AOGs for Events



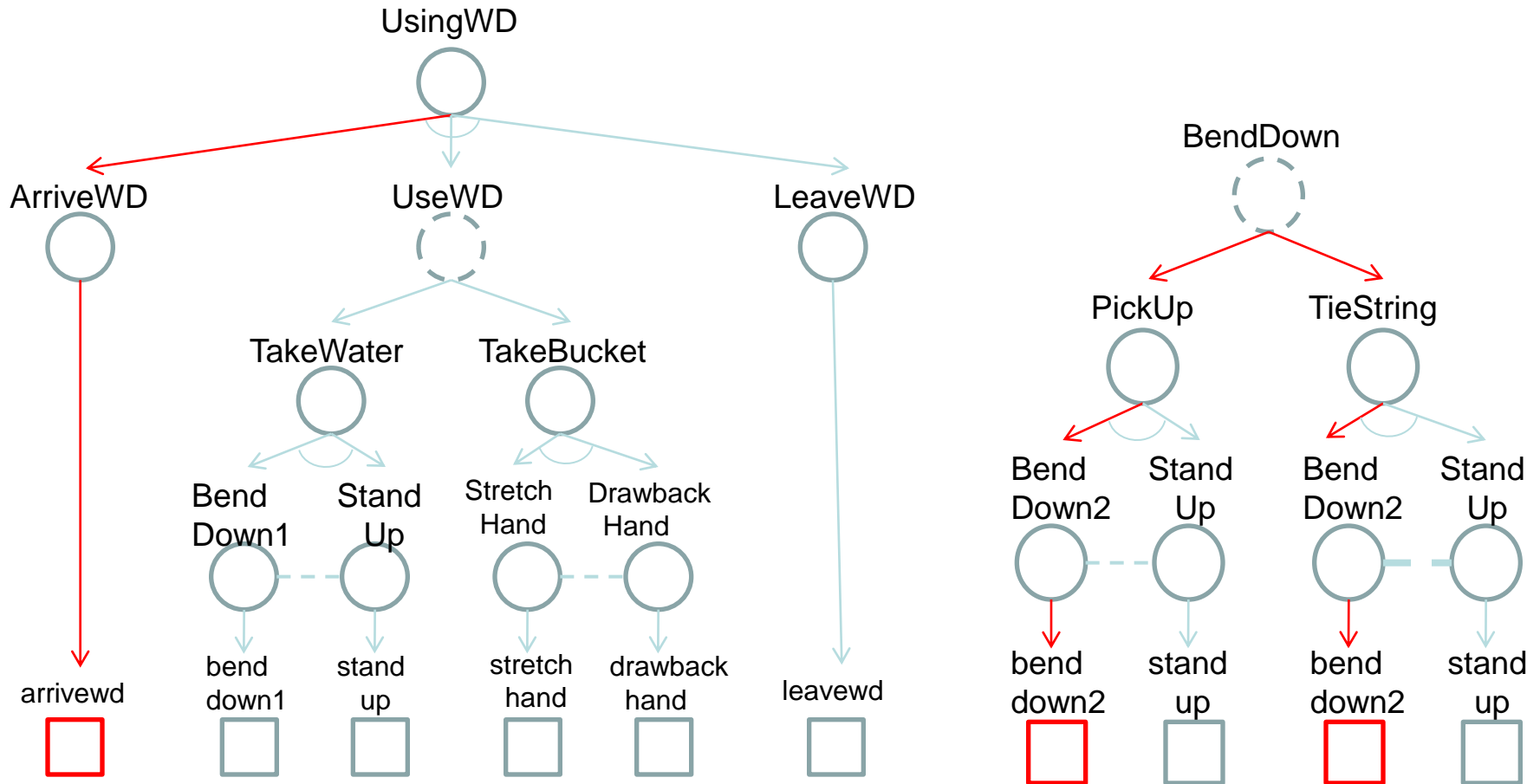
The parsing process



The parsing process

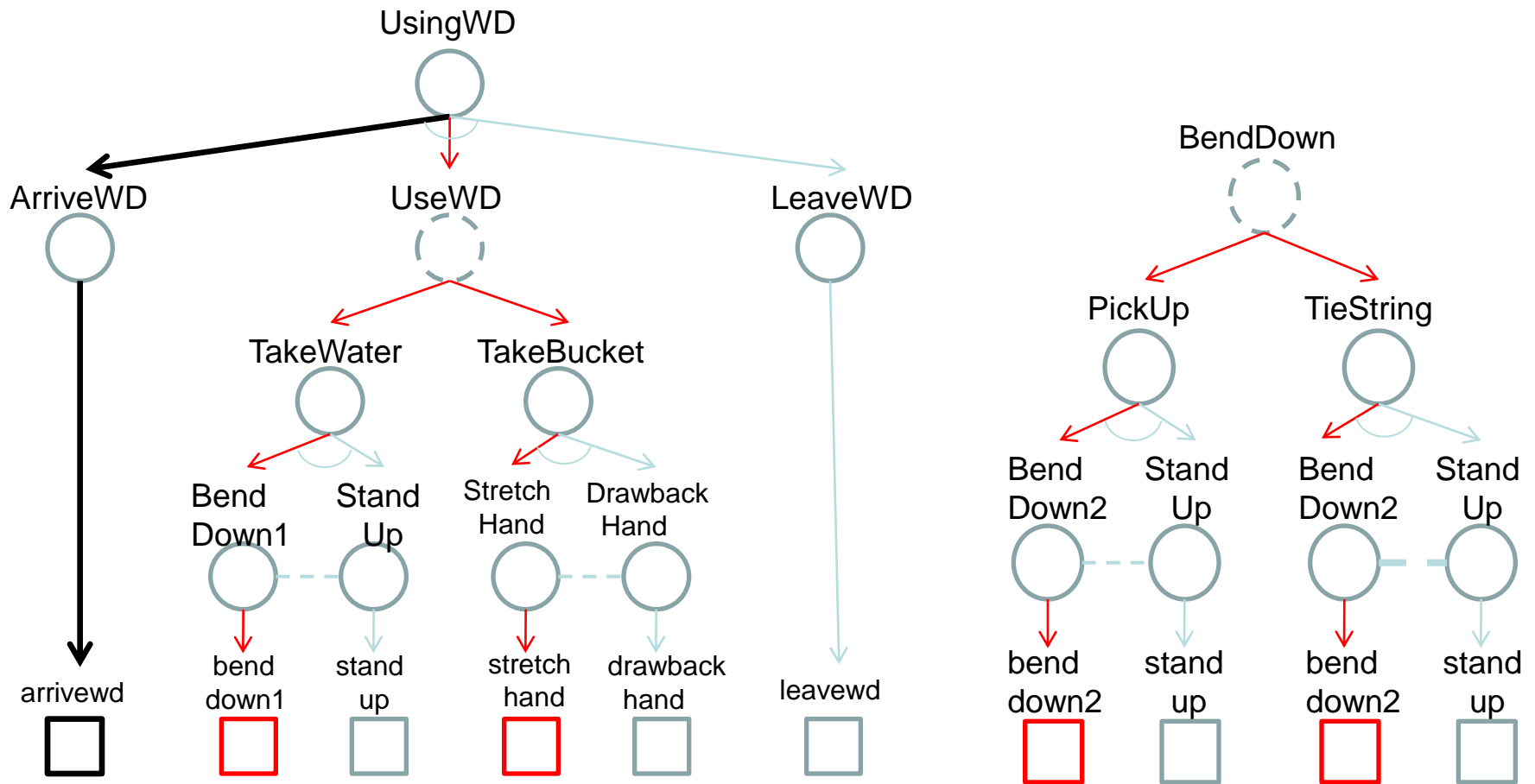


The parsing process



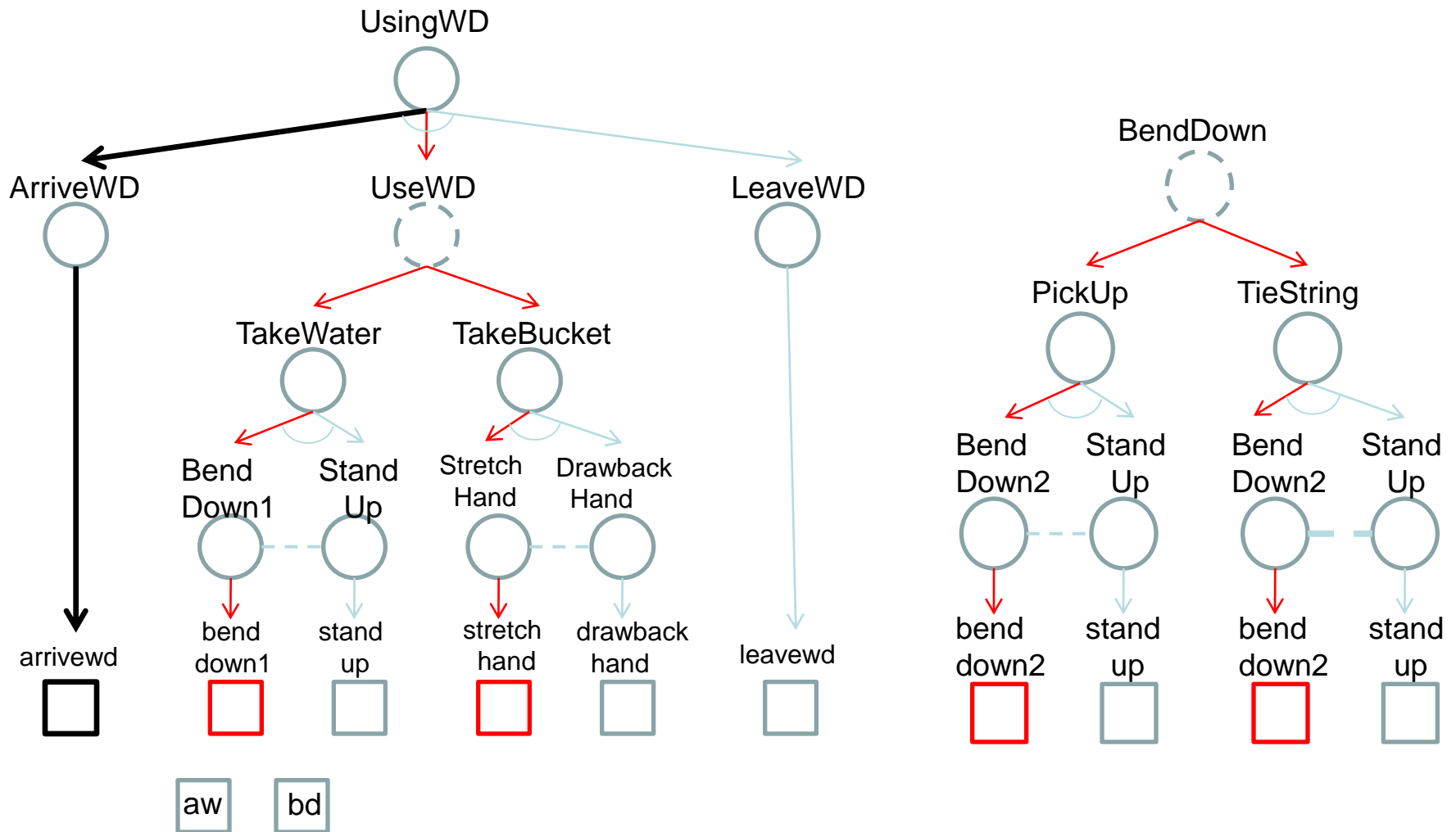
aw

The parsing process

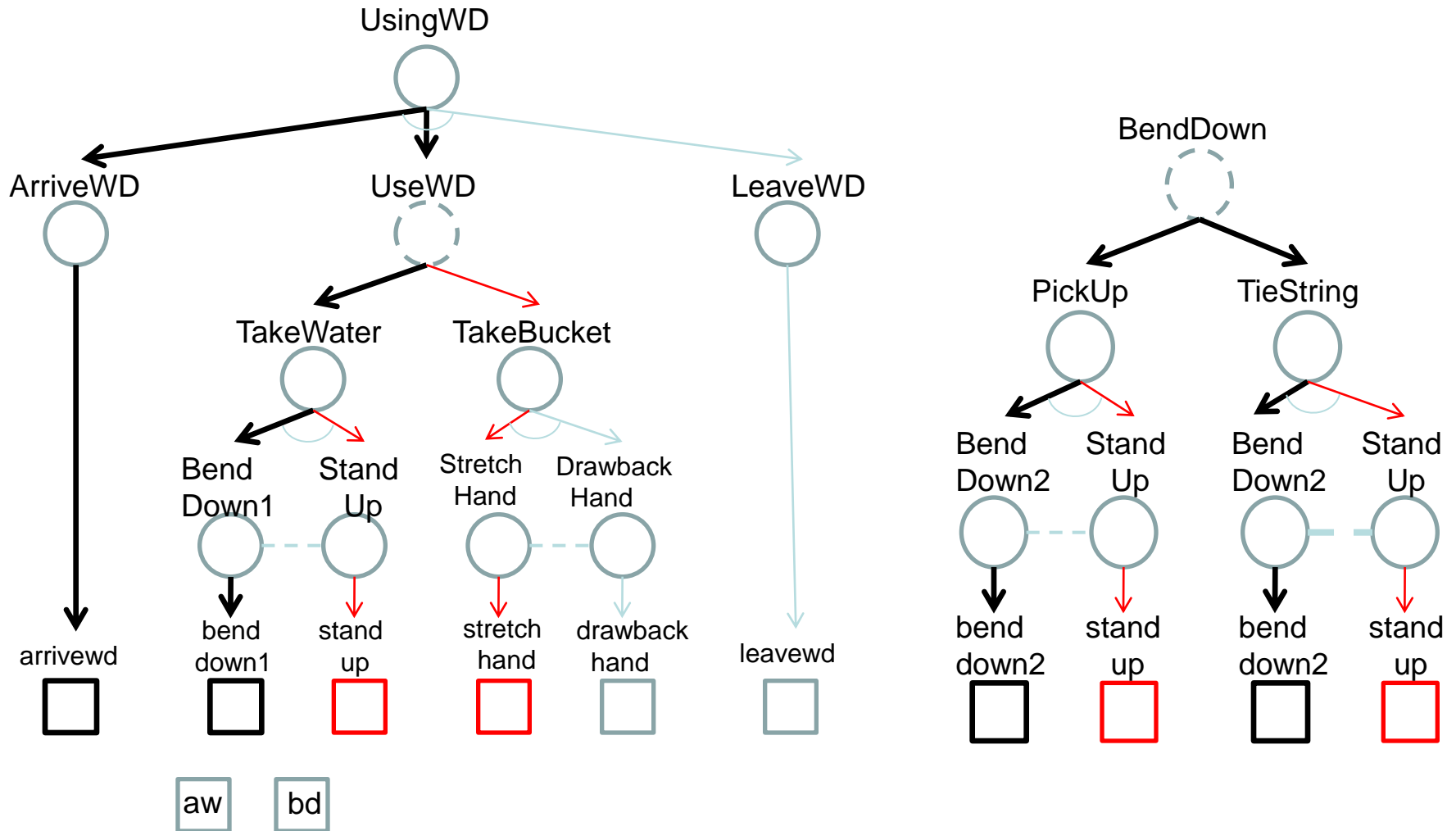


aw

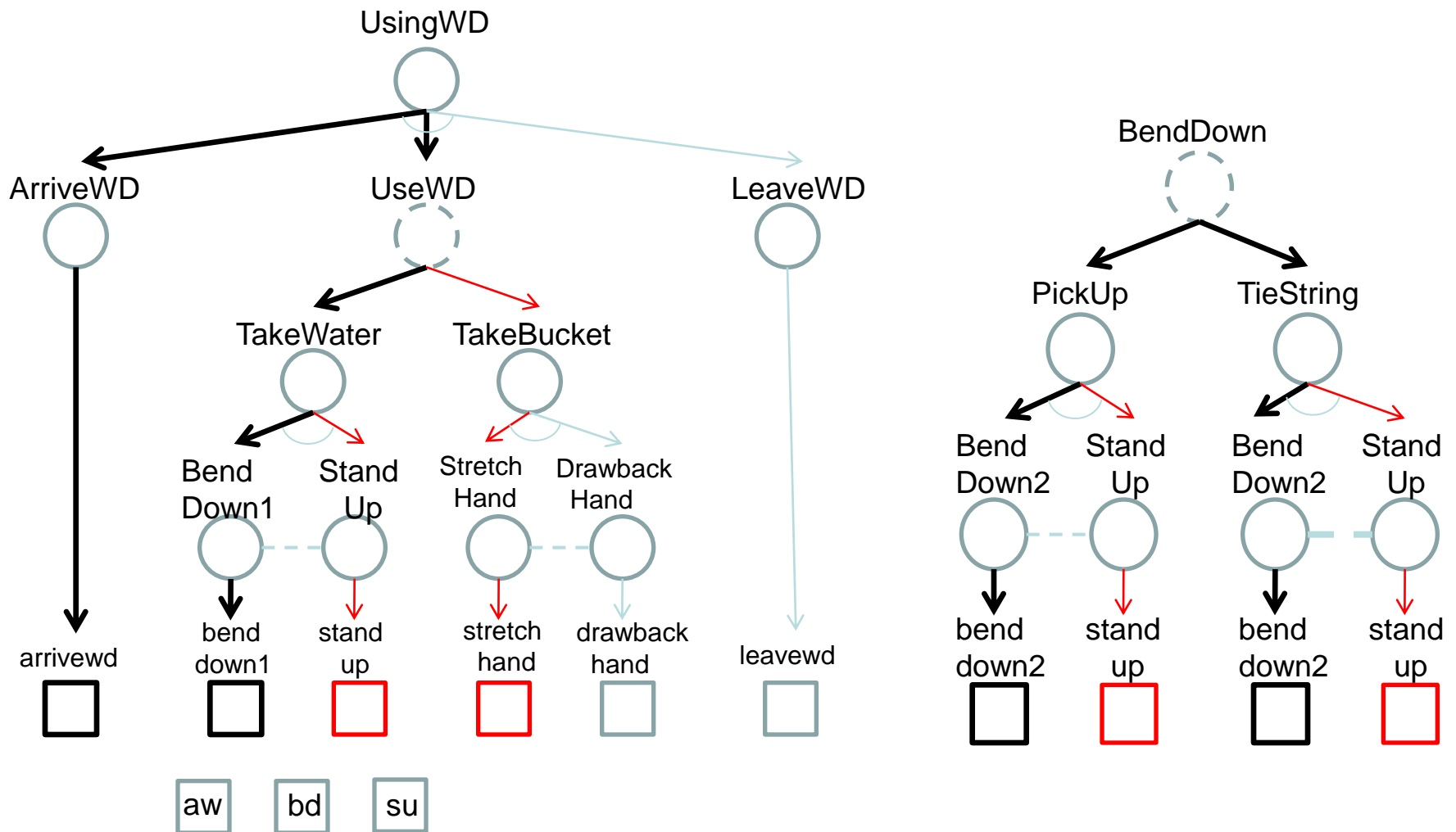
The parsing process



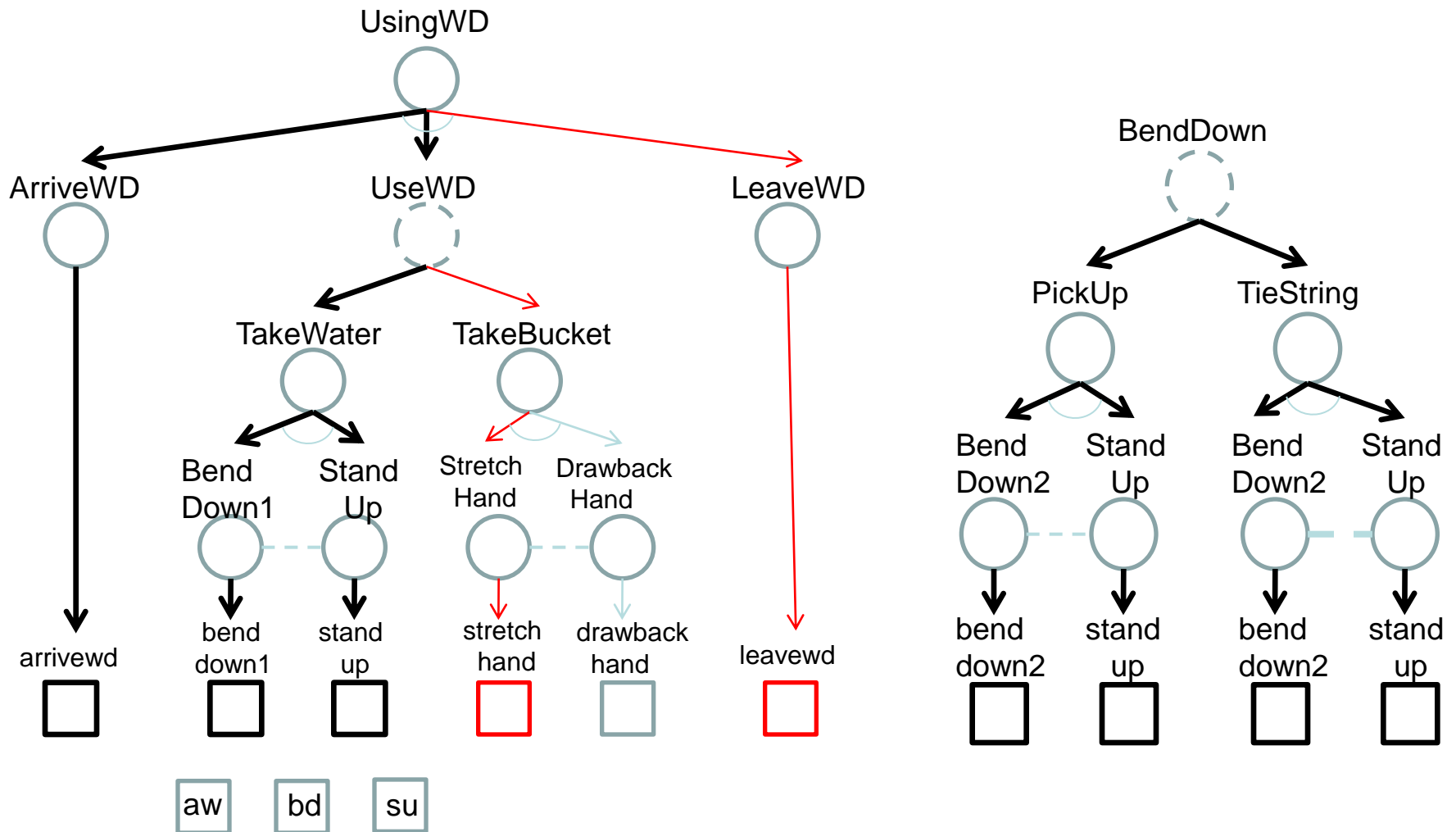
The parsing process



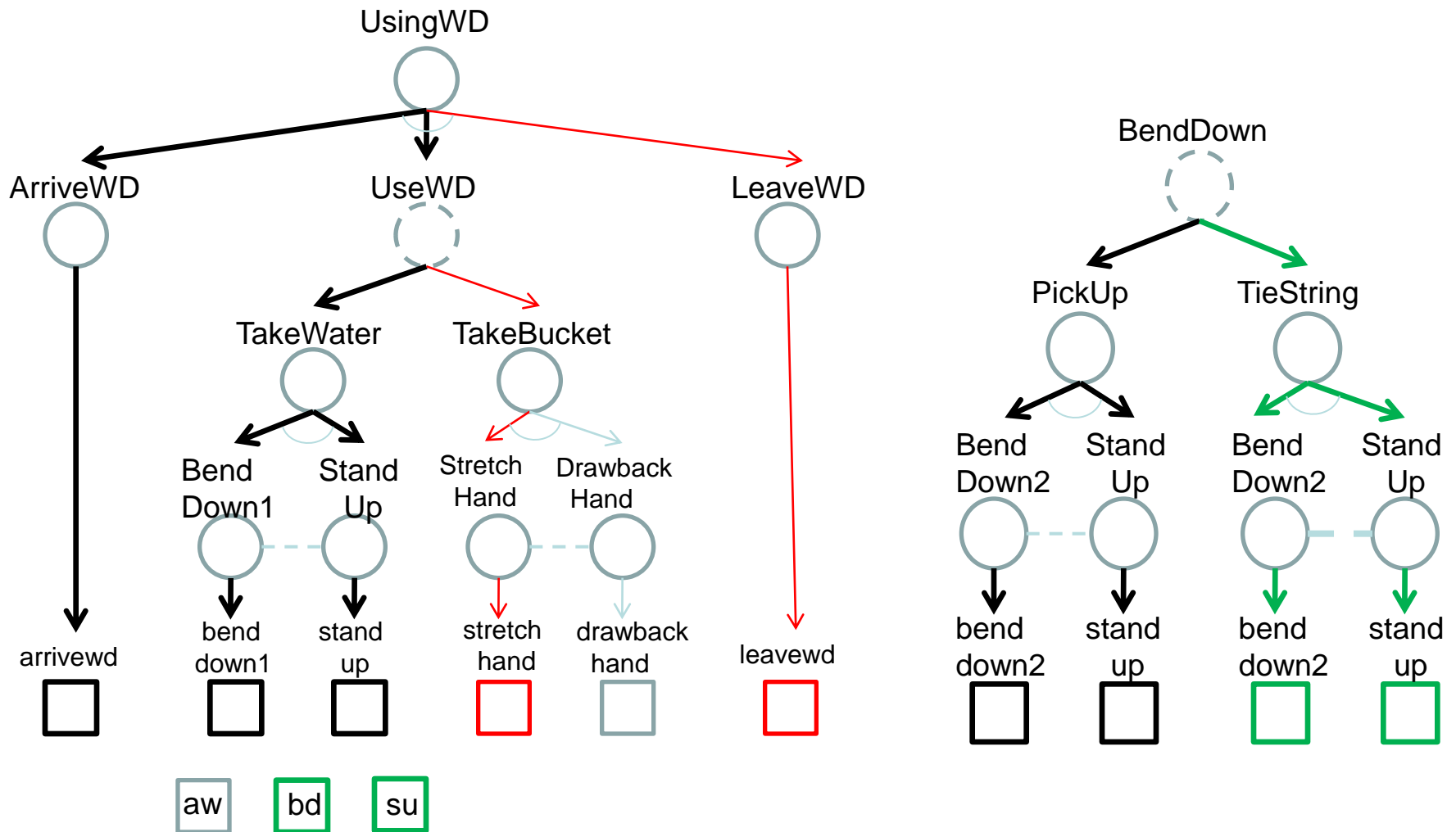
The parsing process



