---

**Algorithm 1:** Unification Networks

---

**Input:** Invariant example $G$, variableness network $\psi$, example $K$, features network $\phi$, unifying features network $\phi_U$, upstream predictor network $f$

**Output:** Predicted label for example $K$

---

1 **begin** ▷ Unification Network
2     **return** $f \circ g(G, K)$                 ▷ Predictions using Soft Unification $g$

3 **begin** ▷ Soft Unification function $g$
4     **foreach** *symbol $s$* **in** $G$ **do**
5        $\boldsymbol{A}_{s,:} \leftarrow \phi(s)$                ▷ Features of $G$, $\boldsymbol{A} \in \mathbb{R}^{|G| \times d}$
6        $\boldsymbol{B}_{s,:} \leftarrow \phi_U(s)$            ▷ Unifying features of $G$, $\boldsymbol{B} \in \mathbb{R}^{|G| \times d}$
7     **foreach** *symbol $s$* **in** $K$ **do**
8        $\boldsymbol{C}_{s,:} \leftarrow \phi(s)$                ▷ Features of $K$, $\boldsymbol{C} \in \mathbb{R}^{|K| \times d}$
9        $\boldsymbol{D}_{s,:} \leftarrow \phi_U(s)$            ▷ Unifying features of $K$, $\boldsymbol{D} \in \mathbb{R}^{|K| \times d}$
10     **Let** $P = \text{softmax}(\boldsymbol{B}\boldsymbol{D}^T)$        ▷ Attention map over symbols, $\boldsymbol{P} \in \mathbb{R}^{|G| \times |K|}$
11     **Let** $\boldsymbol{E} = \boldsymbol{P}\boldsymbol{C}$             ▷ Attended representations of $G$, $\boldsymbol{E} \in \mathbb{R}^{|G| \times d}$
12     **foreach** *symbol $s$* **in** $G$ **do**
13        $\boldsymbol{U}_{s,:} \leftarrow \psi(s)\boldsymbol{E}_{s,:} + (1 - \psi(s))\boldsymbol{A}_{s,:}$    ▷ Unified representation of $G$ with $K$, $\boldsymbol{U} \in \mathbb{R}^{|G| \times d}$
14     **return** $\boldsymbol{U}$

---

1 Dear reviewers, thank you for your comments. We are pleased with the unanimous consensus on the novelty of our contribution,
2 and recognition by some reviewers of a comprehensive evaluation. We thank Reviewer 3 for their suggestion of an algorithm box
3 to clarify Section 2. We present Algorithm 1 above that defines our approach and the soft unification function $g$ with all the
4 learnable components $\psi, \phi, \phi_U, f$. We compute 2 sets of features (see L87, L90, L refers to lines in paper), then use the unifying
5 features to let each symbol in $G$ attend to a symbol in $K$. Depending on how much each symbol in $G$ is a variable, determined
6 by $\psi$, the representation of $G$ is interpolated between its and $K$'s features. We plan to include the algorithm in Section 2 and
7 simplify the notation, move Figure 2 to Section 3 and add further explanation using the extra page if the paper gets accepted. We
8 hope these changes will resolve clarity issues also raised by Reviewers 1 and 4. Reviewer 3 asks about hyper-parameters and
9 their selection as well as the number examples used for invariances. In Appendix A, second paragraph of each model details the
10 dimensions, layers etc. used. These were selected based on similar previous works on the datasets. We use one example to learn
11 from but have experimented with multiple (L209) of which the predictions are aggregated (e.g. sum), Appendix D Figure 8.

12 Reviewers 1 and 4 point out some baselines were missing in the text. We included them (e.g. DMN, G2N2) only as references
13 in the captions due to space constraints; we will expand their definitions to clarify further. So in response to Reviewer 4, in
14 Table 3 DMN and IMA are baselines without soft-unification and our approach achieves comparable results to DMN with half
15 the data size. We explain weak supervision briefly in L213 and will clarify strong supervision further. To answer Reviewer
16 4, we do not constrain the same symbol variable to attain the same value because it can appear in different contexts. When
17 visualising invariants, the choice of variable symbols is given by the variabless $\psi$ which is the same for same symbols. We leave
18 contextualised variablisation as future work (L119). However, the value assignment can be context dependant (L123, L140).

19 Reviewer 2 points out there is no quantitative analysis of identifying invariants. This was not included because (i) the focus
20 of our evaluation was to solve an upstream task with less examples using our approach whilst maintaining task performance,
21 and (ii) our approach might solve the task without necessarily using the *expected* invariants as shown in Appendix D Figure 11
22 and 12. For completeness, we can add this analysis to our synthetic dataset where we know what are the expected invariances.
23 Reviewer 2 mentions attention weights in Figure 5 are low $\leq 0.3$ (also mentioned by Reviewer 4) which might indicate low
24 confidence. There might be a misunderstanding here as these are the variableness $\psi$ of the symbols, not the attention weights. We
25 purposefully penalise the magnitude of $\psi$ in equation 3 (Sparsity) so we expect them to be low. This is because we want to find
26 the minimum variablisation of G to correctly predict K's label (see L252). In Figures 1 and 6, we threshold $\psi$ for visualisation
27 purposes (L227) although the interpretation of what symbols are variables is not binary, i.e. is it or is it not a variable. This is a
28 soft view of the notion of a variable. Therefore, to answer Reviewer 2's comment on soft-interpolation, this soft view indeed
29 produces varying symbol embeddings as intended. But this does not hurt the interpretability of the attention maps which are
30 instead computed from the unifying features (Algorithm 1, line 10). The interpolation happens after the attention is computed.

31 Reviewer 2 asks about context embeddings to disamguiate unification. This is indeed used in some models through unifying
32 features $\phi_U$ (L123, L140). In other words, in the sentence "$\boldsymbol{X}$:lily is a $\boldsymbol{Y}$:frog" what $\boldsymbol{X}$ unifies with takes into account that it is
33 related to a frog / animal, e.g. Figure 3 pink diamond represents unification RNN (L140). We thank Reviewer 2 for the suggested
34 related work [a] and will cite it; however, we do not claim our approach is more interpretable than Neural Turing Machines
35 (NTM) (L273). We discuss NTMs as an approach that lacks an explicit representation like our invariants (L278) in relation to
36 generalisation. We acknowledge that interpretability can be subjective and biased (L281), for example we might want to think
37 $\boldsymbol{X}$:Mary means *someone* but we do not know how the model uses these representations (L297, L300). We discuss why and why
38 not our results could be interpretable but refrain from claiming that our model provides an explanation like Grad-CAM (L301) or
39 is more interpretable than NTMs.

40 We thank you for your reviews and hope to have addressed your questions. We intend to incorporate all the suggestions and
41 remaining remarks into the paper.