

# PIEClass: Task Weakly-Supervised Text Classification Method with Prompting and Noise-Robust Iterative Ensemble Training

Source: EMNLP 2023  
Advisor: JIA-LING KOH  
Speaker: FAN-CHI-YU  
Date: 2023/02/27

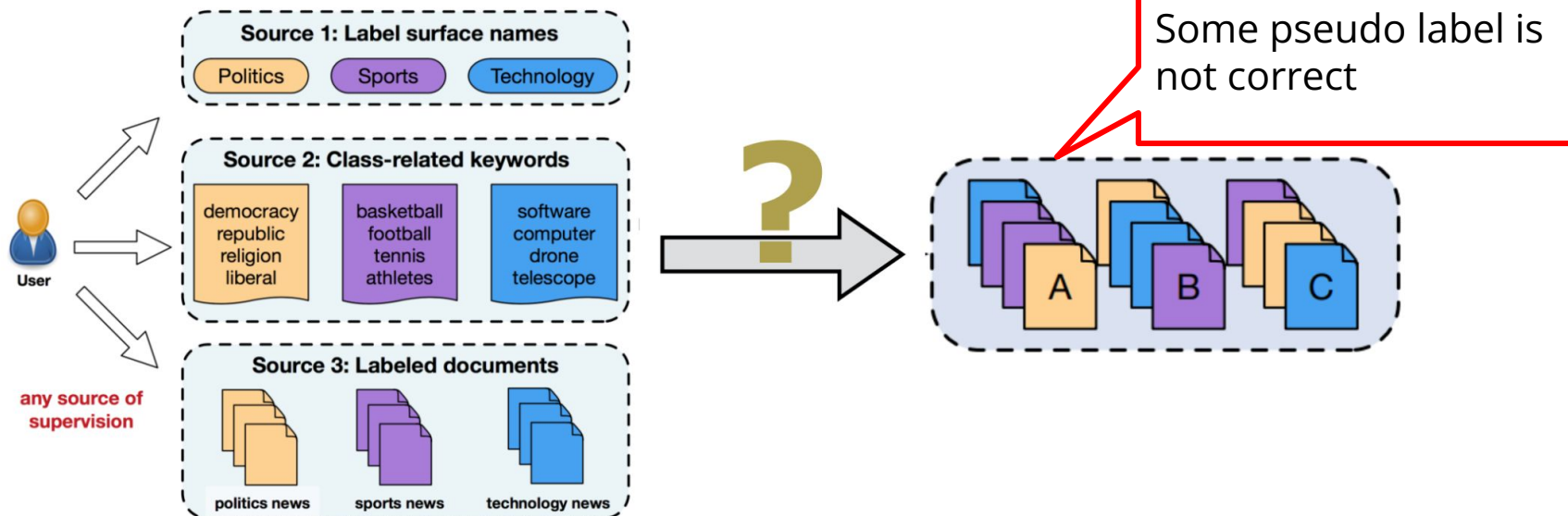
# Outline

- Introduction
- Method
- Experiment
- Conclusion

# Introduction

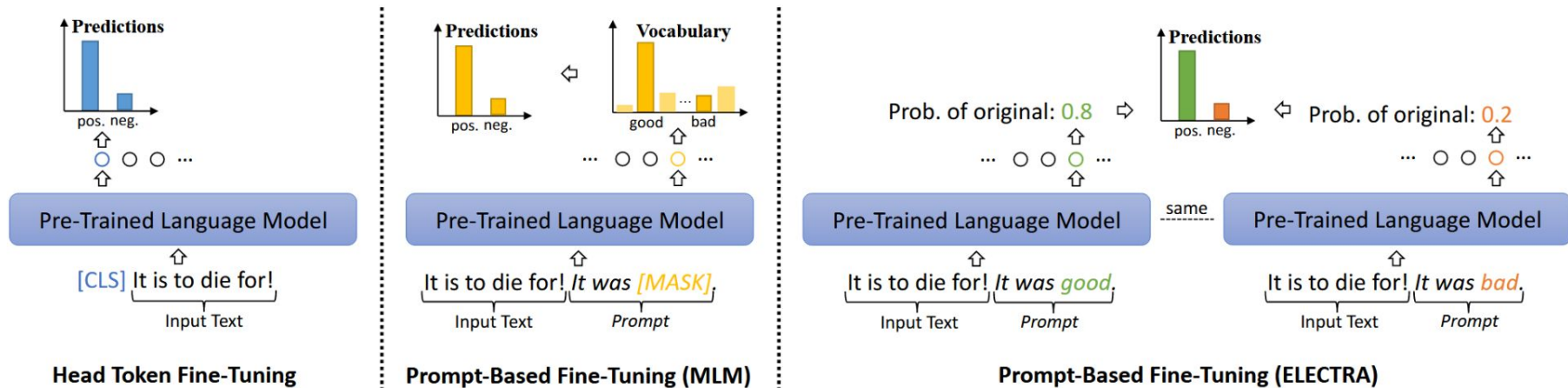
# Weakly-Supervised Text Classification

Any **labeled documents** are not allowed, **surface names** or **limited word-level descriptions** of each category can be used.