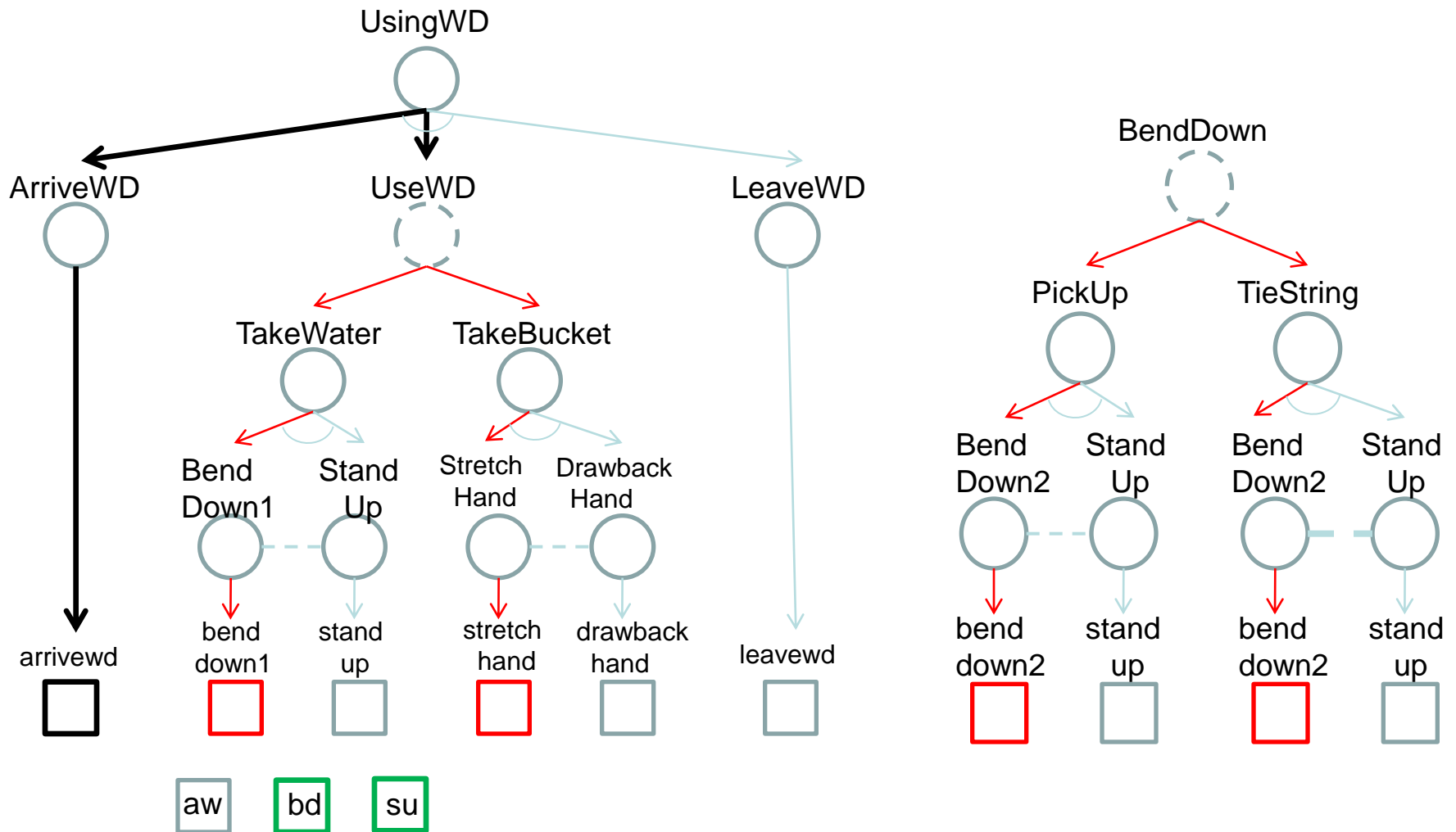
The parsing process



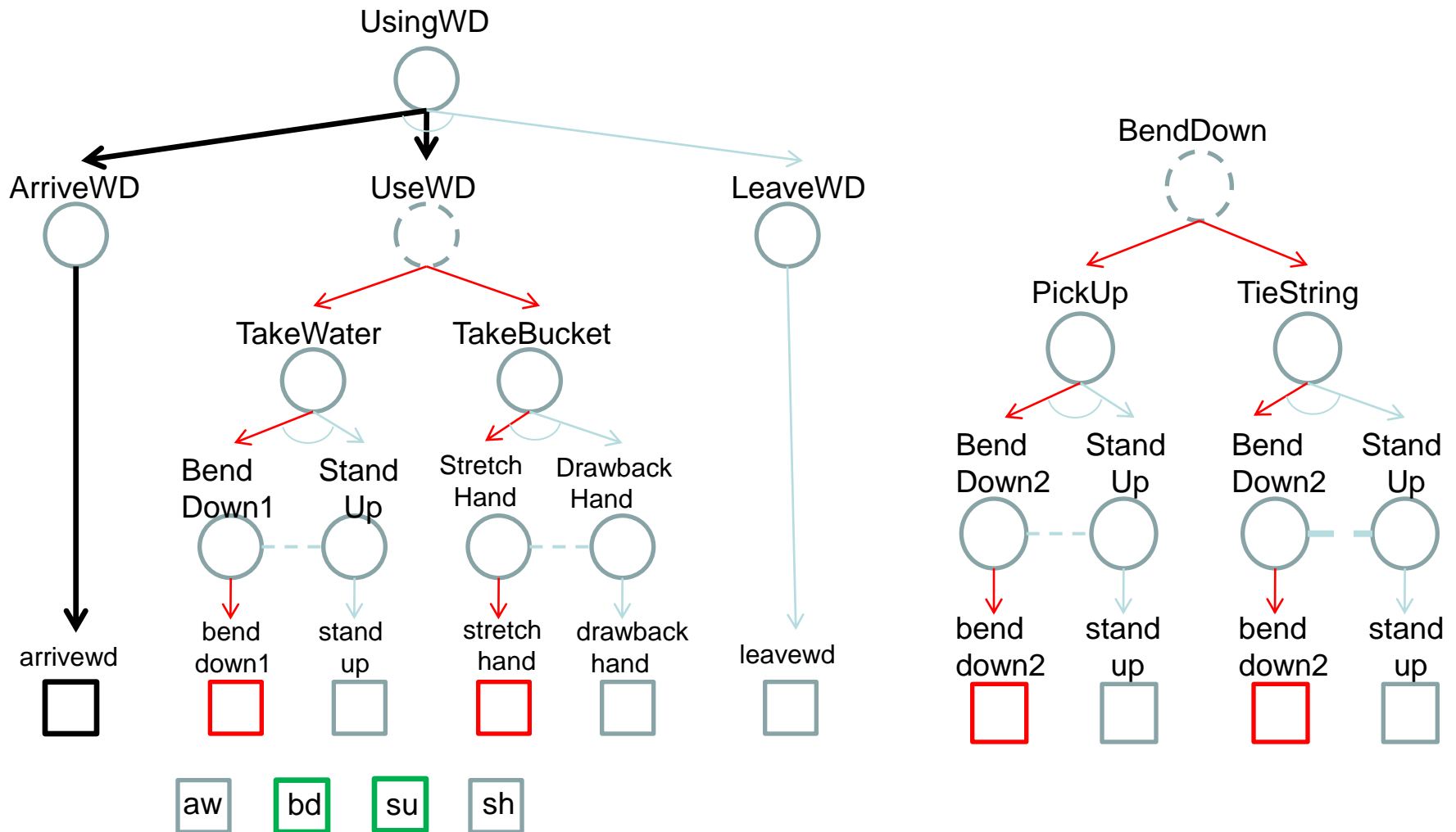
The parsing process



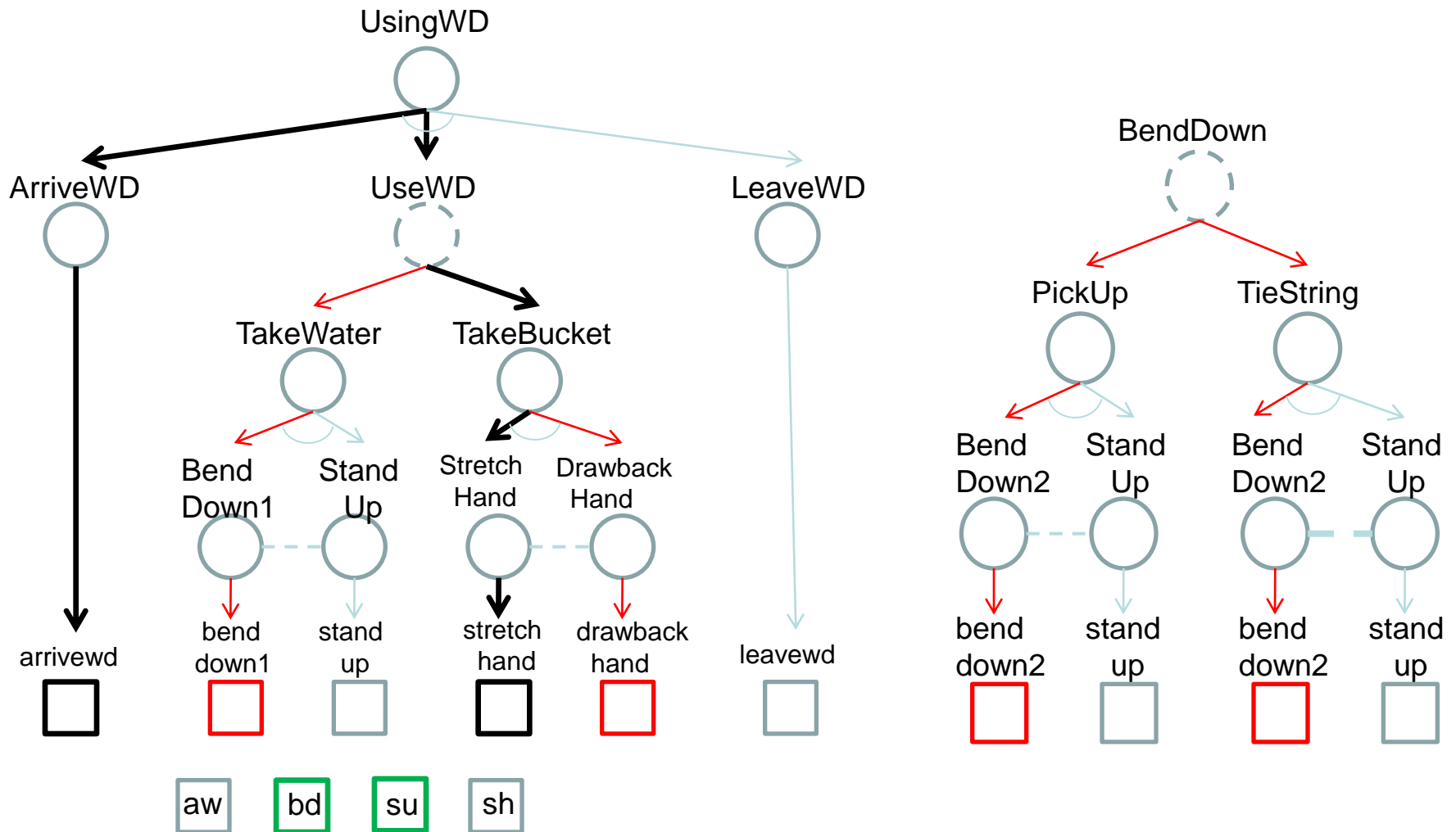
The parsing process



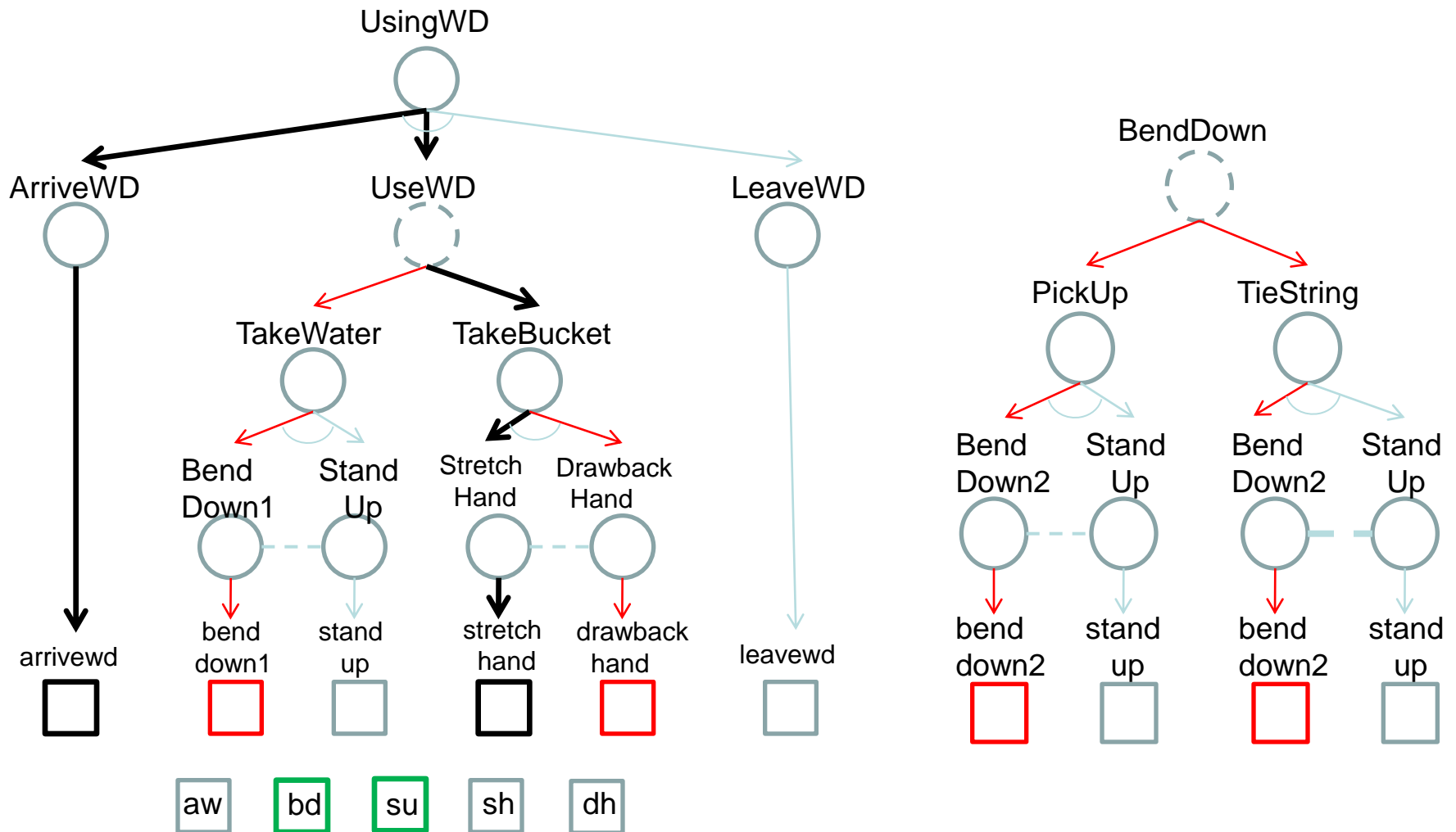
The parsing process



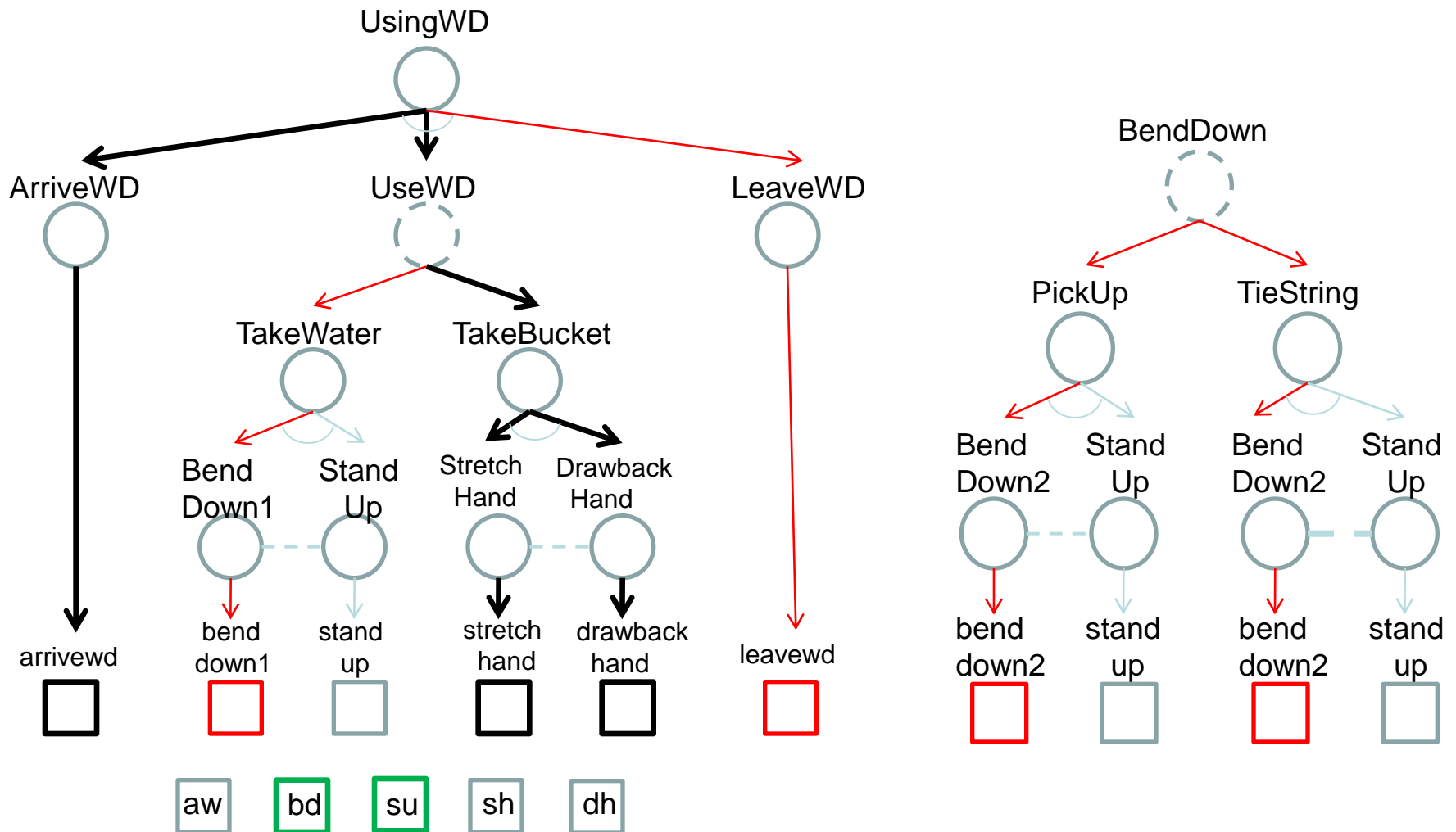
The parsing process



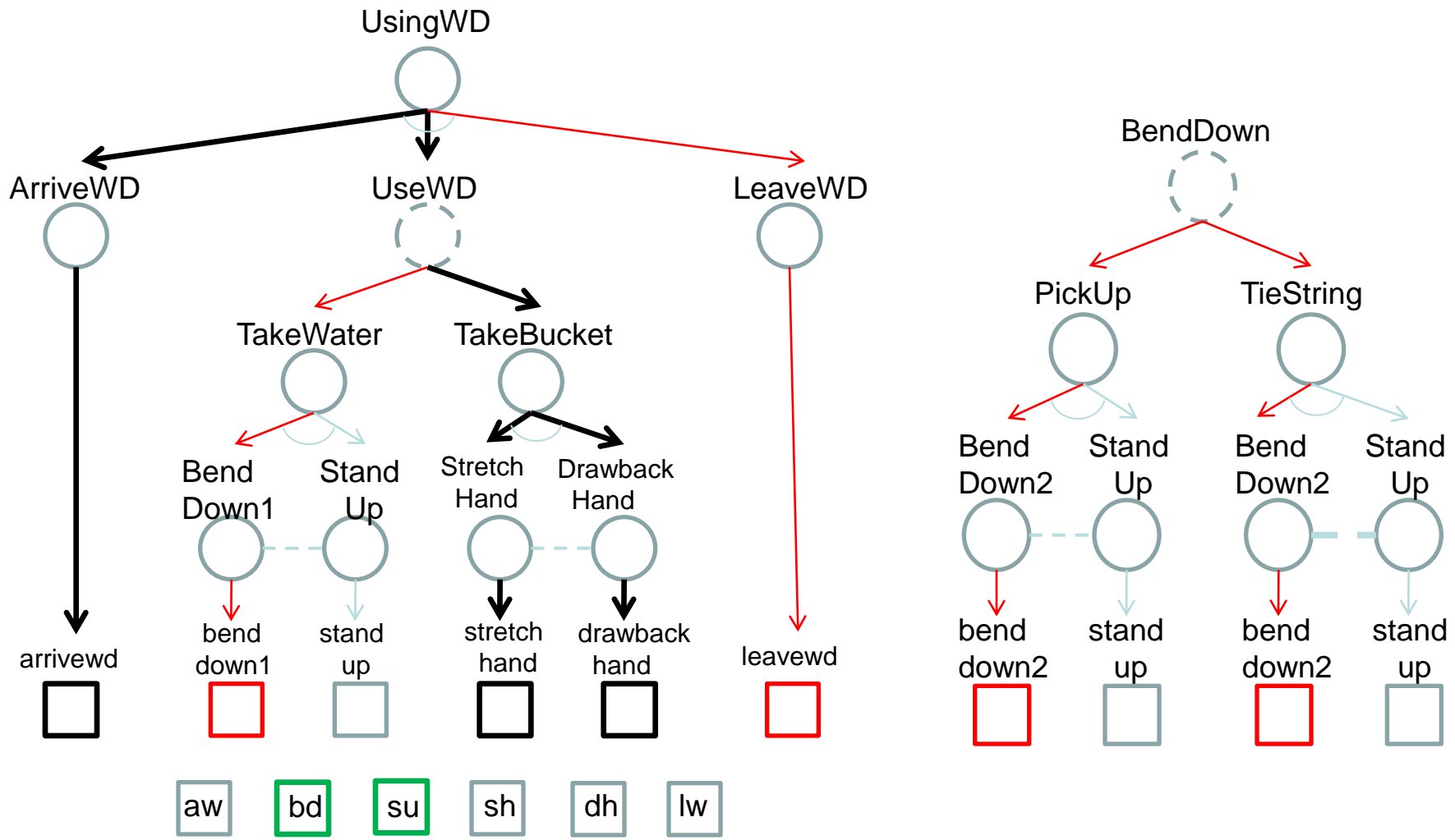
The parsing process



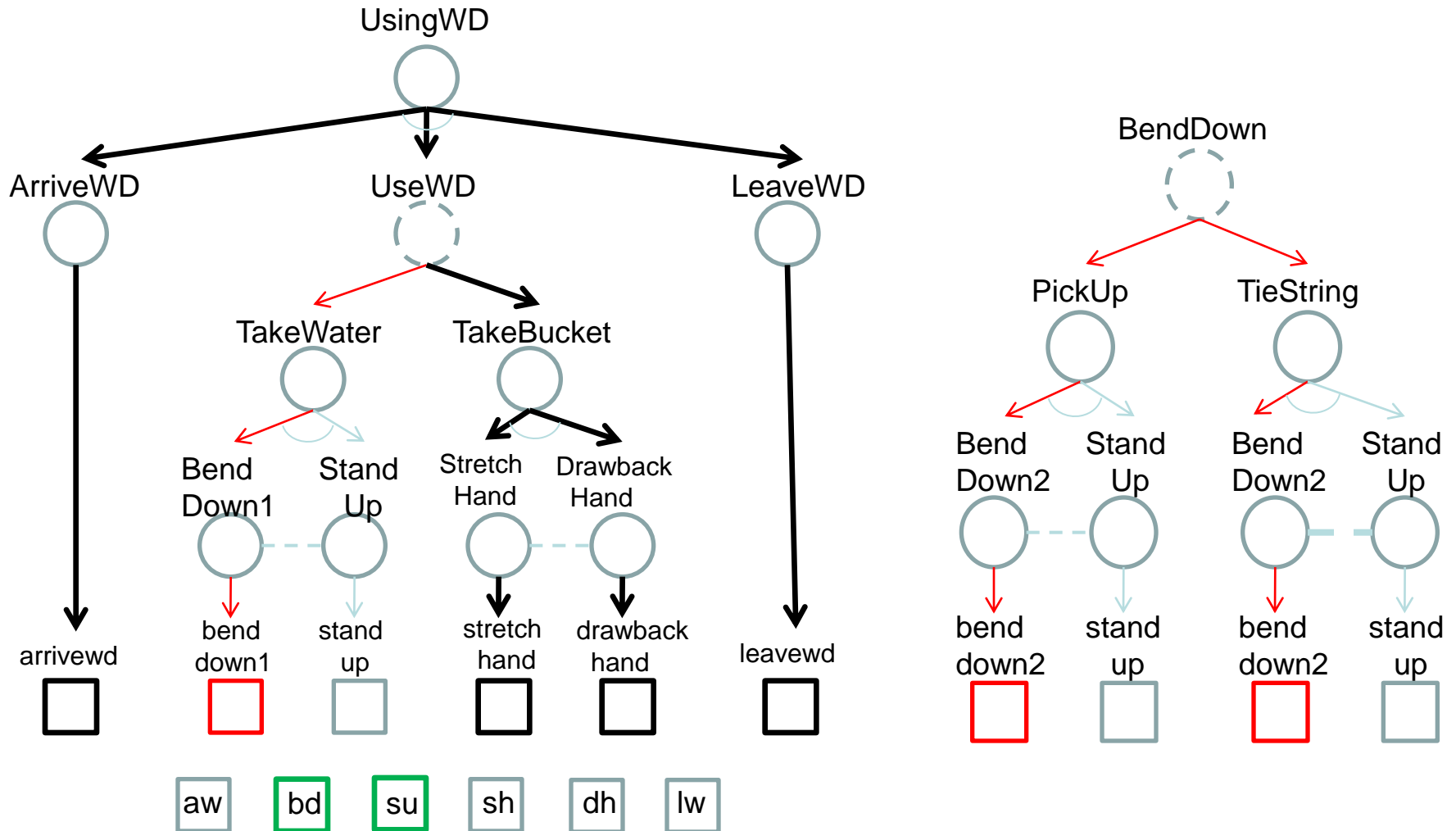
The parsing process



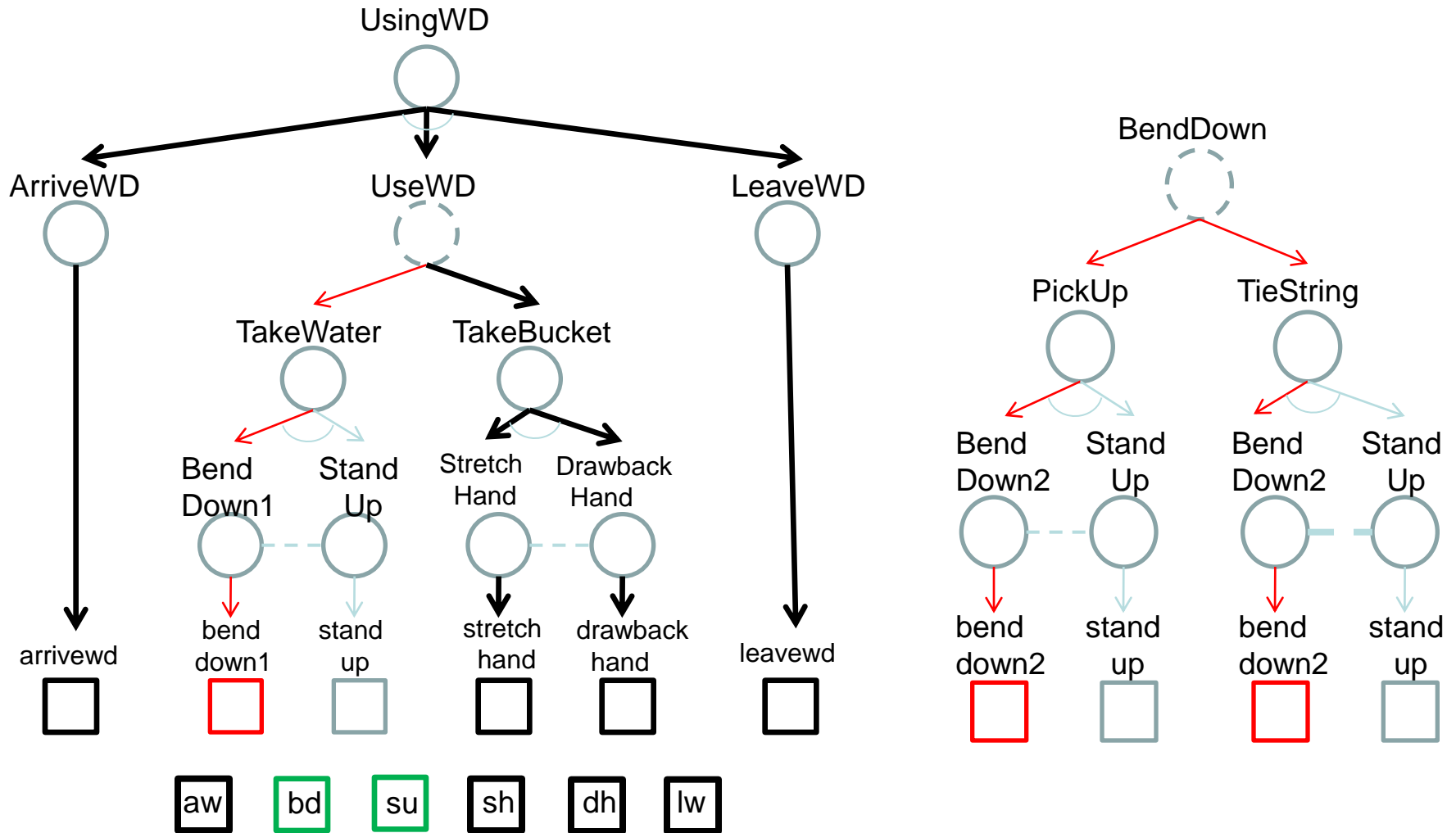
The parsing process



The parsing process



The parsing process



Goal: Recognize events in daily scenes

For example, in an office scene



Work by Mingtao Pei, UCLA

Challenges

1. Events happen over an extended time period

- Variant time-span
- Could be interrupted
- Multiple routes
- Intention and prediction



2. Actions are hard to recognize

- Subtle and similar
- No salient motion/pose at most of the time
- Contextual OBJECTS are key!!



