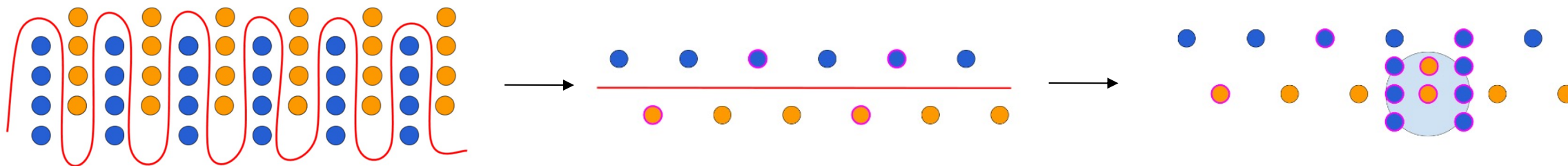
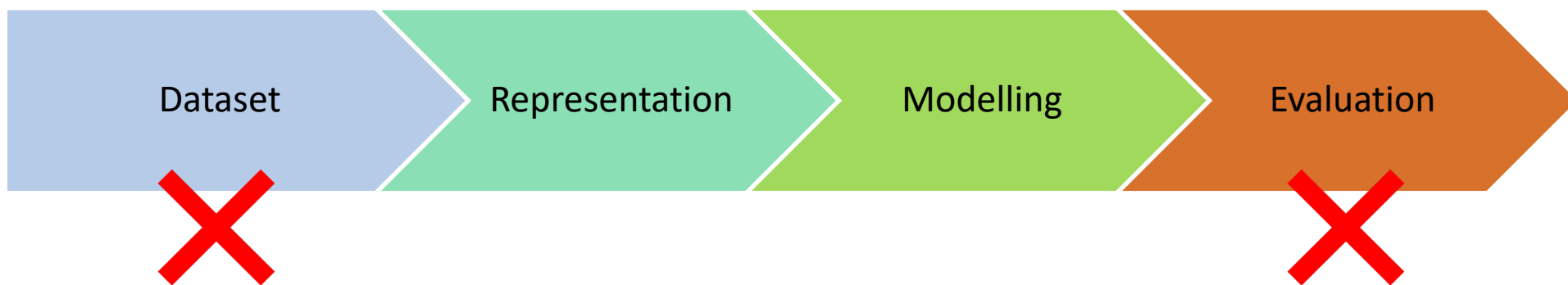


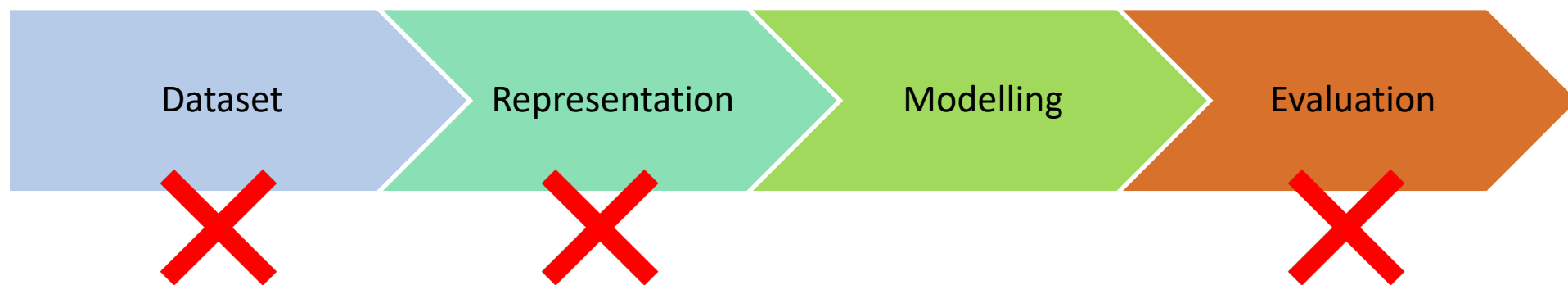


Robustness Improvement in Natural Language Processing

Tao Gui

Fudan University





ORIGINAL

The government made a quick decision

BAE - R

The **MASK** made a quick decision

judge , doctor , captain

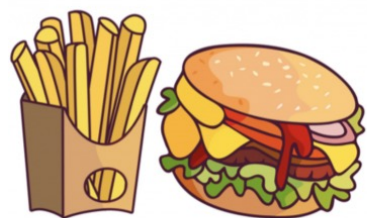
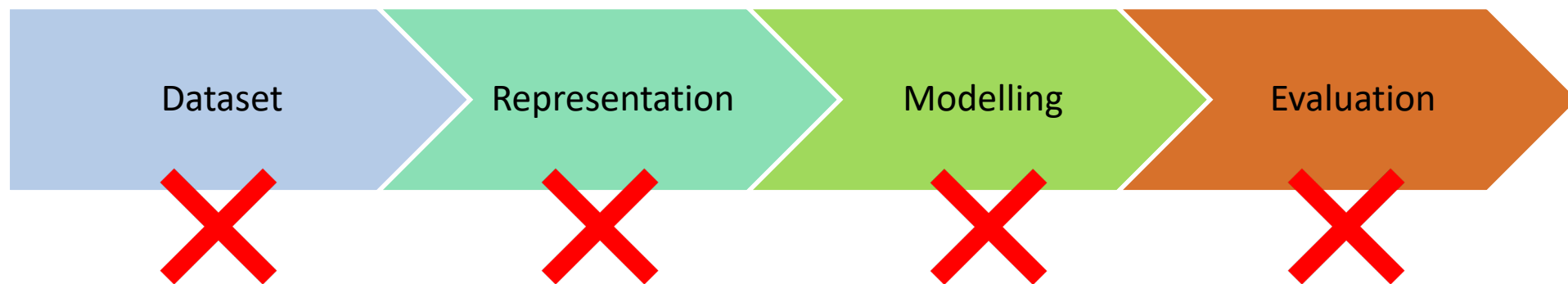
BAE - I

The **MASK** government made a quick decision

state , british , federal

The government **MASK** made a quick decision

officials , then , immediately



Sentiment Analysis Data

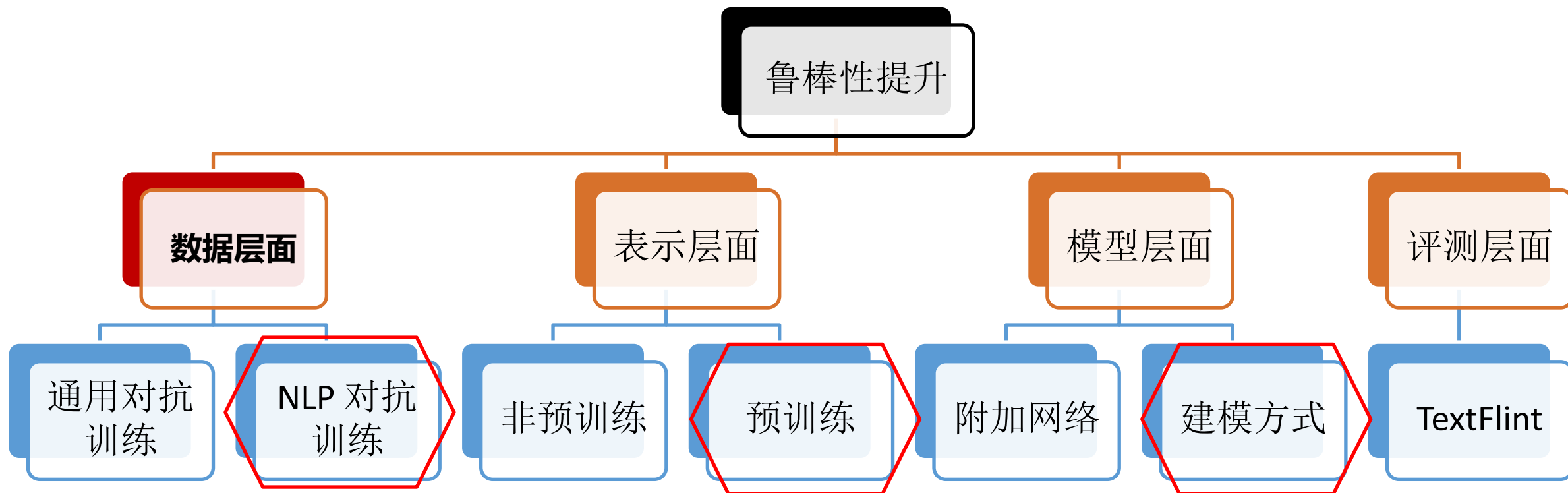
Tasty **burgers**, and crispy **fries**.

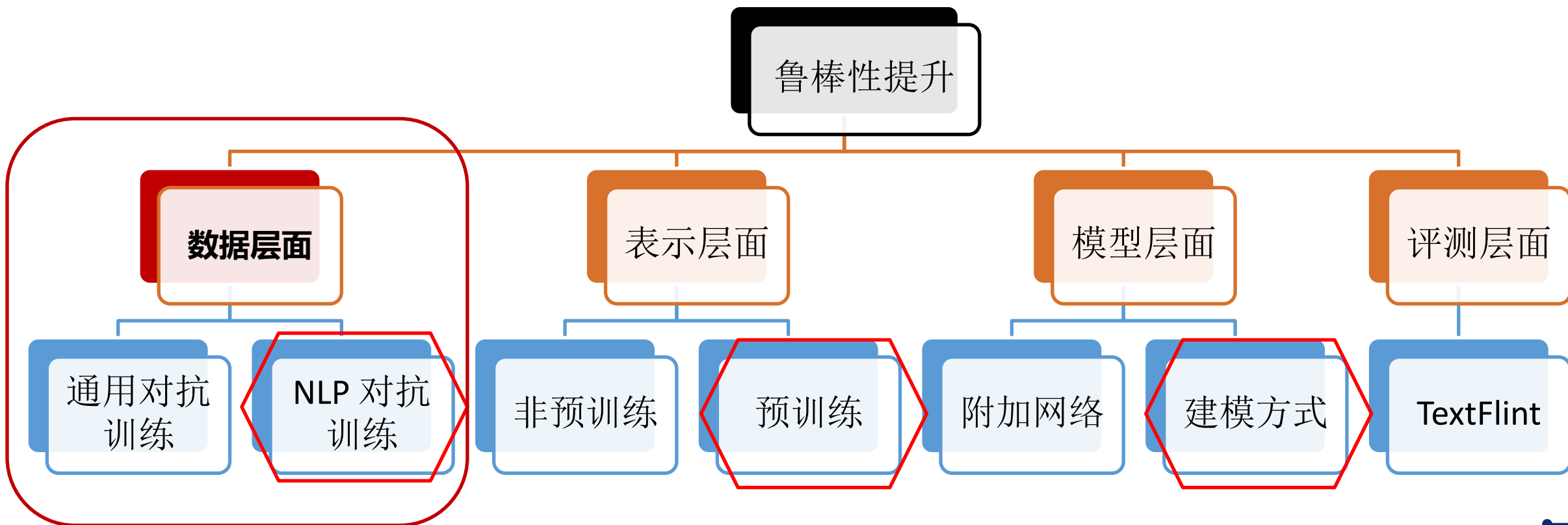
burgers 😊 **fries** 😊 **SA** 😊

Positive/Negative/Neutral



(3) Aspect Sentiment Classification







$$\begin{aligned} \mathbf{x}' &= \mathbf{x} + \eta, f(\mathbf{x}) = y, \mathbf{x} \in \mathbf{X} \\ f(\mathbf{x}') &\neq y \\ \text{or } f(\mathbf{x}') &= y', y' \neq y \end{aligned}$$

$$\min_{\theta} \mathbb{E}_{(\mathbf{Z}, y) \sim \mathcal{D}} \left[\max_{\|\delta\| \leq \epsilon} L(f_{\theta}(\mathbf{X} + \delta), y) \right],$$



x
“panda”
57.7% confidence

+ .007 ×

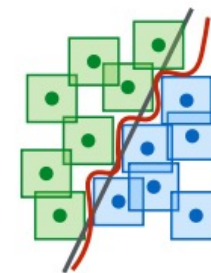
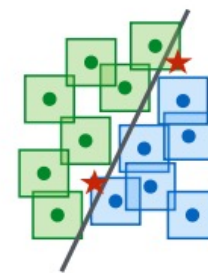
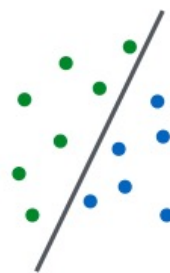


$\text{sign}(\nabla_x J(\theta, x, y))$
“nematode”
8.2% confidence

=



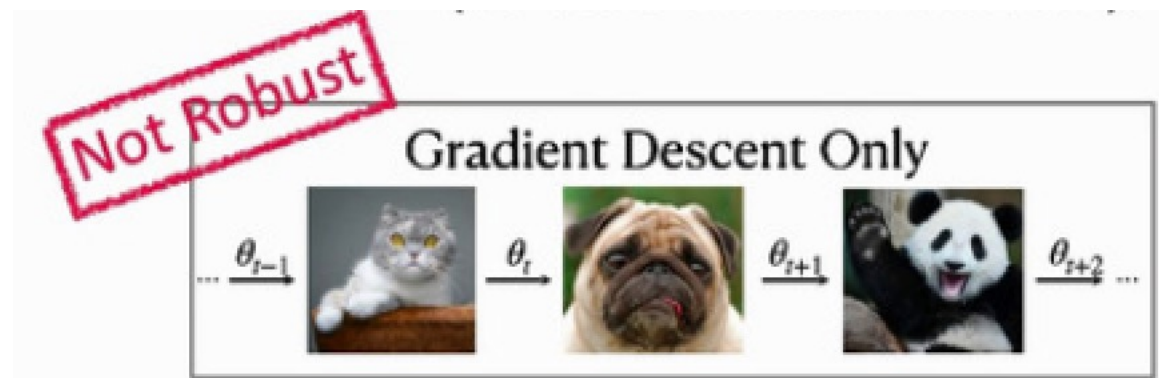
$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
“gibbon”
99.3 % confidence



Adversarial Attack



Adversarial Training



$$\text{FGSM:} \quad \delta = \epsilon \cdot \text{Sign}(g)$$

$$\text{FGM:} \quad \delta = \epsilon \cdot (g / \|g\|_2)$$

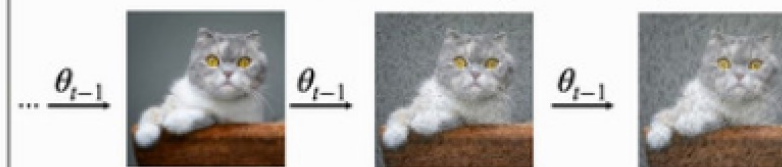


Not Robust

Gradient Descent Only



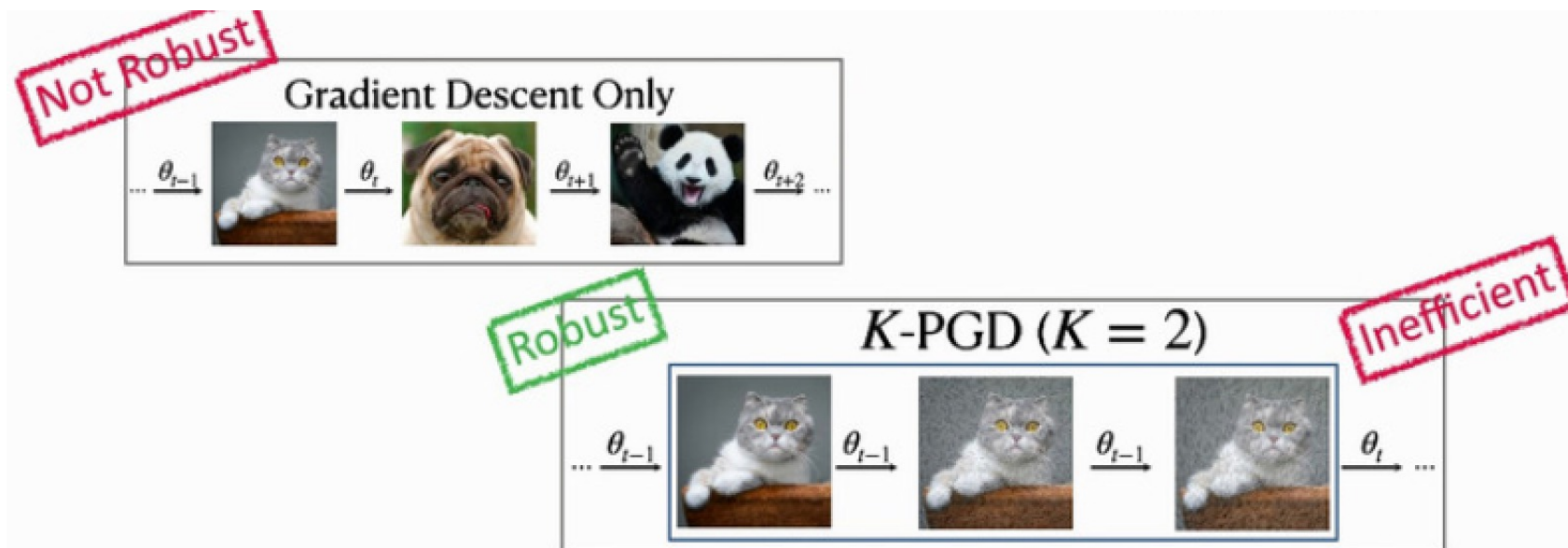
K -PGD ($K = 2$)



$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} \left[\max_{\delta \in \mathcal{S}} L(x + \delta, y; \theta) \right]$$

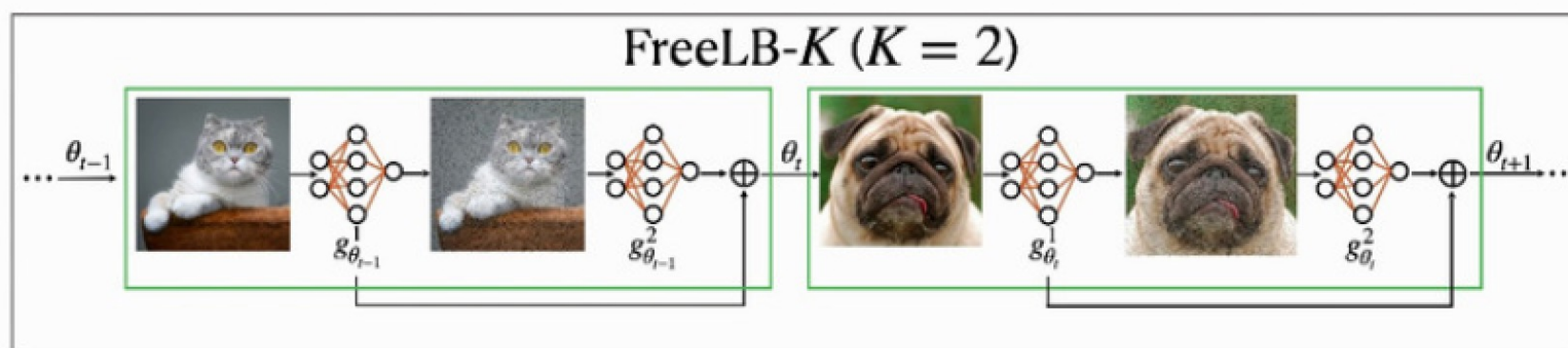
K steps of projected gradient ascent for δ

Followed by one gradient descent on θ



N parameter update takes:

- $(K + 1)N$ forward-backward passes
- KN input updates



Equivalent to enlarging the batch size with adversarial examples

$$\min_{\theta} \mathbb{E}_{(Z,y) \sim \mathcal{D}} \left[\frac{1}{K} \sum_{t=0}^{K-1} \max_{\delta_t \in \mathcal{F}_t} L(f_{\theta}(X + \delta_t), y; \theta) \right]$$

Improved invariance with a set of different transforms

Better generalization (Sokolic et al. 2017, Gong et al. 2020)



AT in Image

Table 1: Validation accuracy and robustness of CIFAR-10 models trained with various methods.