# Fine-Tuning

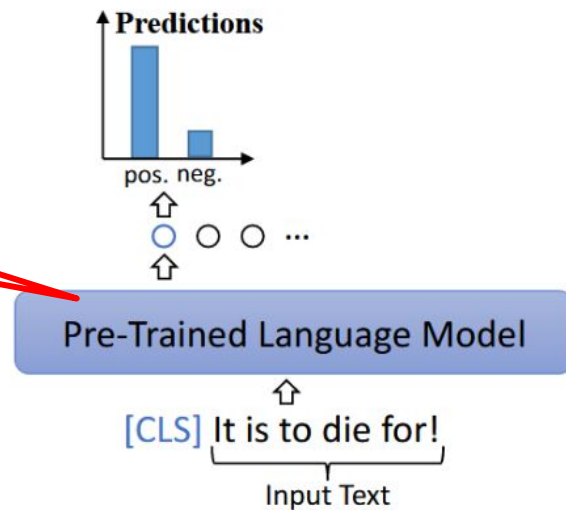
## Type of fine tuning



# Head Token Fine-Tuning

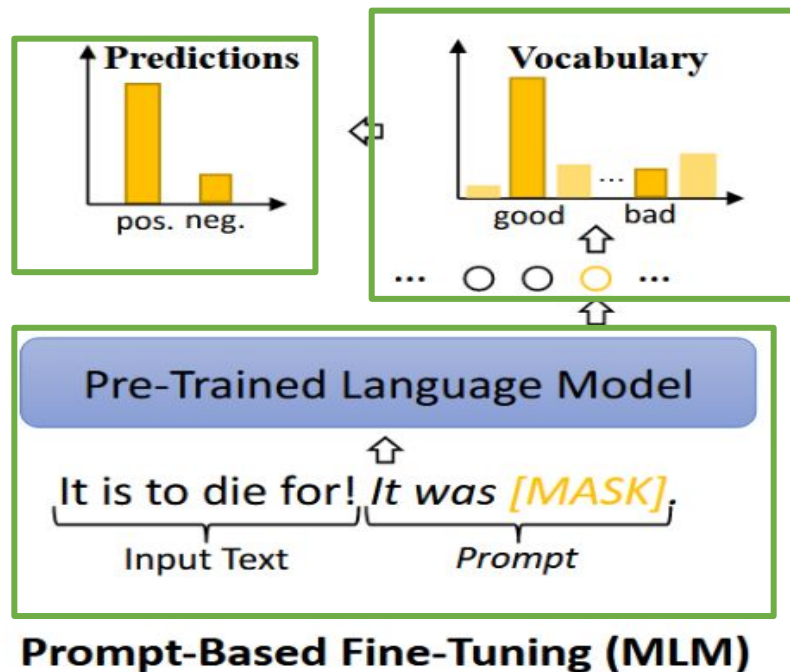
Classifier

$$p(c|d) = \text{Softmax}(g(\mathbf{h}^{\text{CLS}}))$$



**Head Token Fine-Tuning**

# Prompt-Base Fine-Tuning(MLM)



$$\mathcal{T}^{\text{MLM}}(d) = d \text{ It was [MASK]}.