



## Fine-grained Entity Typing via Label Reasoning

**Qing Liu<sup>1,3</sup>, Hongyu Lin<sup>1\*</sup>, Xinyan Xiao<sup>4</sup>, Xianpei Han<sup>1,2\*</sup>, Le Sun<sup>1,2</sup>, Hua Wu<sup>4</sup>**

<sup>1</sup>Chinese Information Processing Laboratory <sup>2</sup>State Key Laboratory of Computer Science

Institute of Software, Chinese Academy of Sciences, Beijing, China

<sup>3</sup>University of Chinese Academy of Sciences, Beijing, China

<sup>4</sup>Baidu Inc., Beijing, China

{liuqing2020, hongyu, xianpei, sunle}@iscas.ac.cn

{xiaoxinyan, wu\_hua}@baidu.com

Code: <https://github.com/loriqing/Label-Reasoning-Network>

—EMNLP2021

2022. 03. 17



gesis  
Leibniz-Institut  
für Sozialwissenschaften



Reported by ChangJiang Hu



# Introduction

What is Fine-grained entity typing(FET)

Given a **candidate entity**(mention) and its **context**,  
predict **a set of possible categories**(Type)

Context: "They were arrested by **FBI agents**."

Mention: **FBI agents**

Type: {organization, administration, force, agent, police}.

Case: Jack robs Mike, Jack is eventually caught.

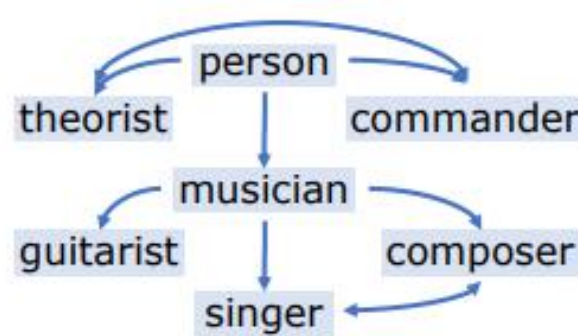
NER:Jack: {person},Mike: {person}

FET:Jack: {person, criminal},Mike: {person, victim}

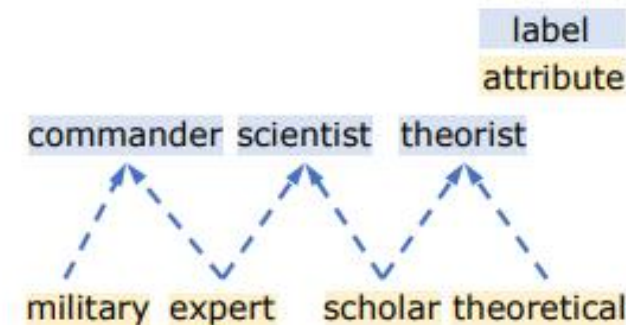
By providing fine-grained semantic labels, FET is **critical** for **entity recognition** and can **benefit** many NLP tasks, such as **relation extraction**, **entity linking** and **question answering**.

# Introduction

First Due to the massive label set, it is impossible to independently recognize each entity label without considering their dependencies

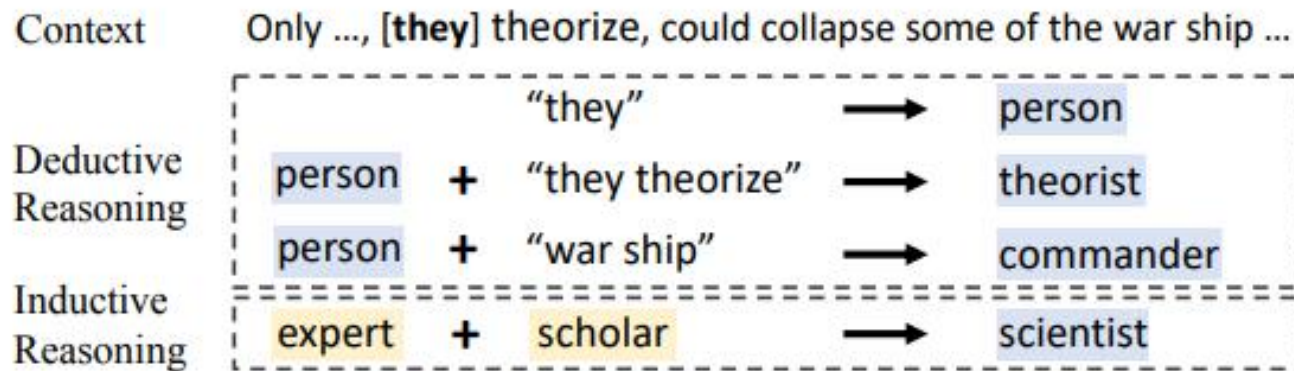


(a) Extrinsic dependency



(b) Intrinsic dependency

Second, because of the fine-grained and large-scale label set, many long tail labels are only provided with several or even no training instances