Use
laptop

Read
book

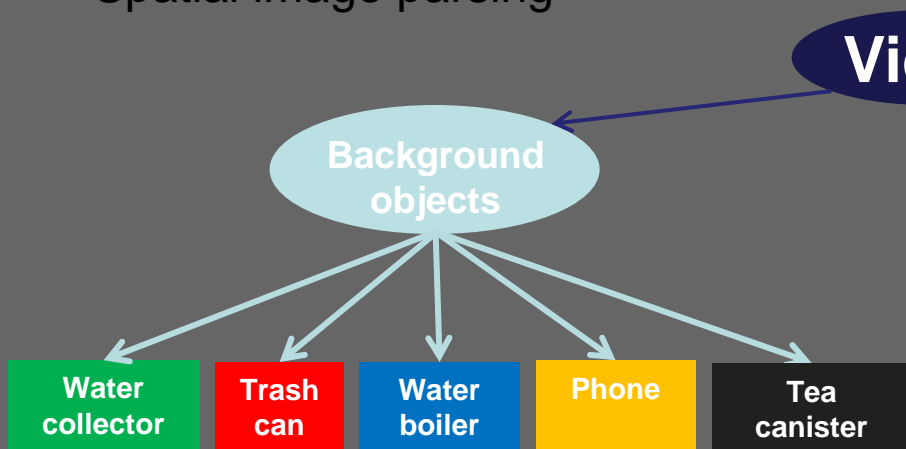
Dump
water

Use
microwave

Overview of our approach

Spatial image parsing

Temporal event parsing

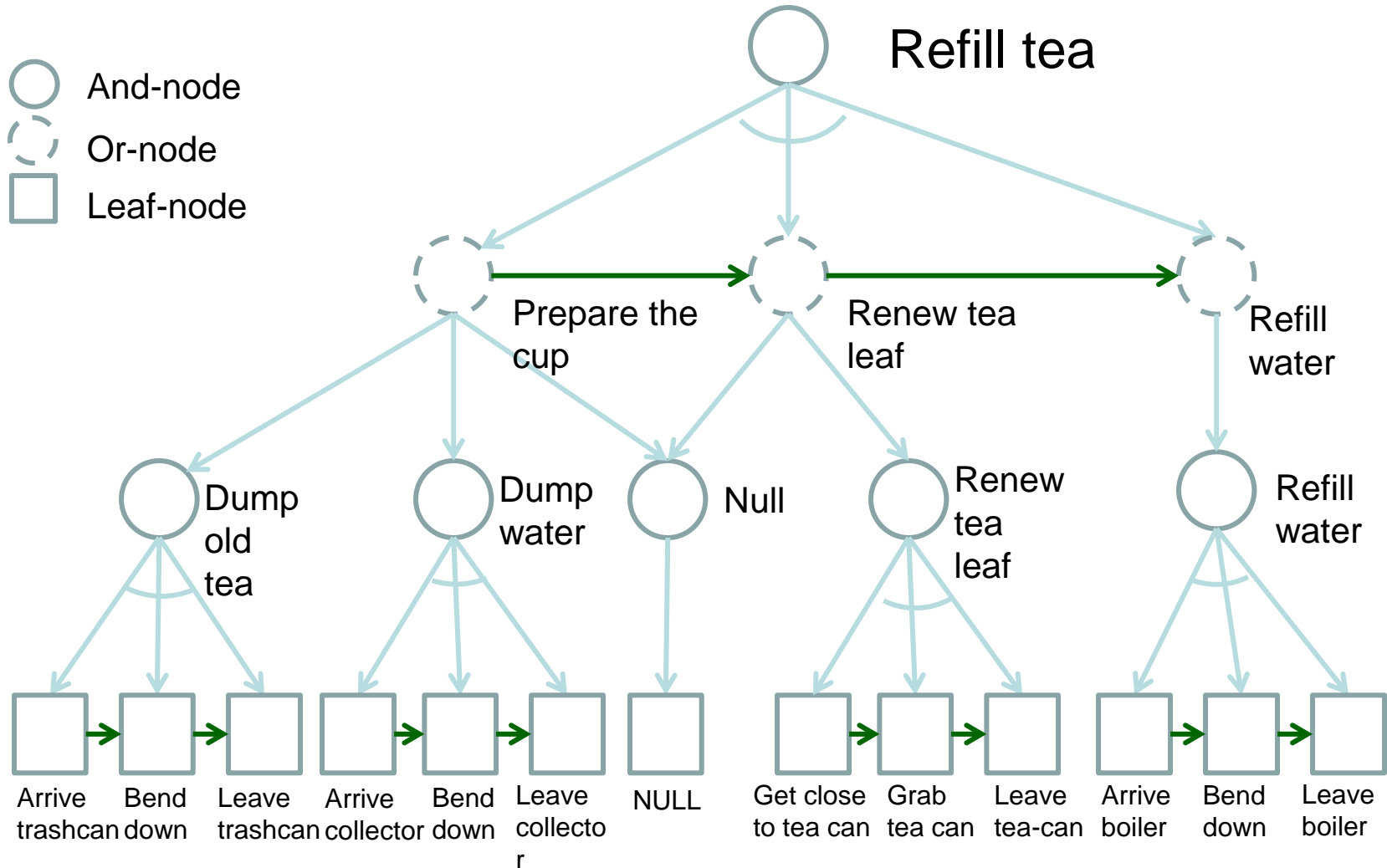


Scene parsing

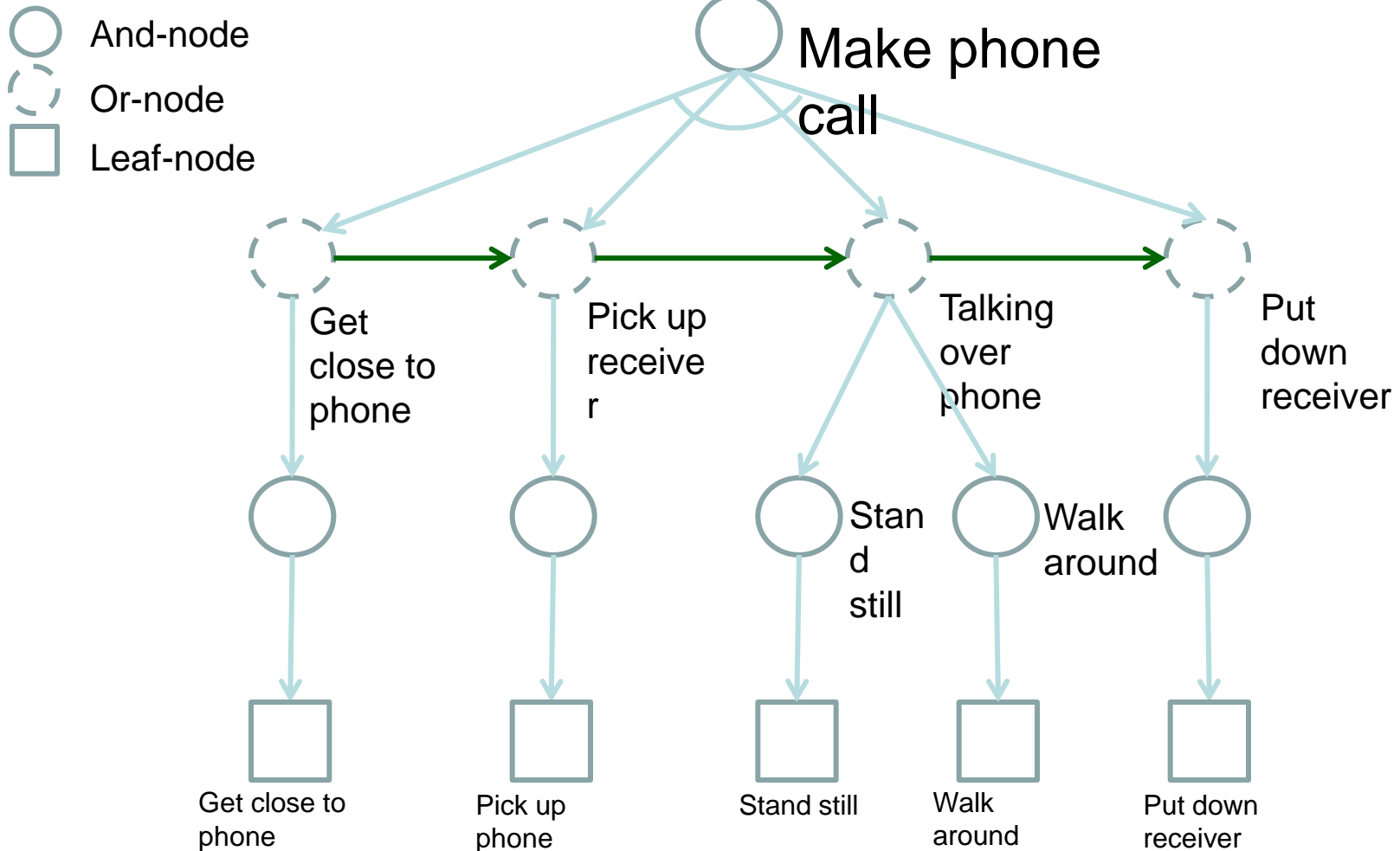
	chair
	desk top
	computer
	paper
	phone
	cup
	tea box
	microwave
	water dispenser
	trash can
	basin
	whiteboard
	floor1
	floor2
	floor
	wall



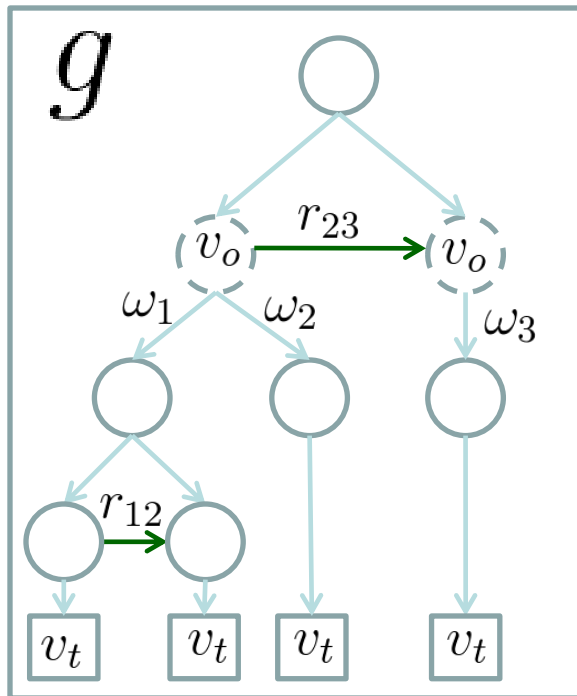
Event as temporal And-Or-Graph



Event as temporal And-Or-Graph



Formulation



$$p(g) = \frac{1}{Z} \exp\{\text{score}(g)\}$$

Grammar

r

Data term

$$\text{score}(g) = \sum_{v_t \in T(g)} \lambda_{v_t} \alpha(v_t) +$$

Or node
Frequenc

$$\sum_{v \in V_o(g)} \lambda_v \omega(v) +$$

Temporal
Relations

$$\sum_{(i,j) \in E(g)} \lambda_{ij} r_{ij}(v_i, v_j)$$

$$\alpha(v_t) = \sum_{i \in \mathcal{F}} \beta_i, h_i(v_t) - \text{dist}(P_{\text{person}}, P_{\text{obj}})$$

