

## APPENDIX

### A. PRUNING METHOD OF CANDIDATE STRATEGIES GENERATION

At a high level, strategies with only few chunks are cost-effective, but may miss the ground truth. On the contrary, candidates with more chunks may improve the accuracy but are costly. However, intuitively, all valid strategies should include highly relevant chunks as input. Therefore, for each attribute  $a_j$  and document  $d_i$ , let  $C_{ij} = [ch_1, ch_2, \dots, ch_l]$  denote an ordered list of  $l$  retrieved chunks ranked by the cosine similarity between their embeddings and  $e_j$  in descending order, and we only consider the concatenation of prefix chunks as a candidate such that highly ranked chunks are always included. To be specific, we use  $C_{ij}^k, k \in [1, l]$  to represent the concatenation of the top- $k$  most relevant chunks, e.g.,  $C_{ij}^1 = ch_1$  and  $C_{ij}^2 = ch_1 \oplus ch_2$ . This approach greatly reduces computational complexity to  $O(l)$  while not sacrificing extraction accuracy much. Finally, considering that non-LLM strategies take the entire document as input,  $S_{ij} = \{\text{OpenIE}(d_i), \text{Codegen}(d_i), \text{PLM}(d_i), \text{LLM}(C_{ij}^1), \dots, \text{LLM}(C_{ij}^l)\}$ .

### BA. PROOF OF THEOREM 0.1

Next, we will show that solving Equation 1 is an NP-hard problem.

$$f^* = \arg \max_{f \in \mathcal{F} = \{f | f: \{r_{ij}\} \rightarrow S_{ij}\}} \frac{\sum_{i=1}^n \sum_{j=1}^m P(v_{ij}^s = v_{ij}^*)}{m \times n} \quad (1)$$

$$s.t. \quad \sum_{i=1}^n \sum_{j=1}^m c_{ij}^s \leq B, s = f(r_{ij})$$

**THEOREM 0.1.** *The problem of solving Equation 1 is NP-hard.*

**PROOF.** The problem of solving Equation 4 can be proven to be NP-hard by a reduction from the Group Knapsack Problem (GKP). The Group Knapsack Problem is defined as follows: given a set of items, which are partitioned into groups  $G$ , and exactly one item must be selected from each group. The goal is to maximize the total value of selected items without exceeding a given budget. To reduce the Group Knapsack Problem to the problem of solving Equation 1, we take each group in GKP as  $r_{ij}$ , where each item in  $G$  corresponds to a candidate strategy  $s \in S_{ij}$ . The value of an item is the estimated accuracy  $P(v_{ij}^s = v_{ij}^*)$  of strategy  $s$ , and its weight is the token cost  $c_{ij}^s$ . The capacity of the knapsack corresponds to the total budget  $B$ . The objective of Equation 1, i.e., selecting one strategy for each  $r_{ij}$  to maximize the total accuracy under the budget  $B$ , directly aligns with the GKP objective. Since GKP is NP-hard, this shows that solving Equation 1 is also NP-hard.