(c) Label reasoning process

# Method

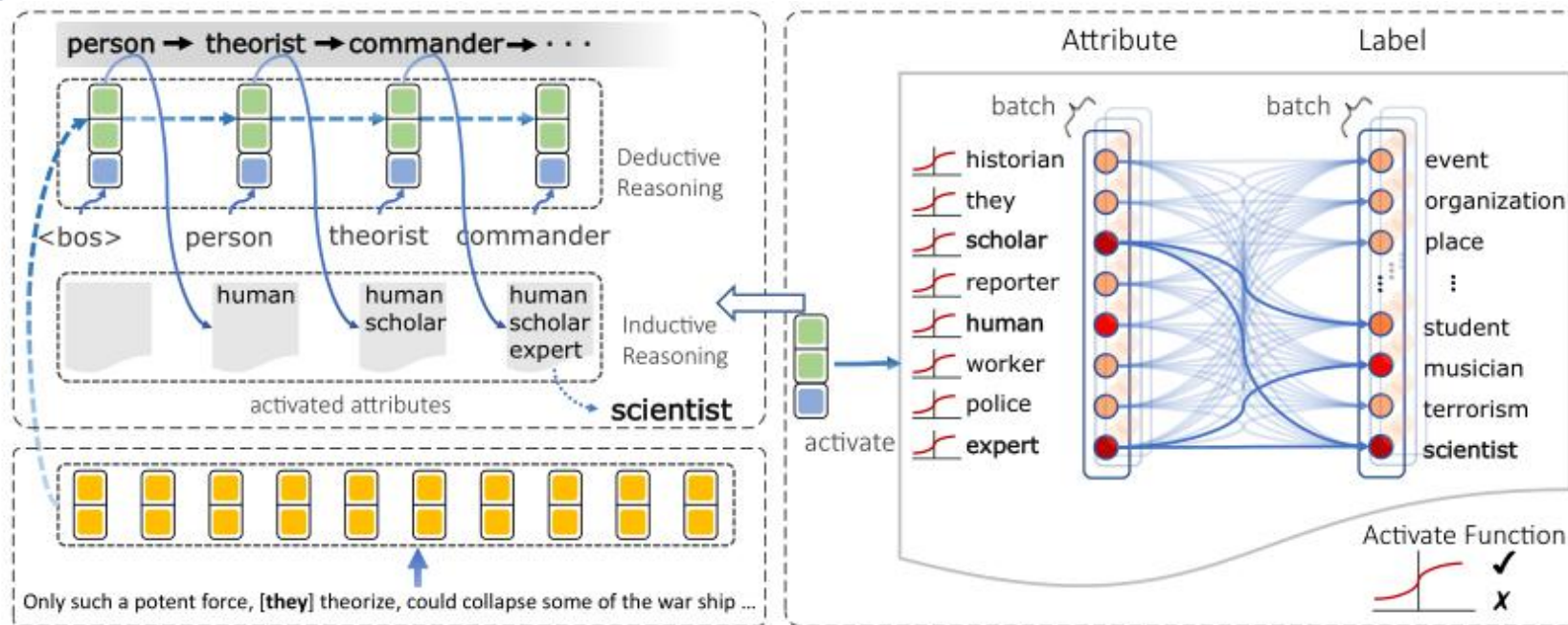
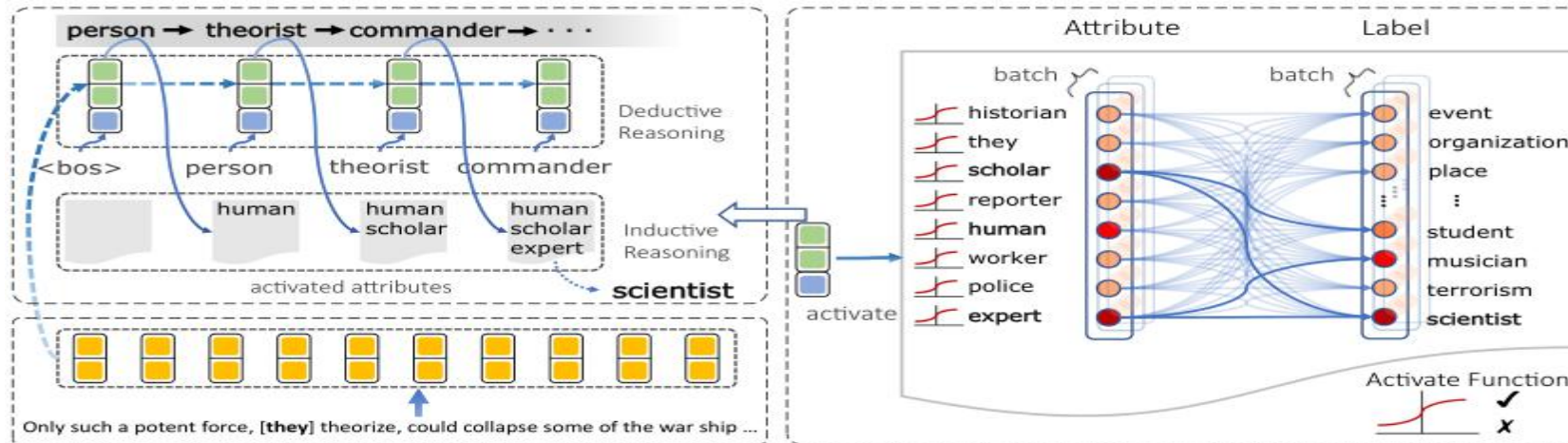