Combing action and contextual object



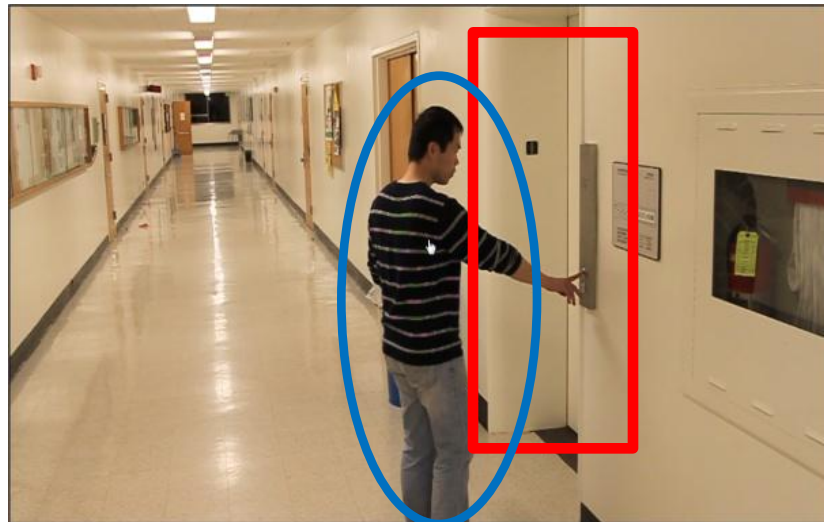
Bend down
action

+

=

Close to
Trash Can

Event1:
Drop waste



Reach out
action

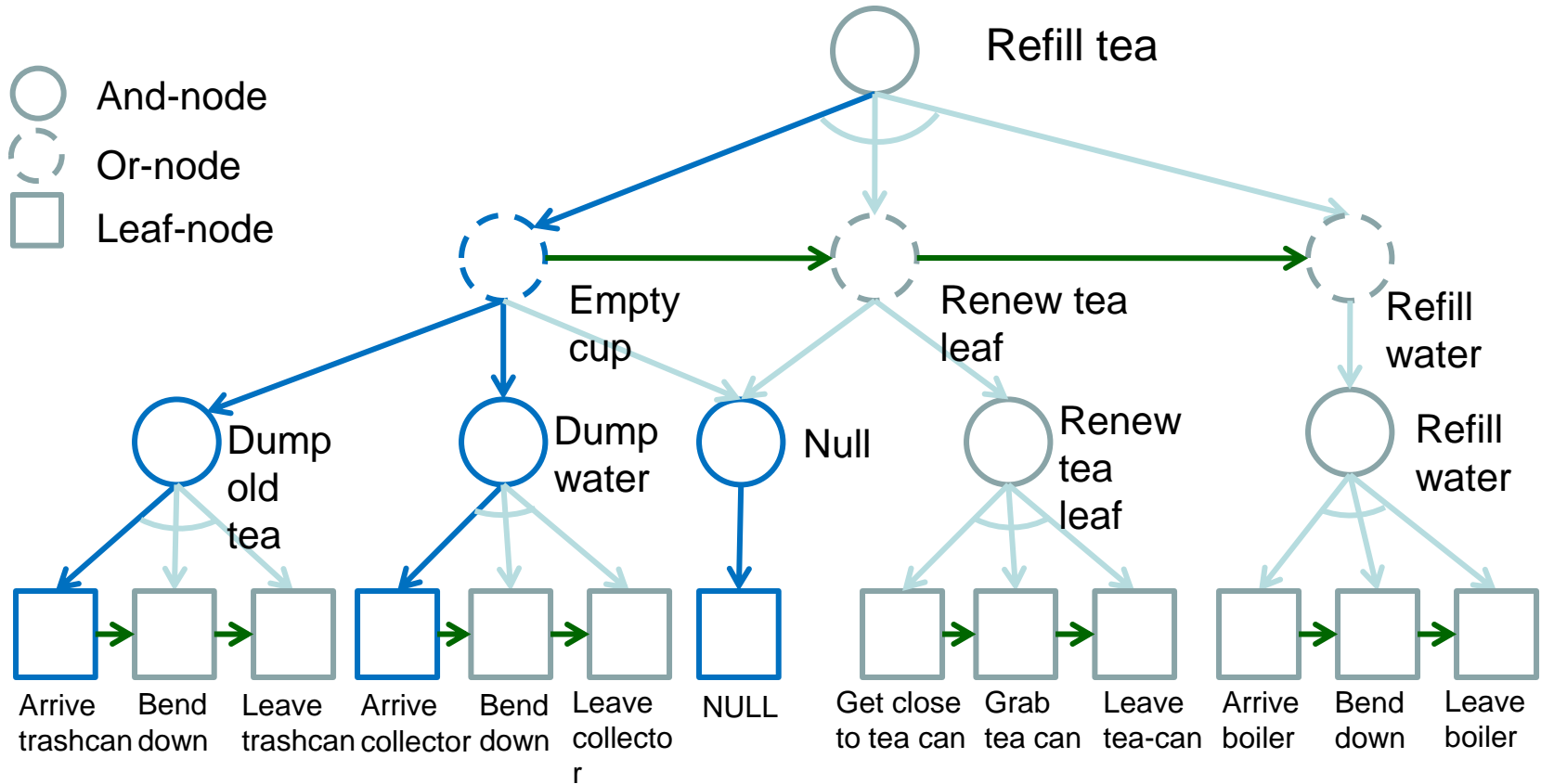
+

=

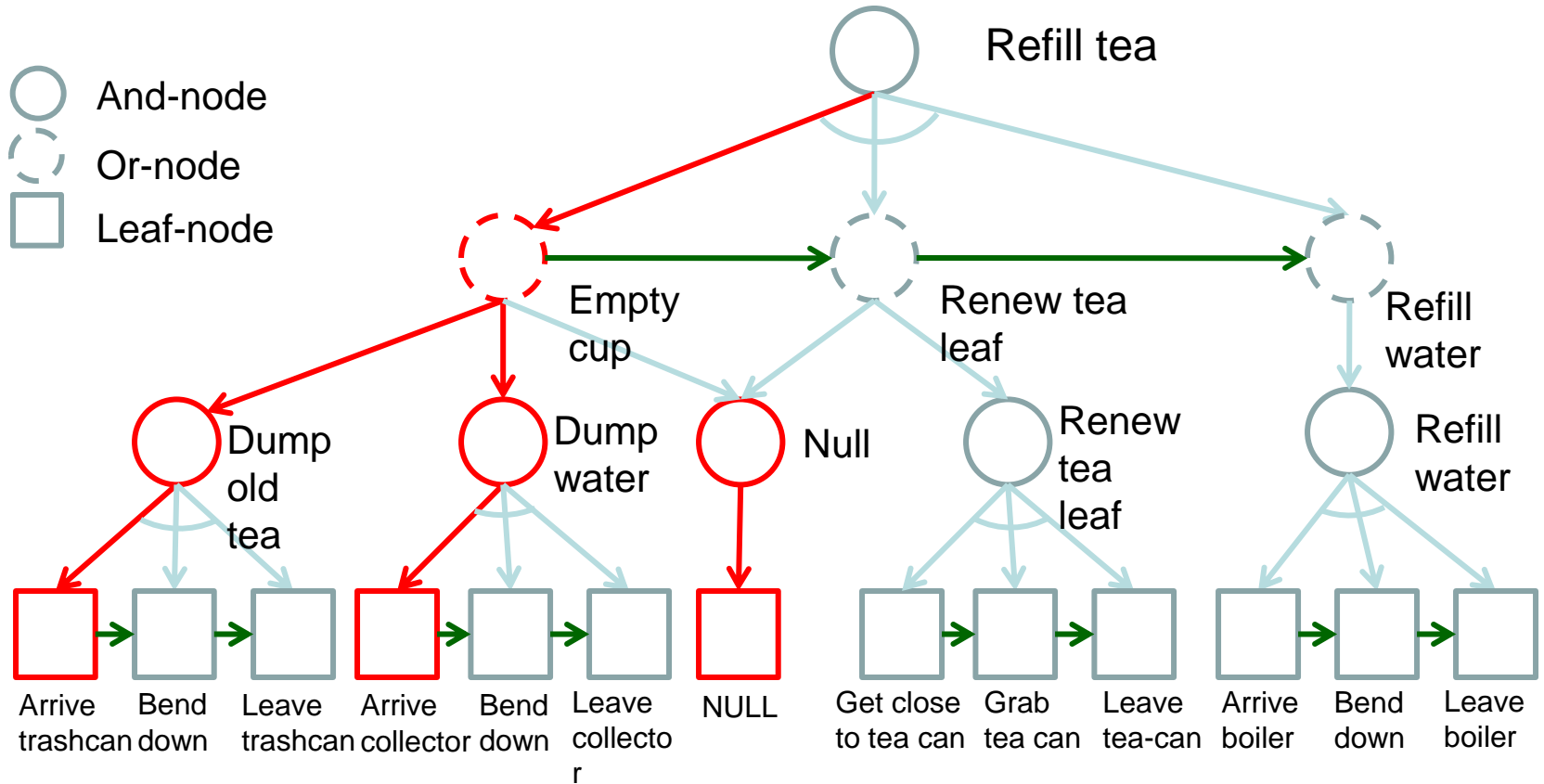
Close to
Elevator

Event2:
Use elevator

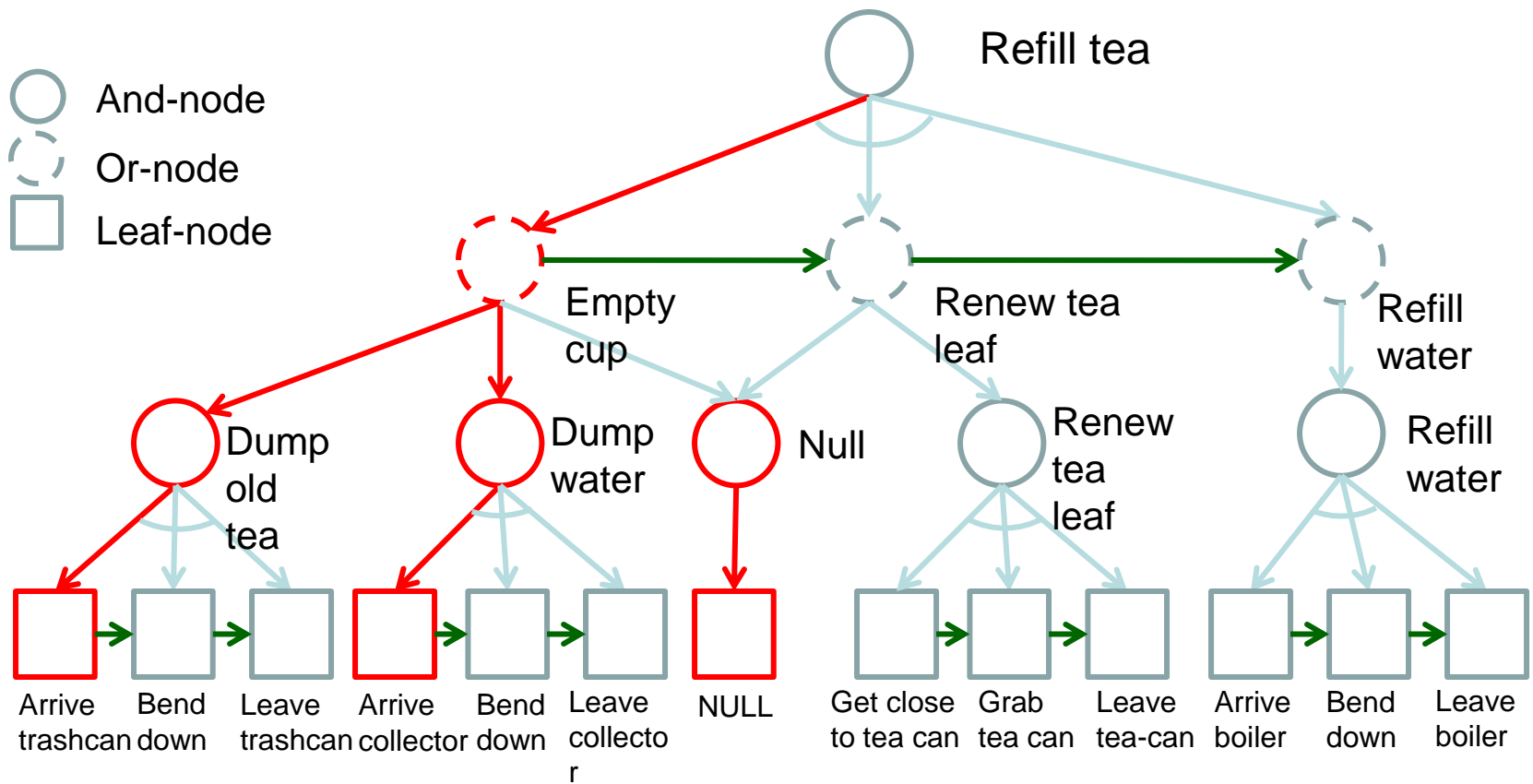
Parsing process (Earley Parser [Earley 1970])



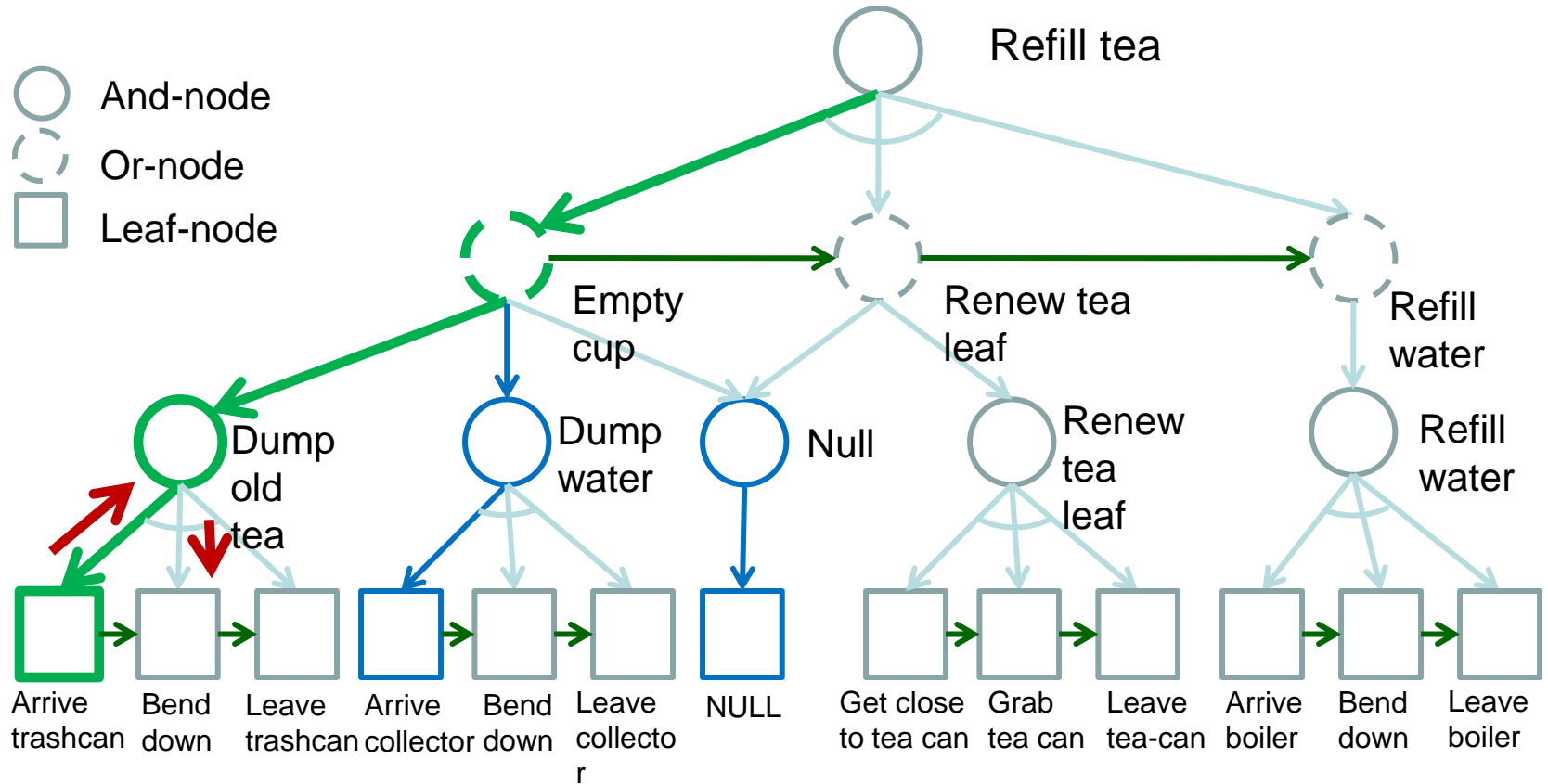
Parsing process (Earley Parser [Earley 1970])



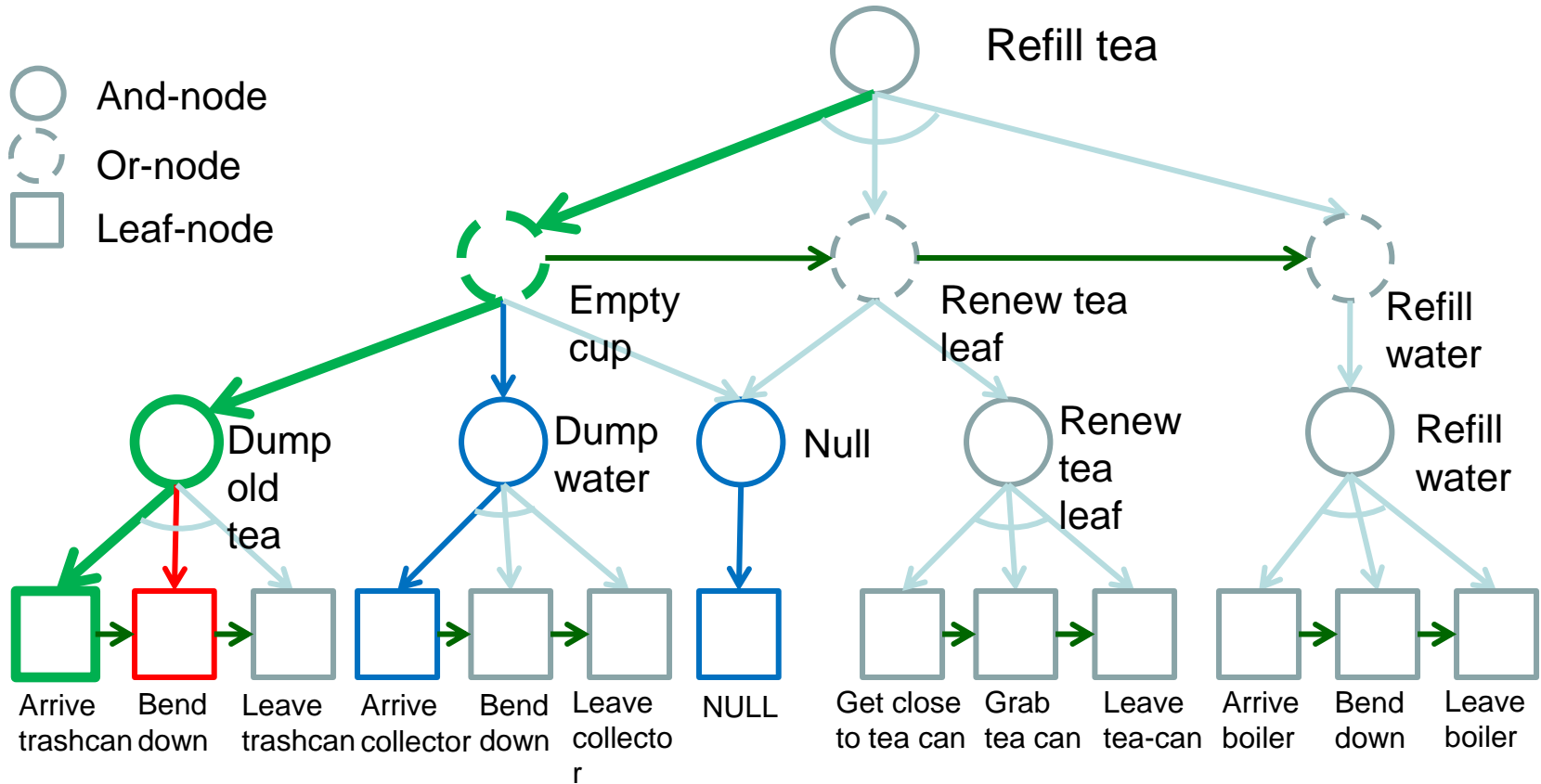
Parsing process (Earley Parser [Earley 1970])



Parsing process (Earley Parser [Earley 1970])

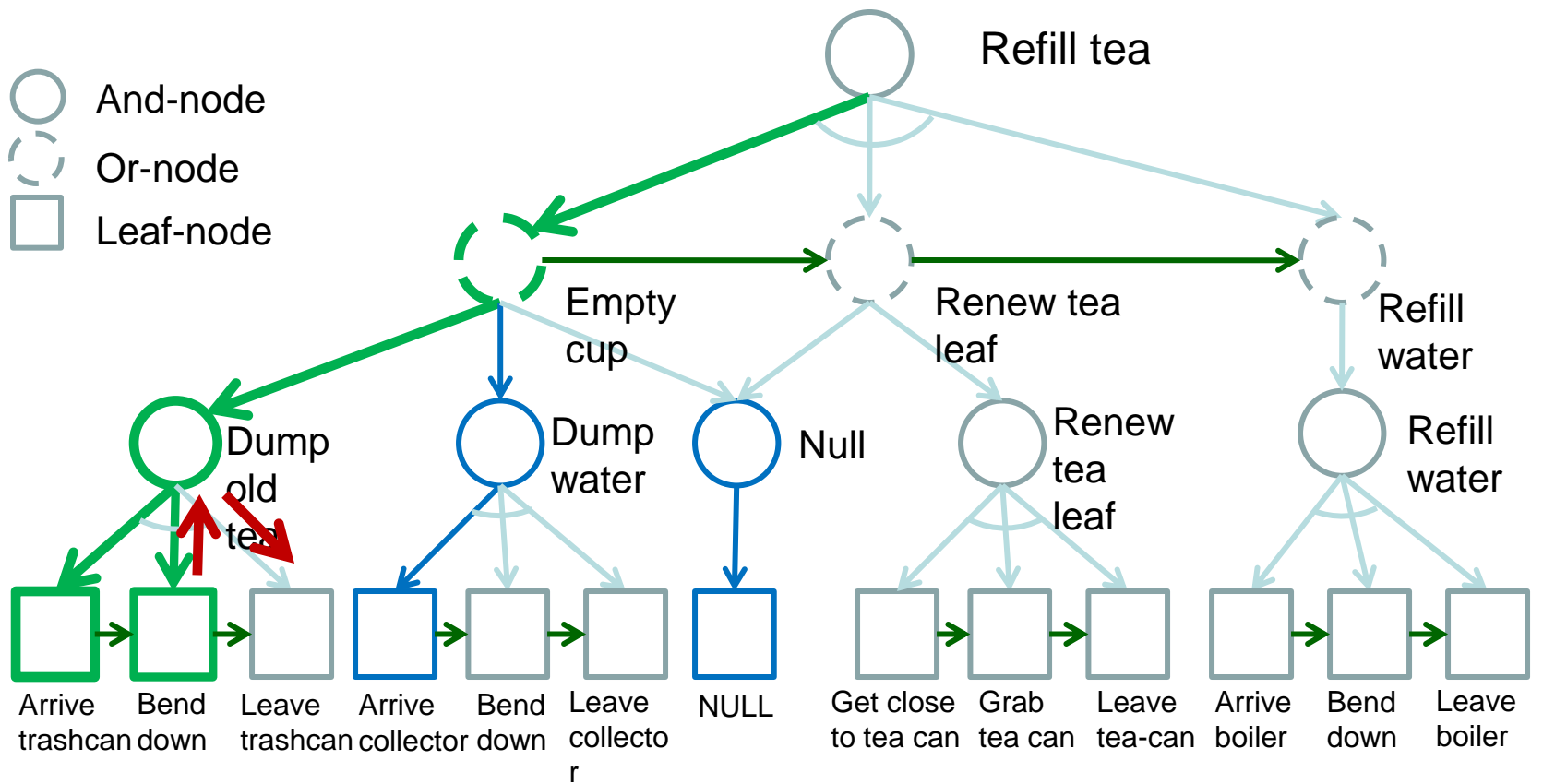


Parsing process (Earley Parser [Earley 1970])