$$

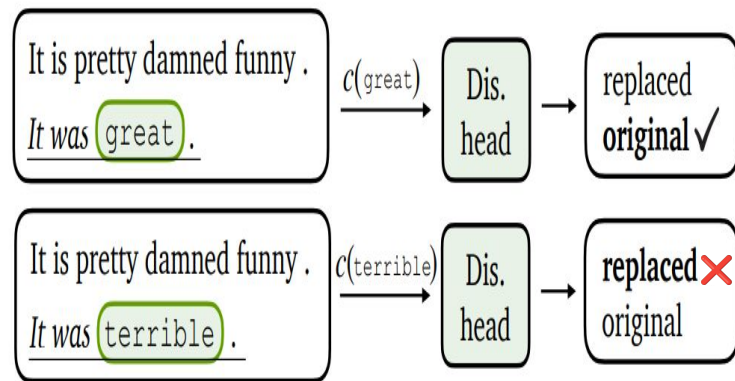
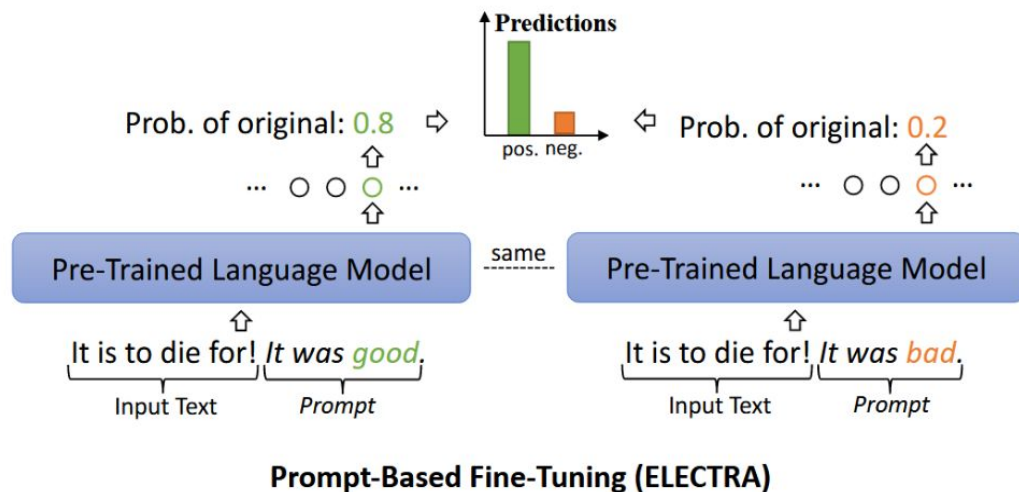


$$p(w|d) = \text{Softmax}(f(\mathbf{h}^{\text{MASK}})). \quad (7)$$



$$p(l(c)|d)$$

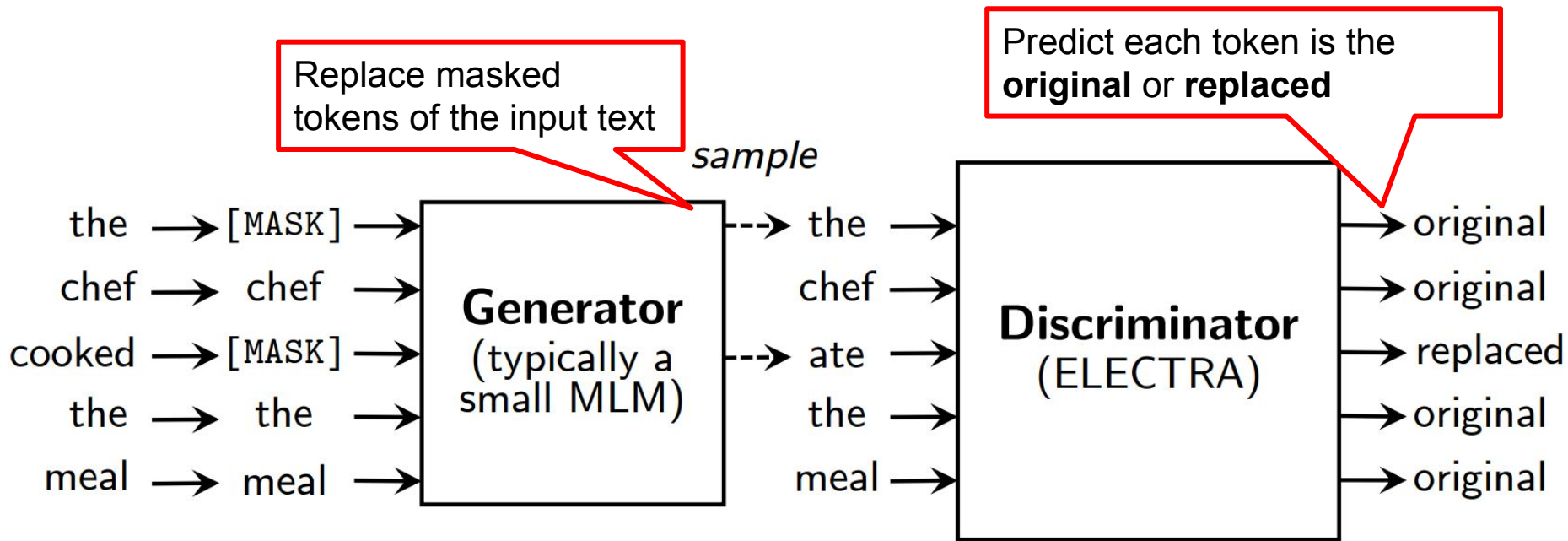
# Prompt-Base Fine-Tuning(ELECTRA)