Training	Evaluated Against					Train Time (min)
	Nat. Images	PGD-20	PGD-100	CW-100	10 restart PGD-20	
Natural	95.01%	0.00%	0.00%	0.00%	0.00%	780
Free $m = 2$	91.45%	33.92%	33.20%	34.57%	33.41%	816
Free $m = 4$	87.83%	41.15%	40.35%	41.96%	40.73%	800
Free $m = 8$	85.96%	46.82%	46.19%	46.60%	46.33%	785
Free $m = 10$	83.94%	46.31%	45.79%	45.86%	45.94%	785
7-PGD trained	87.25%	45.84%	45.29%	46.52%	45.53%	5418

Robustness **improves**
Accuracy **decreases**

AT in Text

Experiment results of different defenders on AGNEWS

Method	Clean%	TextFooler		
		Aua%	Suc%	#Query
Baseline (BERT)	94.5	19.1	79.6	317.4
Adversarial Data Augmentation	94.4	38.6	58.9	404.6
MixADA (Si et al., 2020)	94.3	37.5	60.3	410.7
PGD-K (Madry et al., 2018)	94.7	24.8	73.9	353.5
FreeLB (Zhu et al., 2020)	94.7	31.6	66.7	382.1
TA-VAT (Li and Qiu, 2020)	94.8	31.0	67.3	382.5
InfoBERT (Wang et al., 2020)	95.1	31.8	66.5	369.9
DNE (Zhou et al., 2020)	93.9	28.7	69.8	367.9
ASCC (Dong et al., 2021)	92.3	28.2	69.6	326.5
SAFER (Ye et al., 2020)	94.3	31.8	66.1	350.1
RanMASK (Zeng et al., 2021)	91.7	37.9	58.7	583.4

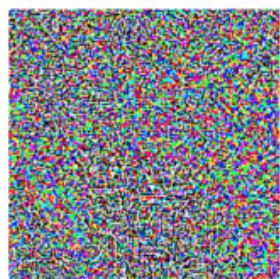
Robustness **improves**
Accuracy **improves**

Problem in current Virtual Adversarial Train


 x

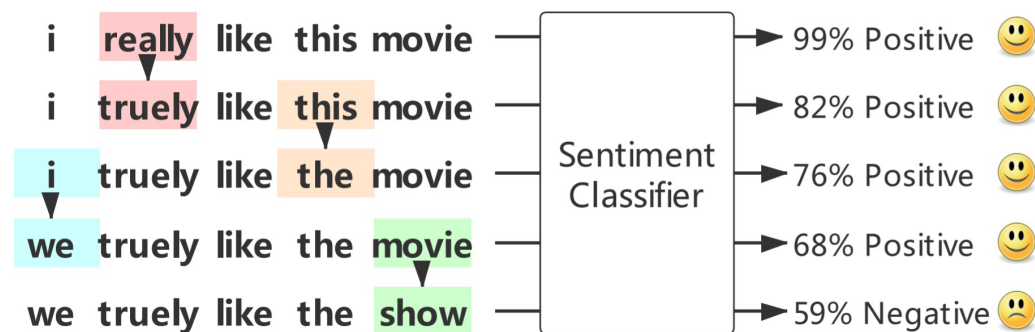
“panda”

57.7% confidence

 $+ .007 \times$

 $\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence





Effectiveness of Adversarial Training

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\|\delta\| \leq \epsilon} \mathcal{L}(f_{\theta}(\mathbf{X} + \delta), y) \right], \quad (3)$$

$$\delta_{t+1} = \prod_{\|\delta\|_F \leq \epsilon} \left(\delta_t + \alpha \frac{g(\delta_t)}{\|g(\delta_t)\|_F} \right), \quad (4)$$

$$\|\delta_{t+1}\| \leq \|\delta_t\| + \left\| \alpha \frac{g(\delta_t)}{\|g(\delta_t)\|_F} \right\| \leq \|\delta_t\| + \alpha. \quad (5)$$

$$\begin{aligned} \|\delta_t\| &\leq \|\delta_{t-1}\| + \alpha \leq \|\delta_{t-2}\| + 2\alpha \\ &\leq \dots \leq \|\delta_1\| + (t-1) * \alpha \leq t\alpha, \end{aligned} \quad (6)$$

Method	Norm ϵ	Clean%	TextFooler		BERT-Attack	
			Aua%	Suc%	Aua%	Suc%
PGD-K	0.01	94.9	21.8	77.0	31.5	66.8
	0.1	95.3	43.6	54.3	45.1	52.7
	1	95.2	45.2	55.3	45.3	52.4
	w/o	95.2	45.2	55.3	45.3	52.4
FreeLB	0.01	95.4	30.5	68.0	43.6	54.3
	0.1	95.5	36.1	62.2	40.0	58.1
	1	94.9	45.8	51.7	42.5	55.2
	w/o	94.9	45.8	51.7	42.5	55.2



Experiments

Method	Clean%	TextFooler			TextBugger			BERT-Attack		
		Aua%	Suc%	#Query	Aua%	Suc%	#Query	Aua%	Suc%	#Query
Baseline (BERT)	94.5	19.1	79.6	317.4	23.5	75.0	320.6	27.2	71.0	338.8
Adversarial Data Augmentation	94.4	38.6	58.9	404.6	43.3	53.9	418.3	42.9	54.5	407.0
MixADA (Si et al., 2020)	94.3	37.5	60.3	410.7	36.4	61.4	423.5	39.1	58.6	408.4
PGD-K (Madry et al., 2018)	94.7	24.8	73.9	353.5	26.7	71.9	367.1	39.4	58.5	399.3
FreeLB (Zhu et al., 2020)	94.7	31.6	66.7	382.1	32.9	65.4	390.6	43.9	53.8	417.1
TA-VAT (Li and Qiu, 2020)	94.8	31.0	67.3	382.5	34.2	63.9	415.2	45.0	52.5	436.9
InfoBERT (Wang et al., 2020)	95.1	31.8	66.5	369.9	36.3	61.8	391.6	42.4	55.3	392.8
DNE (Zhou et al., 2020)	93.9	28.7	69.8	367.9	28.2	70.3	377.6	42.4	55.5	470.1
ASCC (Dong et al., 2021)	92.3	28.2	69.6	326.5	37.0	60.1	307.4	32.7	64.7	337.1
SAFER (Ye et al., 2020)	94.3	31.8	66.1	350.1	41.2	56.1	398.8	39.3	58.2	373.5
RanMASK (Zeng et al., 2021)	91.7	37.9	58.7	583.4	45.0	50.9	626.8	49.5	46.1	661.8
FreeLB++	95.1	51.5	46.0	419.1	55.9	41.4	416.9	41.8	56.2	386.1

Table 2: The experiment results of different defenders on AGNEWS, where all models are trained on BERT. The best performance is marked in **bold**. FreeLB++ not only achieves best defense performance under both TextBugger and TextFooler, but also improves **Clean%**. Although RanMASK has also achieved significant defense performance, it drops a lot in **Clean%**.



Problem in current Virtual Adversarial Train

- 1. Randomization Problem
- random initialize at step zero

At step t :

generate perturbations

$$\delta_{t+1} = \prod_{\|\delta_t\|_F \leq \epsilon} \frac{(\delta_t + \alpha g(\delta_t))}{\|g(\delta_t)\|_F} \quad (2)$$

constrain using L2norm

- 2. Constraint Problem
- Frobenius Norm on X

$$g(\delta_t) = \nabla_{\delta} L(f_{\theta}(X + \delta_t), y) \quad (3)$$

- $X = [v_0, v_1, \dots, v_n \dots]$ is output of a sequence



Method

- 1. Global Perturbation Vocabulary
- 2. Token-Level Constraint





1. Perturbation Generation

Algorithm 1 Token-Aware Virtual Adversarial Training

Require: Training Samples $S = \{(X = [w_0, \dots, w_i, \dots], y)\}$, perturbation bound ϵ , initialize bound σ adversarial steps K , adversarial step size α , model parameter θ

```

1:  $V \in \mathbb{R}^{N \times D} \leftarrow \frac{1}{\sqrt{D}} U(-\sigma, \sigma)$  // Initialize perturbation vocabulary  $V$ 
2: for epoch = 1,  $\dots$ , do
3:   for batch  $B \subset S$  do
4:      $\delta_0 \leftarrow \frac{1}{\sqrt{D}} U(-\sigma, \sigma)$ ,  $\eta_0^i \leftarrow V[w_i]$ ,  $g_0 \leftarrow 0$  //Initialize perturbation and gradient of  $\theta$ 
5:     for  $t = 1, \dots, K$  do
6:        $g_t \leftarrow g_{t-1} + \frac{1}{K} \mathbb{E}_{(X,y) \in B} [\nabla_{\theta} L(f_{\theta}(X + \delta_{t-1} + \eta_{t-1}), y)]$  //Accumulate gradients of  $\theta$ 
7:       Update token-level perturbation  $\eta$ :
8:        $g_{\eta}^i \leftarrow \nabla_{\eta^i} L(f_{\theta}((X + \delta_{t-1} + \eta_{t-1}), y))$ 
9:        $\eta_t^i \leftarrow \eta_{t-1}^i + \alpha \cdot g_{\eta}^i / \|g_{\eta}^i\|_F$ 
10:       $\eta_t \leftarrow \prod_{\|\eta\|_F < \epsilon} (\eta_t)$ 
11:      Update instance-level perturbation  $\delta$ :
12:       $g_{\delta} \leftarrow \nabla_{\delta} L(f_{\theta}((X + \delta_{t-1} + \eta_{t-1}), y))$ 
13:       $\delta_t \leftarrow \prod_{\|\delta\|_F < \epsilon} (\delta_{t-1} + \alpha \cdot g_{\delta} / \|g_{\delta}\|_F)$ 
14:    end for
15:     $V[w_i] \leftarrow \eta_K^i$  //Update perturbation vocabulary  $V$ 
16:     $\theta \leftarrow \theta - g_K$  //Update model parameter  $\theta$ 
17:  end for
18: end for

```



2. Global Perturbation Vocabulary

Algorithm 1 Token-Aware Virtual Adversarial Training

Require: Training Samples $S = \{(X = [w_0, \dots, w_i, \dots], y)\}$, perturbation bound ϵ , initialize bound σ adversarial steps K , adversarial step size α , model parameter θ

```

1:  $V \in \mathbb{R}^{N \times D} \leftarrow \frac{1}{\sqrt{D}} U(-\sigma, \sigma)$  // Initialize perturbation vocabulary  $V$ 
2: for epoch = 1, ..., do
3:   for batch  $B \subset S$  do
4:      $\delta_0 \leftarrow \frac{1}{\sqrt{D}} U(-\sigma, \sigma)$ ,  $\eta_0^i \leftarrow V[w_i]$ ,  $g_0 \leftarrow 0$  // Initialize perturbation and gradient of  $\theta$ 
5:     for  $t = 1, \dots, K$  do
6:        $g_t \leftarrow g_{t-1} + \frac{1}{K} \mathbb{E}_{(X,y) \in B} [\nabla_{\theta} L(f_{\theta}(X + \delta_{t-1} + \eta_{t-1}), y)]$  // Accumulate gradients of  $\theta$ 
7:       Update token-level perturbation  $\eta$ :
8:        $g_{\eta}^i \leftarrow \nabla_{\eta^i} L(f_{\theta}((X + \delta_{t-1} + \eta_{t-1}), y))$ 
9:        $\eta_t^i \leftarrow \eta_{t-1}^i + \alpha \cdot g_{\eta}^i / \|g_{\eta}^i\|_F$ 
10:       $\eta_t \leftarrow \prod_{\|\eta\|_F < \epsilon} (\eta_t)$ 
11:      Update instance-level perturbation  $\delta$ :
12:       $g_{\delta} \leftarrow \nabla_{\delta} L(f_{\theta}((X + \delta_{t-1} + \eta_{t-1}), y))$ 
13:       $\delta_t \leftarrow \prod_{\|\delta\|_F < \epsilon} (\delta_{t-1} + \alpha \cdot g_{\delta} / \|g_{\delta}\|_F)$ 
14:    end for
15:     $V[w_i] \leftarrow \eta_K^i$  // Update perturbation vocabulary  $V$ 
16:     $\theta \leftarrow \theta - g_K$  // Update model parameter  $\theta$ 
17:  end for
18: end for

```





3. Token-Level Constraint

Algorithm 1 Token-Aware Virtual Adversarial Training

Require: Training Samples $S = \{(X = [w_0, \dots, w_i, \dots], y)\}$, perturbation bound ϵ , initialize bound σ adversarial steps K , adversarial step size α , model parameter θ

```

1:  $V \in \mathbb{R}^{N \times D} \leftarrow \frac{1}{\sqrt{D}} U(-\sigma, \sigma)$  // Initialize perturbation vocabulary  $V$ 
2: for epoch = 1,  $\dots$ , do
3:   for batch  $B \subset S$  do
4:      $\delta_0 \leftarrow \frac{1}{\sqrt{D}} U(-\sigma, \sigma)$ ,  $\eta_0^i \leftarrow V[w_i]$ ,  $g_0 \leftarrow 0$  //Initialize perturbation and gradient of  $\theta$ 
5:     for t = 1,  $\dots$ ,  $K$  do
6:        $g_t \leftarrow g_{t-1} + \frac{1}{K} \mathbb{E}_{(X,y) \in B} [\nabla_{\theta} L(f_{\theta}(X + \delta_{t-1} + \eta_{t-1}), y)]$  //Accumulate
7:       Update token-level perturbation  $\eta$ .
8:        $g_{\eta}^i \leftarrow \nabla_{\eta^i} L(f_{\theta}((X + \delta_{t-1} + \eta_{t-1}), y))$ 
9:        $\eta_t^i \leftarrow \eta_{t-1}^i + \alpha \cdot g_{\eta}^i / \|\eta_{t-1}^i\|_F$ 
10:       $\eta_t \leftarrow \prod_{\|\eta\|_F < \epsilon} (\eta_t)$ 
11:      Update instance-level perturbation  $\delta$ :
12:       $g_{\delta} \leftarrow \nabla_{\delta} L(f_{\theta}((X + \delta_{t-1} + \eta_{t-1}), y))$ 
13:       $\delta_t \leftarrow \prod_{\|\delta\|_F < \epsilon} (\delta_{t-1} + \alpha \cdot g_{\delta} / \|g_{\delta}\|_F)$ 
14:    end for
15:     $V[w_i] \leftarrow \eta_K^i$  //Update perturbation vocabulary  $V$ 
16:     $\theta \leftarrow \theta - g_K$  //Update model parameter  $\theta$ 
17:  end for
18: end for
  
```

$$n^i = \frac{\|\eta_t^i\|_F}{\max_j (\|\eta_t^j\|_F)}$$



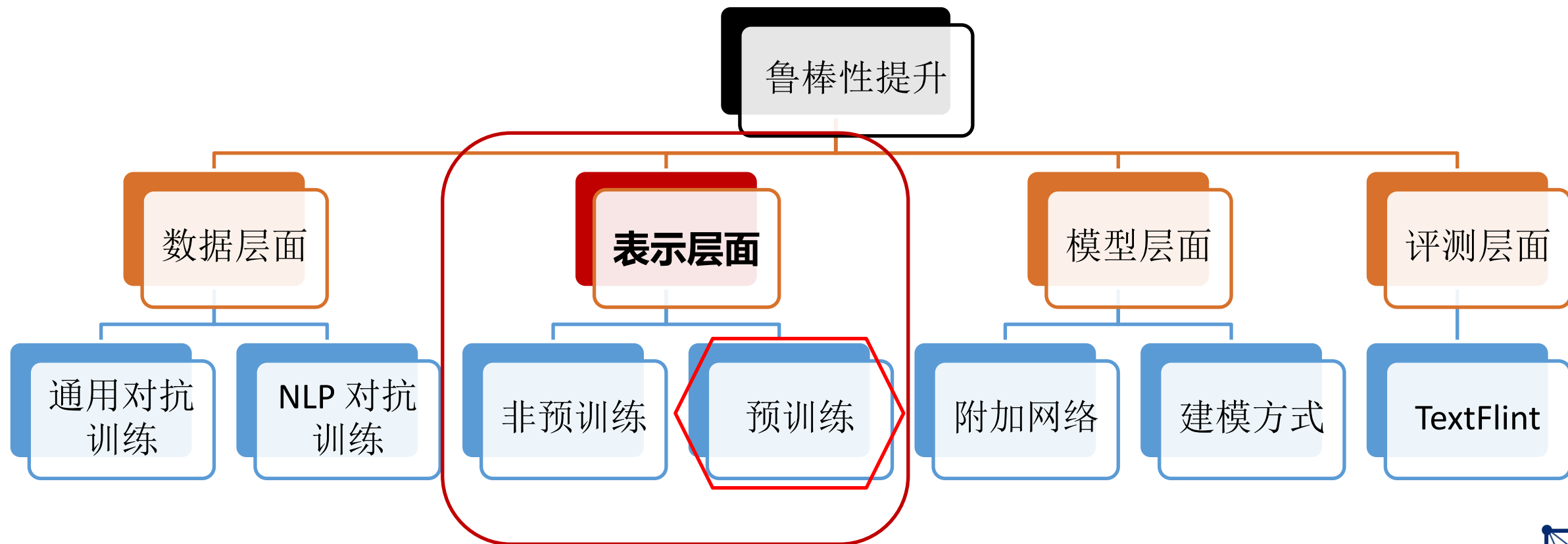
Results

Model	RTE	QNLI	MRPC	CoLA	SST	STS-B	MNLI-m/mm	QQP
	Acc	Acc	Acc/f1	Mcc	Acc	P/S Corr	Acc	Acc/f1
BERT-BASE								
BERT (Devlin et al. 2018)	-	88.4	-/86.7	-	92.7	-	84.4/-	-
BERT-ReImp	63.5	91.1	84.1/89.0	54.7	92.9	89.2/88.8	84.5/84.4	90.9/88.3
FreeAT-ReImp	68.0	91.3	85.0/89.2	57.5	93.2	89.5/89.0	84.9 / 85.0	91.2/88.5
FreeLB-ReImp	70.0	91.5	86.0/90.0	58.9	93.4	89.7/89.2	85.3 / 85.5	91.4/88.6
TA-VAT(ours)	74.0	92.4	88.0/91.6	62.0	93.7	90.0/89.6	85.7 / 85.8	91.6/88.9
ALBERT-xxlarge-v2								
ALBERT-xxlarge-v2(Lan et al. 2019)	89.2	95.3	-/90.9	71.4	96.9(96.5)	93.0/-	90.8/-	92.2/-
FreeLB(Zhu et al. 2020)	89.9	95.6*	-/92.4	73.1	97.0	93.2/-	90.9/-	92.5/-
TA-VAT(ours)	90.3	95.7	-/ 93.4	74.1	96.8	93.4/-	91.1/-	92.6 /-

Table 1: Evaluation results on the development set of GLUE benchmark. QNLI* in FreeLB is formed as pairwise ranking task.

Model	RTE	QNLI	MRPC	CoLA	SST	STS-B	MNLI-m/mm	QQP
	Acc	Acc	Acc/f1	Mcc	Acc	P/S Corr	Acc	Acc/f1
BERT-BASE(Devlin et al. 2018)	66.4	90.5	88.9/84.8	52.1	93.5	87.1/85.8	84.6/83.4	71.2/89.2
FreeLB(Zhu et al. 2020)	70.1	91.8*	88.1/83.5	54.5	93.6	87.7/86.7	85.7/84.6	72.7/89.6
TA-VAT(ours)	71.0	91.7	88.9/84.5	55.9	94.5	86.8/85.7	85.2/ 84.7	72.8/89.5

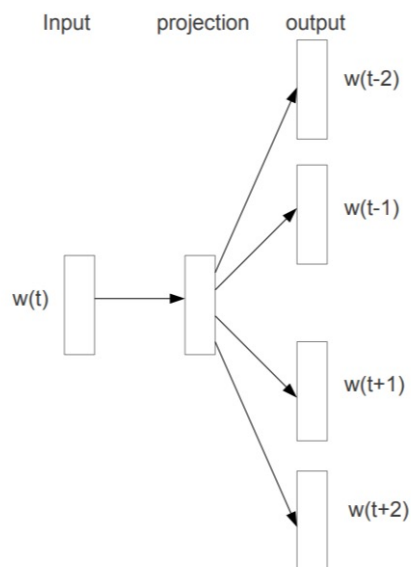
Table 2: Evaluation results on the test set of GLUE benchmark. Results use the evaluation server on GLUE website. QNLI* in FreeLB is formed as pairwise ranking task.