Figure 2: Overview of the process for LRN which contains an encoder, a deductive reasoning-based decoder and an inductive reasoning-based decoder. The figure shows: at step 1, the label *person* is predicted by deductive reasoning, and the attribute *human* is activated; at step 3, the label *scientist* is generated by inductive reasoning.

For encoding, we form the input instance  $X$  as “[CLS],  $x_1, \dots, [E_1], m_1, \dots, m_k, [E_2], \dots, x_n$ ” where  $[E_1], [E_2]$  are entity markers,  $m$  is mention word and  $x$  is context word. We then feed  $X$  to BERT and obtain the source hidden state  $\mathcal{H} = \{h_1, \dots, h_n\}$ . Finally, the hidden vector of [CLS] token is used as sentence embedding  $\mathbf{g}$

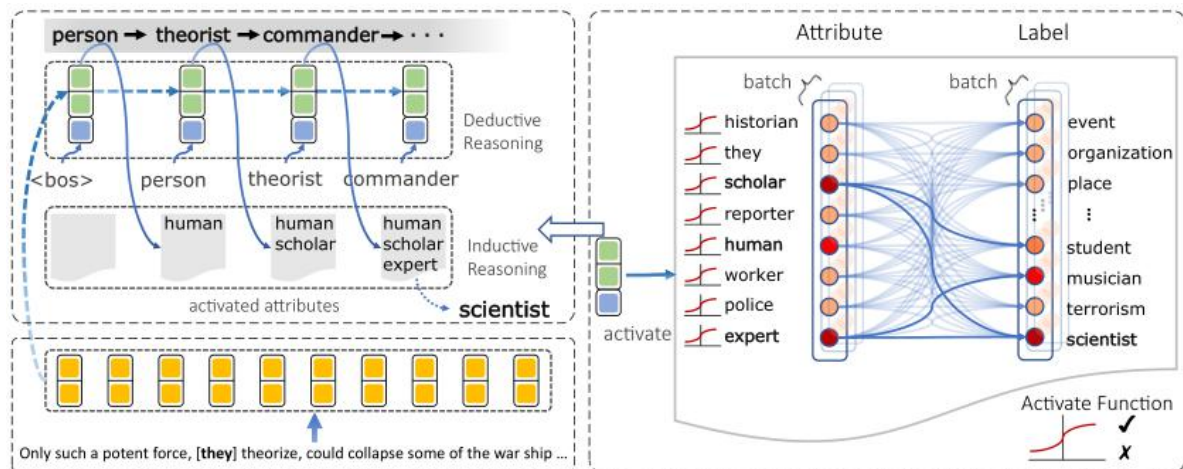
# Introduction

## Deductive Reasoning for Extrinsic Dependencies



Concretely, we utilize a **LSTM-based auto-regressive network** as decoder and obtain the hidden state of decoder  $S = \{s_0, \dots, s_k\}$ , where  $k$  is the number of predicted labels. We first initialize  $s_0$  using sentence embedding  $\mathbf{g}$ , then at each time step, **two attention mechanisms** – contextual attention and premise attention, are designed to capture context and label information for next prediction.

# Method



## Label Prediction

$$\mathbf{m}_t = [\mathbf{c}_t + \mathbf{g}; \mathbf{u}_t + \mathbf{s}_t]$$

$$\mathbf{o}_t = \mathbf{W}_o \mathbf{m}_t \quad (8)$$

$$\mathbf{y}_t = \text{softmax}(\mathbf{o}_t + \mathbf{I}_t) \quad (9)$$

$$(\mathbf{I}_t)_i = \begin{cases} -\text{inf} & , l_i \in \mathcal{Y}_{t-1}^* \\ 1 & , \text{otherwise} \end{cases} \quad (10)$$

$$\mathbf{s}_t = \text{LSTM}(\mathbf{s}_{t-1}, \mathbf{W}_b \mathbf{y}_{t-1})$$