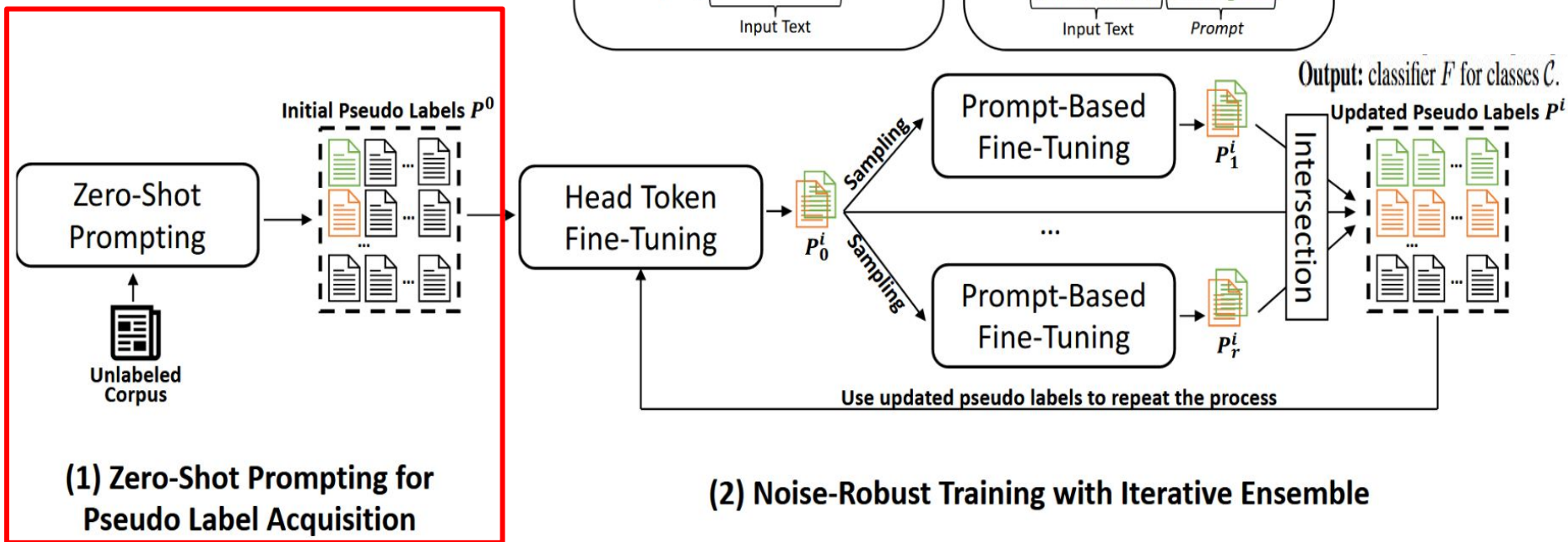
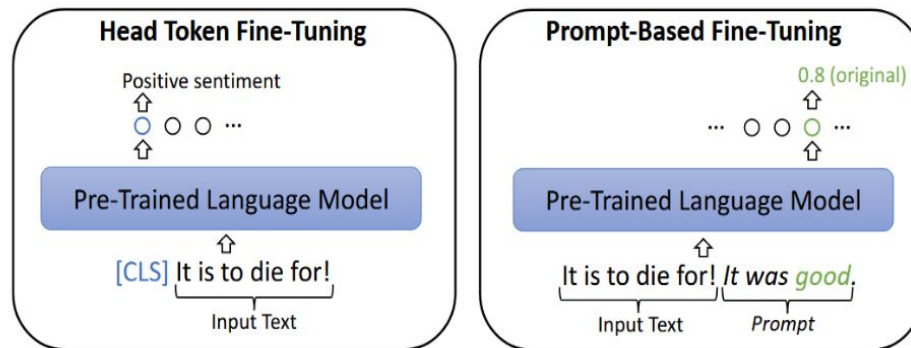
# ELECTRA Pre-train model