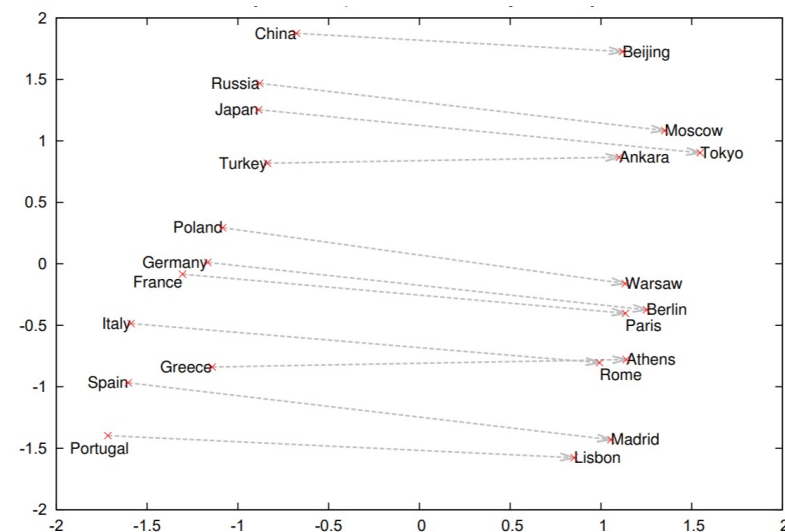


➤ Word Embedding (unk 未见词)

- Glove / Word2Vec / Random Init



Word2Vec (skip-gram)



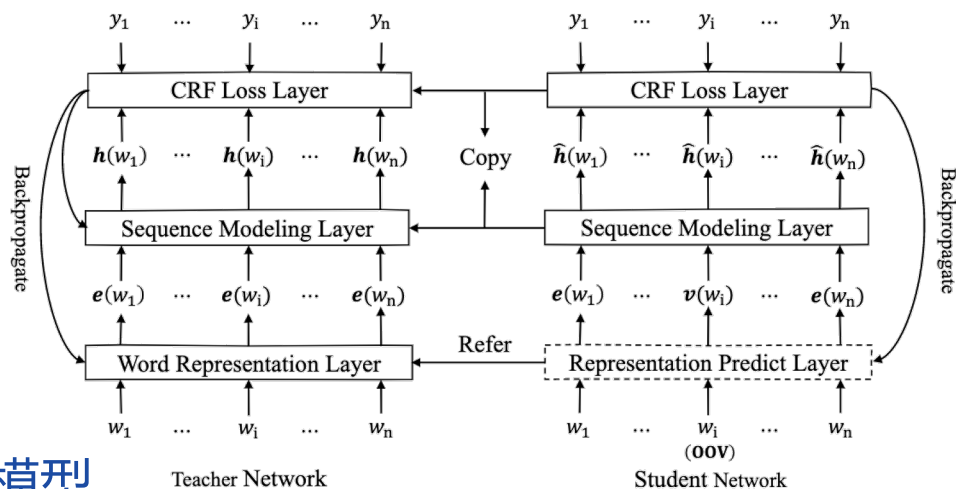
Good Representation :
Country and Capital Vectors

学习 OOV 词在序列标注任务中的表示

未登录词

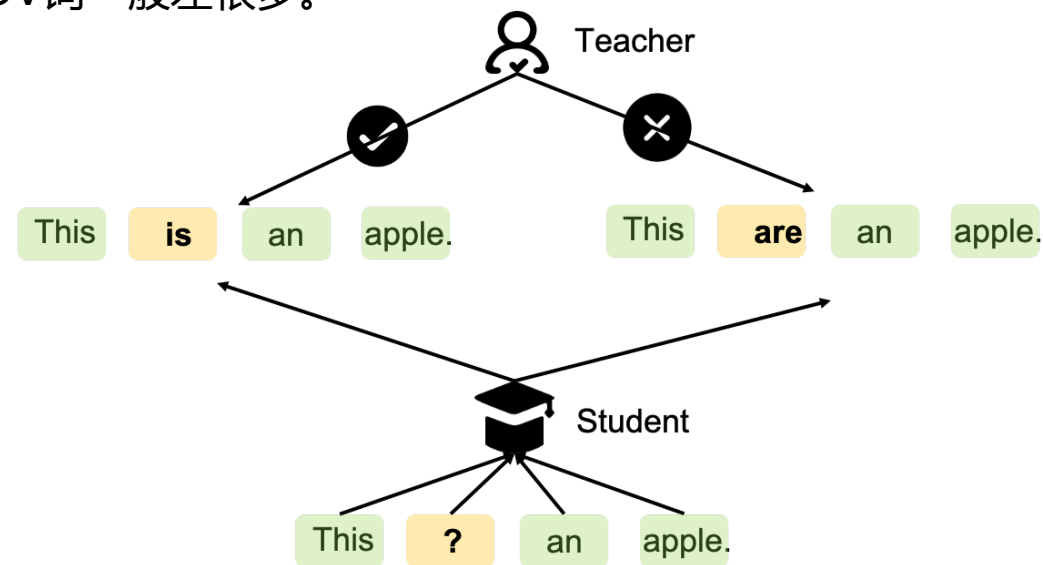
研究动机

- 在测试集中经常会出现一些在训练集中没有出现过的OOV词。
- 在序列标注任务中，模型在OOV词上的表现相较于非OOV词一般差很多。



模型

- 训练一个普通的序列标注模型，称为Teacher Network。
- 训练一个OOV词表示模型(Student Network)：利用一个词的上下文和该词的字符序列预测改词的代表，使得预测的词表示在Teacher Network中有好的表现。





Dataset	Dev		Test	
	#OOV	OOV Rate	#OOV	OOV Rate
NER				
CoNLL02-Spanish	2,216	50.91%	1,544	43.38%
CoNLL02-Dutch	1,819	69.53%	2,564	65.05%
Twitter-English	1,266	79.15%	4,131	79.13%
CoNLL03-German	3,928	81.27%	2,685	73.10%

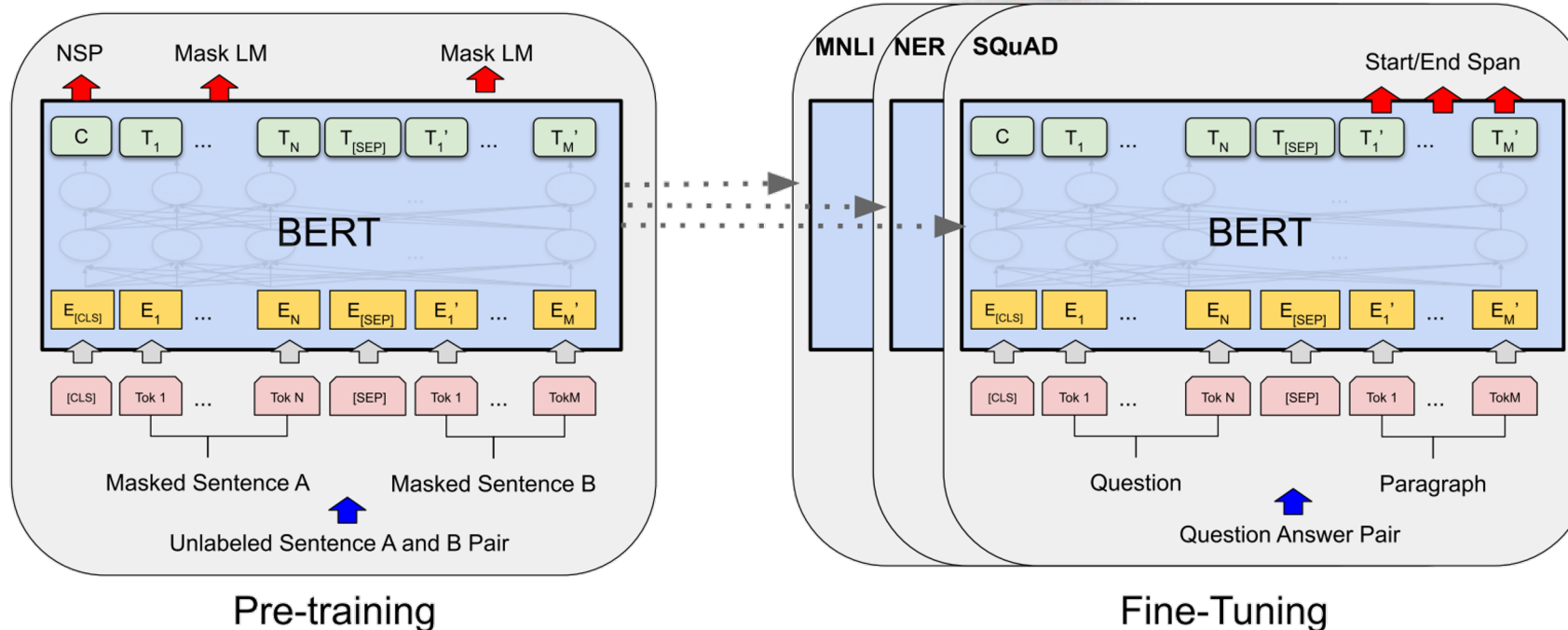
数据集 OOV 统计

Arch	Model	CoNLL02-Spanish		CoNLL02-Dutch		Twitter-English		CoNLL03-German	
		Dev	Test	Dev	Test	Dev	Test	Dev	Test
LSTM	RandomUNK	69.36	72.06	64.23	64.08	56.88	56.38	55.92	56.89
	SingleUNK	68.79	71.59	67.83	66.39	56.82	56.39	59.69	60.16
	[Lazaridou <i>et al.</i> , 2017]	68.61	69.08	65.99	65.43	47.72	47.20	47.87	49.17
	[Khodak <i>et al.</i> , 2018]	68.74	69.53	66.34	65.70	48.22	47.28	47.97	49.33
	[Schick and Schütze, 2018]	70.84	72.88	68.88	67.51	59.18	57.21	55.83	58.42
	[Akbik <i>et al.</i> , 2018]	61.78	64.06	60.49	62.09	49.68	50.22	55.06	53.01
	Proposed	73.91	74.63	70.33	70.12	60.14	58.32	60.55	61.79

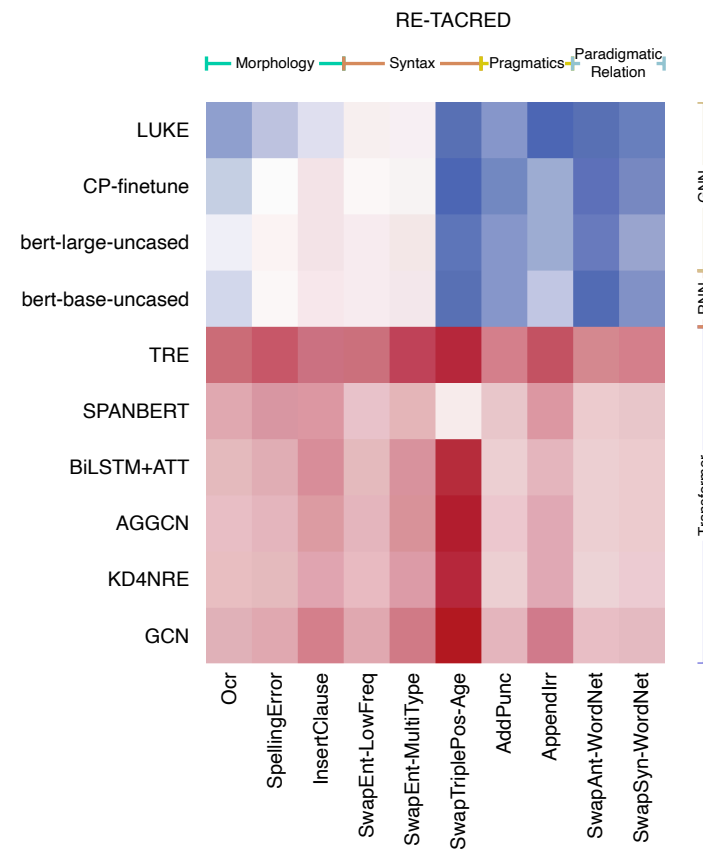
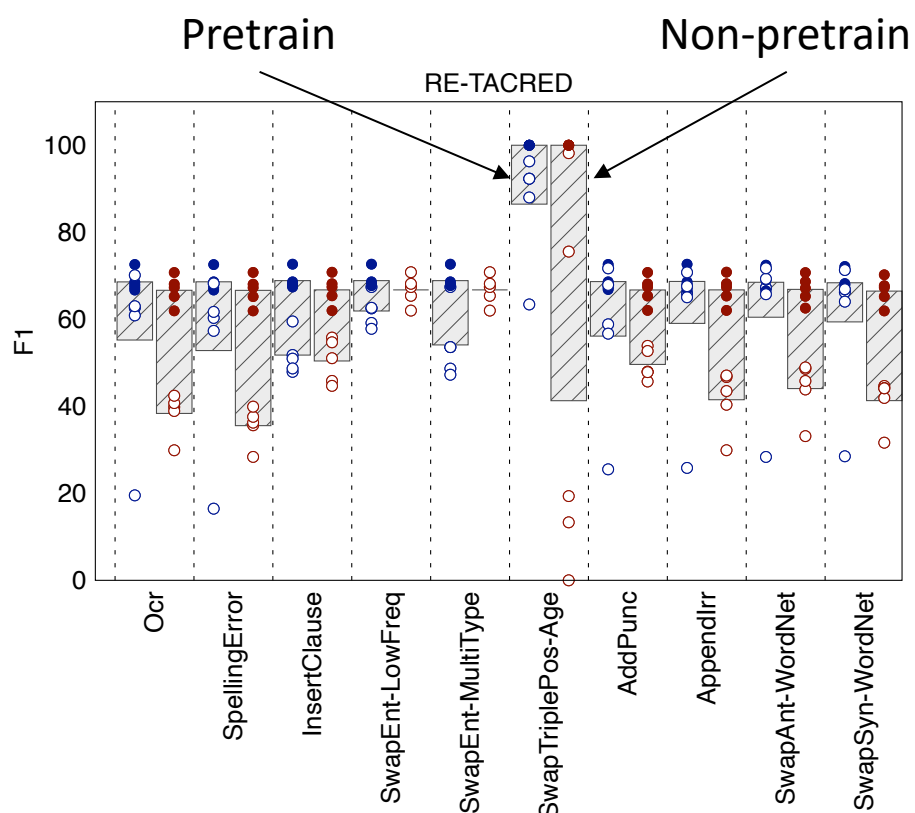
面对 OOV 鲁棒性

➤ Pretrained Contextual Embedding (contextual)

- ELMo / BERT (我在苹果公司工作/我吃了苹果)



BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding
Deep contextualized word representations



预训练与非预训练鲁棒性对比

Language Models Need Knowledge

- **PLMs perform poorly on entity recognition**

- Contextualized PLMs achieved small improvements on entity & semantic related tasks compared with non-contextualized methods. ([Tenney et al.](#))

- BPE tokenization breaks entities

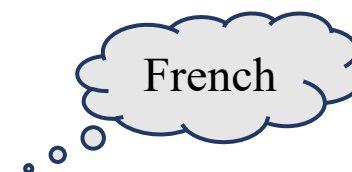
The native language of Jean Mara ##is is French.

- Surface form-based reasoning

	original BERT	E-BERT- replace	E-BERT- concat	ERNIE	Know- Bert
Jean Marais	French	French	French	french	french
Daniel Ceccaldi	Italian	French	French	french	italian
Orane Demazis	Albanian	French	French	french	french
Sylvia Lopez	Spanish	French	Spanish	spanish	spanish
Annick Alane	English	French	French	english	english

**BERT does not know *Daniel Ceccaldi* (as an entity) at all.
It just think *Daniel Ceccaldi* looks like an Italian name.**

The native language of Jean Marais is [MASK] .



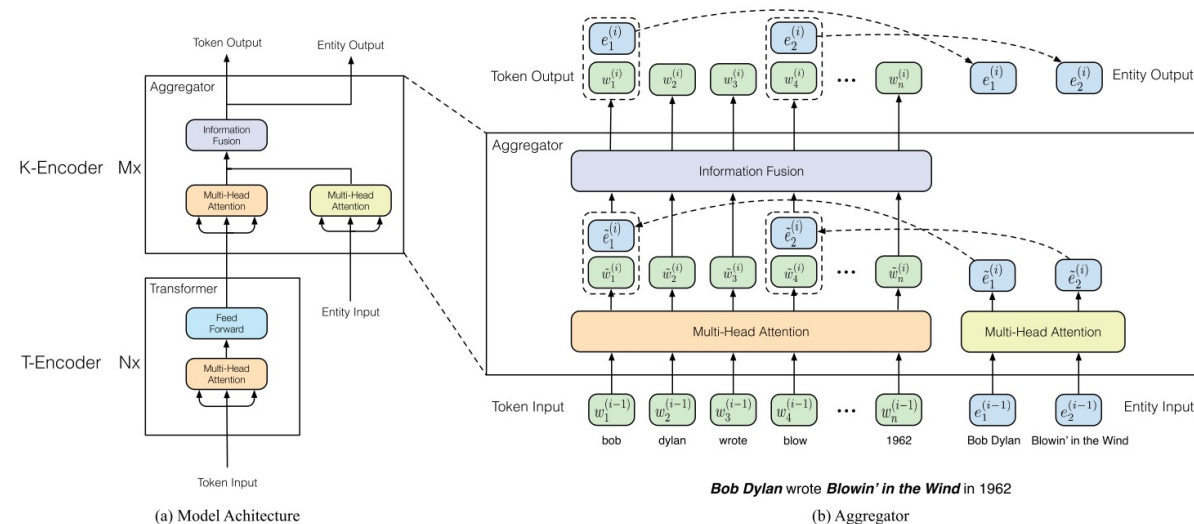
Injecting Knowledge into PLMs

- **Injecting entity embeddings**

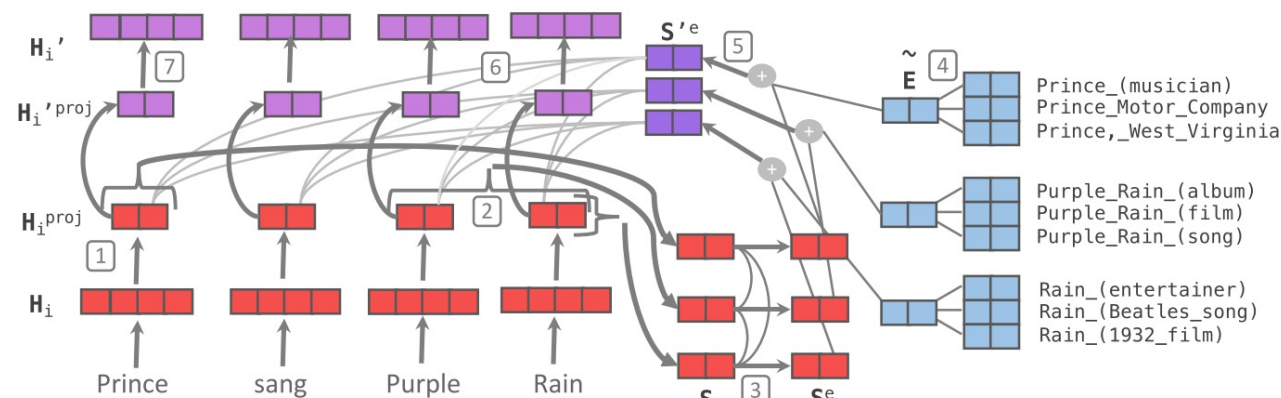
- ERNIE, KnowBERT, K-BERT, etc.
- The entity embeddings are NOT
 - Jointly learned along with PLM
 - Contextualized

- **Knowledge as supervision**

- WKLM, etc.



ERNIE: Enhanced Language Representation with Informative Entities
<https://arxiv.org/abs/1905.07129>



Knowledge Enhanced Contextual Word Representations
<https://arxiv.org/abs/1909.04164>

Model	Acc.	Macro	Micro
NFGEC (Attentive)	54.53	74.76	71.58
NFGEC (LSTM)	55.60	75.15	71.73
BERT	52.04	75.16	71.63
ERNIE	57.19	76.51	73.39

Table 2: Results of various models on F1GER (%).

