



# A Noise-tolerant Differentiable Learning Approach for Single Occurrence Regular Expression with Interleaving

Rongzhen Ye<sup>1</sup>, Tianqu Zhuang<sup>1</sup>, Hai Wan<sup>1,\*</sup>, Jianfeng Du<sup>2,\*</sup>, Weilin Luo<sup>1</sup>, Pingjia Liang<sup>1</sup>

<sup>1</sup>Sun Yat-sen University, Guangzhou, China <sup>2</sup>Guangdong University of Foreign Studies, China



## Motivation

**Problem.** Learning *single occurrence regular expression with interleaving* (SOIRE) in *full expressive power* from a set of text strings with *noise*.

**Challenge.**

1. *heavy computation in searching* to get the full expressive power
2. *wrong search bias* resulting from noisy data

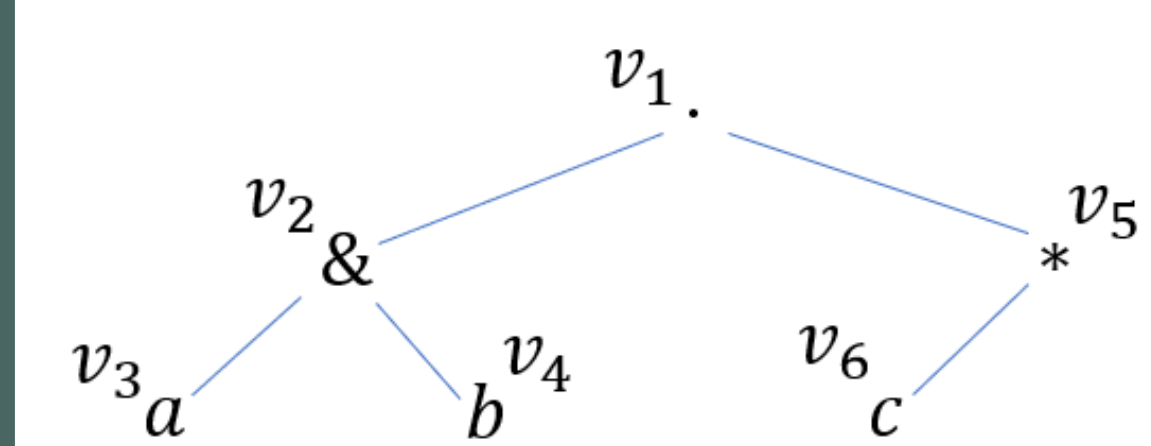
## Contributions

1. We have proposed a noise-tolerant differentiable learning approach SOIREDL. We first train the neural network that simulates SOIRE matching to classify the given strings, and then interpret the target SOIRE from the parameters.
2. Theoretically, the faithful encodings learnt by SOIREDL one-to-one correspond to SOIREs for a bounded size.
3. Experimental results have demonstrated higher performance compared with the SOTA approaches.

## Filter Matching

Filter matching for a SOIRE  $r$  and a string  $s$  is to check if  $r$  matches  $filter(s, \alpha(r))$ , where  $\alpha(r)$  denotes the set of symbol in  $r$ , and function  $filter(s, V)$  returns a string that only retains symbols in  $V$ , where  $V \subseteq \Sigma$ . Let  $g_{i,j}^t \in \{0, 1\}$  denote whether  $r^t$  matches  $filter(s_{i,j}, \alpha(r^t))$ , where  $s_{i,j}$  denotes the substring of  $s$  from  $i$  to  $j$ , and where  $s_{1,0} = \epsilon$  specially.

**Example 1.** For  $(a\&b)c^*$  and  $s = dbac$ ,



- *filter matching is to check if  $(a\&b)c^*$  matches  $filter(dbac, \{a, b, c\}) = bac$ , as  $\alpha((a\&b)c^*) = \{a, b, c\}$ .*
- $g_{1,2}^2$  denotes if  $a\&b$  ( $r^2$ ) matches  $ba$  and  $g_{1,2}^2 = 1$ .

## SOIRETM

**Theorem 1.** Given a SOIRE  $r$  and a string  $s$ ,  $r \models s$  iff  $filter(s, \alpha(r)) = s$  and  $r \models filter(s, \alpha(r))$ .

**The steps of SOIRETM:**

1. build the syntax tree of  $r$
2. check if  $filter(s, \alpha(r)) = s$
3. calculate  $g_{i,j}^t$  from shorter substrings to longer ones and from bottom to top of the syntax tree (dynamic programming)
4. return  $g_{1,|s|}^1$  (filter matching)