$$(7) \quad \mathcal{S} = \{s_0, \dots, s_k\}$$

## Contextual Attention

$$e_{ti} = \mathbf{v}_c^T \tanh(\mathbf{W}_c \mathbf{s}_t + \mathbf{U}_c \mathbf{h}_i)$$

(1)

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{i=1}^n \exp(e_{ti})}$$

(2)

$$\mathbf{c}_t = \sum_{i=1}^n \alpha_{ti} \mathbf{h}_i \quad (3)$$

## Premise Attention

$$e_{tj} = \mathbf{v}_p^T \tanh(\mathbf{W}_p \mathbf{s}_t + \mathbf{U}_p \mathbf{s}_j)$$

(4)

$$\alpha_{tj} = \frac{\exp(e_{tj})}{\sum_{j=0}^{t-1} \exp(e_{tj})}$$

(5)

$$\mathbf{u}_t = \sum_{j=0}^{t-1} \alpha_{tj} \mathbf{s}_j \quad (6)$$

where  $\mathcal{Y}_{t-1}^*$  is the predicted labels before step  $t$  and  $l_i$  is the  $i^{\text{th}}$  label in label set  $L$ . The label with maximum value in  $\mathbf{y}_t$  is generated and used as the input for the next time step until  $[EOS]$  is generated.

# Method

## Inductive Reasoning for Intrinsic Dependencies

### BAG Construction

BAG  $\mathbf{g} = \{V, E\}$

V contain **attribute nodes**  $V_a$  and **label nodes**  $V_l$

In local BAG, we collect attributes in two ways:

(1) We **mask the entity mention** in the sentence, and predict the [MASK] token using masked language, and the **non-stop words** whose prediction scores **greater than** a confidence threshold  $\theta_c$  will be used as attributes — we denote them as **context attributes**.

(2) We **directly** segment the **entity mention** into words using Stanza2, and all **non-stop words** are used as attributes — we denote them as **entity attributes**.

Figure 3 shows several attribute examples. Given attributes, we compute the attribute-label relatedness (i.e. E in  $\mathbf{g}$ ) using the cosine similarity between their GloVe embeddings.