Cast the word prediction problem into a binary classification problem



# Method

Two fine-tuning strategies for pre-trained language model

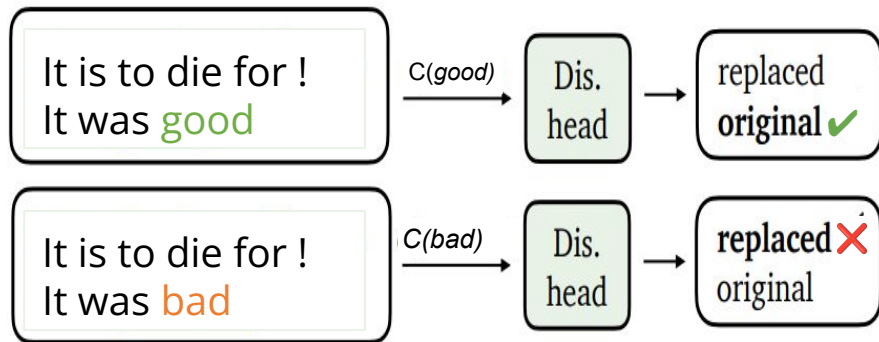


# Zero-Shot Prompting for Pseudo Label Acquisition

Construct input with the template

$$\mathcal{T}^{\text{ELECTRA}}(d, \text{good}) = d \text{ It was } \underline{\text{good}}.$$

$$\mathcal{T}^{\text{ELECTRA}}(d, \text{bad}) = d \text{ It was } \underline{\text{bad}}.$$



$\mathcal{T}(d, l(c)) \leftarrow$  Construct input with the template;

# Zero-Shot Prompting for Pseudo Label Acquisition

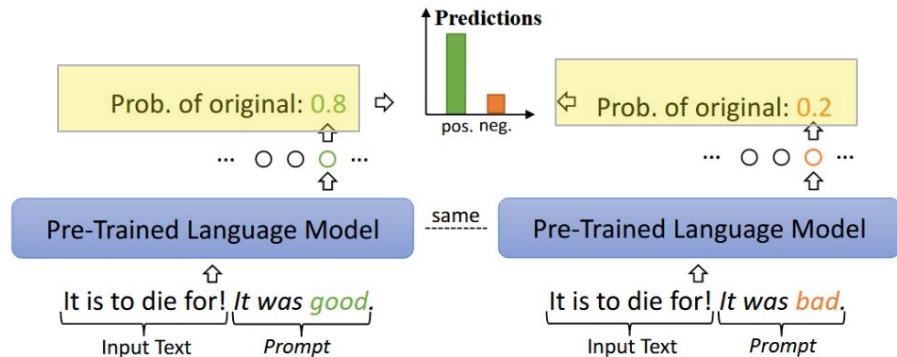
$$p(l(c)|d) = \text{Sigmoid}(f(h^{l(c)})), \quad (1)$$



$$p(c|d) = \frac{p(l(c)|d)}{\sum_{c' \in \mathcal{C}} p(l(c')|d)}. \quad (2)$$



$$P^0 = \text{top}k(p(c|d))$$



Prompt-Based Fine-Tuning (ELECTRA)