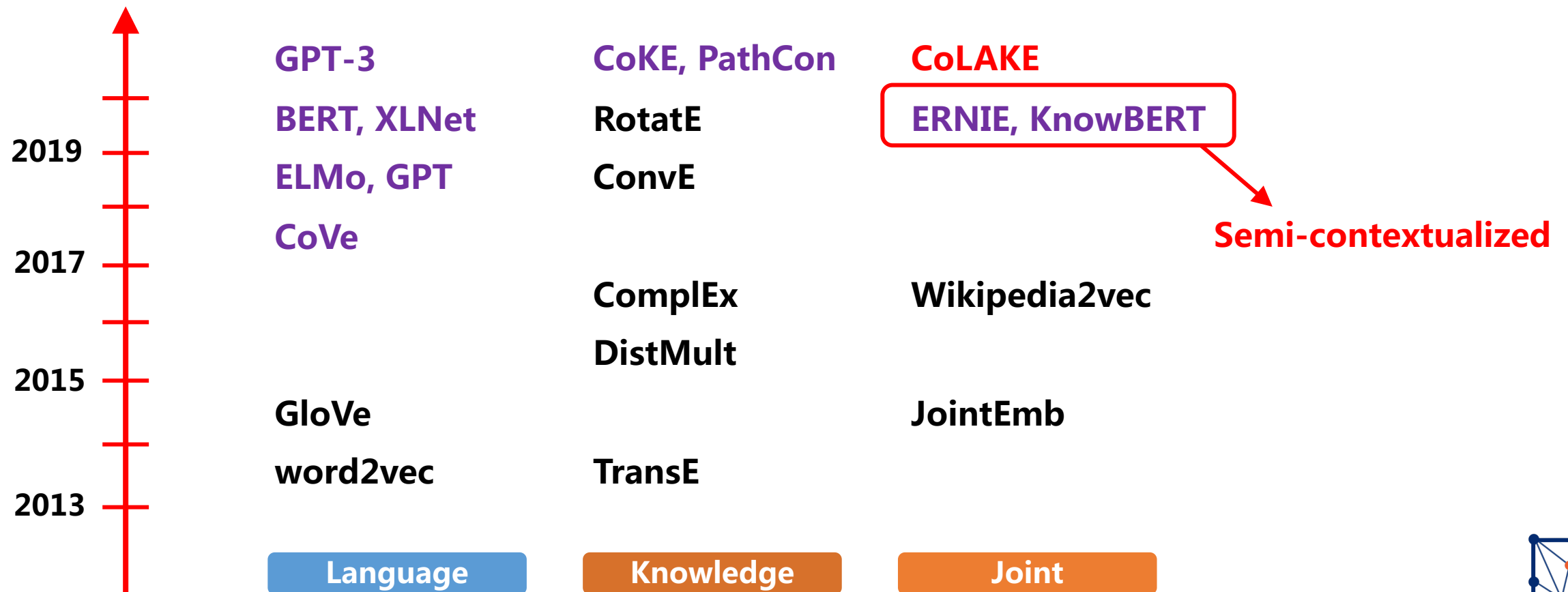
Model	FewRel			TACRED		
	P	R	F1	P	R	F1
CNN	69.51	69.64	69.35	70.30	54.20	61.20
PA-LSTM	-	-	-	65.70	64.50	65.10
C-GCN	-	-	-	69.90	63.30	66.40
BERT	85.05	85.11	84.89	67.23	64.81	66.00
ERNIE	88.49	88.44	88.32	69.97	66.08	67.97

Table 5: Results of various models on FewRel and TACRED (%).

ERNIE 在实体相关的任务上性能更好

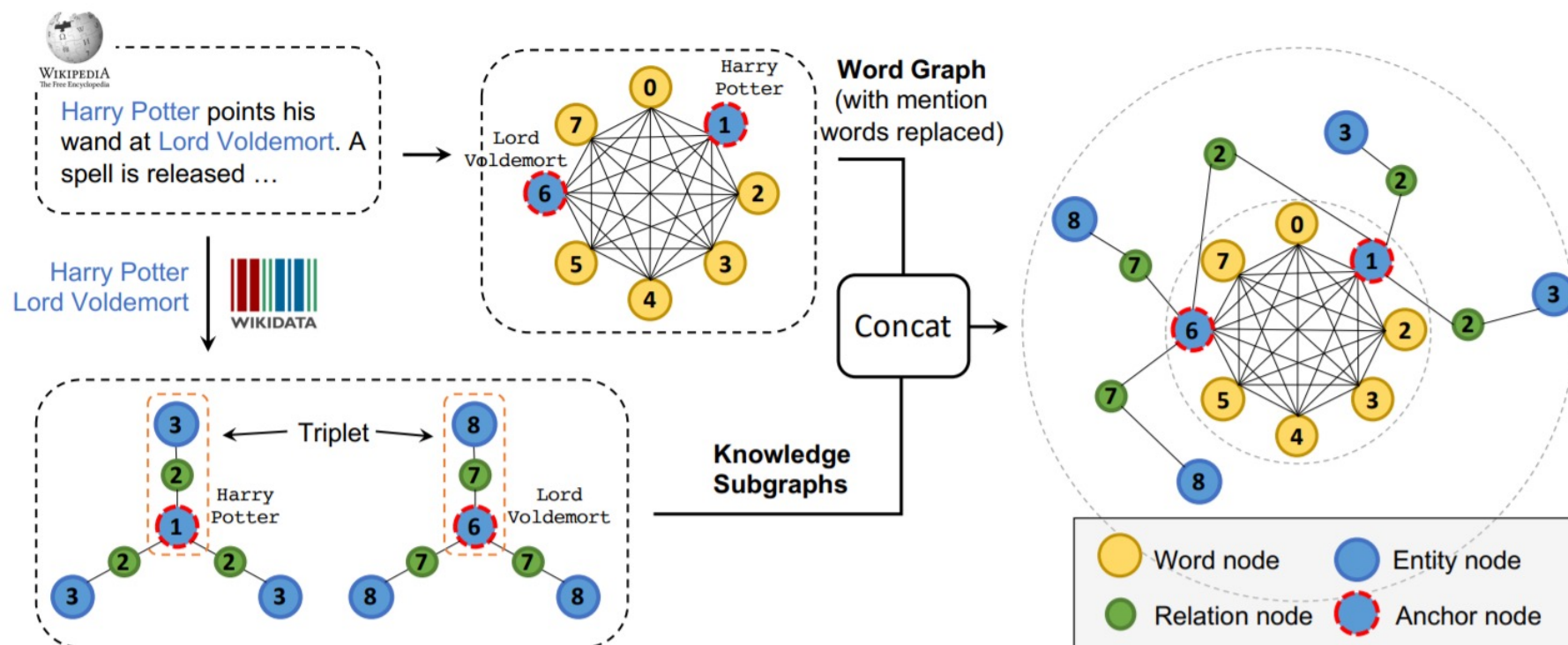
Representation in Language and Knowledge

- Combine the success of both sides -- CoLAKE





- Word graph + Knowledge subgraph





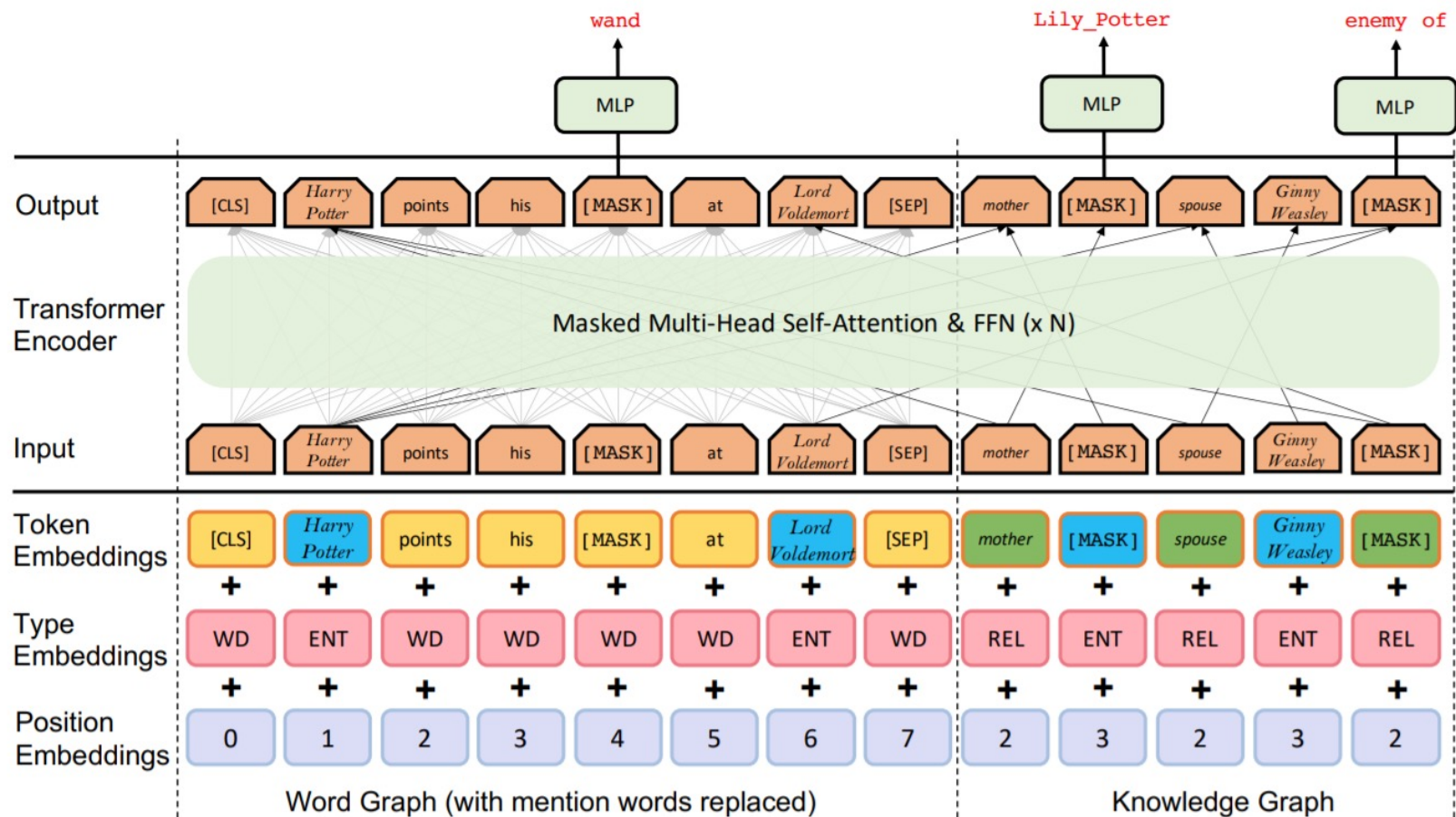
• Modify Transformer for word-knowledge graph

Word-knowledge graph is a **positional heterogeneous graph**.

Position
Embedding

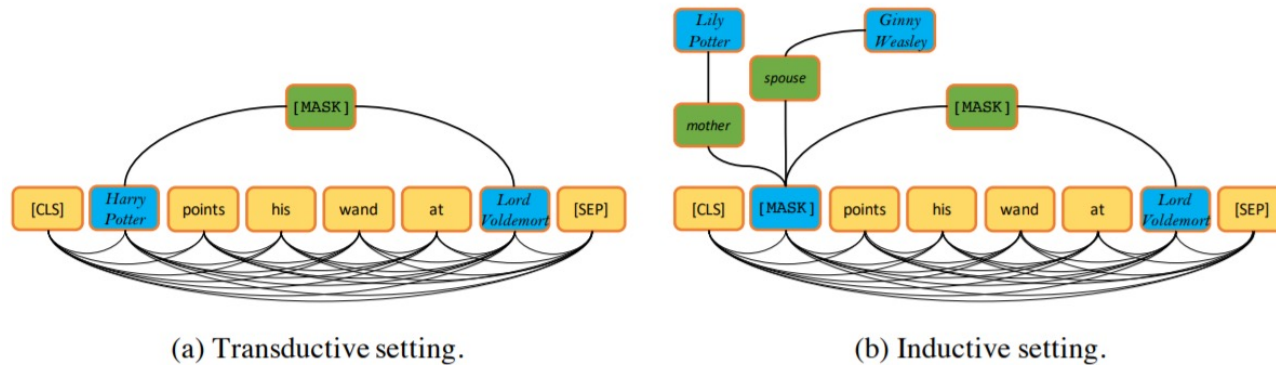
Type
Embedding

Masked
Self-Attention





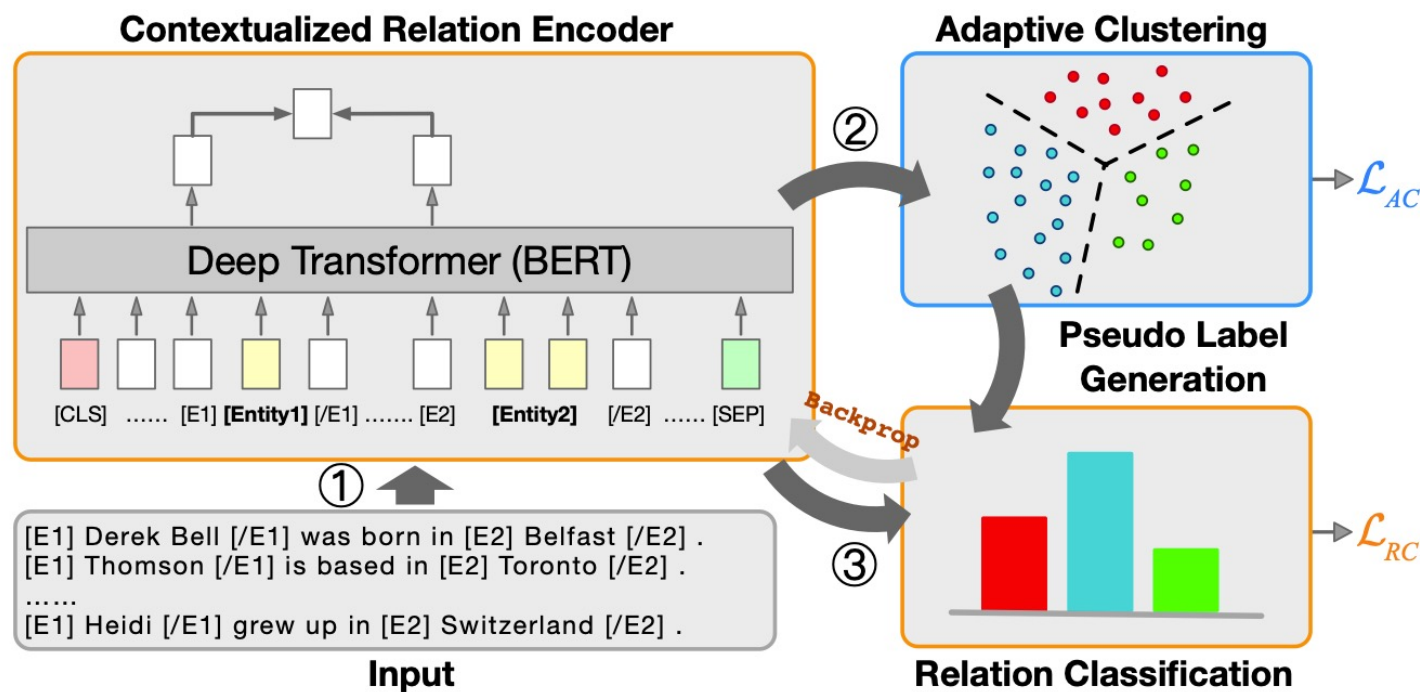
- **Synthetic graph task**
 - Word-knowledge graph completion

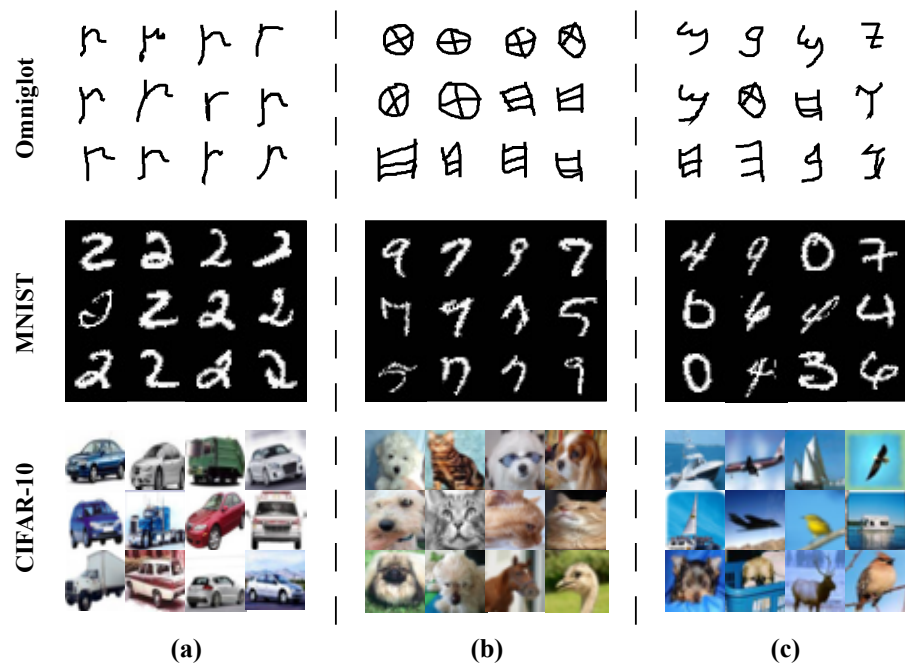


Model	MR ↓	MRR	HITS@1	HITS@3	HITS@10
Transductive setting					
TransE (Bordes et al., 2013)	15.97	67.30	60.28	70.96	79.75
DistMult (Yang et al., 2015)	27.09	60.56	48.66	69.69	79.61
ComplEx (Trouillon et al., 2016)	26.73	61.09	49.80	70.64	79.78
RotatE (Sun et al., 2019)	30.36	70.90	64.74	74.89	81.05
CoLAKE	2.03	82.48	72.14	92.19	98.58
Inductive setting					
DKRL (Xie et al., 2016)	168.21	8.18	5.03	7.28	14.13
CoLAKE	31.01	28.10	15.69	30.28	58.05

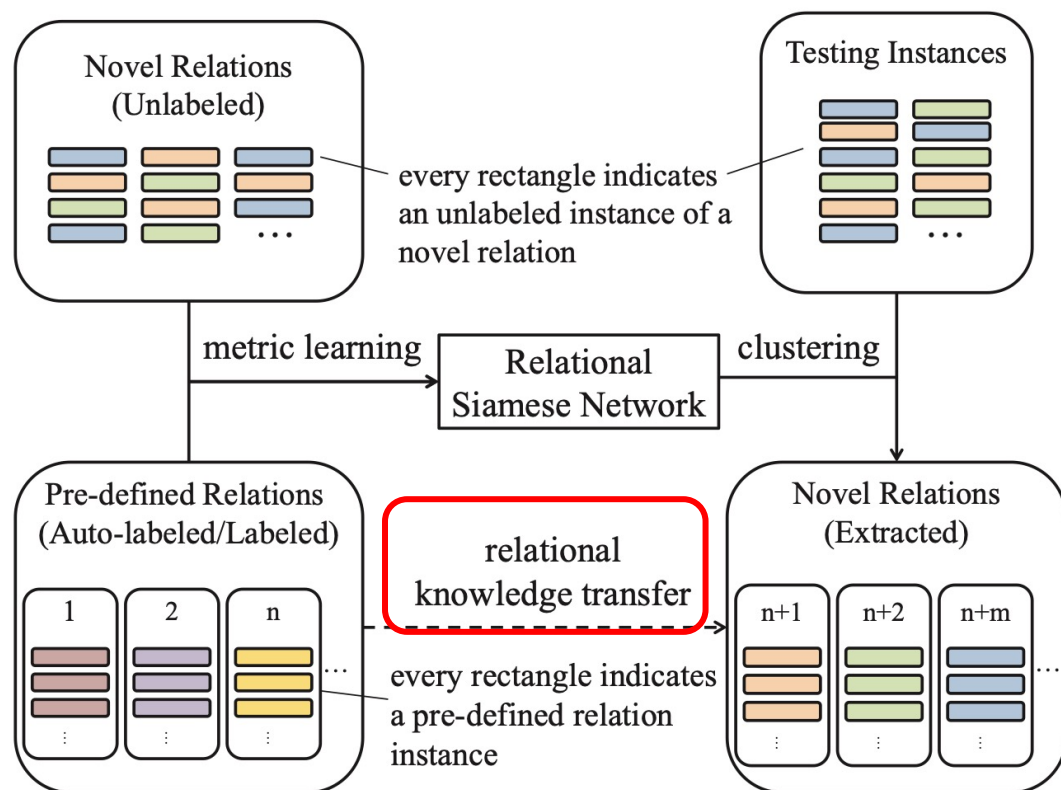
Results on word-knowledge graph completion