**Theorem 2.** Given a SOIRE  $r$  and a string  $s$ ,  $r \models s$  iff  $SOIRETM(r, s) = 1$ .

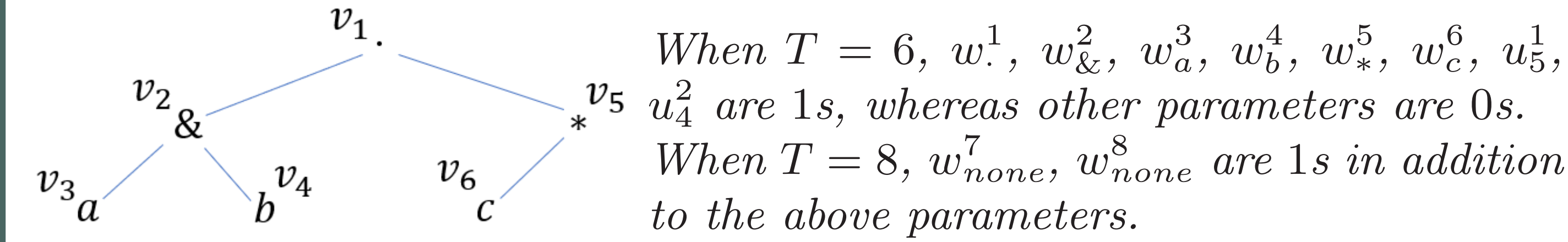
## From SOIRETM to Neural Network

The trainable parameters  $\theta = (w, u)$ :

- $w \in [0, 1]^{T \times |\mathbb{B}|}$ ,  $w_a^t$  denotes the probability of vertex  $t$  representing a symbol in  $\Sigma$  or an ordinary operator or the none operator.
- $u \in [0, 1]^{T \times T}$ ,  $u_{t'}^t$  denotes the probability of vertex  $t$  choosing vertex  $t'$  as its right child.

where  $\mathbb{B} = \Sigma \cup \{?, *, +, \cdot, \&, |, \text{none}\}$ ,  $T$  the bounded size of the target SOIRE.

**Example 2.**



**Four main parts in conversion.**

- $\rho_a^t$ , the differentiable version of  $\alpha(r^t)$ . ( $\sigma_{01}(x) = \min(\max(x, 0), 1)$ )

$$\rho_a^t = \sigma_{01}(w_a^t + \sum_{o \in \{?, *, +, \cdot, \&, | \}} w_o^t \rho_a^{t+1} + \sum_{o \in \{ \cdot, \&, | \}} w_o^t \sum_{t'=t+2}^T u_{t'}^t \rho_a^{t'})$$

- $flag_{i,j}^{t,t'}$ , the probability that there does not exist a symbol occurring in both  $s_{i,j}$  and  $\alpha(r^t)$  but not occurring in  $\alpha(r^{t'})$ .

$$flag_{i,j}^{t,t'} = 1 - \sigma_{01}(\sum_{a \in \Sigma} \sigma_{01}(1[a \in s_{i,j}] + (\rho_a^t - \rho_a^{t'}) - 1))$$

- $g_{i,j}^t$ , whether  $r^t$  matches  $filter(s_{i,j}, \alpha(r^t))$

$$g_{i,j}^t = \sum_{a \in \Sigma} w_a^t \cdot 1[filter(s_{i,j}, a) = a]$$

$$+ \sum_{o \in \{?, *, + \}} w_o^t p_{i,j}^t(o) + \sum_{o \in \{ \cdot, \&, | \}} w_o^t \sum_{t'=t+2}^T u_{t'}^t p_{i,j}^t(o, t')$$

- return value of SOIRETM, combining  $filter(s, \alpha(r)) = s$  and  $g_{1,|s|}^1$ .

$$\hat{y} = g_{1,|s|}^1 - \max_{a \in \Sigma} \sigma_{01}(1[a \in s] - \rho_a^1)$$

Objective function:  $\frac{1}{2}(\hat{y} - y)^2$ , where  $y$  is ground-truth label for  $r$  matching  $s$ .

## Faithful Encoding

**Definition 1** (Faithful encoding). An encoding  $\theta = (w, u)$  of SOIREs with length  $T$  is said to be faithful if it satisfies all the following conditions:

1.  $\forall 1 \leq t \leq T, w^t$  is a one-hot vector.
2.  $\forall 1 \leq t \leq T, u^t$  is either a one-hot vector or an all-zero vector.
3.  $\forall 1 \leq t \leq T, \sum_{t'=t+2}^T u_{t'}^t + \sum_{a \in \Sigma \cup \{?, *, +, \cdot, \&, |, \text{none}\}} w_a^t = 1$ .
4.  $\forall 1 \leq t \leq T-1, w_{none}^{t+1} - w_{none}^t \geq 0$ .
5.  $\forall 2 \leq t \leq T, \sum_{a \in \{?, *, +, \cdot, \&, | \}} w_a^{t-1} + \sum_{t'=1}^{t-2} u_{t'}^{t-1} + w_{none}^t = 1$ .
6.  $\forall 3 \leq t \leq T, \forall 1 \leq p \leq t-2, (t-1-p)u_t^p + \sum_{p'=p+1}^{t-1} u_{t'}^{p-1} \sum_{t'=t+1}^T u_{t'}^{p'} \leq t-1-p$ .
7.  $\forall a \in \Sigma, \sum_{t=1}^T w_a^t \leq 1$ .