$$p(l(\bar{c})|d) \leftarrow \text{Prompt } E \text{ with Eq. (1)}$$

# Zero-Shot Prompting for Pseudo Label Acquisition

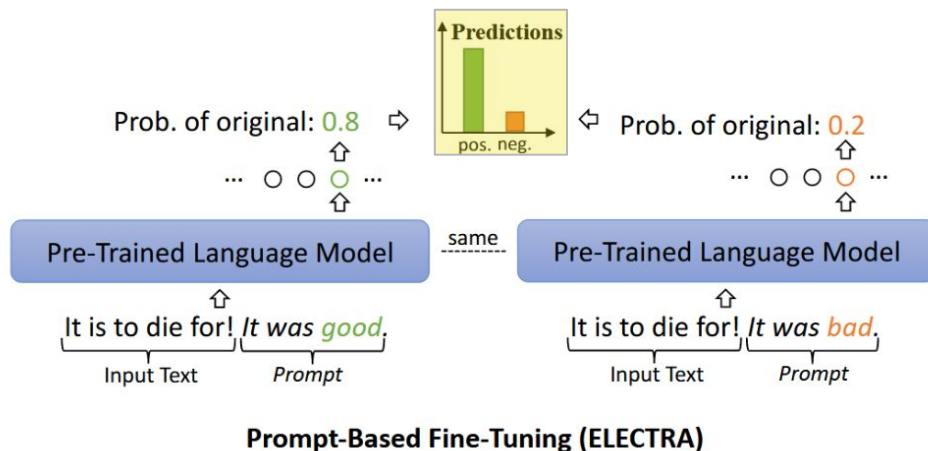
$$p(l(c)|d) = \text{Sigmoid}(f(\mathbf{h}^{l(c)})), \quad (1)$$



$$p(c|d) = \frac{p(l(c)|d)}{\sum_{c' \in \mathcal{C}} p(l(c')|d)}. \quad (2)$$



$$P^0 = \text{topk}(p(c|d))$$



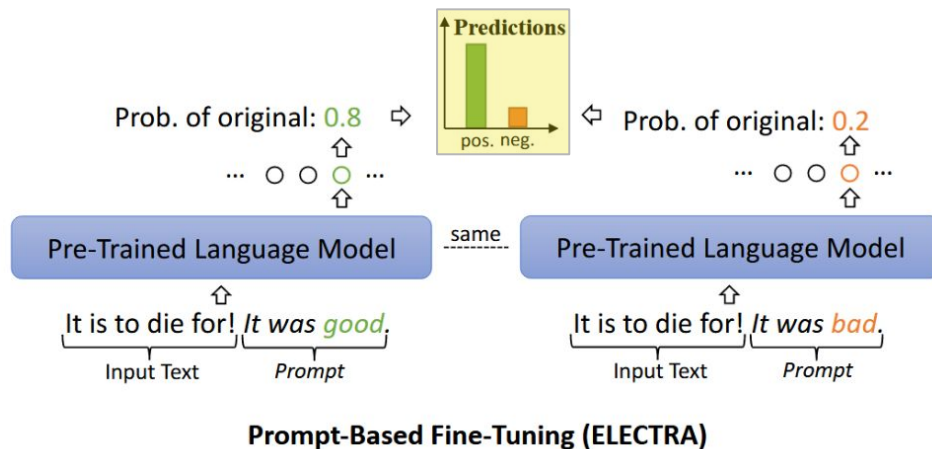
$$p(c|d) \leftarrow \text{Eq. (2)};$$

# Zero-Shot Prompting for Pseudo Label Acquisition

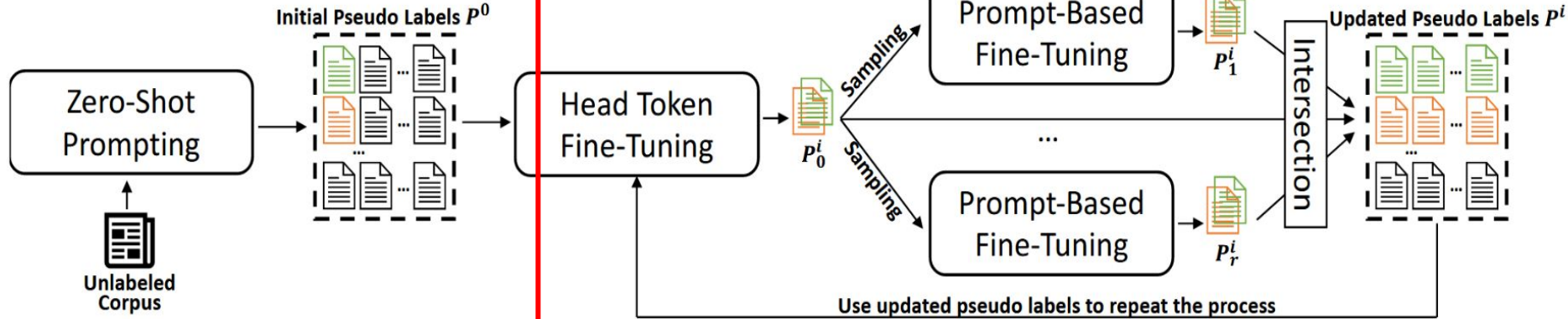
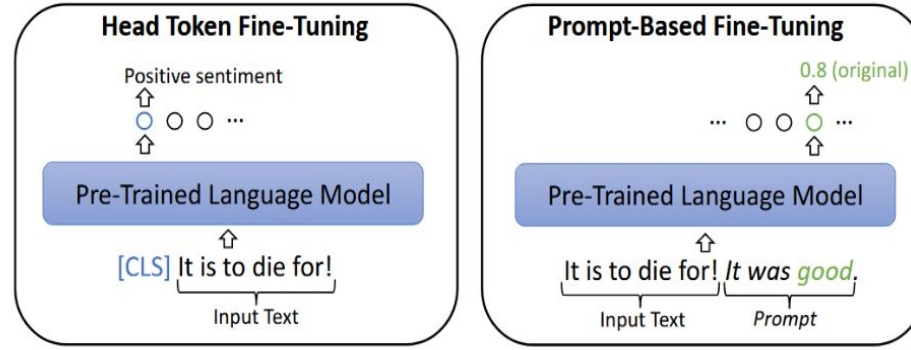
$$p(l(c)|d) = \text{Sigmoid}(f(\mathbf{h}^{l(c)})), \quad (1)$$