2 Representation---Target Oriented

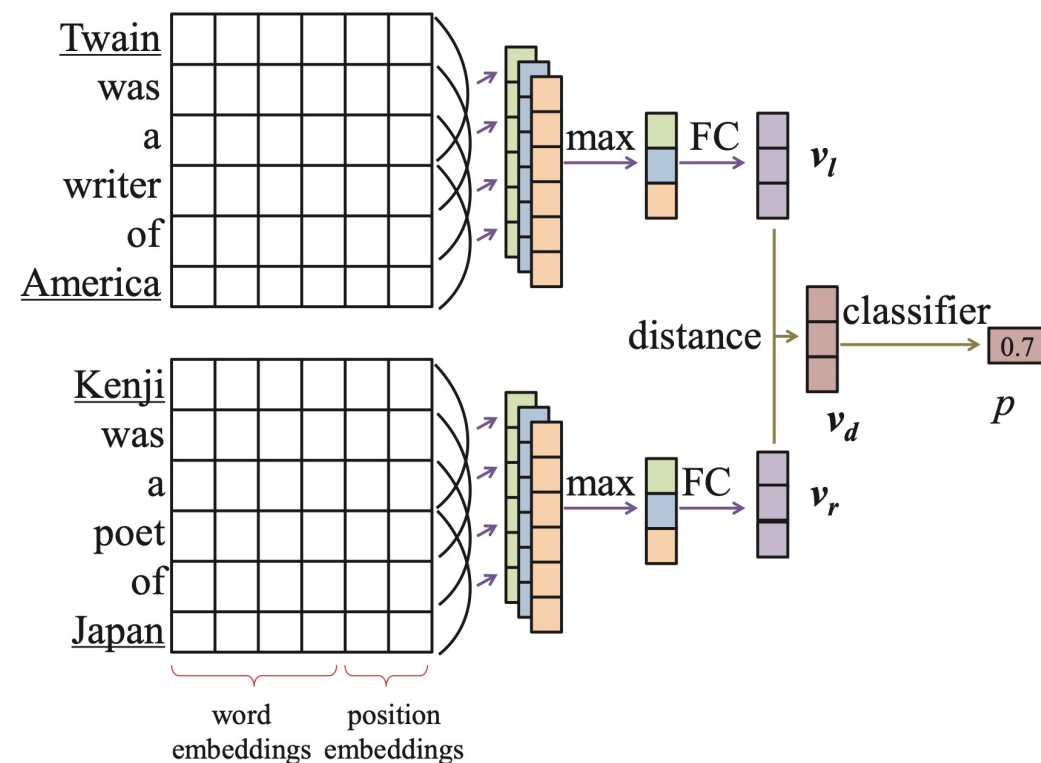




Problem in Unsupervised Clustering



Transfer knowledge from Labeled Data



Relational Siamese Networks

2 Representation---Target Oriented



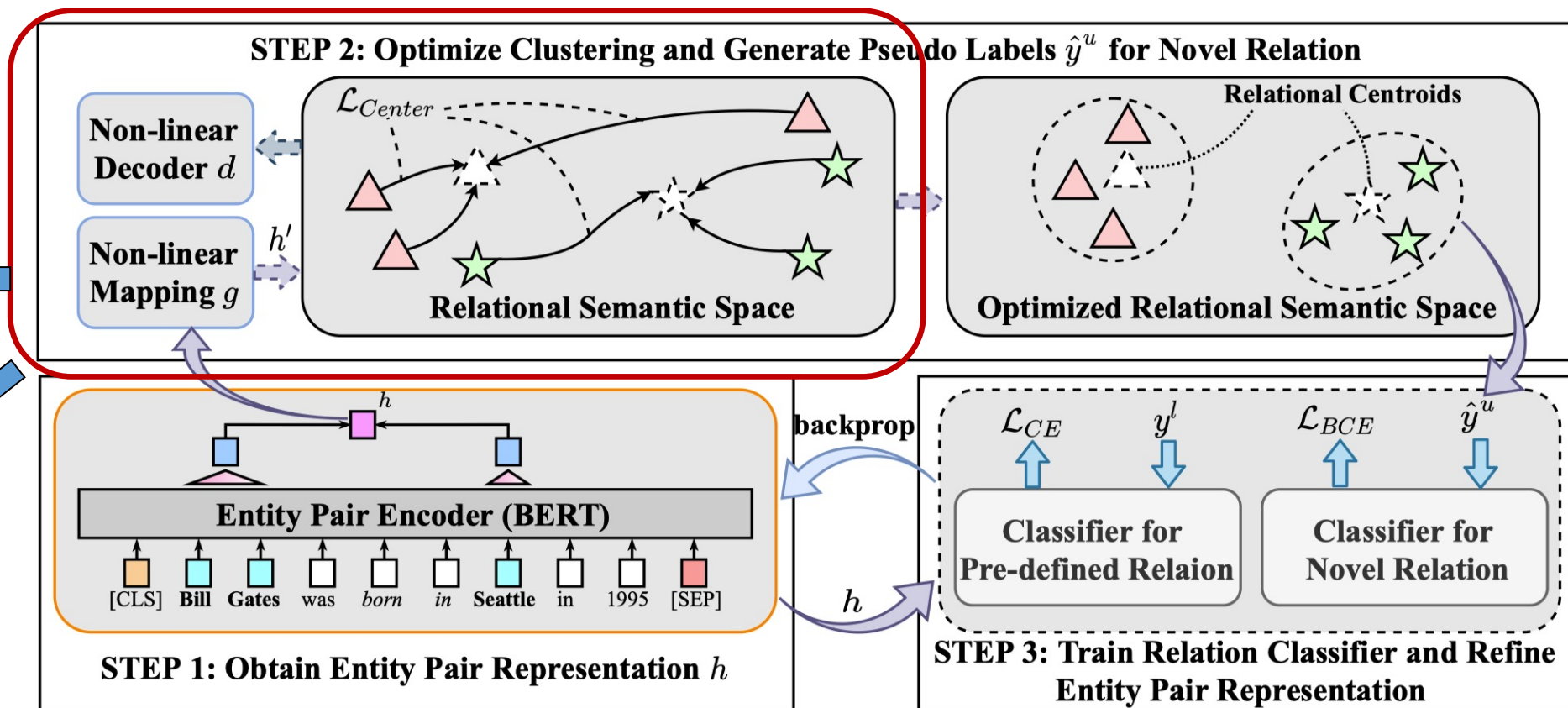
Center Loss for Learning Clustering:

$$\mathcal{L}_{center} = \frac{1}{2} \sum_{i=1}^N \|h_i^{\ell'} - c_{y_i}\|_2^2$$

$$c_r = \frac{1}{|\mathcal{D}_r|} \sum_{i \in \mathcal{D}_r} h_i^{\ell'}$$

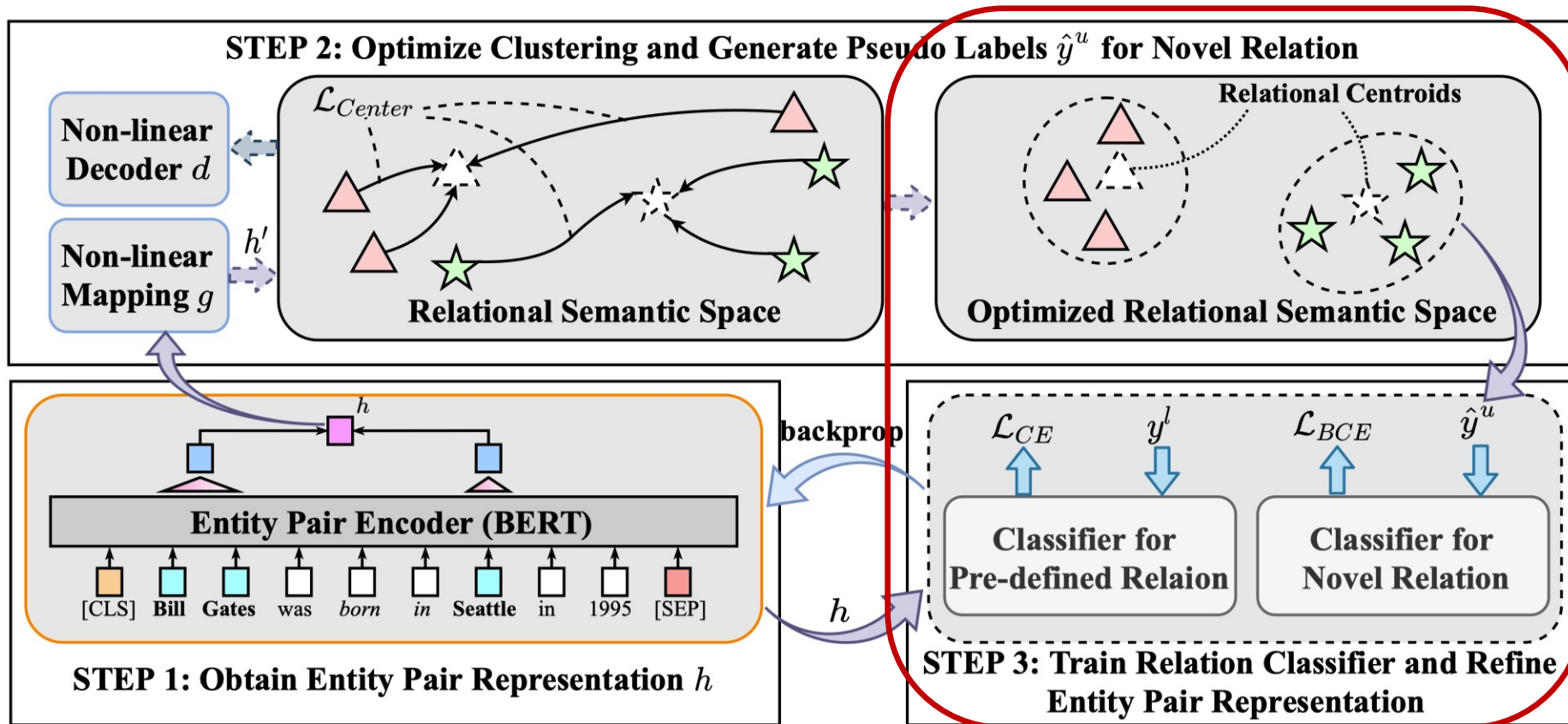
Reconstruction Loss
Avoids Trivial Solution:

$$\mathcal{L}_C = \sum_{i=1}^N \ell(g(h_i'), h_i)$$



Overview of our proposed RoCORE method.

2 Representation---Target Oriented



Kmeans Clustering for Pseudo Labels

$$\hat{y}^u = \text{k-means}(h^{u'})$$

BCE Loss for New Classifier



$$\mathcal{L}_{BCE}(s_i^u, s_j^u) = q_{i,j} \mathcal{L}^+(s_i^u, s_j^u) + (1 - q_{i,j}) \mathcal{L}^-(s_i^u, s_j^u).$$

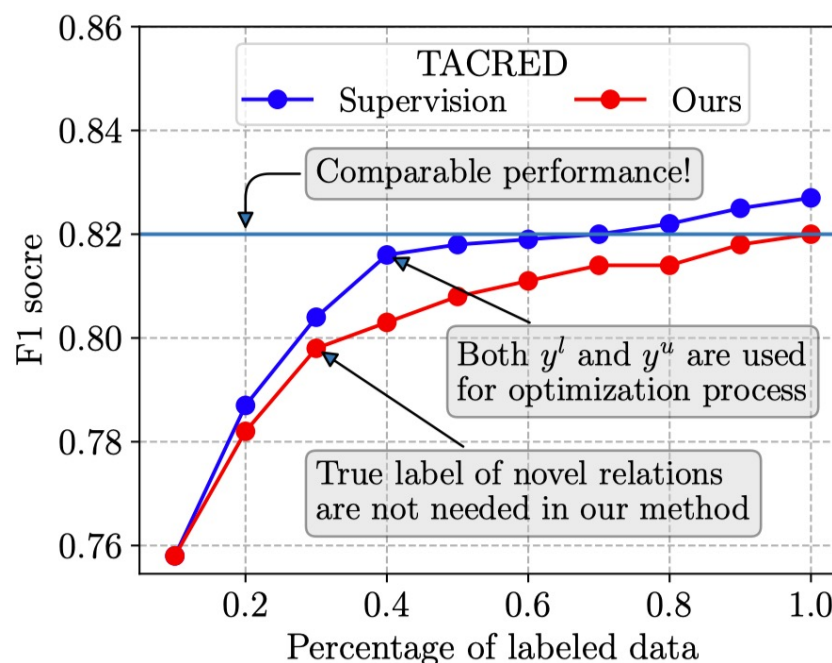
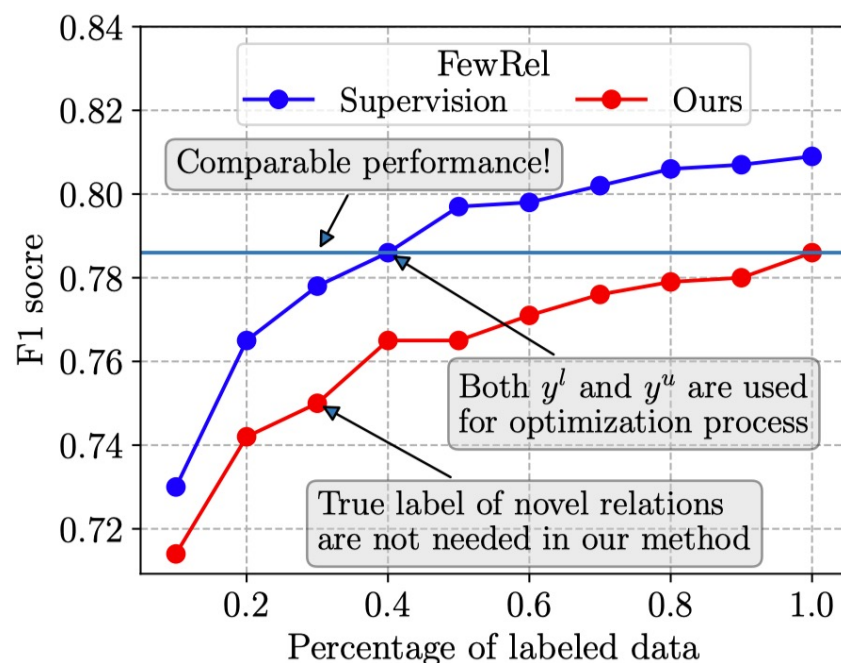
Overview of our proposed RoCORE method.



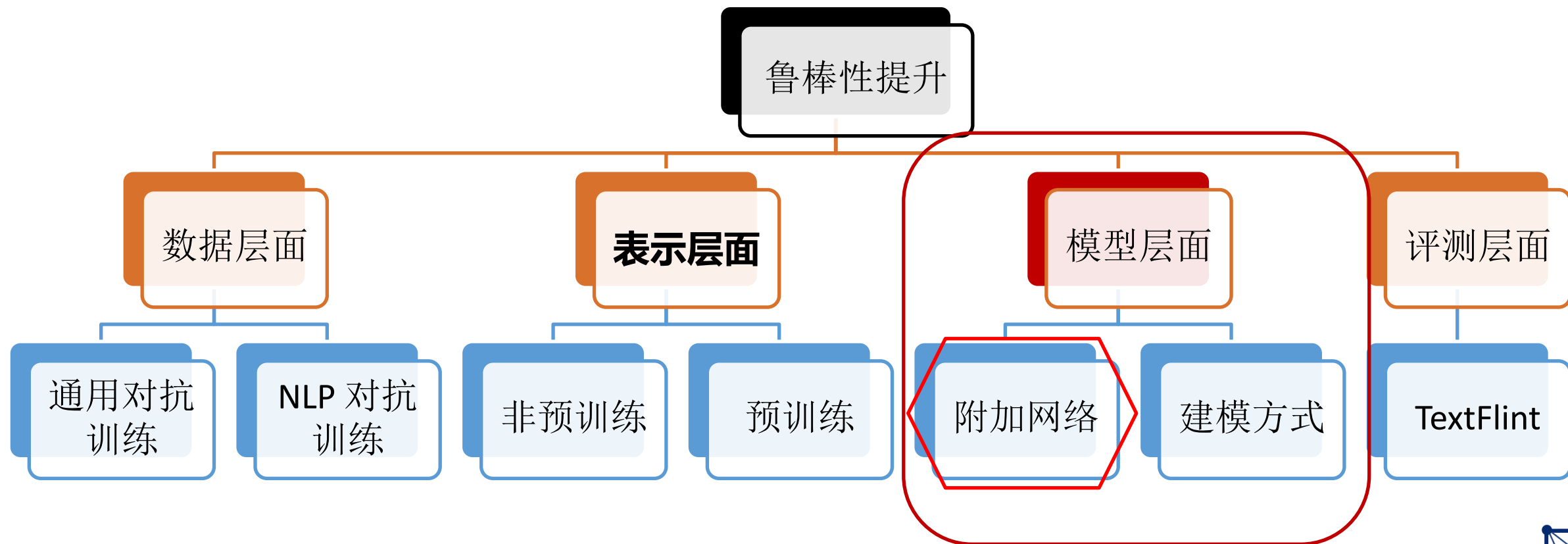
0-63 Labeled
64-79 Unlabeled

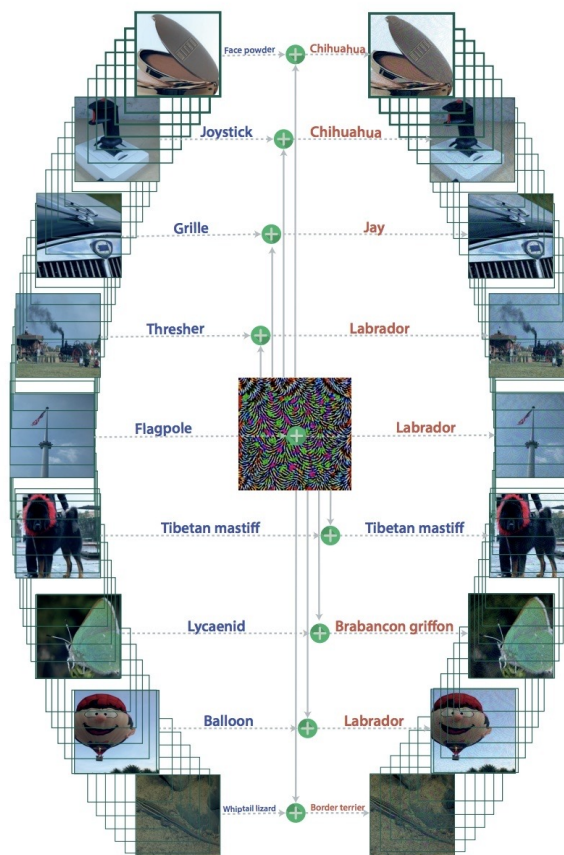
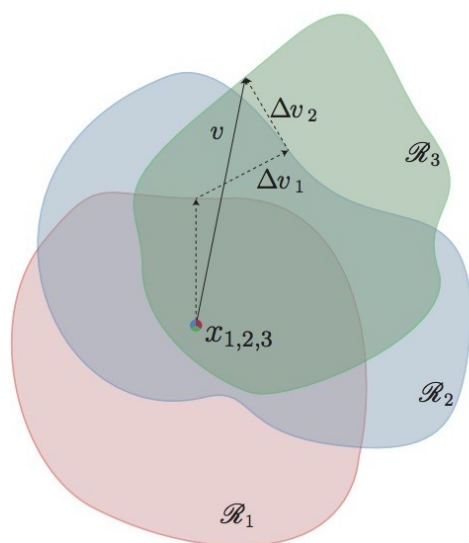
0-30 Labeled
31-40 Unlabeled

