



Figure-r. 1. **Overview of SHIELDAGENT.** (Top) From AI regulations (e.g. EU AI Act) and platform-specific safety policies, SHIELDAGENT first extracts verifiable rules and iteratively refines them to ensure each rule is accurate, concrete, and atomic. It then clusters these rules and assembles them into an action-based safety policy model, associating actions with their corresponding constraints (with weights learned from real or simulated data). (Bottom) During inference, SHIELDAGENT retrieves relevant rule circuits w.r.t. the invoked action and performs action verification. By referencing existing workflows from a hybrid memory module, it first generates a step-by-step shielding plan with operations supported by a comprehensive tool library to assign truth values for all predicates, then produces executable code to perform formal verification for actions. Finally, it runs probabilistic inference in the rule circuits to provide a safety label and explanation and reports violated rules.

Algorithm 1 SHIELDAGENT Inference Procedure

Require: Interaction history $\mathcal{H}_{<i} = \{(o_j, a_j) \mid j \in [1, i-1]\}$ from the target agent; Current observation o_i ; Agent output a_i ; Safety policy model $\mathcal{G}_{\text{ASPM}} = (\mathcal{P}, \mathcal{R}, \pi_\theta)$; Safety threshold ϵ .

- 1: $p_a \leftarrow \text{EXTRACT}(a_i)$ ▷ Extract *action* predicates
- 2: $\mathcal{C}_{\theta_a}^{p_a} = (\mathcal{P}_{p_a}, R_{p_a}, \theta_a) \leftarrow \text{RETRIEVE}(p_a, \mathcal{G}_{\text{ASPM}})$
- 3: $\mathcal{V}_s = \{p_s^i : v_s^i\} \leftarrow \emptyset$ ▷ Initialize predicate-value map
- 4: **for each** rule $r = [\mathcal{P}_r, T_r, \phi_r, t_r] \in R_{p_a}$ **do**
- 5: $\mathcal{W}_r \leftarrow \text{RETRIEVEWORKFLOW}(r, p_a)$
- 6: **while** $\exists p_s \in \mathcal{P}_r$ s.t. $\mathcal{V}_s[p_s]$ is not assigned **do**
- 7: $A_s \leftarrow \text{PLAN}(\mathcal{W}_r, r, \mathcal{P}_r)$ ▷ Generate an action plan with shielding operations (e.g., SEARCH, CHECK)
- 8: **for each** step t_s^i in action plan A_s **do**
- 9: $o_s^i \leftarrow \text{EXECUTE}(t_s^i, \mathcal{H}_{<i}, o_i)$ ▷ Get step result
- 10: $\mathcal{V}_s[p_s] \leftarrow \text{PARSE}(o_s^i), p_s \in \mathcal{P}_r$ ▷ Attempt to assign a truth value to any unassigned predicates
- 11: **end for**
- 12: **end while**
- 13: $l_r \leftarrow \text{VERIFY}(r, \mathcal{V}_s)$ ▷ Run formal verification
- 14: **end for**
- 15: $\epsilon_s \leftarrow P_\theta(\mu_{p_a=1}) - P_\theta(\mu_{p_a=0})$ ▷ Calculate safety condition via Eq. (3)
- 16: **if** $\epsilon_s \geq \epsilon$ **then**
- 17: $l_s \leftarrow 1$ ▷ Action p_a is safe
- 18: **else**
- 19: $l_s \leftarrow 0$ ▷ Action p_a is unsafe
- 20: **end if**
- 21: **return** (l_s, V_s, T_s) ▷ Return safety label, violated rules, textual explanation