$$p(c|d) = \frac{p(l(c)|d)}{\sum_{c' \in \mathcal{C}} p(l(c')|d)}. \quad (2)$$

$$P^0 = \text{top}k(p(c|d))$$



Two fine-tuning strategies for pre-trained language model

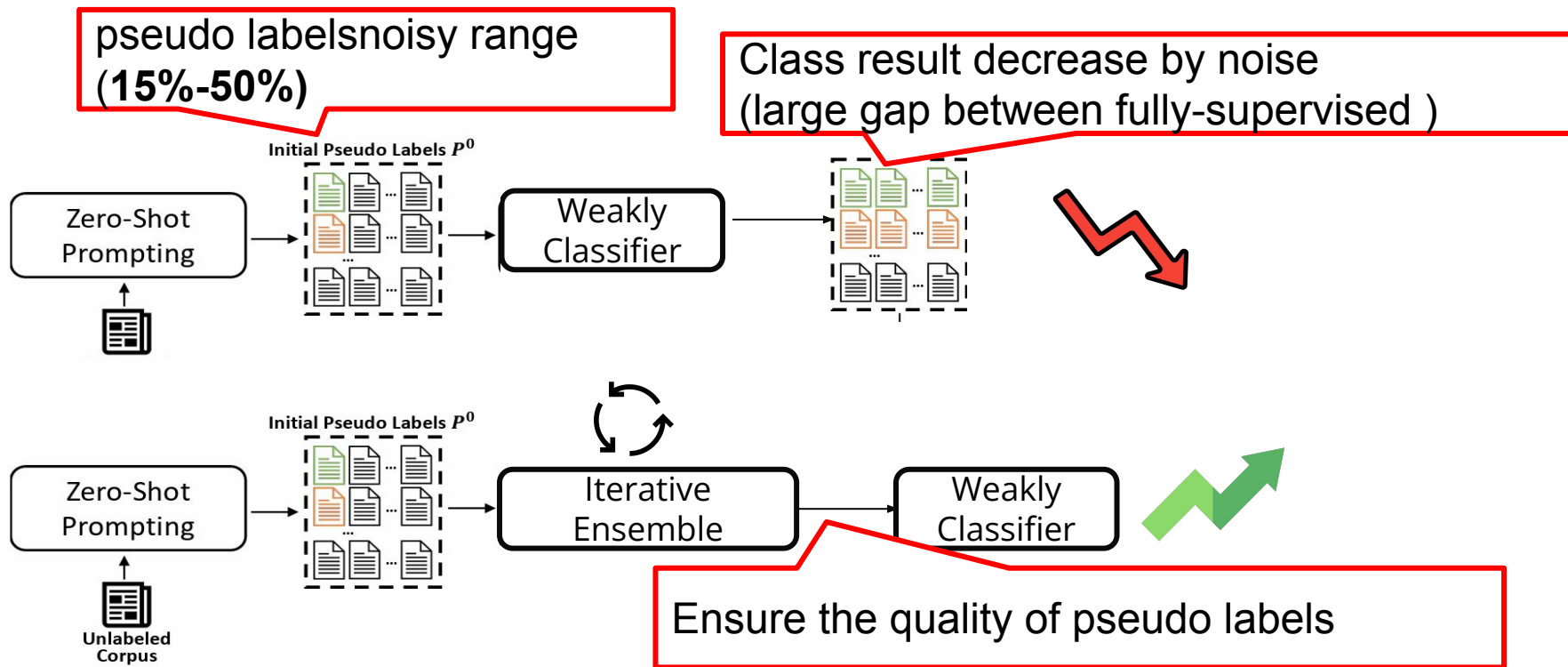


(1) Zero-Shot Prompting for Pseudo Label Acquisition

(2) Noise-Robust Training with Iterative Ensemble



# Noise-Robust Training with Iterative Ensemble



# Noise-Robust Training with Iterative Ensemble

for  $i \leftarrow 1$  to  $T$  do

$F_0^i \leftarrow$  Head token fine-tuning using  $P^{i-1}$ ;

$P_0^i \leftarrow$  Select top  $t_i$  predictions by  $F_0^i$ ;