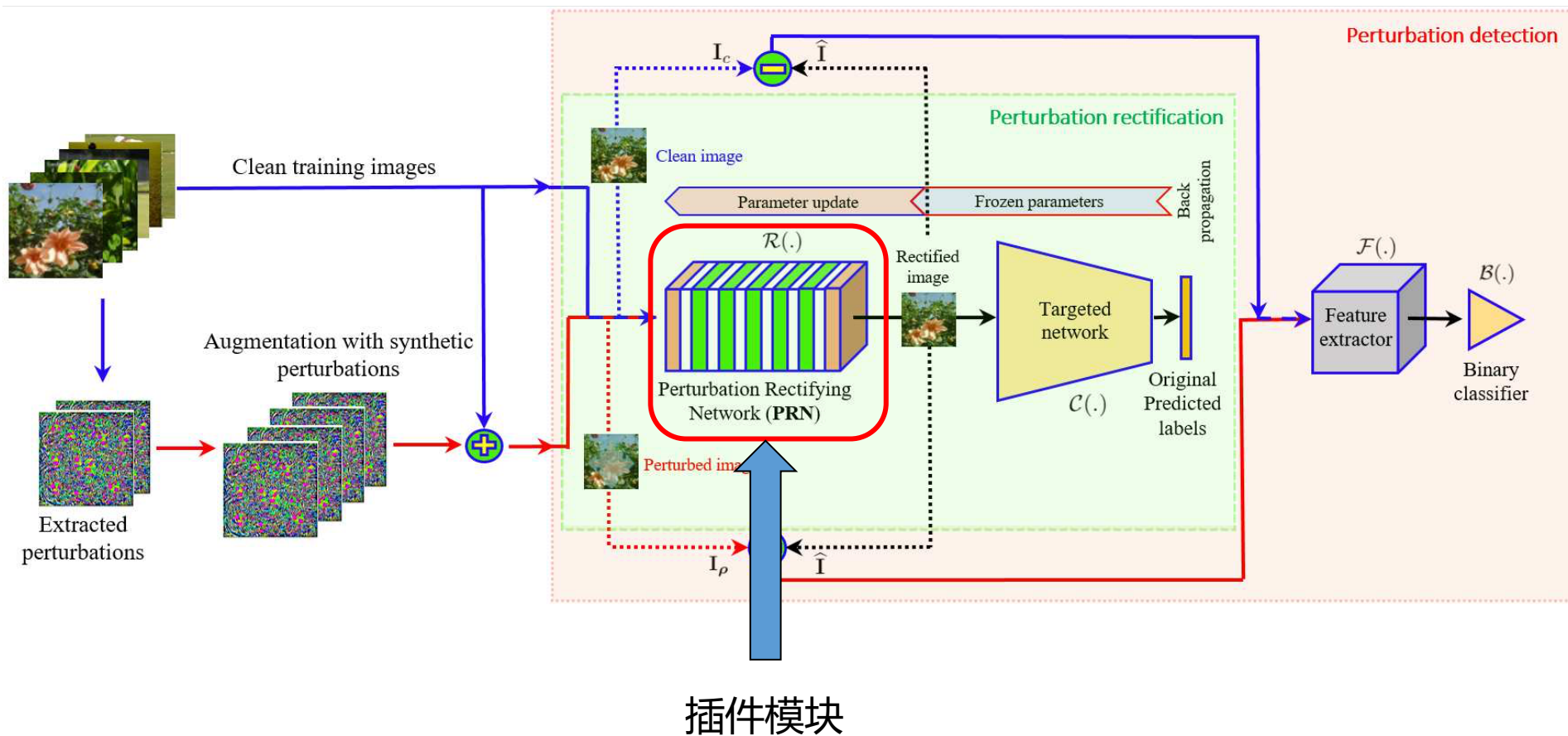
Method	FewRel			TACRED			FewRel→TACRED			TACRED→FewRel		
	P	R	F_1	P	R	F_1	P	R	F_1	P	R	F_1
VAE	0.179	0.697	0.285	0.234	0.606	0.338	0.234	0.606	0.338	0.179	0.697	0.285
RW-HAC	0.318	0.460	0.376	0.534	0.568	0.551	0.534	0.568	0.551	0.318	0.460	0.376
RSN	0.496	0.734	0.592	0.628	0.634	0.631	0.398	0.617	0.484	0.233	0.567	0.331
RSN+BERT	0.585	0.899	0.709	0.795	0.878	0.834	0.379	0.843	0.523	0.262	0.898	0.406
RoCORE	0.749 ₁₂	0.846 ₀₉	0.794 ₀₉	0.871 ₄₂	0.849 ₃₅	0.860 ₃₄	0.616 ₃₃	0.599 ₅₇	0.607 ₄₁	0.687 ₃₅	0.766 ₄₆	0.724 ₂₅
 In Domain						 Cross Domain						

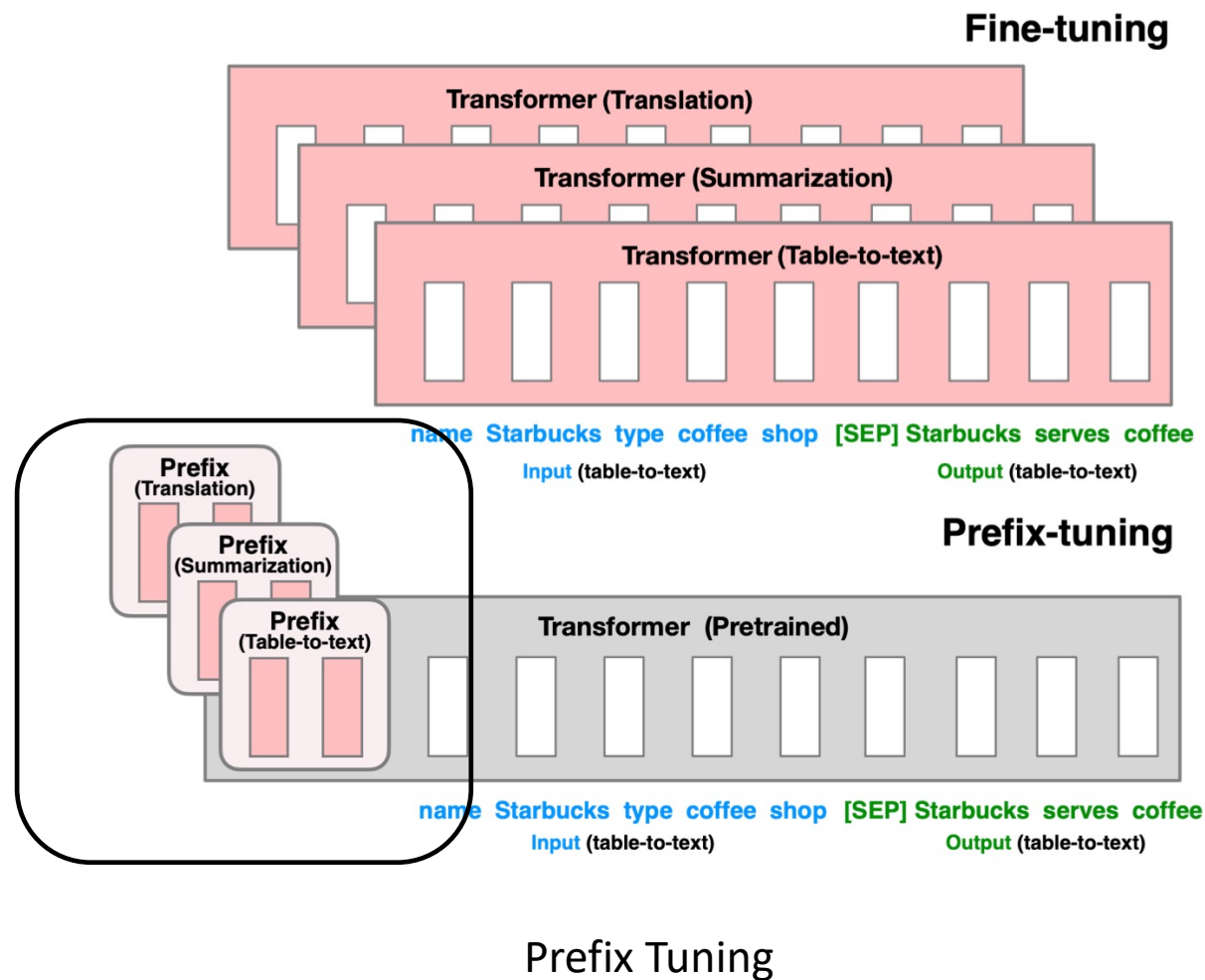


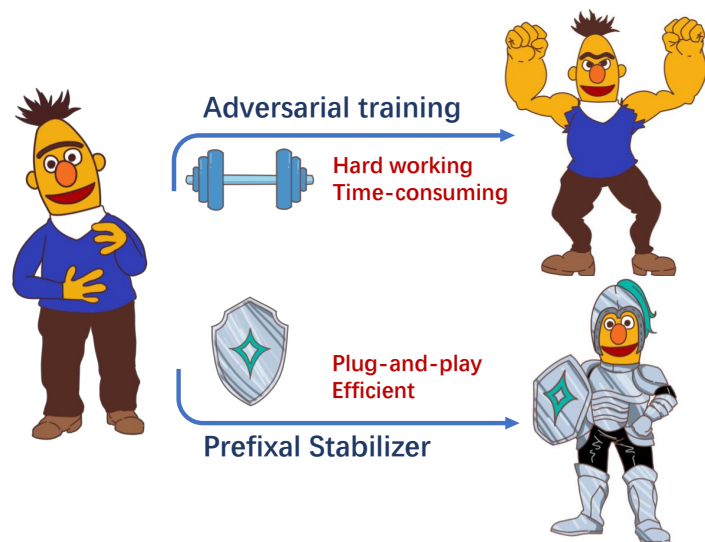
Model performance with different amounts of labeled data.



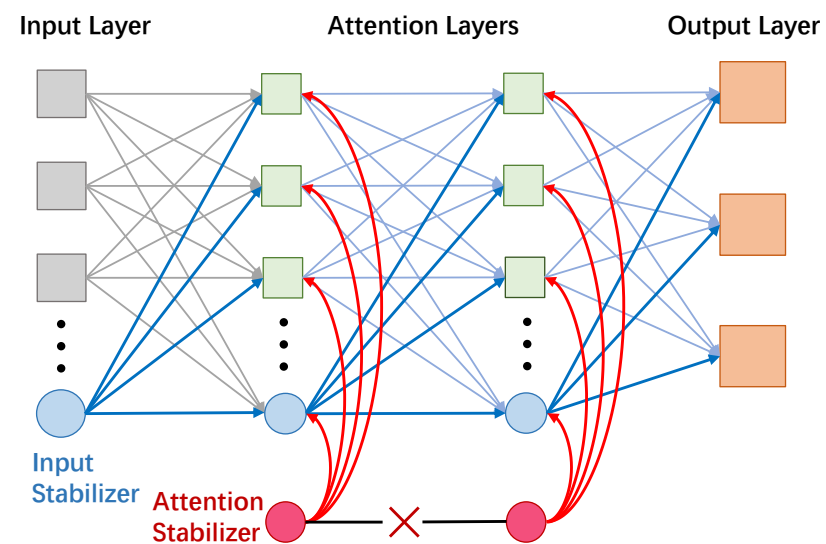








Adversarial training VS Prefixal Stabilizer



Architecture of Prefixal Stabilizer

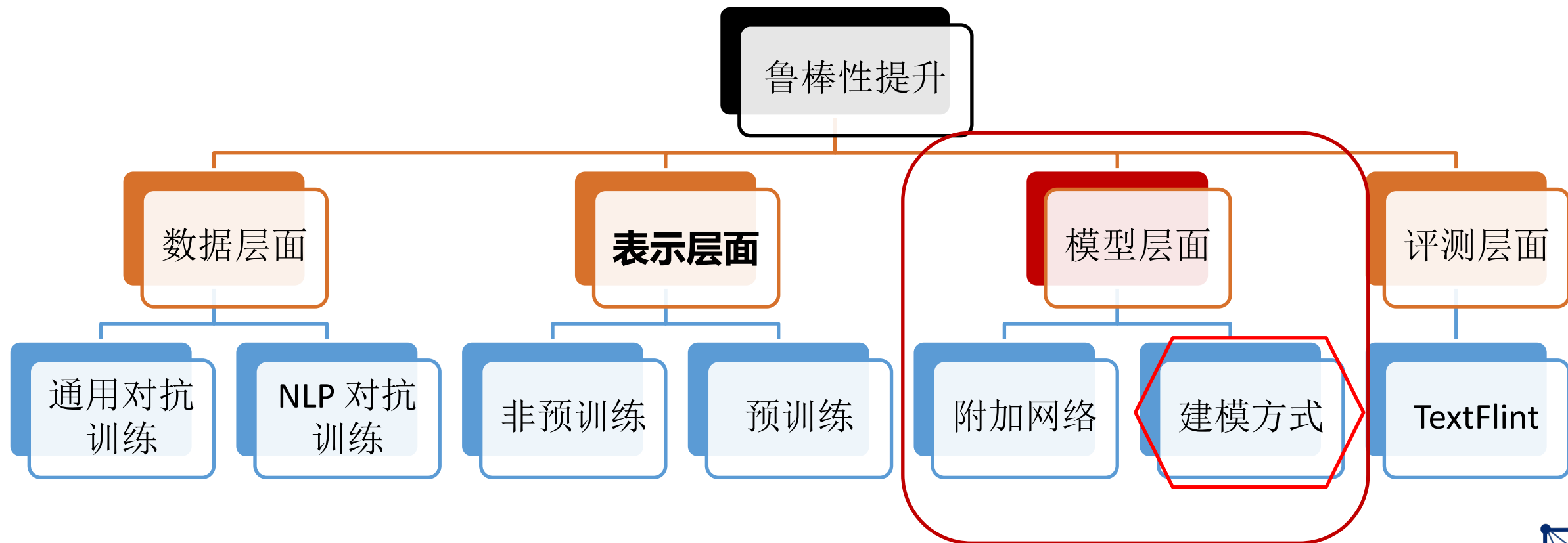


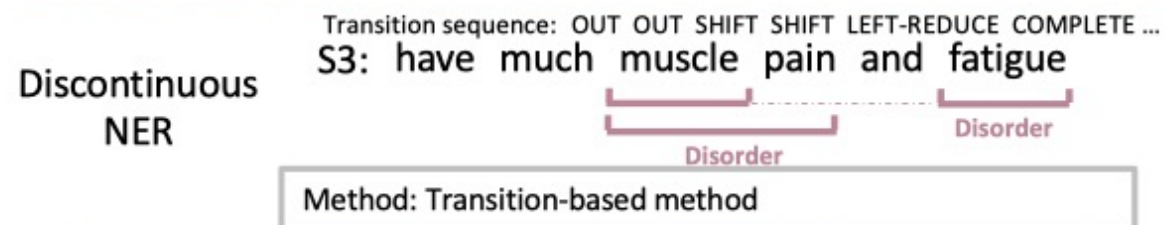
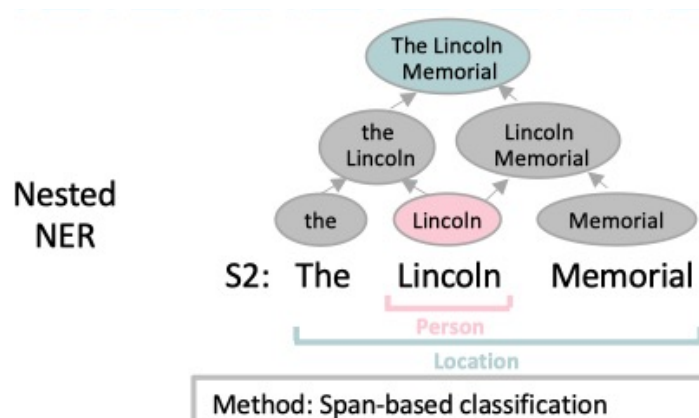
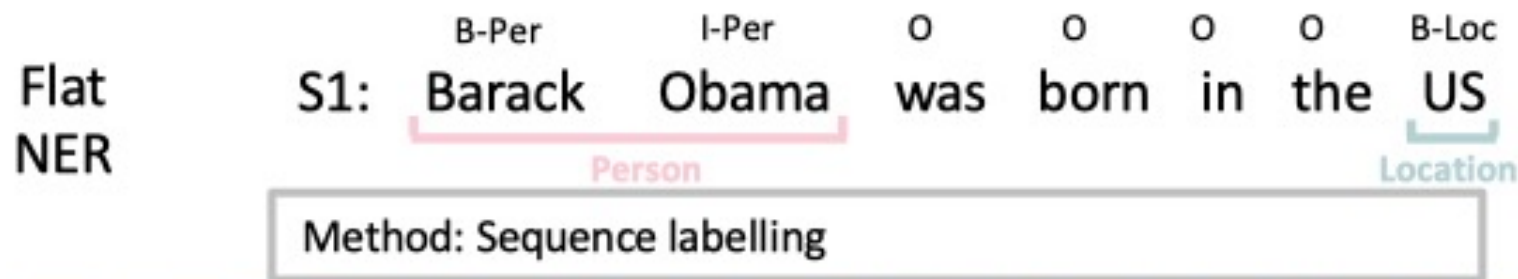
Methods	Clean%	TextFooler			PWWS			TextBugger		
		Aua%	Suc%	#Query	Aua%	Suc%	#Query	Aua%	Suc%	#Query
Baseline (BERT)	94.2	9.5	89.9	1007.8	20.0	78.8	2654.2	6.0	93.6	1038.7
PGD (Madry et al. 2017)	94.6	25.5	73.0	1416.0	35.0	68.3	2741.5	23.0	75.9	1625.5
FreeLB (Zhu et al. 2019)	94.7	18.0	81.0	1160.7	25.5	73.1	2699.4	19.0	79.9	1445.0
TAVAT (Li and Qiu 2020)	94.3	14.0	85.2	1117.7	14.5	84.6	2625.4	12.0	87.3	1122.6
InfoBERT (Wang et al. 2020)	94.5	23.5	75.1	1374.7	27.5	70.9	2655.9	21.0	77.8	1420.6
MixADA (Si et al. 2021)	93.9	10.5	88.8	1023.4	23.5	75.0	2656.6	6.5	93.1	1077.3
Stabilizer-PGD	93.9	33.5	64.3	1626.6	39.0	58.5	2777.1	32.5	65.4	1845.3
Stabilizer-FreeLB	94.1	24.0	74.5	1300.3	32.5	65.5	2772.6	23.0	75.6	1472.7

Table 2: Experiment results of different defenders on IMDB, where all models are trained on BERT. The best performance is marked in bold.

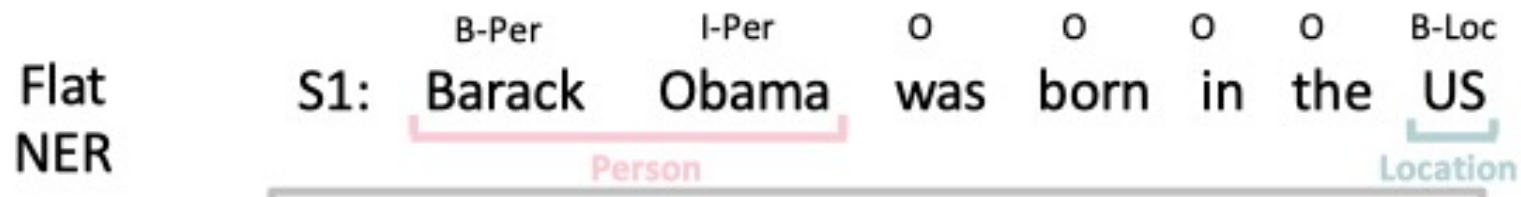
Methods	SST-2 to IMDB			
	Clean%	Aua%	Suc%	#Query
BERT (SST-2)	90.8	2.0	97.8	564.1
PGD (SST-2)	89.9	20.0	77.8	827.8
FreeLB (SST-2)	89.2	22.0	75.3	874.2
Stabilizer-PGD (SST-2) +BERT (IMDB)	93.7	36.0	61.6	1664.5
Stabilizer-FreeLB (SST-2) +BERT (IMDB)	93.8	34.0	63.7	1536.0

Table 3: Experimental results when the models **trained on SST-2** is transferred to **IMDB for testing**. BERT (fine-tuned on IMDB) equipped with our module (trained on SST-2) has improved robustness to a large extent and has little damage to cleaning accuracy.