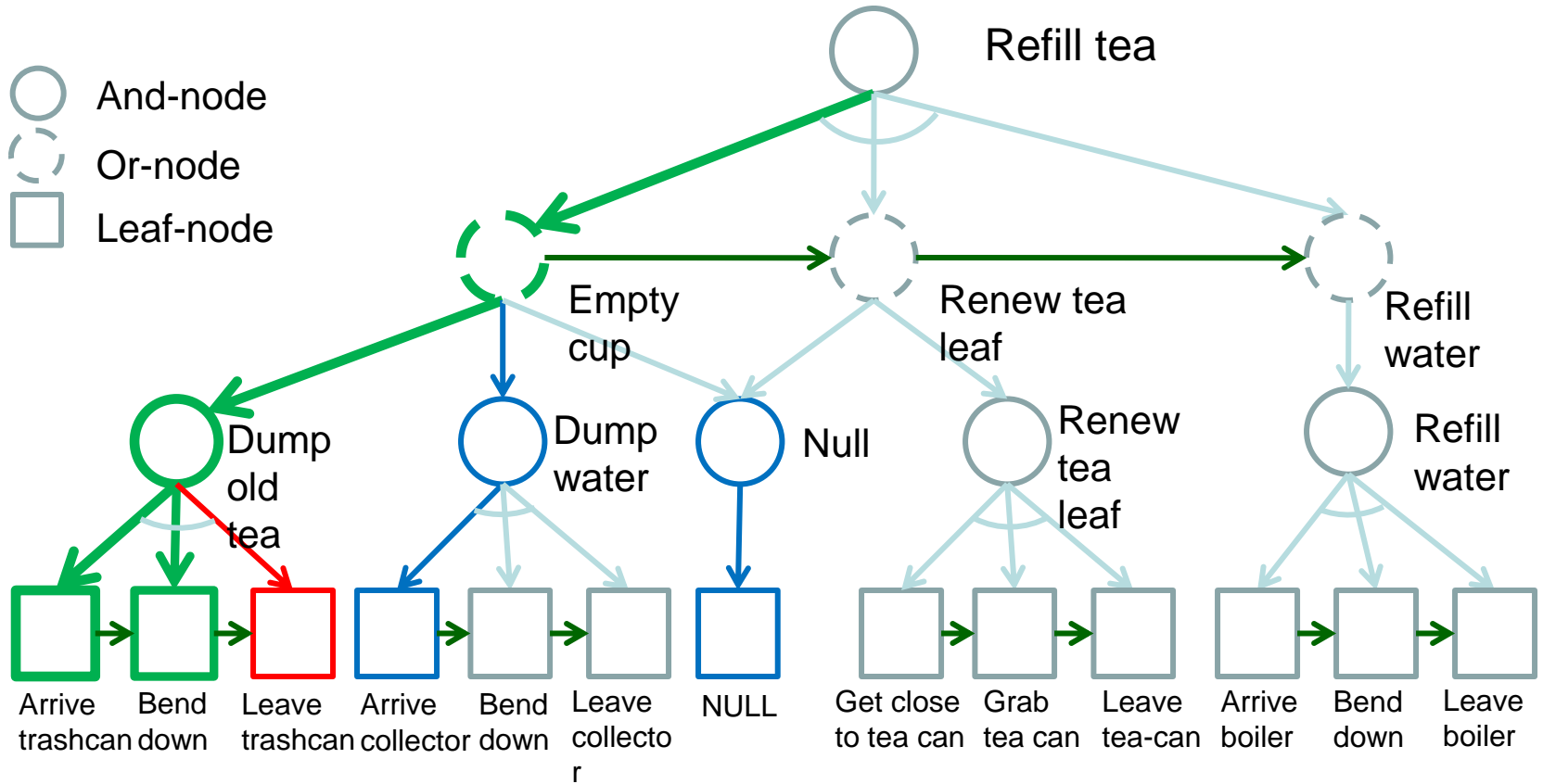


Parsing process (Earley Parser [Earley 1970])

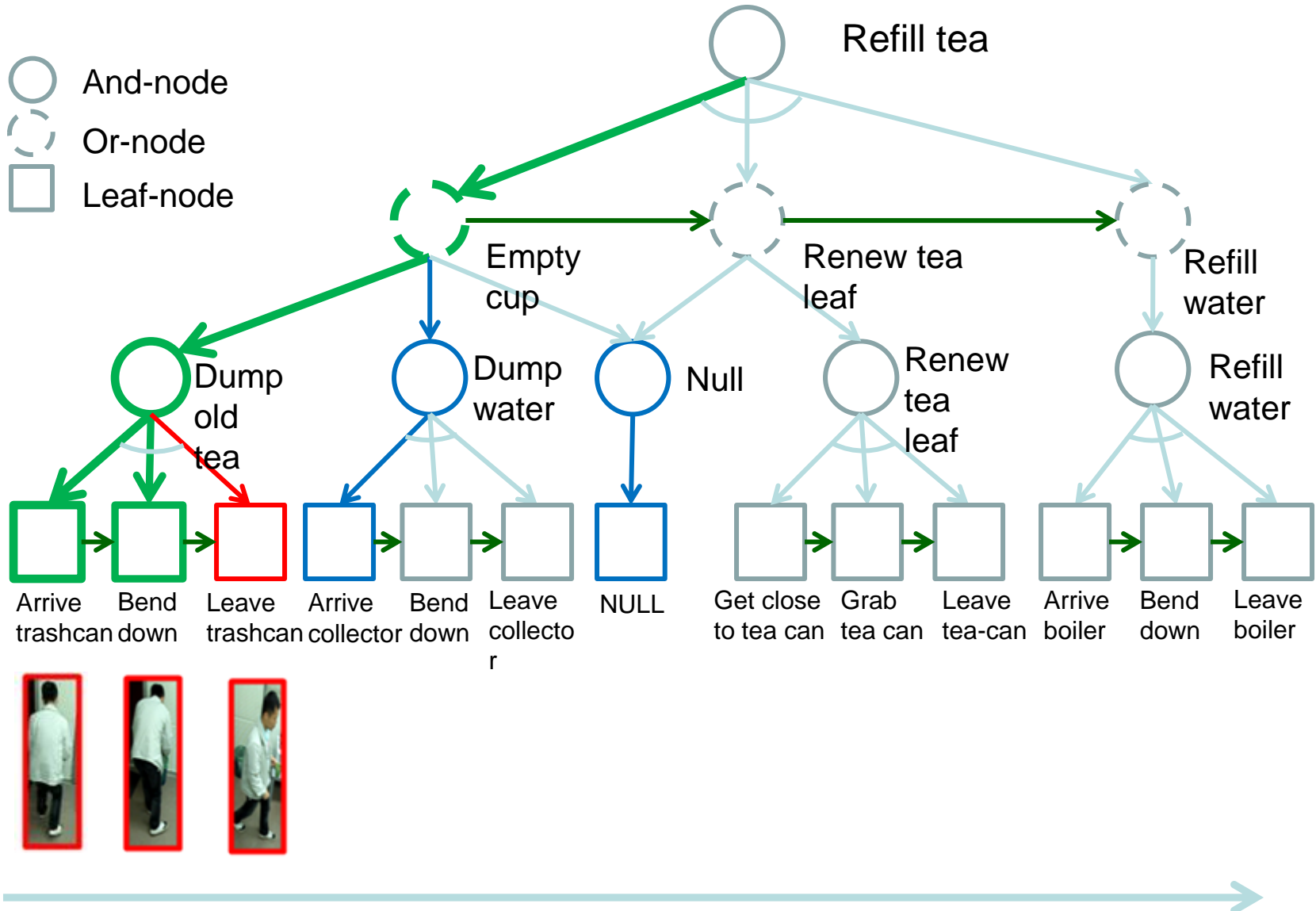
- And-node
- Or-node
- Leaf-node



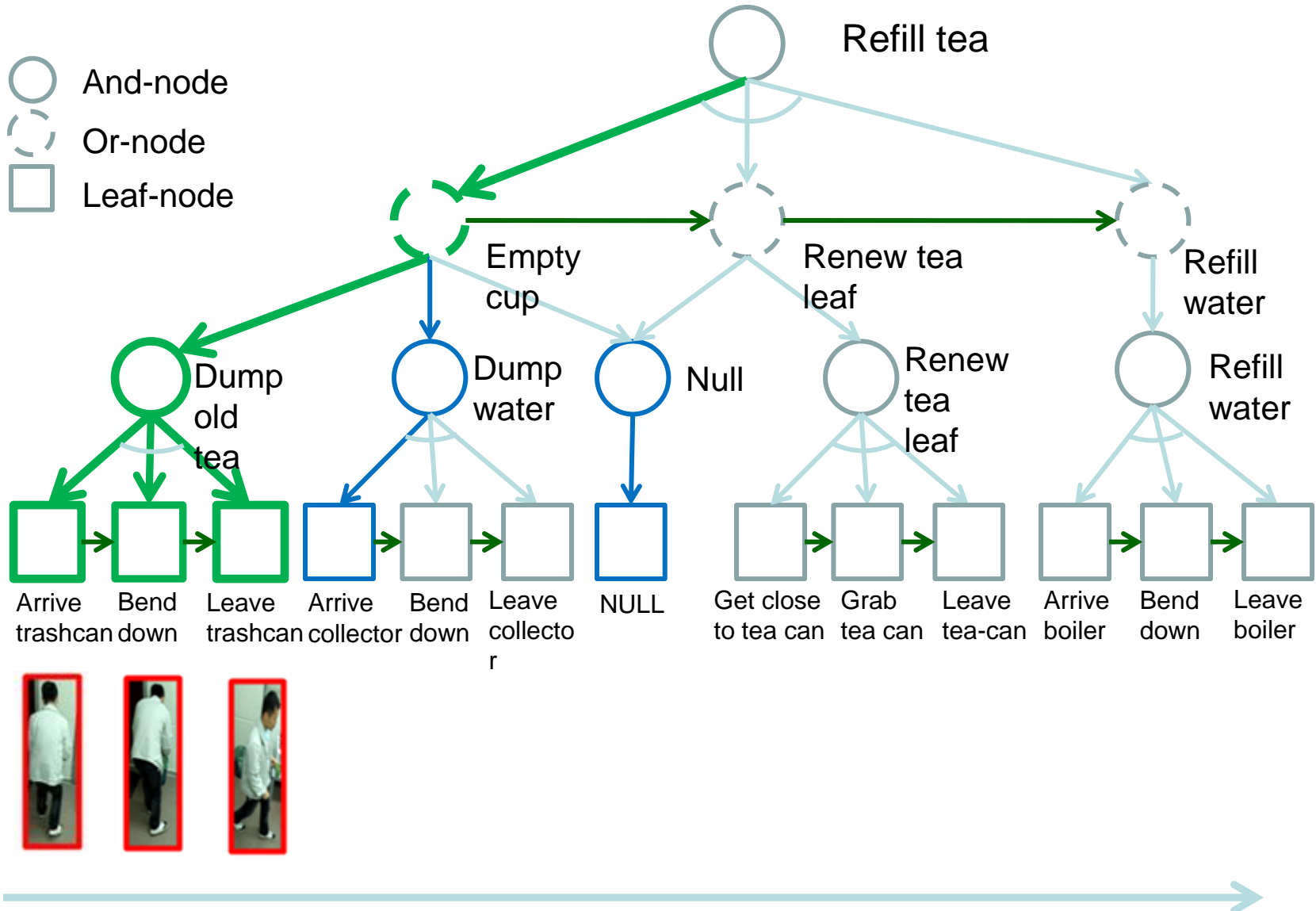
Parsing process (Earley Parser [Earley 1970])



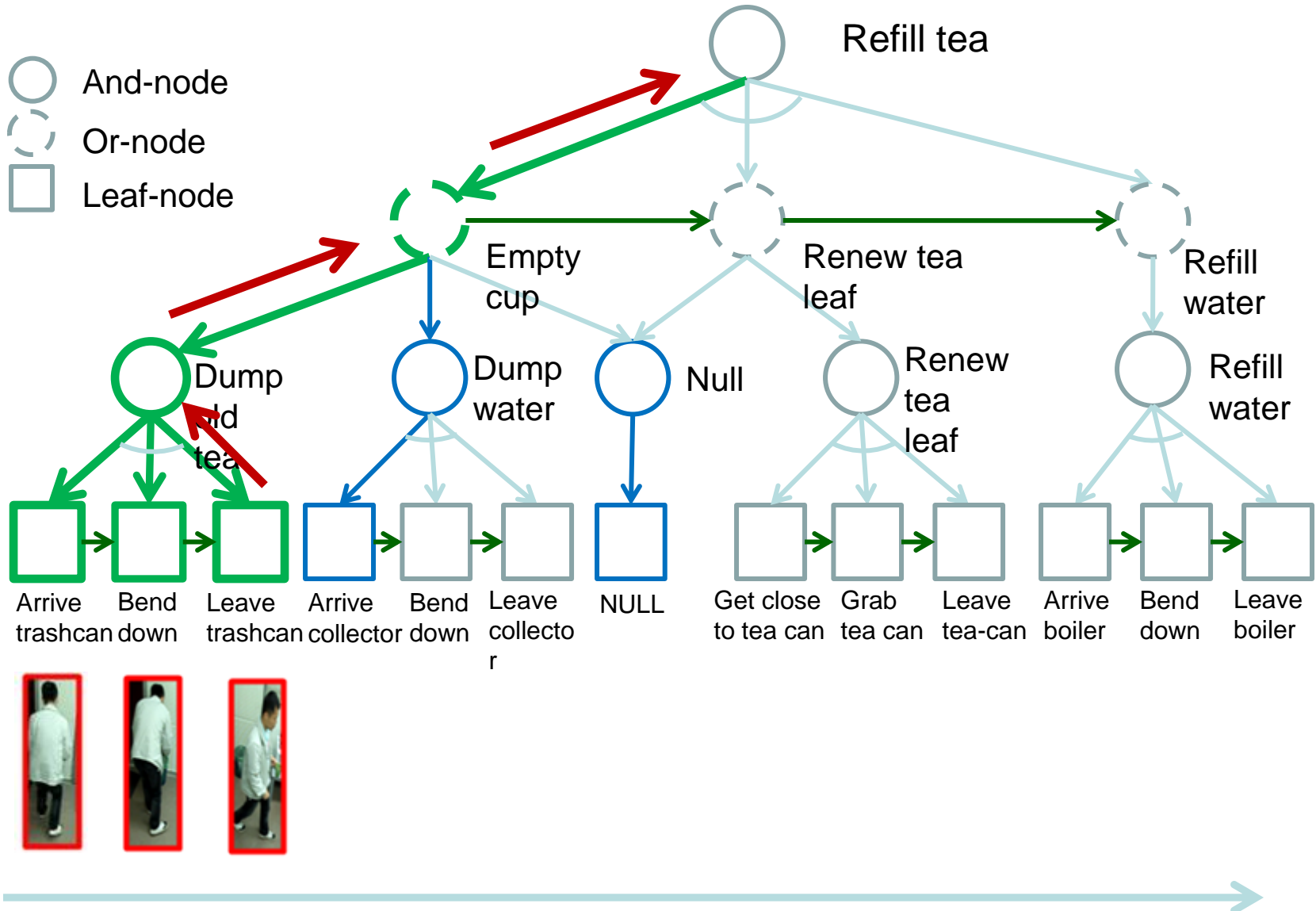
Parsing process (Earley Parser [Earley 1970])



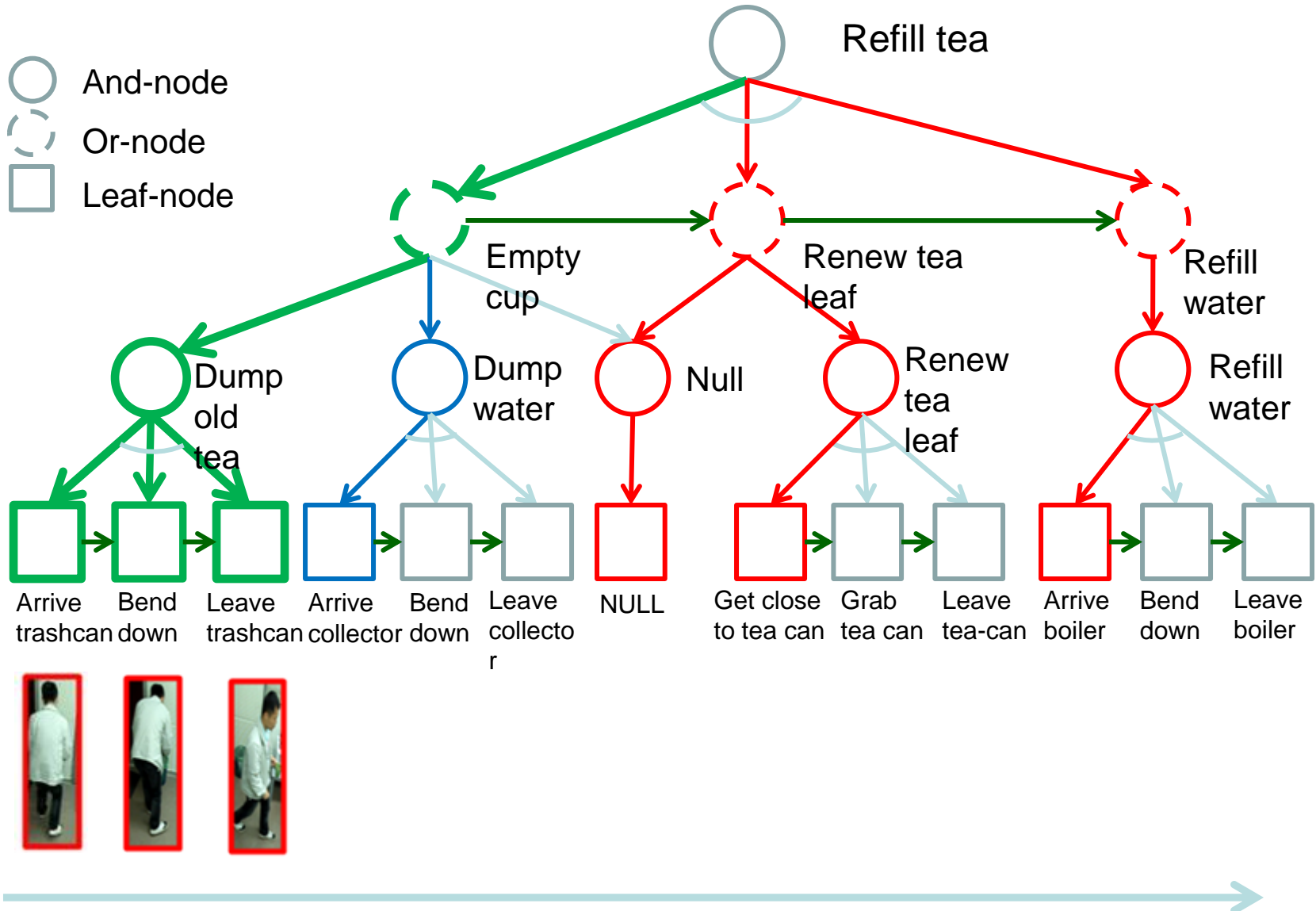
Parsing process (Earley Parser [Earley 1970])



