# APPENDIX

# A PROOFS

LEMMA 1 *(Correctness) The* GROUP-COVERAGE *algorithm successfully identifies if a group* $\mathbf{g}$ *is covered or not, i.e., if there are at least* $\tau$ *instances of* $\mathbf{g}$ *in the data set* $\mathcal{D}$.

PROOF. Each set query with a yes answer contains at least one object belonging to the group $\mathbf{g}$. Using this, the algorithm maintains a lower bound $cnt$ on the $|\mathbf{g}|$ in $\mathcal{D}$. That is, $cnt \leq |\mathbf{g}|$. The algorithm returns true when $cnt = \tau$. When $cnt = \tau$, $\mathbf{g}$ is covered, because $\tau = cnt \leq |\mathbf{g}|$. The algorithm returns false when the queue is empty and $cnt < \tau$. For sets with yes answers, the algorithm divides the set in two halves, unless the set size is 1. As a result, when the queue, all the set questions with yes answers have had size 1, otherwise the queue would have not empty. Therefore, $|\mathbf{g}| = cnt < \tau$, meaning that $\mathbf{g}$ is uncovered. □

# B PSEUDO-CODES

---

**Algorithm 5** PARTITION & LABEL

---

1: **function** PARTITION($\mathcal{D}, \mathcal{G}, n$)
2:     Let $Q$ = an empty queue
3:     **for** $i \leftarrow 0$ to $N$ with step size $n$ **do**:
4:         $t \leftarrow \{t_i, \cdots, t_j\}$
5:         $Q.add(t)$
6:     Let $S$ = an empty set
7:     **while** $Q$ is not empty **do**
8:         $T \leftarrow Q.del\_top()$
9:         $(i, j) \leftarrow (T.b\_index, T.e\_index)$
10:        ans $\leftarrow$ASKQUESTION($\{t_i, \cdots, t_j\}, \mathbf{g}'$)
11:        **if** ans=no **then**
12:            $S$.ADD($\{t_i, \cdots, t_j\}$)
13:        **else**
14:            **if** $j > i$ /*if setsize>1*/ **then**
15:                $T_1 \leftarrow \{t_i, t_{\lfloor \frac{i+j}{2} \rfloor}\}$
16:                $T_2 \leftarrow (t_{\lfloor \frac{i+j}{2} \rfloor + 1}, t_j)$
17:                $Q.add(T_1); Q.add(T_2)$
18:    **return** $S$
19: **function** LABEL($\mathcal{D}, \mathcal{G}, \tau$)
20:     $cnt \leftarrow 0$
21:     **for** $t \in \mathcal{G}$ **do**
22:         $l \leftarrow$POINTQUERY($t$)
23:         **if** $l \neq \mathbf{g}$ **then** $\mathcal{G}$.REMOVE($t$)
24:         **else** $cnt \leftarrow cnt + 1$
25:         **if** $cnt \geq \tau$ **then**
26:             **break**
27:     **return** $\mathcal{G}$

---

**Algorithm 6** LABEL SAMPLES & AGGREGATE

---

1: **function** LABELSAMPLES($\mathcal{D}, \tau, c = 2$)
2:     **for** $c\tau$ random samples $t$ from $\mathcal{D}$ **do**
3:         $l \leftarrow$ POINTQUERY($t$)
4:         $\mathcal{L}.add(\langle t, l \rangle); \mathcal{D}.remove(t)$
5:     **return** $\mathcal{D}, \mathcal{L}$
6: **function** AGGREGATE($\mathcal{L}, N, \tau, \mathbb{G}, mult$)
7:     sort $\mathbb{G}$ based on $\mathcal{L}$.COUNT(g), $\mathbf{g} \in \mathbb{G}$, ascending
8:     $sum \leftarrow 0; \mathcal{G} \leftarrow \{\}$
9:     **for** $\mathbf{g} \in \mathbb{G}$ **do**
10:        $E_{\mathbf{g}} \leftarrow \frac{\mathcal{L}.\text{COUNT}(\mathbf{g})}{|\mathcal{L}|} \times N$
11:        **if:** $mult = $ true **then** $\mathcal{G} \leftarrow \mathcal{G} \mid \mathcal{G}$.parent= g.parent
12:        **if:**$sum + E_{\mathbf{g}} < \tau$ **then** $\mathcal{G}$.ADD(g); $sum \leftarrow sum + E_{\mathbf{g}}$
13:        **else:** $\mathbb{G}_{agg}$.add($\mathcal{G}$) ; $\mathcal{G} \leftarrow \{\mathbf{g}\}; sum \leftarrow E_{\mathbf{g}}$
14:    **return** $\mathbb{G}_{agg}$.add($\mathcal{G}$)

---

**Algorithm 7** BASE-COVERAGE($\mathcal{D}, \tau, \mathbf{g}$)

---

**Input:** Dataset $\mathcal{D}$, coverage threshold $\tau$, and target group $\mathbf{g}$
**Output:** Coverage of group $\mathbf{g}$
1: $cnt \leftarrow 0$
2: **for** $\forall t \in \mathcal{D}$ **do**
3:     ans $\leftarrow$ASKQUESTION($t, \mathbf{g}$)
4:     **if** $ans =$**true then** $cnt \leftarrow cnt + 1$
5:     **if** $cnt = \tau$ **then return true** // covered
6: **return false** //uncovered