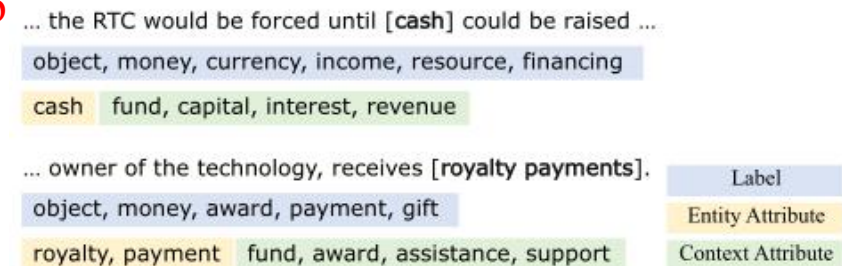
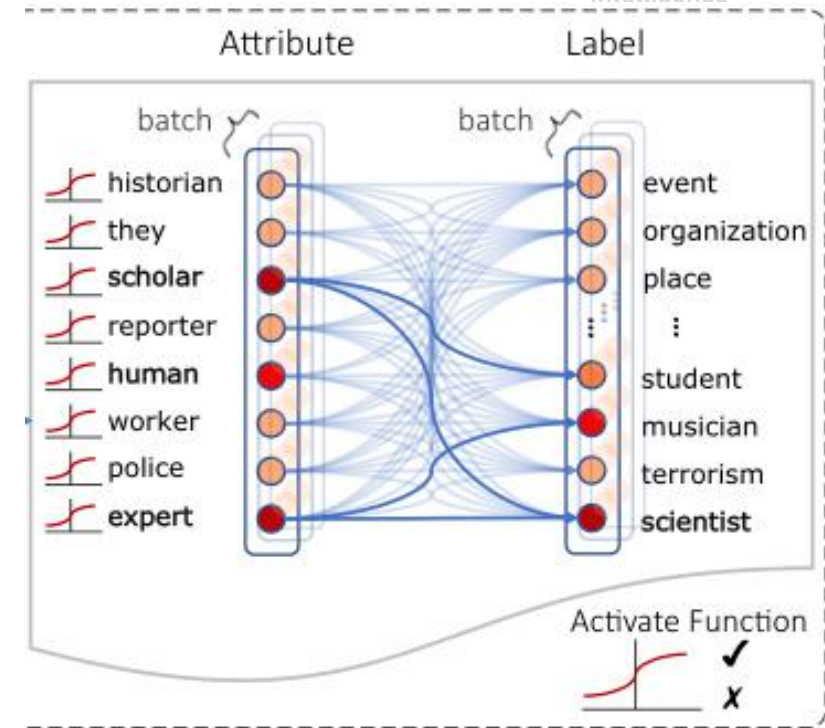
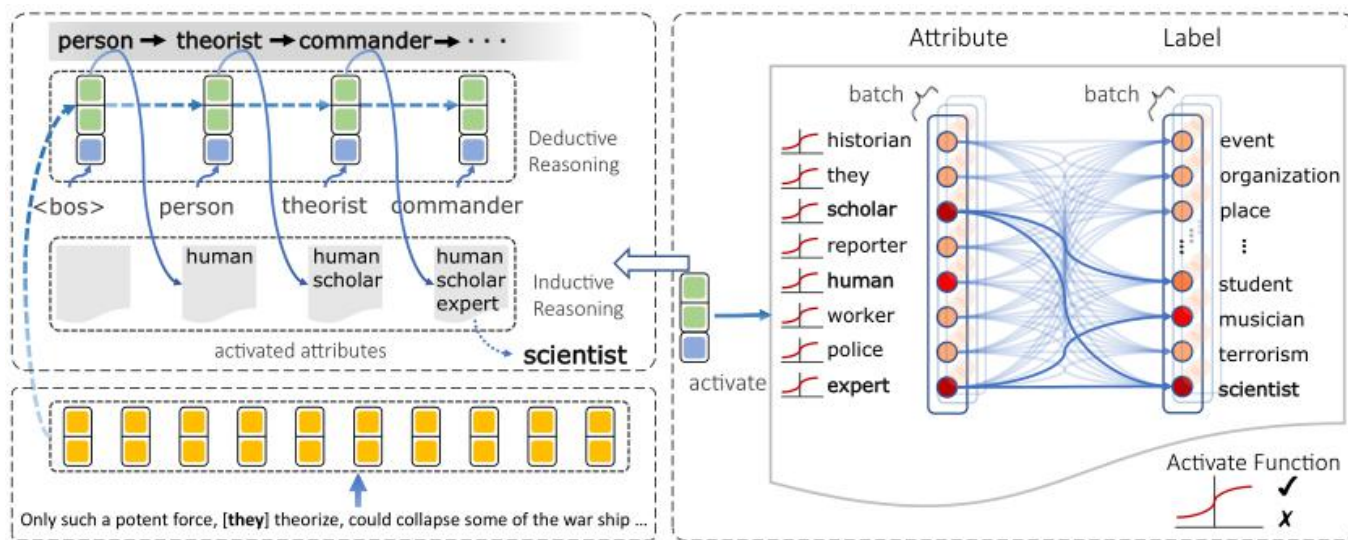


Figure 3: Examples of attributes.