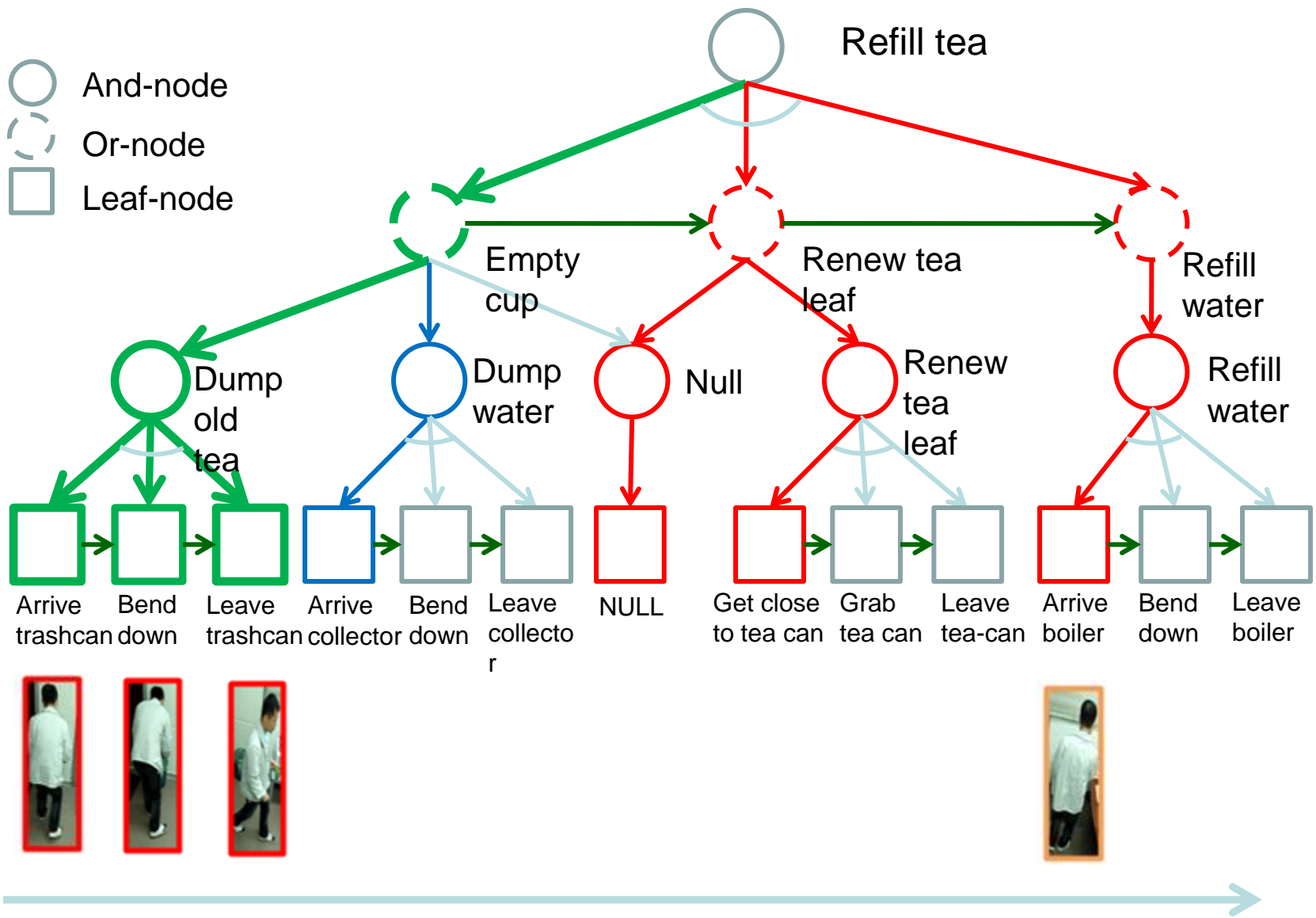
Parsing process (Earley Parser [Earley 1970])



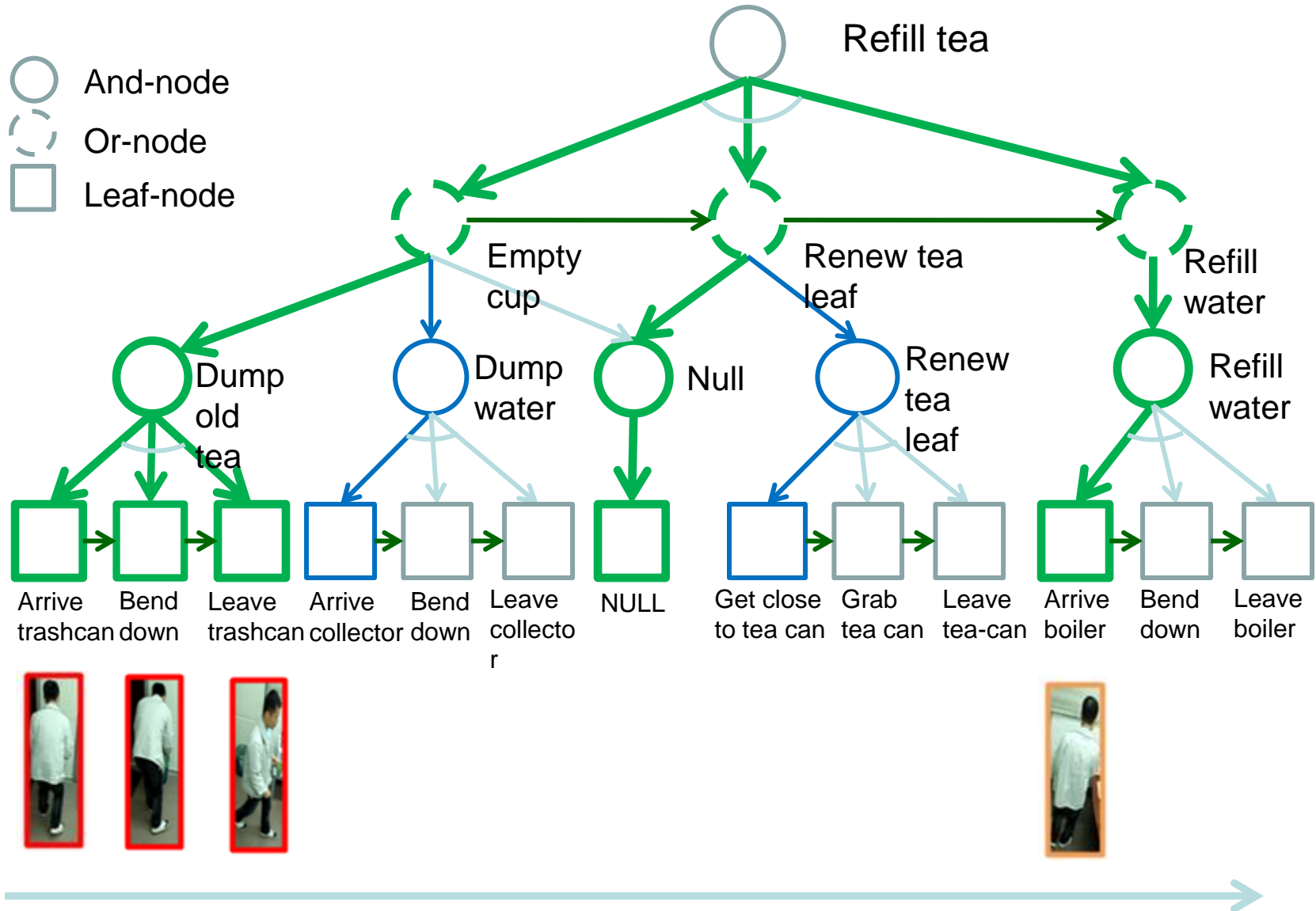
Parsing process (Earley Parser [Earley 1970])



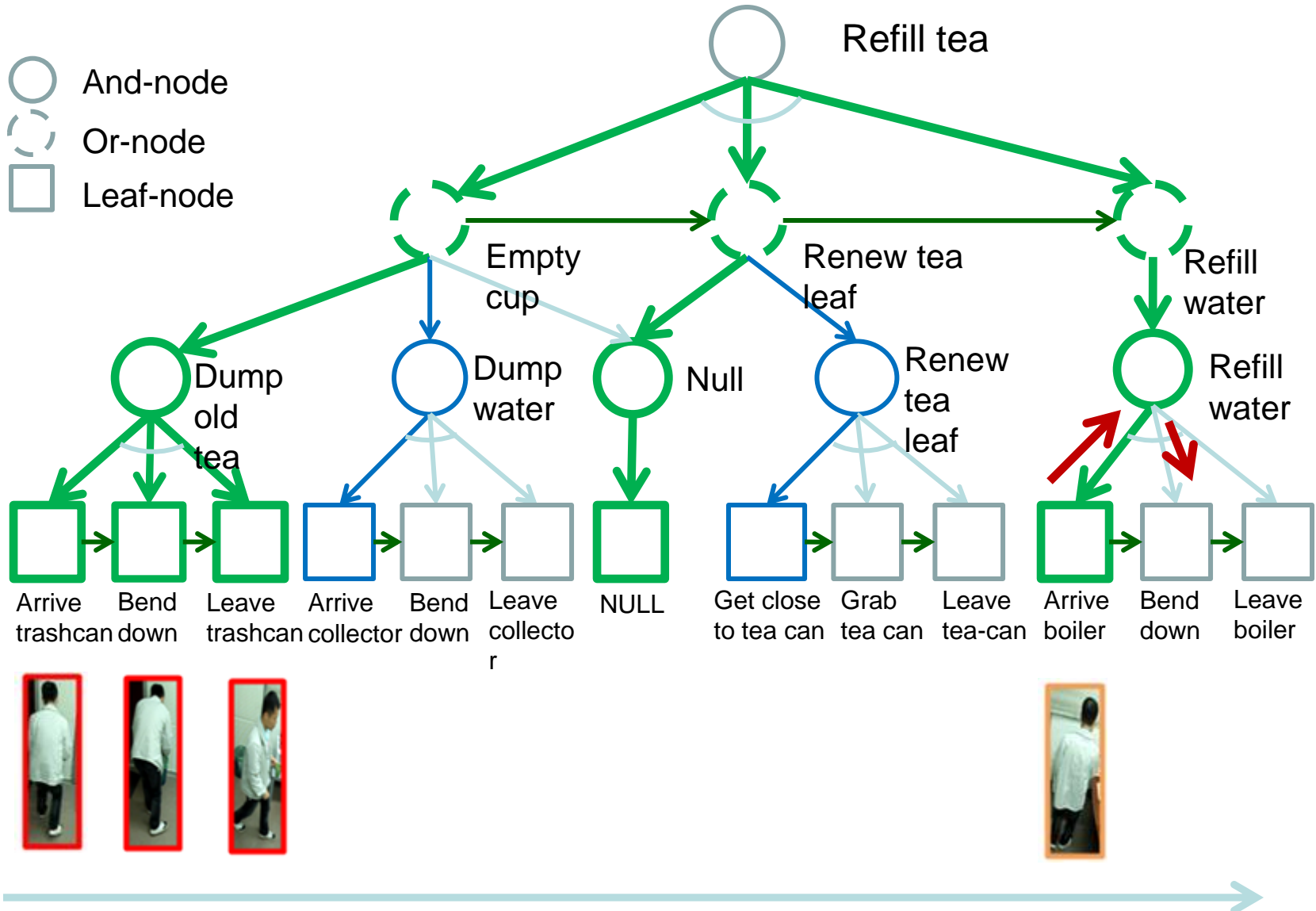
Parsing process (Earley Parser [Earley 1970])



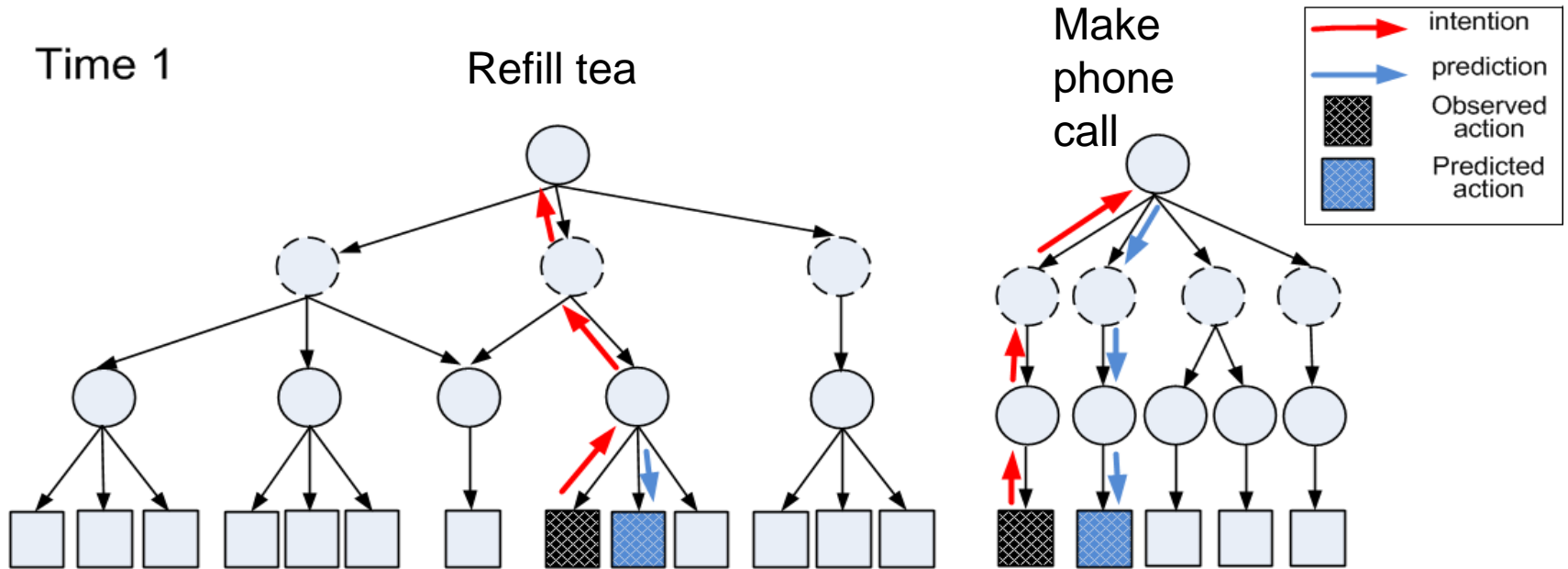
Parsing process (Earley Parser [Earley 1970])



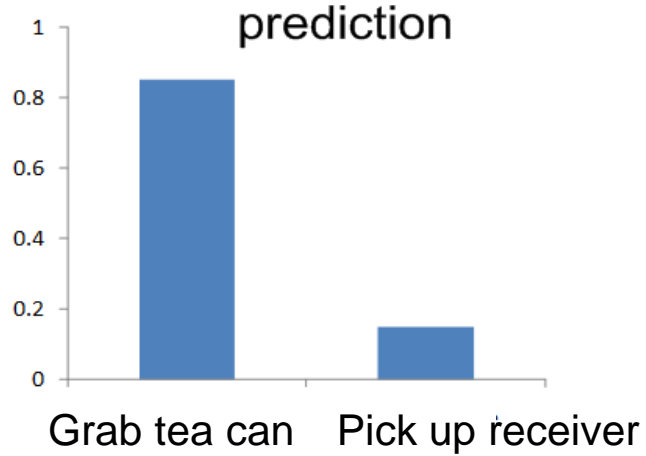
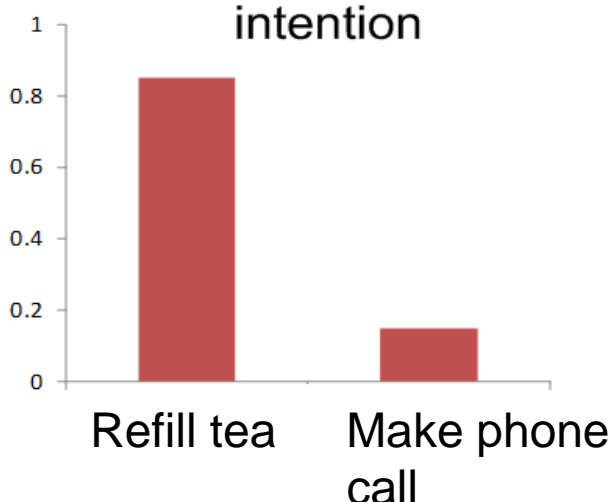
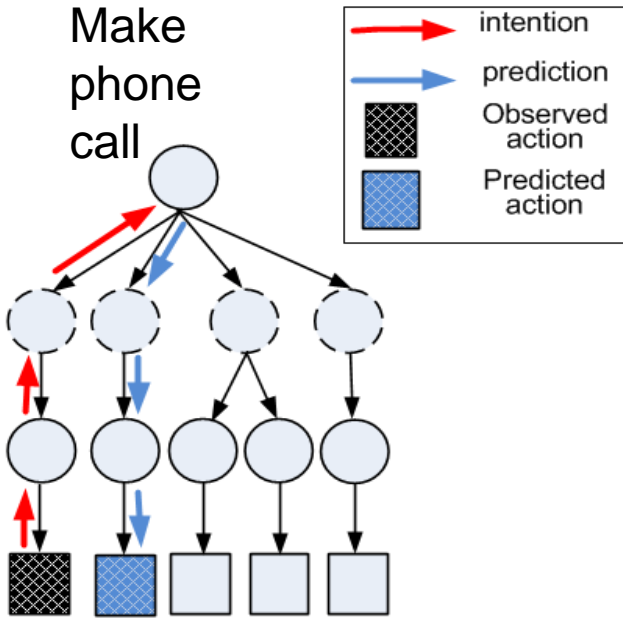
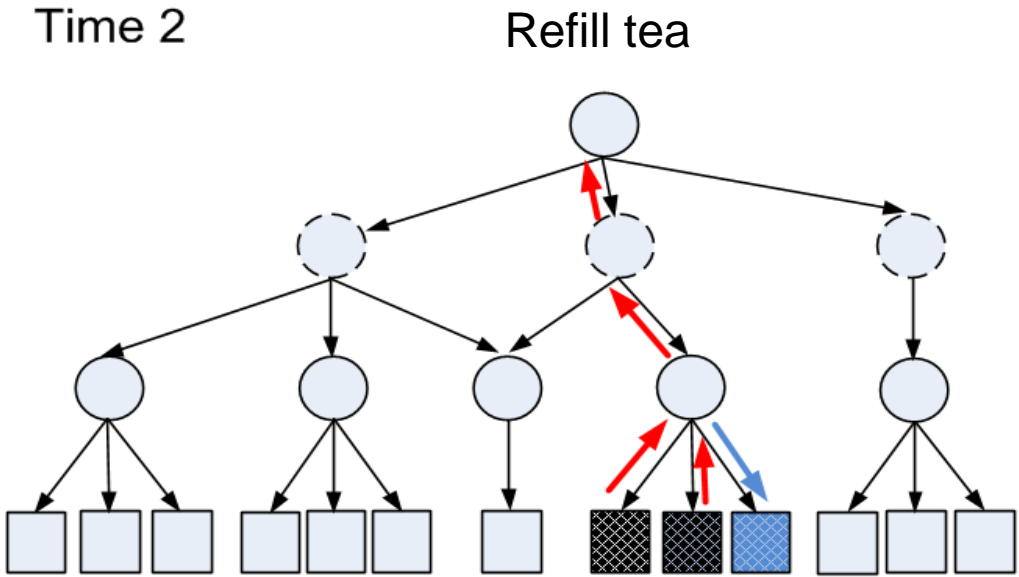
Parsing process (Earley Parser [Earley 1970])