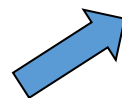
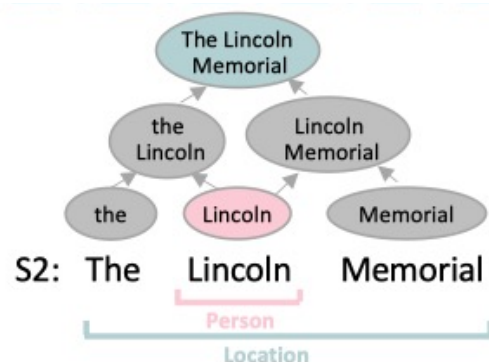


Three kinds of named entities

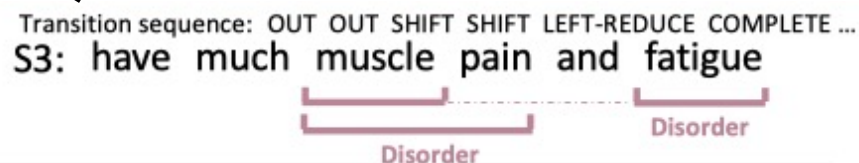


Real Application

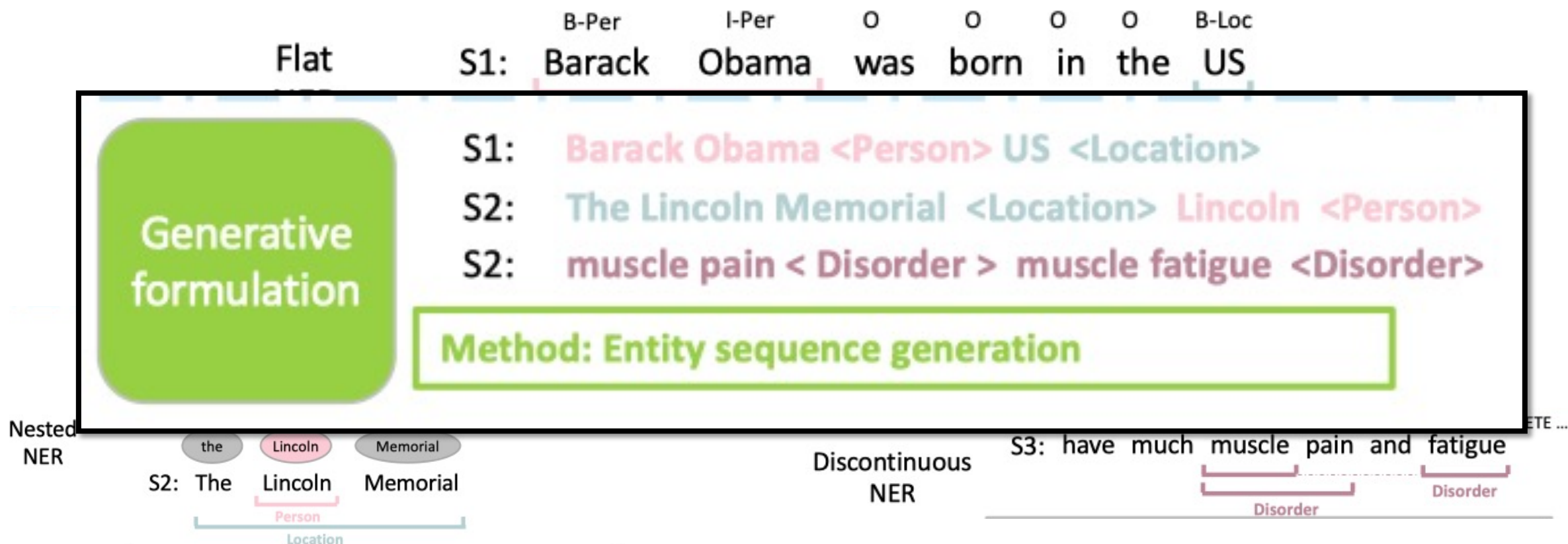
Nested NER



Discontinuous NER



The different formulations make it hard to solve all NER tasks in a unified method

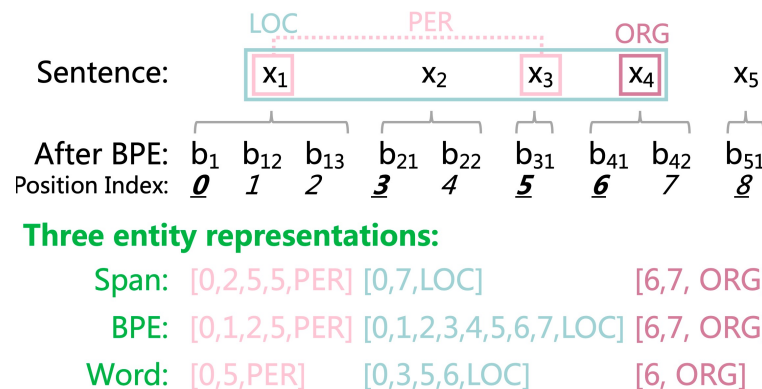
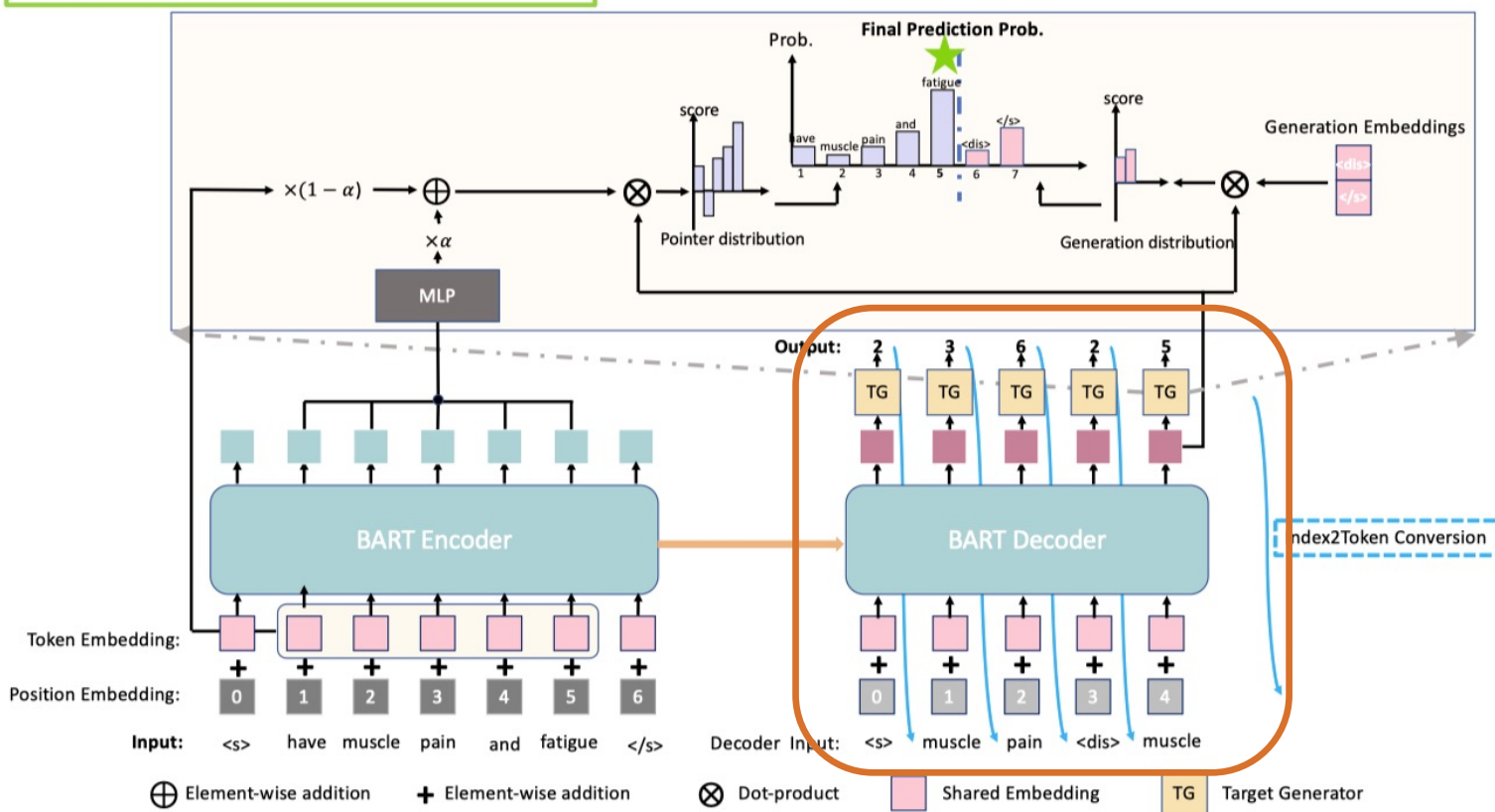


The different formulations make it hard to solve all NER tasks in a unified method



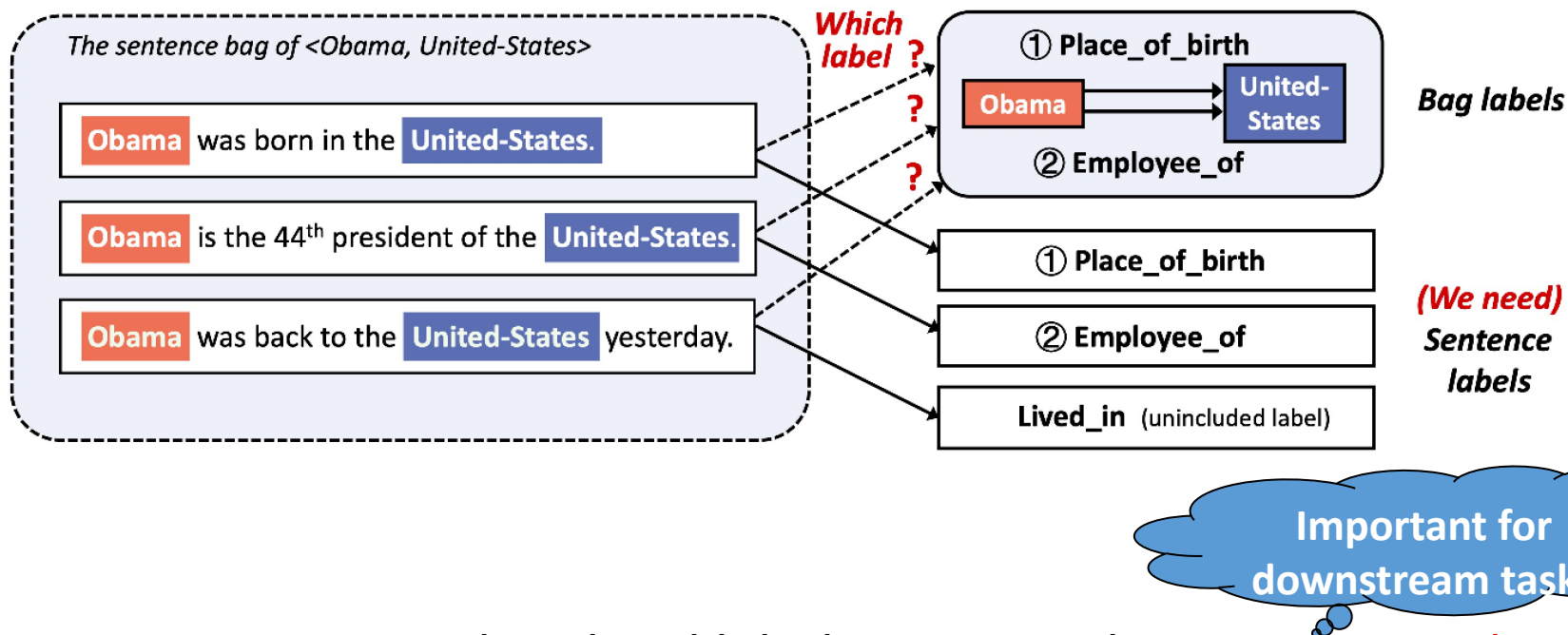
Input: <s> have muscle pain and fatigue </s>

Output: 2 3 6 2 5 </s>



Three kinds of entity representations

Figure 2: Model structure used in our method. "<s>" and "</s>" are the predefined start-of-sentence and end-of-sentence tokens in BART. We assume there is only one entity tag "<dis>", therefore, only this token and the "</s>" token need to be generated, other tokens can be generated by the pointer.



1. Given bag-level labels, can we obtain **sentence-level labels**?
2. Sentence bag contains **correct** labels, **incorrect** labels, and **unincluded** labels.
3. Previous positive learning framework **cannot distinguish** noisy data.

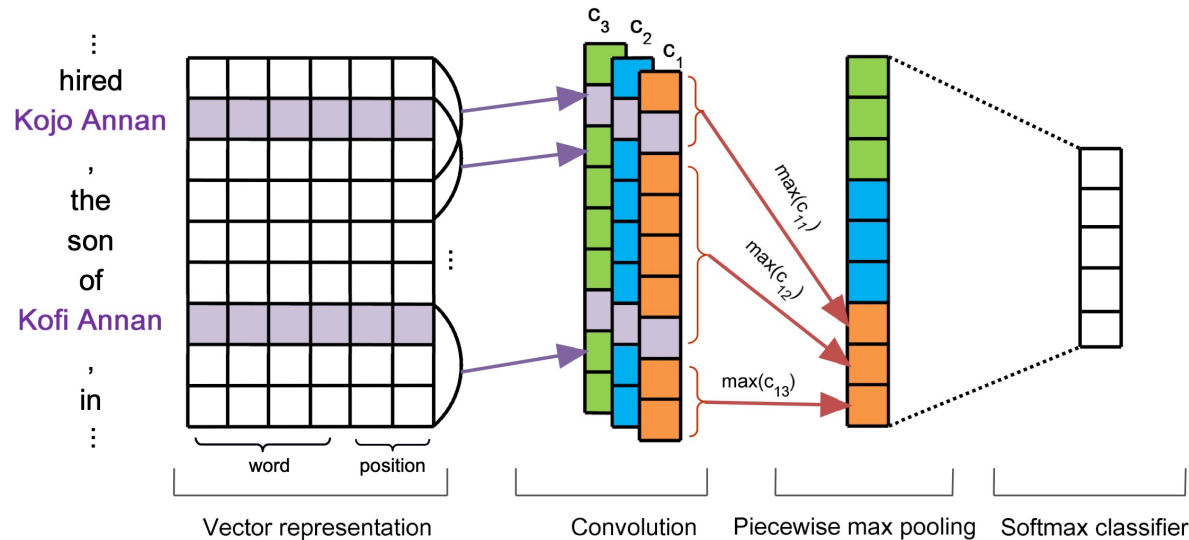


Figure 3: The architecture of PCNNs (better viewed in color) used for distant supervised relation extraction, illustrating the procedure for handling one instance of a bag and predicting the relation between *Kojo Annan* and *Kofi Annan*.

Algorithm 1 Multi-instance learning

- 1: Initialize θ . Partition the bags into mini-batches of size b_s .
- 2: Randomly choose a mini-batch, and feed the bags into the network one by one.
- 3: Find the j -th instance m_i^j ($1 \leq i \leq b_s$) in each bag according to Eq. (9).
- 4: Update θ based on the gradients of m_i^j ($1 \leq i \leq b_s$) via Adadelta.
- 5: Repeat steps 2-4 until either convergence or the maximum number of epochs is reached.

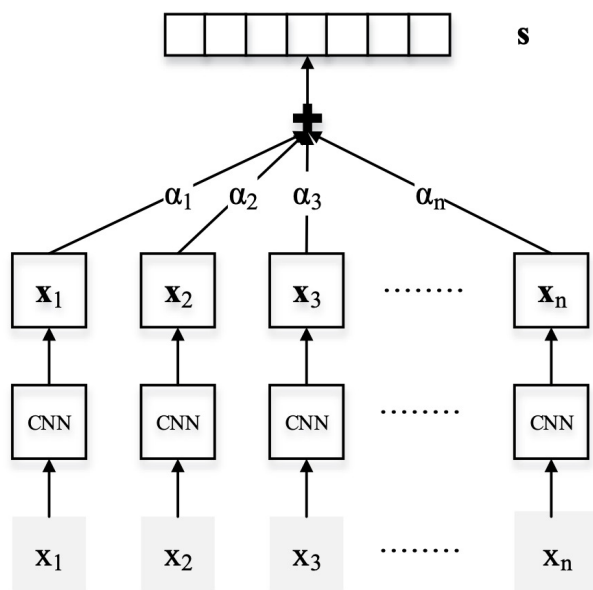


Figure 1: The architecture of sentence-level attention-based CNN, where x_i and \mathbf{x}_i indicate the original sentence for an entity pair and its corresponding sentence representation, α_i is the weight given by sentence-level attention, and s indicates the representation of the sentence set.

Attention

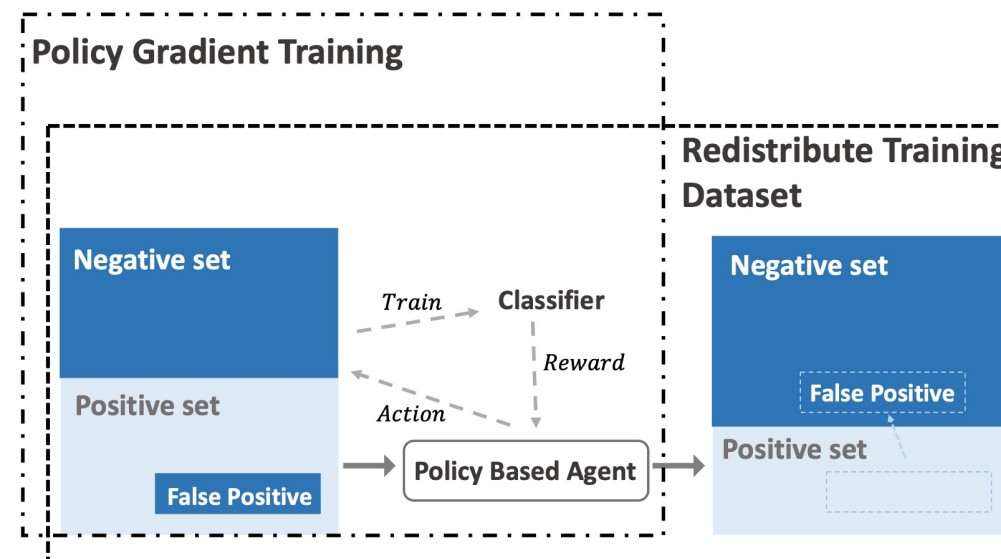


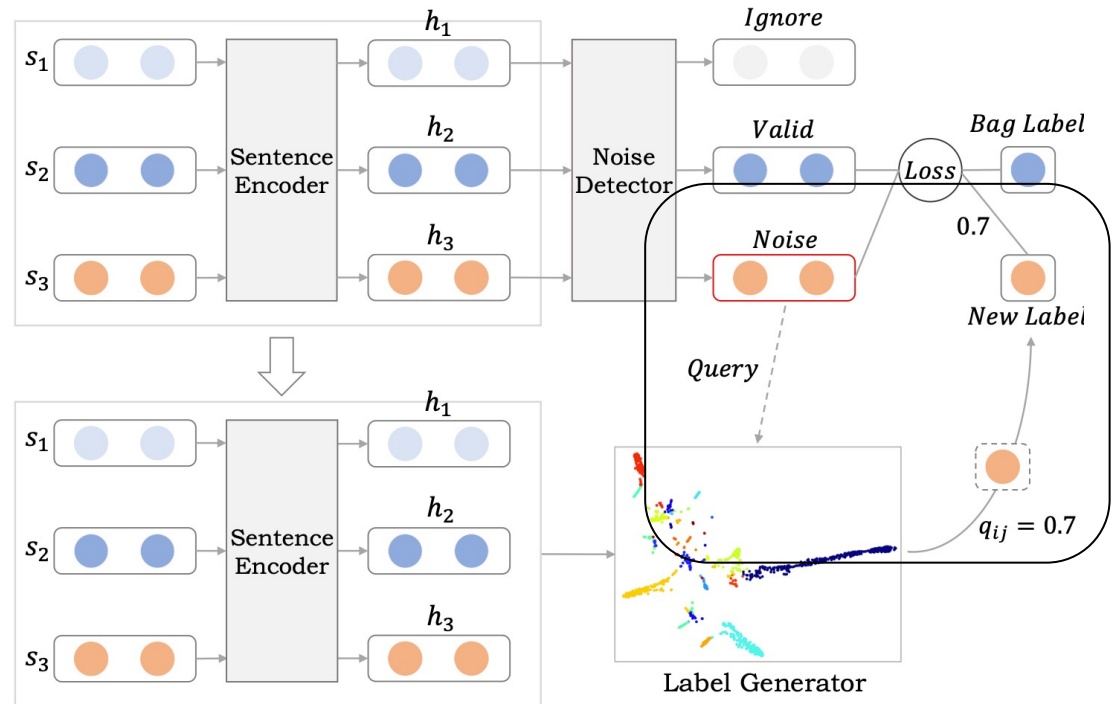
Figure 1: Our deep reinforcement learning framework aims at dynamically recognizing false positive samples, and moving them from the positive set to the negative set during distant supervision.