Algorithm 2 ASPM Structure Optimization

Require: Predicate set $\mathcal{P} = \{\mathcal{P}_a, \mathcal{P}_s\}$; Rule set $\mathcal{R} = \{\mathcal{R}_a, \mathcal{R}_p\}$; Embedding model \mathcal{E} ; Clustering algorithm \mathcal{C} ; Refinement budget N_b ; Max iterations M_{it} ; Surrogate LLM; Graph $G = (\mathcal{P}, E)$ with initial edge weights E .

```
1: Initialize vagueness score for each predicate  $\mathcal{V}_p, p \in \mathcal{P}$ 
2:  $\mathcal{V}_r = \max\{\mathcal{V}_{p_1}, \dots, \mathcal{V}_{p_{|\mathcal{P}_r|}}\}, \mathcal{P}_r \subseteq \mathcal{P}$  ▷ Compute vagueness score for each rule
3: Initialize a max-heap  $\mathcal{U} \leftarrow \{(\mathcal{V}_r, r) \mid r \in \mathcal{R}\}$ 
4:  $n \leftarrow 0$  ▷ Count how many refinements have been done
5: for  $m = 1$  to  $M_{it}$  do
6:   changed  $\leftarrow$  false ▷ Tracks if any update occurred in this iteration
7:   while  $\mathcal{U} \neq \emptyset \wedge n \leq N_b$  do
8:      $(\_, r) \leftarrow \text{HeapPop}(\mathcal{U})$  ▷ Pop the most vague rule
9:     if  $\text{LLM\_verifiable}(r) = \text{false}$  then
10:       $r_{\text{new}} \leftarrow \text{LLM\_refine}(r, \mathcal{P}_r)$  ▷ Refine rule  $r$  to be verifiable; update its predicates if needed
11:      Update  $\mathcal{R}$ : replace  $r$  with  $r_{\text{new}}$ 
12:      Update  $\mathcal{P}$ : if  $r_{\text{new}}$  introduces or revises predicates
13:      Recompute  $\mathcal{V}_p$  for any changed predicate  $p$  in  $r_{\text{new}}$ 
14:      Recompute  $\mathcal{V}_{r_{\text{new}}} = \max\{\mathcal{V}_p \mid p \in \mathcal{P}_{r_{\text{new}}}\}$ 
15:      Push  $(\mathcal{V}_{r_{\text{new}}}, r_{\text{new}})$  into  $\mathcal{U}$ 
16:       $n \leftarrow n + 1$ 
17:      changed  $\leftarrow$  true
18:   end if
19: end while
20:  $\mathcal{K} \leftarrow \mathcal{C}(G)$  ▷ Cluster predicates in  $G$  to prune redundancy
21: for each cluster  $C \in \mathcal{K}$  do
22:    $p_{\text{merged}} \leftarrow \text{LLM\_merge}(C, \mathcal{R})$  ▷ Merge similar predicates/rules in  $C$  if beneficial
23:   if  $p_{\text{merged}} \neq \emptyset$  then
24:     Update  $G$ : add  $p_{\text{merged}}$ , remove predicates in  $C$ 
25:     Update  $\mathcal{R}$  to replace references of predicates in  $C$  with  $p_{\text{merged}}$ 
26:     Recompute  $\mathcal{V}_{p_{\text{merged}}}$  and any affected  $\mathcal{V}_r$ 
27:     Push updated rules into  $\mathcal{U}$  by their new  $\mathcal{V}_r$ 
28:     changed  $\leftarrow$  true
29:   end if
30: end for
31: if changed = false then
32:   break ▷ No more refinements or merges
33: end if
34: end for
35: return ASPM  $\mathcal{G}_{\text{ASPM}}$  with optimized structure and randomized weights
```

Algorithm 3 ASPM TRAINING PIPELINE

Require: Rule set \mathcal{R} ; state predicates \mathcal{P}_s and action predicates \mathcal{P}_a ; similarity threshold θ ; number of clusters k .

- 1: $A \in \{0, 1\}^{|\mathcal{P}_s| \times |\mathcal{P}_s|} \leftarrow \mathbf{0}$ ▷ Initialize adjacency matrix
- 2: $A_{ij} \leftarrow 1$ if (p_s^i, p_s^j) co-occur in any rule OR $\text{cosSim}(\text{emb}(p_s^i), \text{emb}(p_s^j)) \geq \theta$; else 0. ▷ Build adjacency matrix
- 3: $\text{labels} \leftarrow \text{SPECTRALCLUSTERING}(A, k)$ ▷ Cluster the state predicates into k groups
- 4: **for** $\ell = 1$ **to** k **do**
- 5: $C_p^\ell \leftarrow \{p_s \mid \text{labels}[p_s] = \ell\}$ ▷ Form predicate clusters C_p
- 6: **end for**
- 7: **for each** pair (p_s^i, p_s^j) **that co-occur do**
- 8: **if** $\text{labels}[p_s^i] \neq \text{labels}[p_s^j]$ **then**
- 9: $C_p^\ell \leftarrow C_p^\ell \cup C_p^m$ s.t. $p_s^i \in C_p^\ell, p_s^j \in C_p^m$ ▷ If two co-occurring predicates appear in different clusters, merge them
- 10: **end if**
- 11: **end for**
- 12: **for** $\ell = 1$ **to** k' **do**
- 13: $C_r^\ell \leftarrow \{r_s \mid p_s \in C_p^\ell\}$ ▷ Group rules which share state predicates in the same cluster
- 14: **end for**
- 15: $\mathcal{G}_{\text{ASPM}} \leftarrow \emptyset$ ▷ Initialize ASPM as an empty dictionary with actions as keys
- 16: **for each** $p_a \in \mathcal{P}_a$ **do**
- 17: **for each** rule cluster $C_r^\ell \in \mathcal{C}_r$ **do**
- 18: **for each** rule $r \in C_r^\ell$ **do**
- 19: **if** $p_a^r \in r$ **then**
- 20: $\mathcal{G}_{\text{ASPM}}[p_a] = \mathcal{G}_{\text{ASPM}}[p_a] \cup C_r^\ell$ ▷ Associate action circuits with any relevant rule clusters
- 21: **break**
- 22: **end if**
- 23: **end for**
- 24: **end for**
- 25: **end for**
- 26: **for each** action circuit $\mathcal{C}_{\theta_a}^{p_a}$ **do**
- 27: **for each** rule $r \in \mathcal{C}_{\theta_a}^{p_a}$ **do**
- 28: Initialize rule weight θ_r randomly
- 29: **end for**
- 30: **for** epoch = 1 **to** max epochs **do**
- 31: **for** $i = 1$ **to** N **do**
- 32: Compute $P_\theta(\mu_{p_a=1}^{(i)})$ and $P_\theta(\mu_{p_a=0}^{(i)})$ ▷ Run probabilistic inference to obtain corresponding safety probabilities via Eq. (1)
- 33: Compute loss $\mathcal{L}(\theta)$ ▷ Calculate loss w.r.t. the groundtruth labels via Eq. (2)
- 34: Update θ using gradient descent
- 35: **end for**
- 36: **end for**
- 37: **end for**
- 38: **return** Action-based safety policy model $\mathcal{G}_{\text{ASPM}}$ with trained weights