Intention and prediction

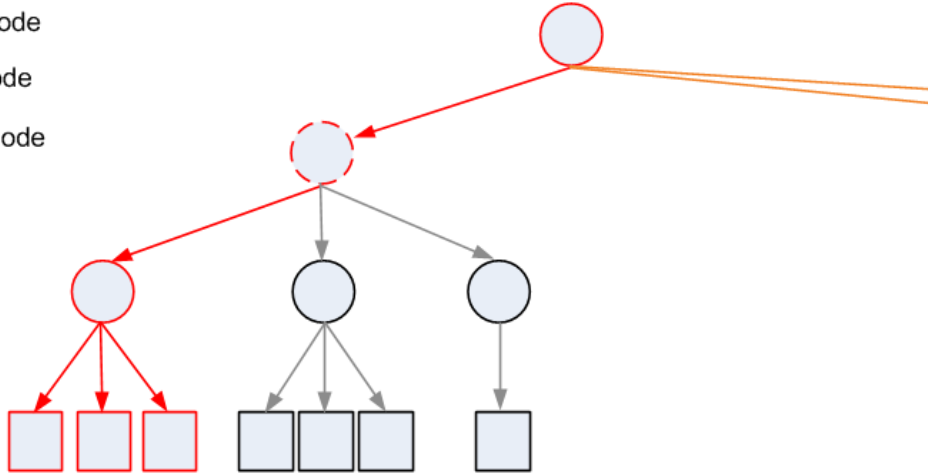


Intention and prediction



Handle event interruption




- And-Node
- Or-Node
- Leaf-Node

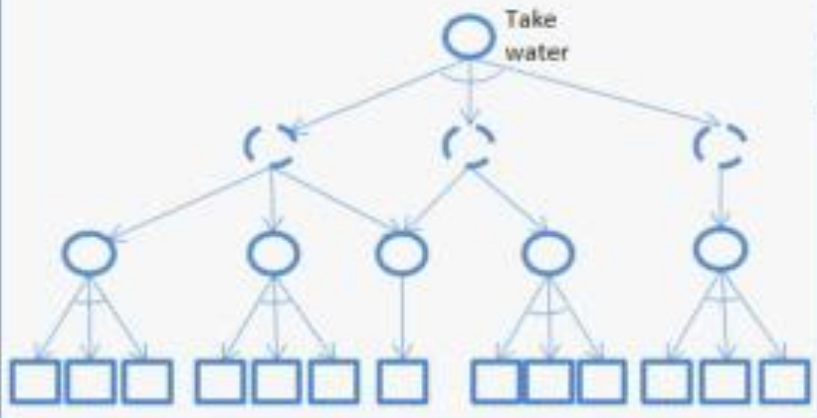
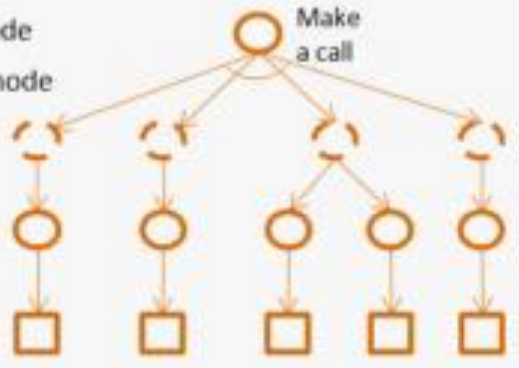


- First Partial parse tree of take water
- Parse tree of take a phone
- Second Partial parse tree of take water

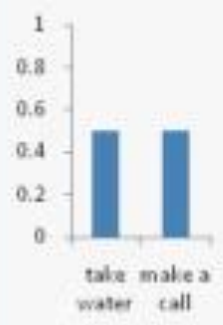
Observed
Data



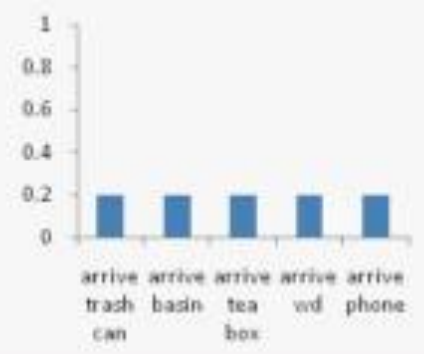
-  And-node
-  Or-node
-  Leaf-node



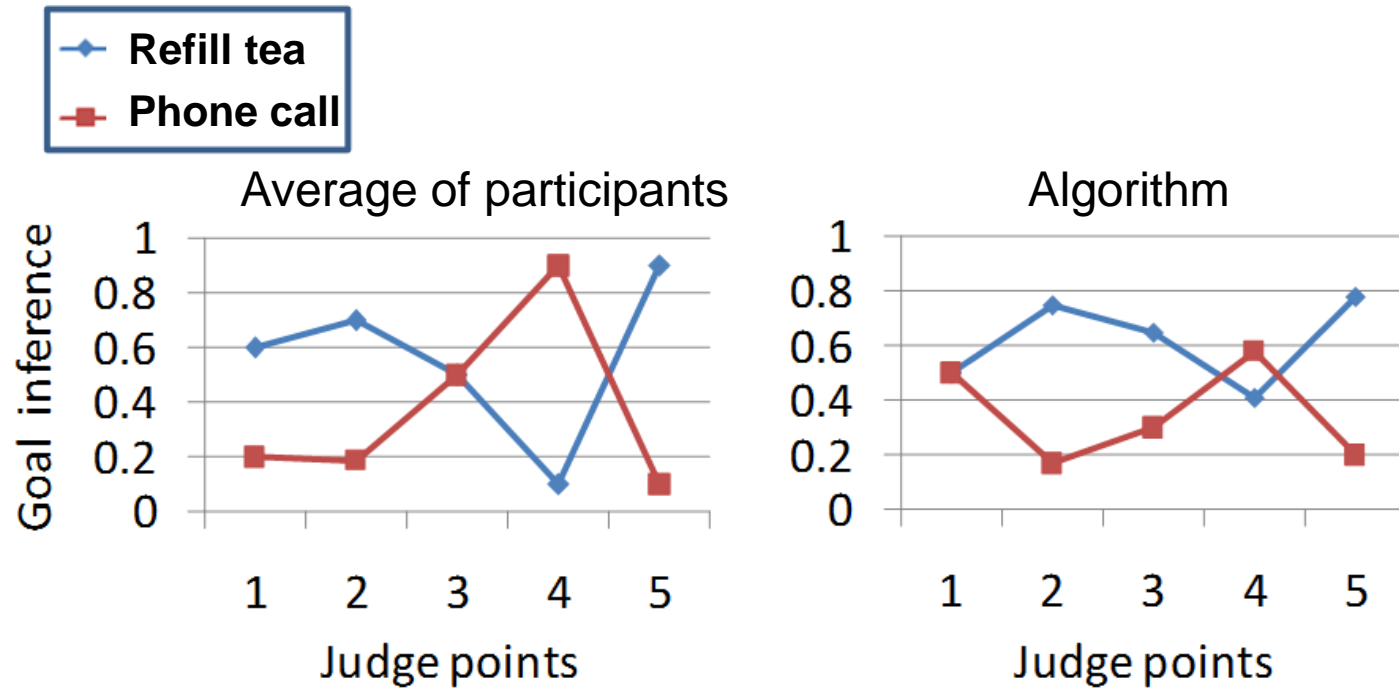
Intention



Prediction

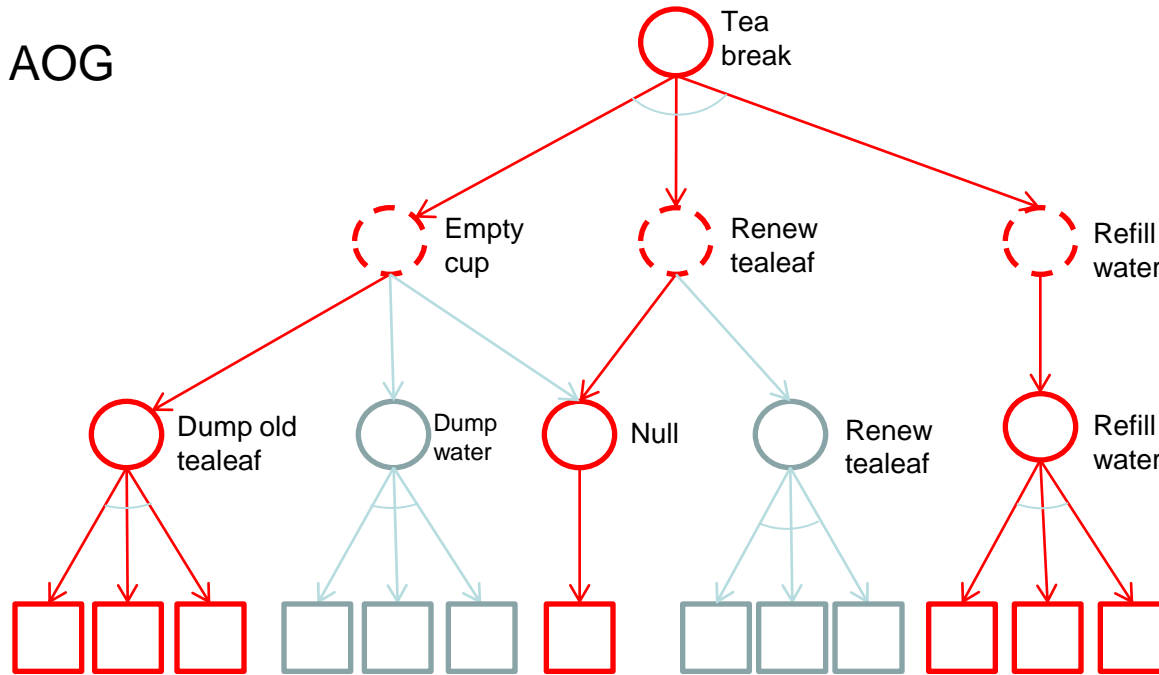


Comparison with human prediction [Baker, Saxe and Tenenbaum 2009]



Synthesize new events by sampling the AoG

Sample from AOG

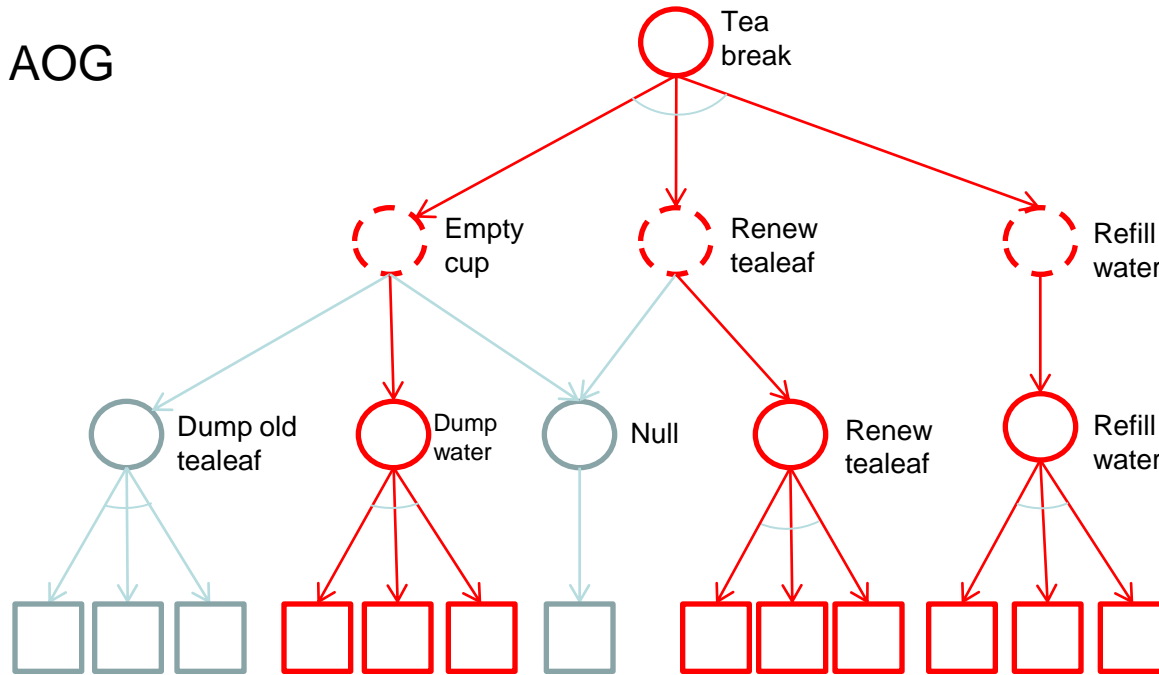


Synthesized Event



Event synthesis

Sample from AOG

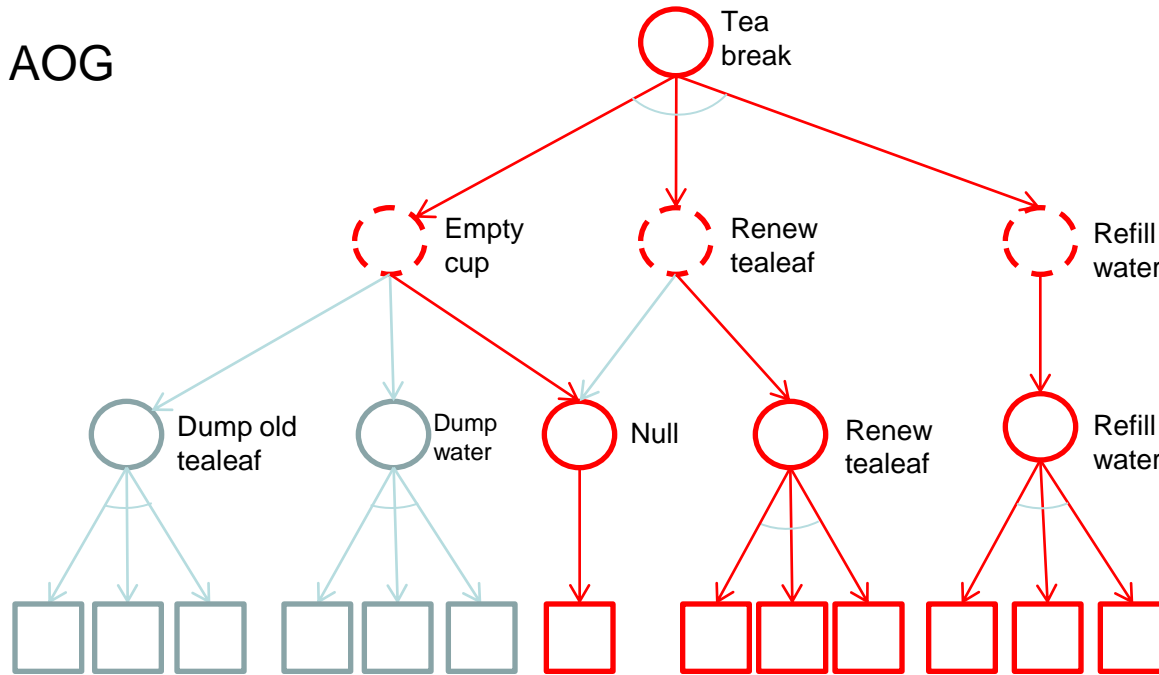


Synthesized
Event



Event synthesis

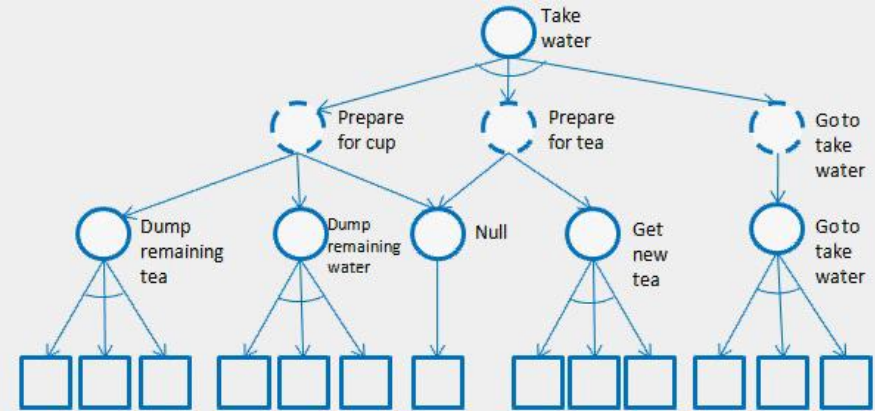
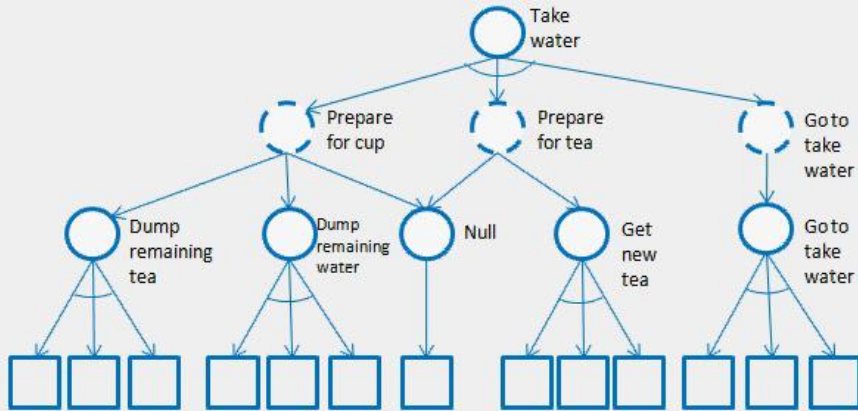
Sample from AOG



Synthesized
Event

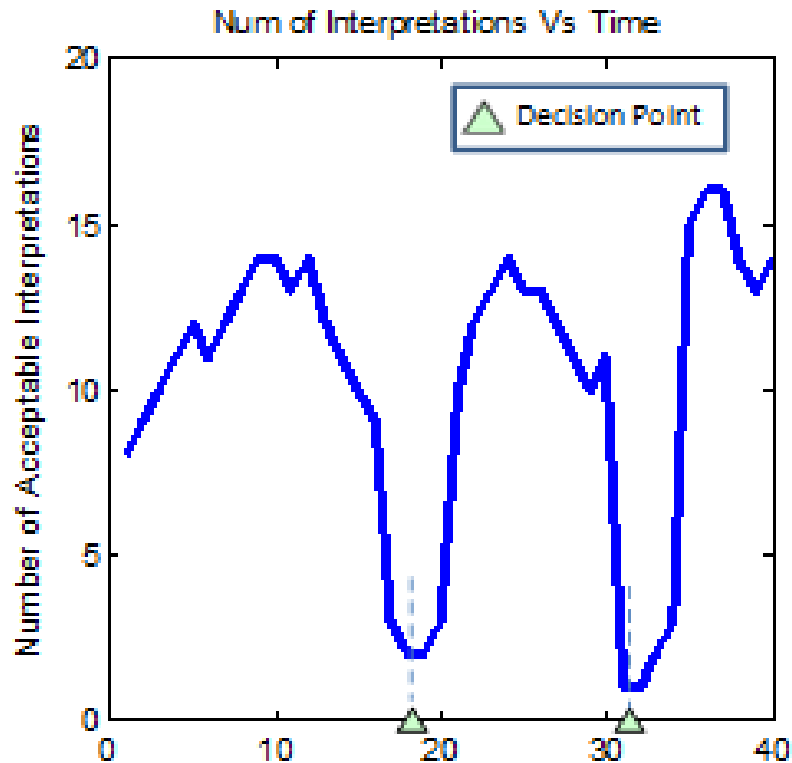


Event synthesise examples*



* We use stored foreground windows and a background frame to synthesize new videos

Computation complexity of parsing



- Initially the number of interpretations above a threshold grows rapidly over time.
- At certain decisive moments, i.e. when informative actions are observed, large number of unlikely interpretation drops below the threshold and hence is pruned.



Pickup phone



Reach water boiler