# Method

Reasoning over BAG



$$\text{score}_{V_a}^{(i)} = \text{ReLU}(\text{sim}(\mathbf{W}_s \mathbf{s}_t, \mathbf{W}_a V_a^{(i)})) \quad (11)$$

$$\text{score}_{V_l}^{(j)} = \sum_{i=1}^{n_a} \text{score}_{V_a}^{(i)} E_{ij} \quad (12)$$

where  $n_a$  is the number of attributes,  $V_l^{(j)}$  is the  $j$ th label nodes and  $E_{ij}$  is the weight between them. Finally a label will be generated if its activation score is greater than a similarity threshold  $\theta_s$ .



# Method

## Learning

**Set Prediction Loss.** In FET, cross entropy loss is not appropriate because the prediction results is a label set, i.e.,  $\{y_1^*, y_2^*, y_3^*\}$  and  $\{y_3^*, y_2^*, y_1^*\}$  should have the same loss. Therefore we measure the similarity of two label set using the bipartite matching loss (Sui et al., 2020). Given the golden label set  $\mathcal{Y} = \{y_1, \dots, y_m\}$  and generated label set  $\mathcal{Y}^* = \{y_1^*, \dots, y_m^*\}$ , the matching loss  $\mathcal{L}(ij)_S$  of  $y_i$  and  $y_j^*$  is calculated by 13, then we use the Hungarian Algorithm (Kuhn, 1955) to get the specific order of golden label set as  $\tilde{\mathcal{Y}} = \{\tilde{y}_1, \dots, \tilde{y}_m\}$  to obtain minimum matching loss  $\mathcal{L}_S$ :

$$\mathcal{L}(ij)_S = \text{CE}(y_i, y_j^*) \quad (13)$$

$$\mathcal{L}_S = \text{CE}(\tilde{\mathcal{Y}}, \mathcal{Y}^*) \quad (14)$$

where CE is cross-entropy.

## Joint Entity and Relation Extraction with Set Prediction Networks

Dianbo Sui<sup>♡</sup> ♦ Yubo Chen<sup>♡</sup> Kang Liu<sup>♡</sup> ♦ Jun Zhao<sup>♡</sup> ♦  
Xiangrong Zeng ♦ Shengping Liu ♦

<sup>♡</sup> National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, China  
♦ University of Chinese Academy of Sciences, Beijing, China  
♦ Beijing Unisound Information Technology Co., Ltd, Beijing, China  
{dianbo.sui, yubo.chen, kliu, jzhao}@nlpr.ia.ac.cn, {zengxiangrong, liushengping}@unisound.com

$$\begin{aligned} \mathcal{L}(\mathbf{Y}, \hat{\mathbf{Y}}) = & \sum_{i=1}^m \{ -\log \mathbf{P}_{\pi^*(i)}^r(r_i) \\ & + \mathbb{1}_{\{r_i \neq \emptyset\}} [ -\log \mathbf{P}_{\pi^*(i)}^{s-start}(s_i^{start}) \\ & - \log \mathbf{P}_{\pi^*(i)}^{s-end}(s_i^{end}) \\ & - \log \mathbf{P}_{\pi^*(i)}^{o-start}(o_i^{start}) \\ & - \log \mathbf{P}_{\pi^*(i)}^{o-end}(o_i^{end}) ] \} \end{aligned}$$



# Experiments

**BAG Loss.** To make the model activate labels correctly, we add a supervisory loss to the bipartite attribute graph to active correct labels:

$$\mathcal{L}_A = - \sum_{j=1}^{|L|} score_{V_i}^{(j)} * y_j \quad (15)$$

$$y_j = \begin{cases} 1 & , v_j \in \mathcal{Y} \\ -1 & , v_j \notin \mathcal{Y} \end{cases} \quad (16)$$

