

# Recurrent Interaction Network for Jointly Extracting Entities and Classifying Relations

**Kai Sun**  
BDBC and SKLSDE  
Beihang University, China  
sunkai@buaa.edu.cn

**Richong Zhang\***  
BDBC and SKLSDE  
Beihang University, China  
zhangrc@act.buaa.edu.cn

**Samuel Mensah**  
BDBC and SKLSDE  
Beihang University, China  
samensah@buaa.edu.cn

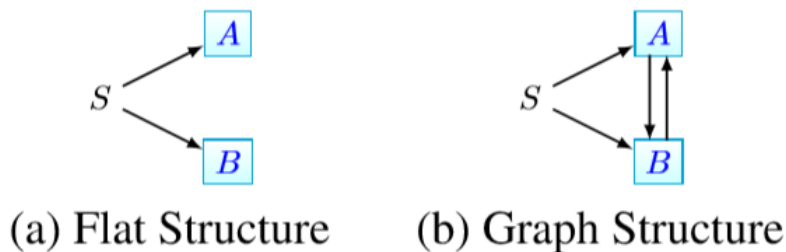
**Yongyi Mao**  
School of EECS  
University of Ottawa, Canada  
ymao@uottawa.ca

**Xudong Liu**  
BDBC and SKLSDE  
Beihang University, Beijing, China  
liuxd@act.buaa.edu.cn

EMNLP2020

这篇做的也是joint，但是感觉思路不错，跟其他千篇一律的不相同。

之前做joint可能是类似pipeline然后loss sum起来一起训，或者是multi-task learning的方式share底层编码结果。但是这样做的后果就是NER跟RE两个task是implicit的学习到对方的feature，这篇paper的核心就是我要explicit来做，我要让他们的feature直接交互。



思路是下面这样的。

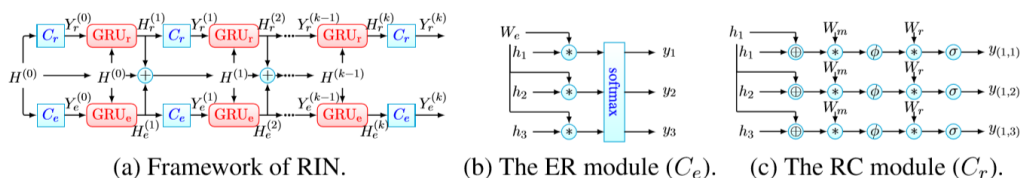


Figure 2: (a) The framework of RIN. (b) The entity recognition module. (c) The relation classification module. In (b) and (c), we use a toy example of shared features  $H = \{h_1, h_2, h_3\}$  to demonstrate the entity prediction for word  $w_i$  and relation prediction for all pairs  $(w_1, w_1), (w_1, w_2), (w_1, w_3)$ .  $+$ ,  $\oplus$ ,  $*$ ,  $\phi$ , and  $\sigma$  denote a summation operator, a concatenation operator, a matrix multiplication, relu activation function and sigmoid function respectively.

首先是通过encoder拿到整个句子的这个编码，相当于H是一个矩阵，每一列是

对应位置token的一个feature vector。

之后它被送进NER模块Ce以及RE模块Cr做预测，拿到了Y矩阵。

对于NER模块，它的Y矩阵每个列向量是预测出来的这个token在BIOES上的distributional probability。

对于RE模块，它的Y矩阵每个列向量是这个token与其他位置token的relation distribution的max pooling结果。

$$y_i = \text{maxpool} (Y_r(w_i)),$$

之后是各自通过GRU来进行更新，注意每个位置都是自己更新，跟其他位置无关，自己只管自己!!!