Reinforcement Learning



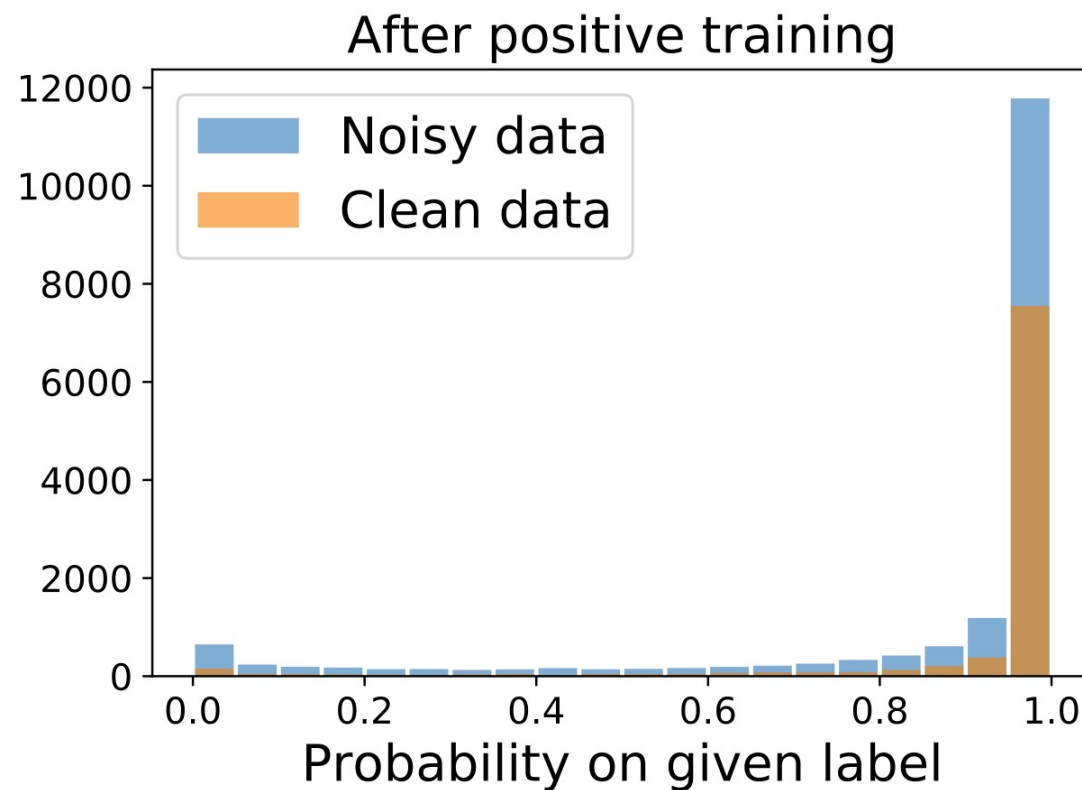
	Sentence	Bag Label	Noise?	Correct Label
Bag	#1: Barack Obama was born in the United States .		Yes	born in
	#2: Barack Obama was the first African American to be elected to the president of the United States .	president of	No	president of
	#3: Barack Obama served as the 44th president of the United States from 2009 to 2017.		No	president of

Table 1: An example of sentence-bag annotated by distant supervision. “Yes” and “No” indicate whether or not each sentence is a noisy sentence. “Correct Label” means the true relationship between the entity pair expressed in each sentence.





$$\mathcal{L}_{PT}(f, y^*) = - \sum_{k=1}^C y_k \log p_k$$



Positive Training 范式很难区分出干净数据与噪音数据



Positive Training



"is a **dog**"



Correct
Labels

"is a **cat**"

"is a **bird**"

"is an **elephant**"

"is a **chicken**"

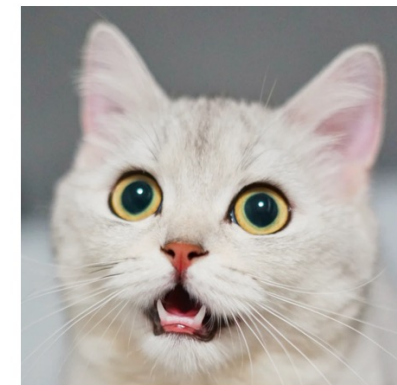
"is a **pig**"



Incorrect
labels

$$\mathcal{L}_{PT}(f, y^*) = - \sum_{k=1}^C y_k \log p_k$$

Negative Training



Correct
Labels

"is not a **dog**"

"is not a **bird**"

"is not an **elephant**"

"is not a **chicken**"

"is not a **pig**"

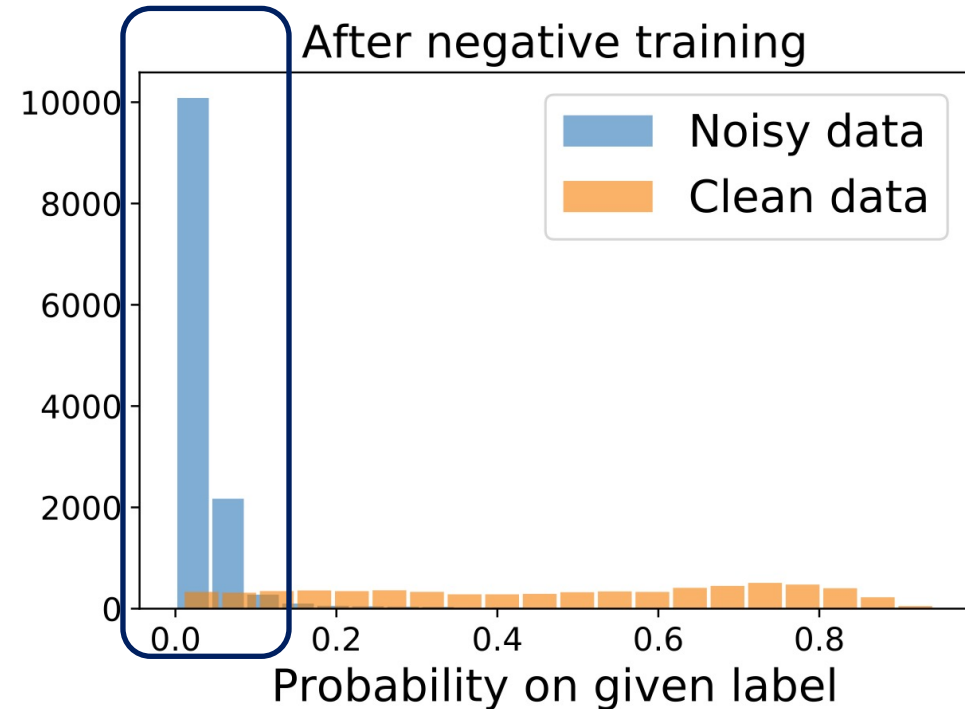
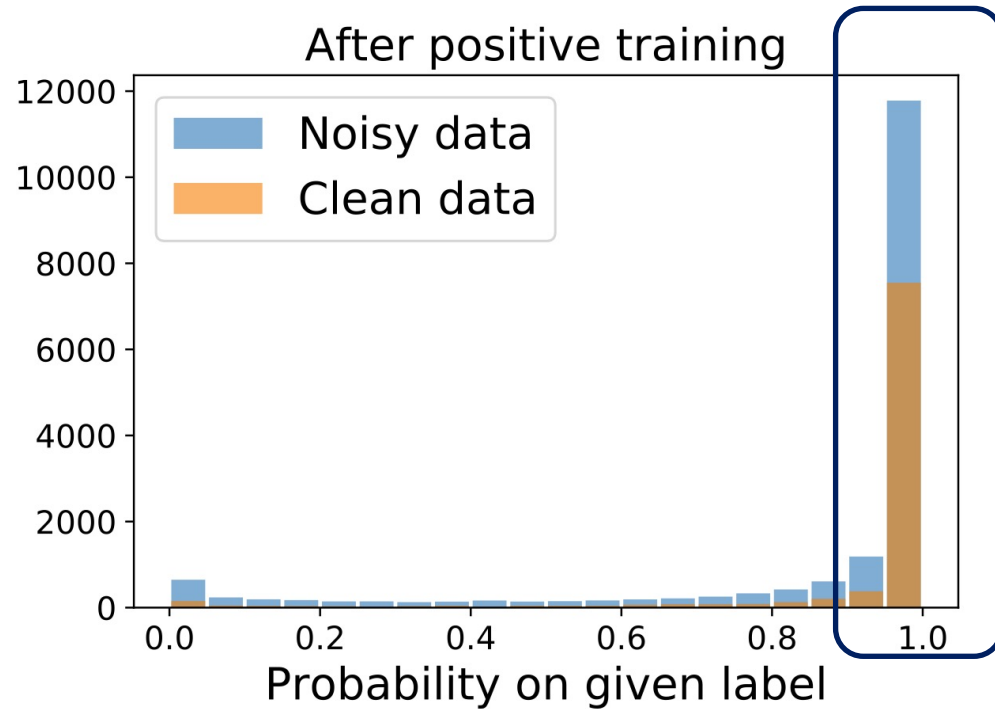


Incorrect
labels

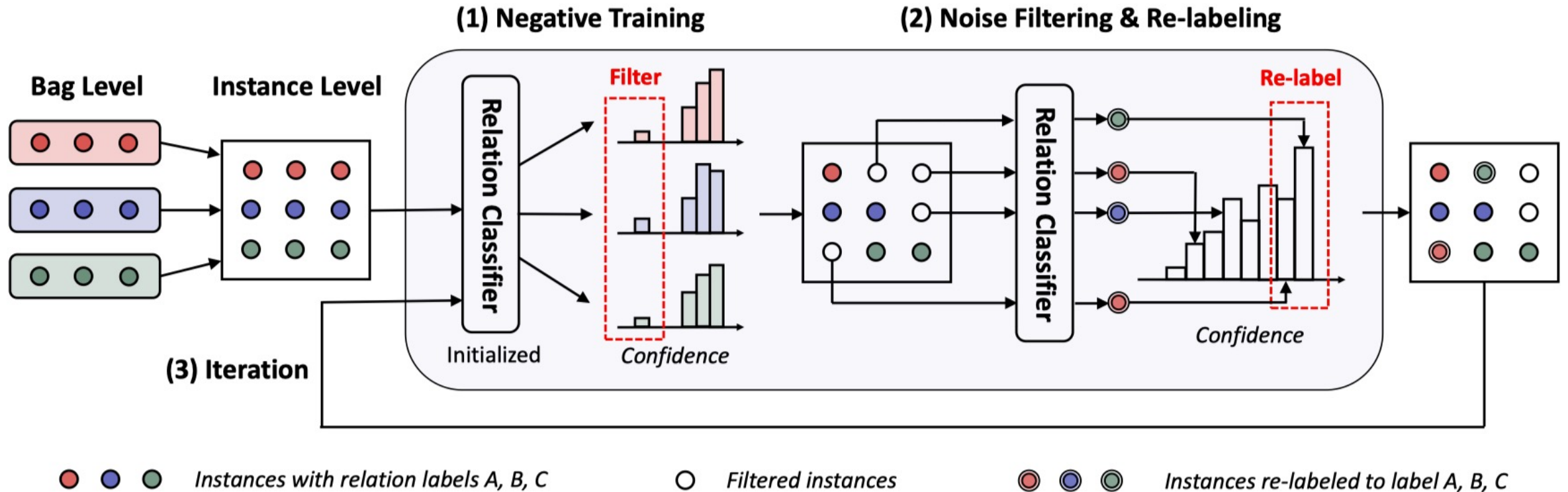


"is not a **cat**"

$$\mathcal{L}_{NT}(f, y^*) = - \sum_{k=1}^C \bar{y}_k \log(1 - p_k)$$

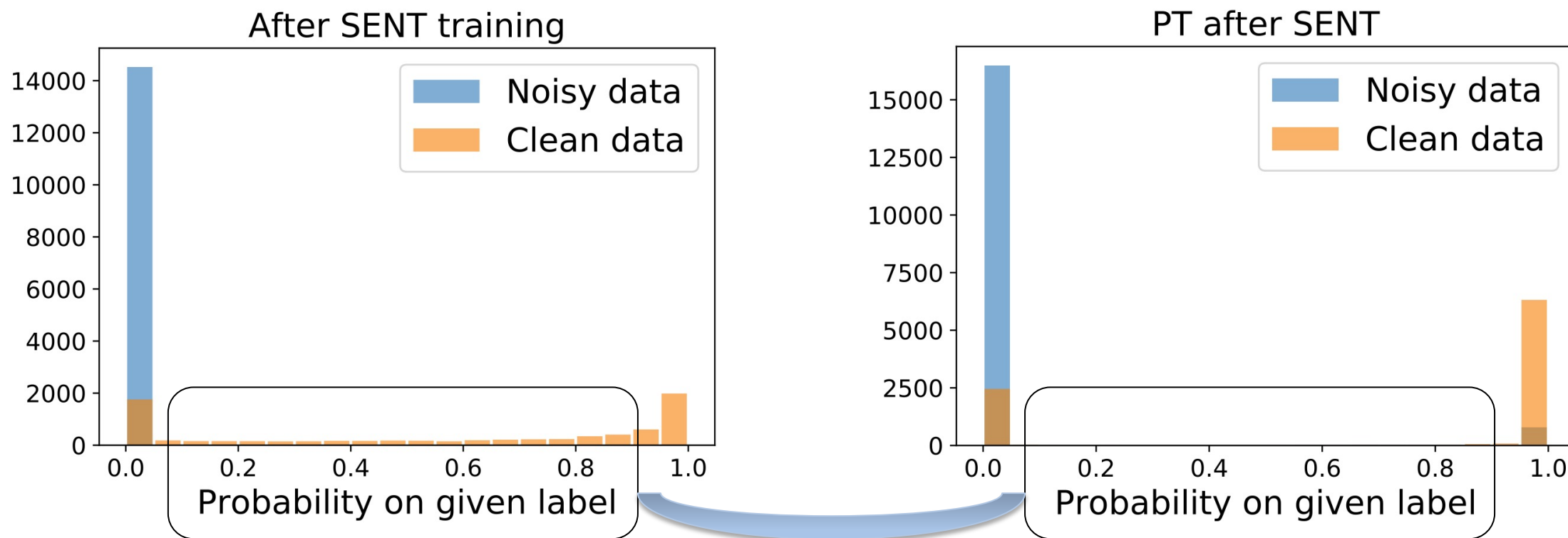


Comparison between positive and negative training



An overview of the proposed framework, SENT, for sentence-level distant RE.

- (1) **Negative training** for separating the noisy data from the training data
- (2) **Noise-filtering** and **re-labeling**
- (3) **Iterative training** to further boost the performance.



- (1) After SENT training, the clean and noisy data are further separated
- (2) PT after SENT helps improve the convergence of the clean data



Datasets		NYT-10	noisy-TACRED
Label		24	41
Train	inst.	371461	68124
	posi.	110518	13012
	noise	unknown	20586
Dev	inst.	2379	22631
	posi.	337	5436
Test	inst.	2164	15509
	posi.	323	3325

Table 1: Statistics of datasets.

Method	Dev			Test		
	Prec.	Rec.	F1	Prec.	Rec.	F1
CNN(Zeng et al., 2014)	38.32	65.22	48.28	35.75	64.54	46.01
PCNN(Zeng et al., 2015)	36.09	63.66	46.07	36.06	64.86	46.35
BiLSTM(Zhang et al., 2015)	36.71	66.46	47.29	35.52	67.41	46.53
BiLSTM+ATT(Zhang et al., 2017)	37.59	64.91	47.61	34.93	65.18	45.48
BERT(Devlin et al., 2019)	34.78	65.17	45.35	36.19	70.44	47.81
BiLSTM+BERT(Devlin et al., 2019)	36.09	73.17	48.34	33.23	72.70	45.61
PCNN+SelATT(Lin et al., 2016)	46.01	30.43	36.64	45.41	30.03	36.15
PCNN+RA_BAG_ATT(Ye and Ling, 2019)	49.84	46.90	48.33	56.76	50.60	53.50
CNN+RL ₁ (Qin et al., 2018)	37.71	52.66	43.95	39.41	61.61	48.07
CNN+RL ₂ (Feng et al., 2018)	40.00	59.17	47.73	40.23	63.78	49.34
ARNOR(Jia et al., 2019)	62.45	58.51	60.36	65.23	56.79	60.90
SENT (BiLSTM)	66.44	56.97	61.34	71.80	59.13	64.86
SENT (BiLSTM+BERT)	69.43	63.72	66.45	75.78	63.82	69.29

Table 2: Main results on sentence-level evaluation. Compared baselines include normal RE model (the first part of the table), and models for distant RE (the second part of the table)





323 Samples including 200 incorrect samples

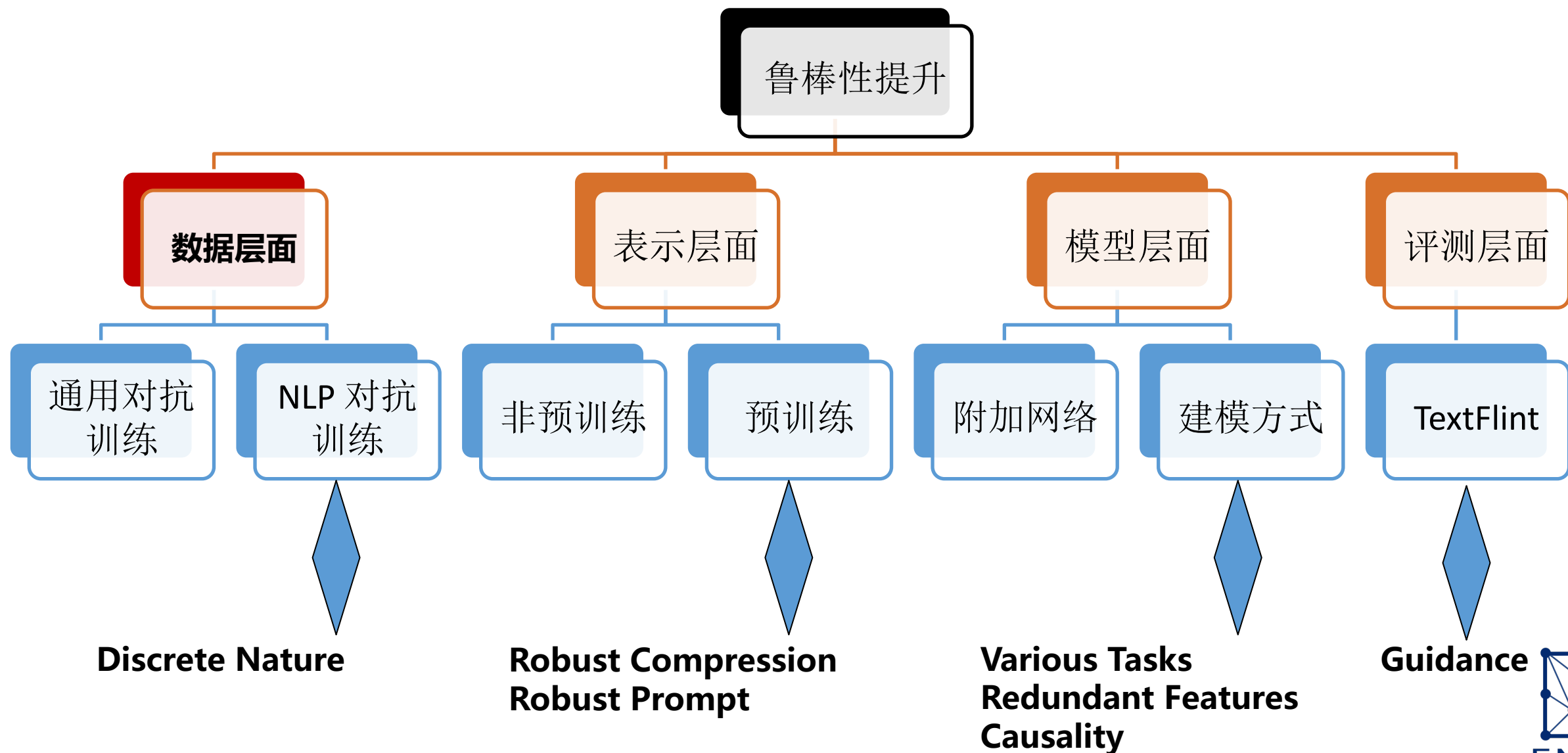
Noise Reduction	Prec.	Rec.	F1
CNN+RL ₂	40.58	96.31	57.10
ARNOR	76.37	68.13	72.02
SENT (biLSTM)	80.00	88.46	84.02
SENT (biLSTM+BERT)	84.33	85.67	84.99

Table 3: The noise-filtering effect evaluated on a noise-annotated test set of NYT-10.

13012 correct samples
20586 incorrect samples

	Method	Prec.	Rec.	F1
Clean Data	BiLSTM+ATT	67.7	63.2	65.4
	BiLSTM	61.4	61.7	61.5
Noisy Data	BiLSTM+ATT	32.8	43.8	37.5
	BiLSTM	37.8	45.5	41.3
	SENT (biLSTM)	66.0	52.9	58.7

Table 4: Model performance on clean and noisy-TACRED. When trained on noisy data, the performance of base models degrade dramatically while SENT achieves comparable results with the models trained on clean data.





Thanks for your attention!



FNL P