**Final Loss.** The final loss is a combination of set loss and BAG loss:

$$\mathcal{L} = \mathcal{L}_S + \lambda \mathcal{L}_A \quad (17)$$



# Experiments

Model	P	R	F1
without label dependency			
*Choi et al. (2018)	47.1	24.2	32.0
*ELMo(Onoe and Durrett, 2019)	51.5	33.0	40.2
BERT(Onoe and Durrett, 2019)	51.6	33.0	40.2
BERT[in-house]	55.9	33.0	41.5
with label dependency			
*LABELGCN (Xiong et al., 2019)	50.3	29.2	36.9
LRN w/o IR	<b>61.2</b>	33.5	43.3
LRN	54.5	<b>38.9</b>	<b>45.4</b>

Table 1: Macro P/R/F1 results on Ultra-Fine test set. \* means using augmented data. "without label dependency" methods formulated FET as multi-label classification without considering associations between labels. "with label dependency" methods leveraged associations between labels explicitly or implicitly.

# Experiments

Model	Total			General			Fine			Ultra-Fine		
	P	R	F	P	R	F	P	R	F	P	R	F
*Choi et al. (2018)	48.1	23.2	31.3	60.3	61.6	61.0	40.4	38.4	39.4	42.8	8.8	14.6
†LABELGCN (Xiong et al., 2019)	49.3	28.1	35.8	66.2	68.8	67.5	43.9	40.7	42.2	42.4	14.2	21.3
HY Large (López and Strube, 2020)	43.4	34.2	38.2	61.4	73.9	67.1	35.7	46.6	40.4	36.5	19.9	25.7
*ELMo (Onoe and Durrett, 2019)	50.7	33.1	40.1	66.9	<b>80.7</b>	73.2	41.7	46.2	43.8	45.6	17.4	25.2
BERT (Onoe and Durrett, 2019)	51.6	32.8	40.1	67.4	80.6	73.4	41.6	54.7	47.3	46.3	15.6	23.4
BERT[in-house]	54.1	32.1	40.3	68.8	79.2	73.6	43.8	<b>57.4</b>	49.7	<b>50.7</b>	14.6	22.6
LRN w/o IR	<b>60.7</b>	32.5	42.3	<b>79.3</b>	75.5	<b>77.4</b>	<b>59.6</b>	44.8	51.2	45.7	18.7	26.5
LRN	53.7	<b>38.6</b>	<b>44.9</b>	77.8	76.4	77.1	55.8	50.6	<b>53.0</b>	43.4	<b>26.0</b>	<b>32.5</b>

Table 2: Macro P/R/F1 of each label granularity on Ultra-Fine dev set, and long tail labels are mostly in the ultra-fine layer. \* means using augmented data. † We adapt the results from López and Strube (2020).

Model	Total			General			Fine			Ultra-Fine		
	P	R	F	P	R	F	P	R	F	P	R	F
HY XLarge (López and Strube, 2020)	/	/	/	/	/	69.1	/	/	39.7	/	/	26.1
BERT[in-house]	55.9	33.0	41.5	69.7	<b>81.6</b>	75.2	43.7	<b>56.0</b>	49.1	<b>53.5</b>	15.5	24.0
LRN w/o IR	<b>61.2</b>	33.5	43.3	<b>78.3</b>	76.7	<b>77.5</b>	<b>61.6</b>	44.1	51.4	47.8	19.9	28.1
LRN	54.5	<b>38.9</b>	<b>45.4</b>	77.4	76.7	77.1	58.4	50.4	<b>54.1</b>	43.5	<b>26.4</b>	<b>32.8</b>

Table 3: Macro P/R/F1 of different label granularity on Ultra-Fine test set.



# Experiments

Number of	Category	Prediction	Shot=0			Shot=1			Shot=2		
			Correct	Predicted	Prec.	Correct	Predicted	Prec.	Correct	Predicted	Prec.
BERT[in-house]	293	5683	0	0	/	1	1	100.0%	9	66	13.6%
LRN w/o IR	330	5740	0	0	/	1	3	33.3%	15	28	53.6%
LRN	997	7808	110	218	50.5%	67	252	26.6%	94	276	34.1%

Table 4: Performance of the zero-shot, shot=1 and shot=2 label prediction. "Category" means how many kinds of types are predicted. "Prediction" means how many labels are generated.