## Interpretation

**Theorem 3.** Given a bounded size  $T \in \mathbb{Z}^+$ , prefix notations of SOIREs  $r$  with  $|r| \leq T$  and faithful encodings  $\theta$  with length  $T$  have a one-to-one correspondence, i.e.,  $Enc2Pre(\theta) = PreForm(r)$ .

Based on this correspondence, we apply beam search to find a faithful encoding nearby the learnt encoding and then interpret it to the target SOIRE.

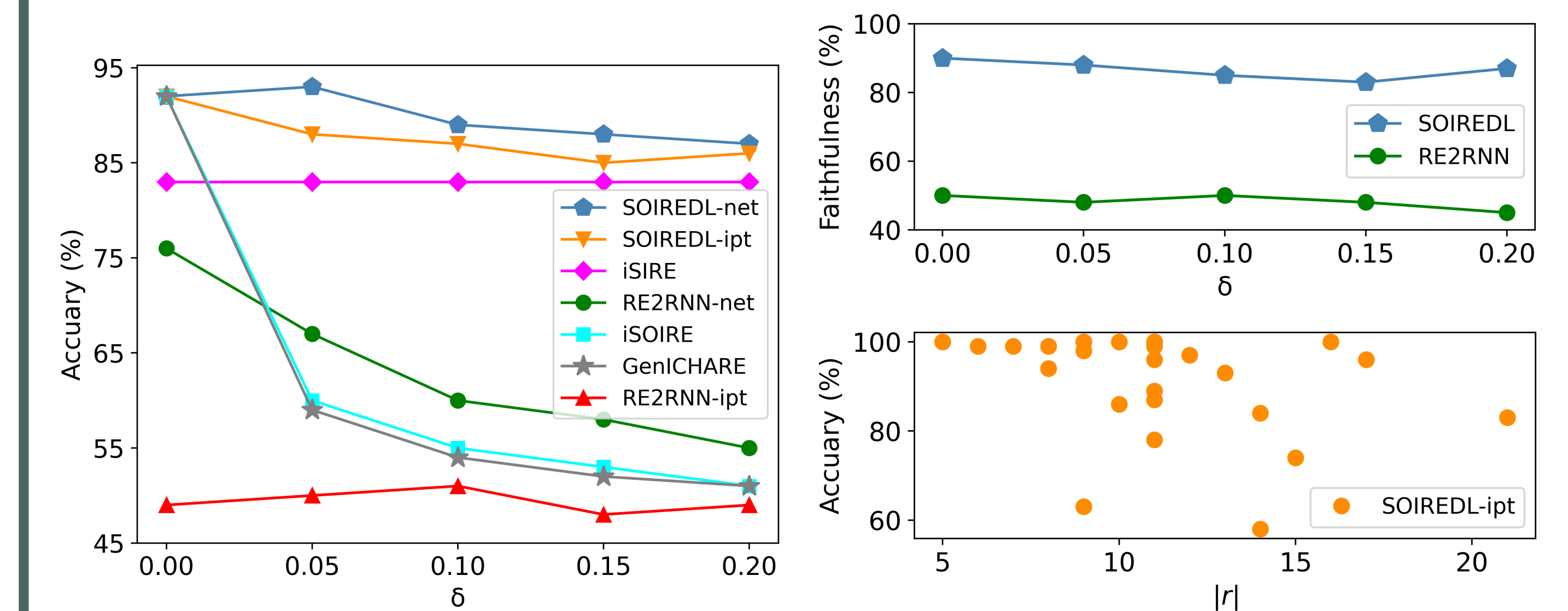
- The interpretation steps are conducted from bottom to top of the syntax tree. The score of each SOIRE is defined as the geometric mean of the probabilities of all operators and symbols.
- calculate the accuracy of each SOIRE in the last step on the training set and pick out the SOIRE with the highest accuracy.

## Result Analysis

**Competitor.**

- Positive strings only: iSOIRE (Li et al. 2019) for RSOIREs, GenICHARE (Zhang et al. 2018) for ICHAREs
- Both positive and negative strings: iSIRE (Li et al. 2020) for SIREs, RE2RNN (Jiang et al. 2020) for automaton, SOIREDL (Ours) for SOIREs

**Comparisons on Noisy Data, Faithfulness and Scalability.**



- SOIREDL gets the SOTA results and is the most robust on noisy data.
- The neural network of SOIREDL and its interpreted SOIRE are more consistent in performing SOIRE matching.
- The accuracy of SOIREDL decreases when the size of the ground-truth SOIRE increases.
  - This may be due to the difficulty for a neural network to capture the long-distance dependency in SOIRE matching.

## Acknowledgments

We thank Kunxun Qi for discussion on the paper and anonymous referees for helpful comments. This paper was supported by the National Natural Science Foundation of China (No. 62276284, 61976232, 61876204), the National Key Research and Development Program of China (No. 2021YFA1000504), Guangdong Basic and Applied Basic Research Foundation (No. 2022A1515011355), Guangzhou Science and Technology Project (No. 202201011699), Guizhou Science Support Project (No. 2022-259), as well as Humanities and Social Science Research Project of Ministry of Education (No. 18YJCZH006).