$$z = \sigma (W_z(h \oplus y))$$

$$u = \sigma (W_u(h \oplus y))$$

$$\check{h} = \tanh (W_o((u * h) \oplus y))$$

$$h_e = (1 - z) * h + z * \check{h}$$

最后就是交叉熵训练了。

奇怪的是好像只用了最后的预测结果来训，中间那些Y并没有给监督信号。

实验是在NYT跟WebNLG上做的，两种评价方式是partial match跟exact match。

partial match只要求relation跟每个entity的head对了就行。

exact match要求relation跟每个entity的head和tail都对。

Evaluation	Model	NYT			WebNLG		
		Prec	Rec	F1	Prec	Rec	F1
Partial Match	OneDecoder	59.4	53.1	56.0	32.2	28.9	30.5
	MultiDecoder	61.0	56.6	58.7	37.7	36.4	37.1
	OrderRL	77.9	67.2	72.1	63.3	59.9	61.6
	RIN <sub>w/o</sub> interaction	83.9 $\pm$ 0.6	83.1 $\pm$ 0.6	83.5 $\pm$ 0.2	84.9 $\pm$ 0.6	86.3 $\pm$ 0.8	85.6 $\pm$ 0.3
	RIN	<b>87.2<math>\pm</math>0.2</b>	<b>87.3<math>\pm</math>0.3</b>	<b>87.3<math>\pm</math>0.1</b>	<b>87.6<math>\pm</math>0.1</b>	<b>87.0<math>\pm</math>0.9</b>	<b>87.3<math>\pm</math>0.4</b>
Exact Match	NovelTagging	62.4	31.7	42.0	52.5	19.3	28.3
	GraphRel <sub>1p</sub>	62.9	57.3	60.0	42.3	39.2	40.7
	GraphRel <sub>2p</sub>	63.9	60.0	61.9	44.7	41.1	42.9
	CopyMLT-One	72.7	69.2	70.9	57.8	60.1	58.9
	CopyMLT-Mul	75.7	68.7	72.0	58.0	54.9	56.4
	RIN <sub>w/o</sub> interaction	77.4 $\pm$ 1.1	76.4 $\pm$ 0.7	76.9 $\pm$ 0.3	75.0 $\pm$ 1.1	73.3 $\pm$ 0.7	74.2 $\pm$ 0.3
	RIN	<b>83.9<math>\pm</math>0.5</b>	<b>85.5<math>\pm</math>0.5</b>	<b>84.7<math>\pm</math>0.4</b>	<b>77.3<math>\pm</math>0.7</b>	<b>76.8<math>\pm</math>1.0</b>	<b>77.0<math>\pm</math>0.2</b>

Table 1: Precision, Recall and F1 performance of different models on the datasets. Results for the compared models are retrieved from their original papers. We report the mean results over five runs and the standard deviation. The best performance is bold-typed.

注意有一个很重要的超参数是GRU的更新次数K，paper也做了探究。

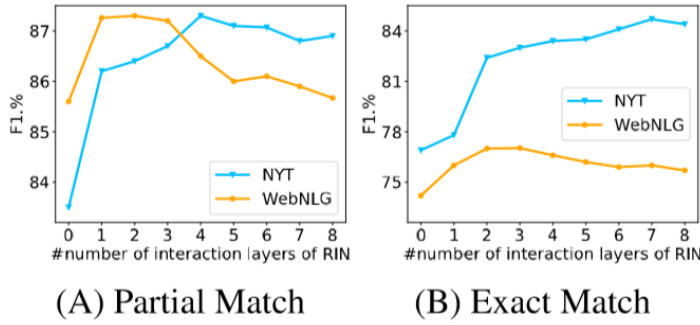


Figure 3: Curves of F1 performance on different number of interaction layers  $K$ .

和一些ablation study

Model	NYT	WebNLG
RIN	84.7	77.0
RIN <sub>w/o</sub> ER	83.9	76.4
RIN <sub>w/o</sub> RC	77.3	76.0
RIN <sub>w/o</sub> interaction	76.9	74.2
RIN <sub>w/o</sub> POS	84.1	76.6

Table 5: F1 performance of different ablation models on the datasets. The Exact Match evaluation is used.

- w/o ER指的是H的更新不考虑ER的GRU output结果
- w/o RC同上，不考虑RC的
- w/o 表示没有interactionn

- w/o POS指的是input没有使用pos embedding