# Experiments

Model	Dev			Test		
	P	R	F	P	R	F
<b>LRN</b>	53.7	38.6	44.9	54.5	38.9	45.4
-PreAtt	53.1	39.3	45.2	52.6	39.5	45.1
-PreAtt-ConAtt	56.3	36.3	44.2	56.4	36.5	44.3
-SetLoss	46.8	40.7	43.5	47.8	40.7	44.0
<b>LRN w/o IR</b>	60.7	32.5	42.3	61.2	33.5	43.3
-PreAtt	54.5	34.2	42.1	55.1	35.0	42.8
-PreAtt-ConAtt	55.2	32.9	41.3	56.2	34.3	42.6
-SetLoss	46.0	37.6	41.4	46.6	37.5	41.6

Table 5: Ablation results on Ultra-Fine dataset: PreAtt denotes premise attention, ConAtt denotes contextual attention, and -SetLoss denotes replacing set prediction loss with cross-entropy loss.

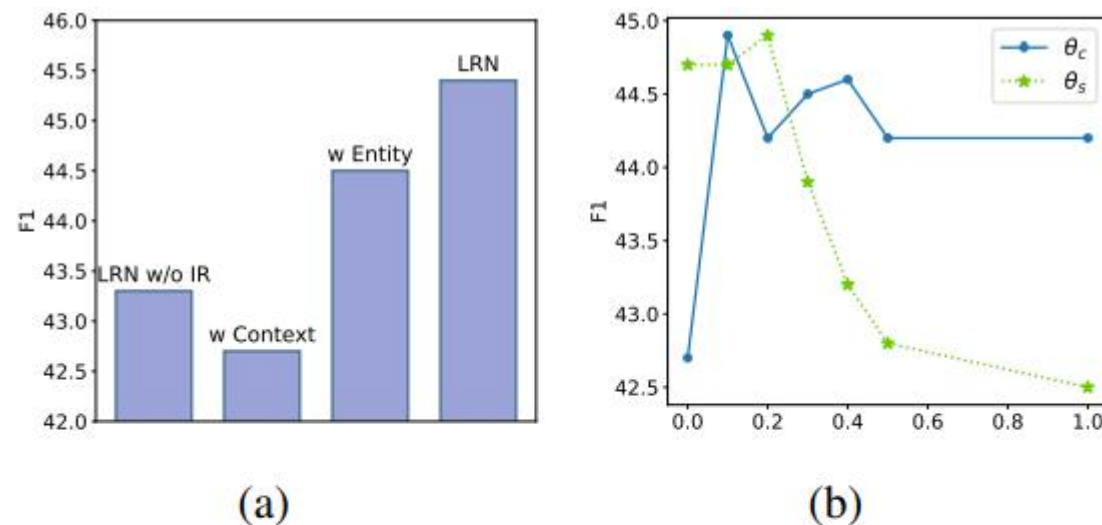


Figure 4: (a) Ablation experiments of context attributes and entity attributes on Ultra-Fine dataset. (b) Performances of different confidence threshold  $\theta_c$  and similarity threshold  $\theta_s$  on dev set.



# Experiments

Encoder	Model	Acc	MaF	MiF
<b>with augmentation</b>				
HYPER	López and Strube (2020)	47.4	75.8	69.4
LSTM	Choi et al. (2018)	59.5	76.8	71.8
	Xiong et al. (2019)	59.6	77.8	72.2
ELMo	*Onoe and Durrett (2019)	64.9	84.5	79.2
	(Lin and Ji, 2019)	63.8	82.9	77.3
BERT	Wang et al. (2020)	61.1	81.8	76.3
	BERT [in-house]	62.2	83.4	78.8
	LRN w/o IR	<b>66.1</b>	<b>84.8</b>	<b>80.1</b>
	LRN	64.5	84.5	79.3
<b>without augmentation</b>				
ELMo	*Onoe and Durrett (2019)	42.7	72.7	66.7
	Chen et al. (2020)	<b>58.7</b>	73.0	68.1
BERT	Onoe and Durrett (2019)	51.8	76.6	69.1
	BERT[in-house]	51.5	76.6	69.7
	LRN w/o IR	55.3	77.3	70.4
	LRN	56.6	<b>77.6</b>	<b>71.8</b>

Table 6: Results on OntoNotes test set. Augmentation is the augmented data created by (Choi et al., 2018) which contains 800K instances and therefore there're little few-shot labels in this setting. And \* indicates using additional features to enhance the label representation.



**Thank you!**