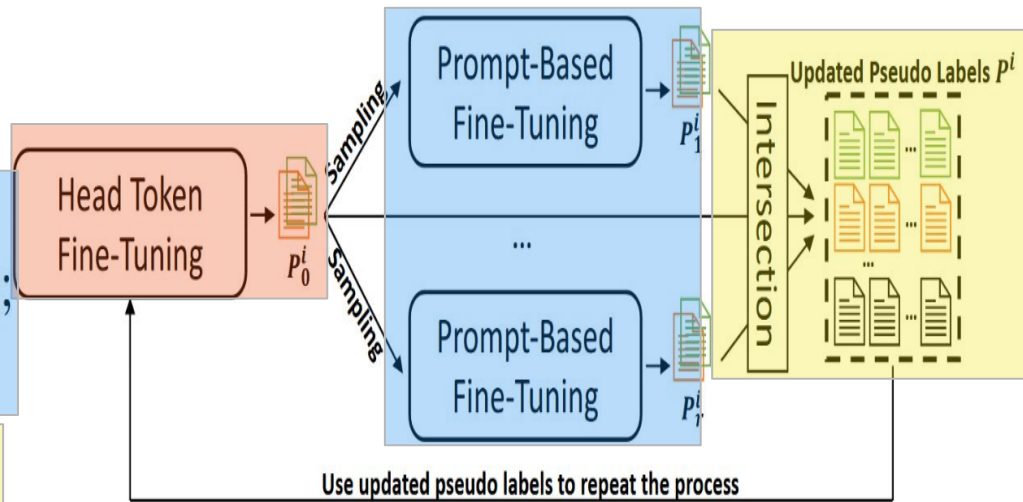
$\mathcal{S} \leftarrow$  Randomly sample  $r$  subsets of  $P_0^i$ ;

for  $S_k \in \mathcal{S}$  do

$F_k^i \leftarrow$  Prompt-based fine-tuning using  $S_k$ ;

$P_k^i \leftarrow$  Select top  $t_i$  percentage by  $F_k^i$ ;

$P^i \leftarrow$  Eq. (4);



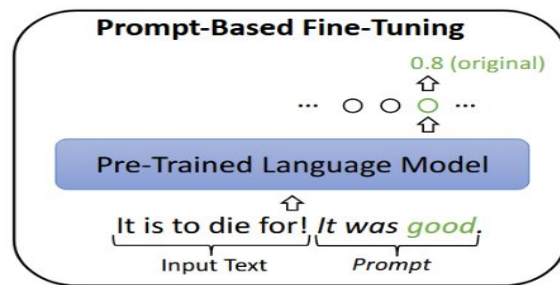
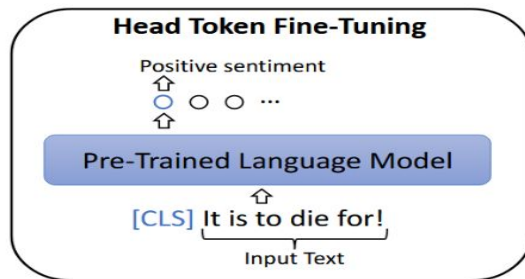
# Noise-Robust Training with Iterative Ensemble

Utilize two PLM fine-tuning methods to ensure the quality of pseudo labels improve the self-training quality

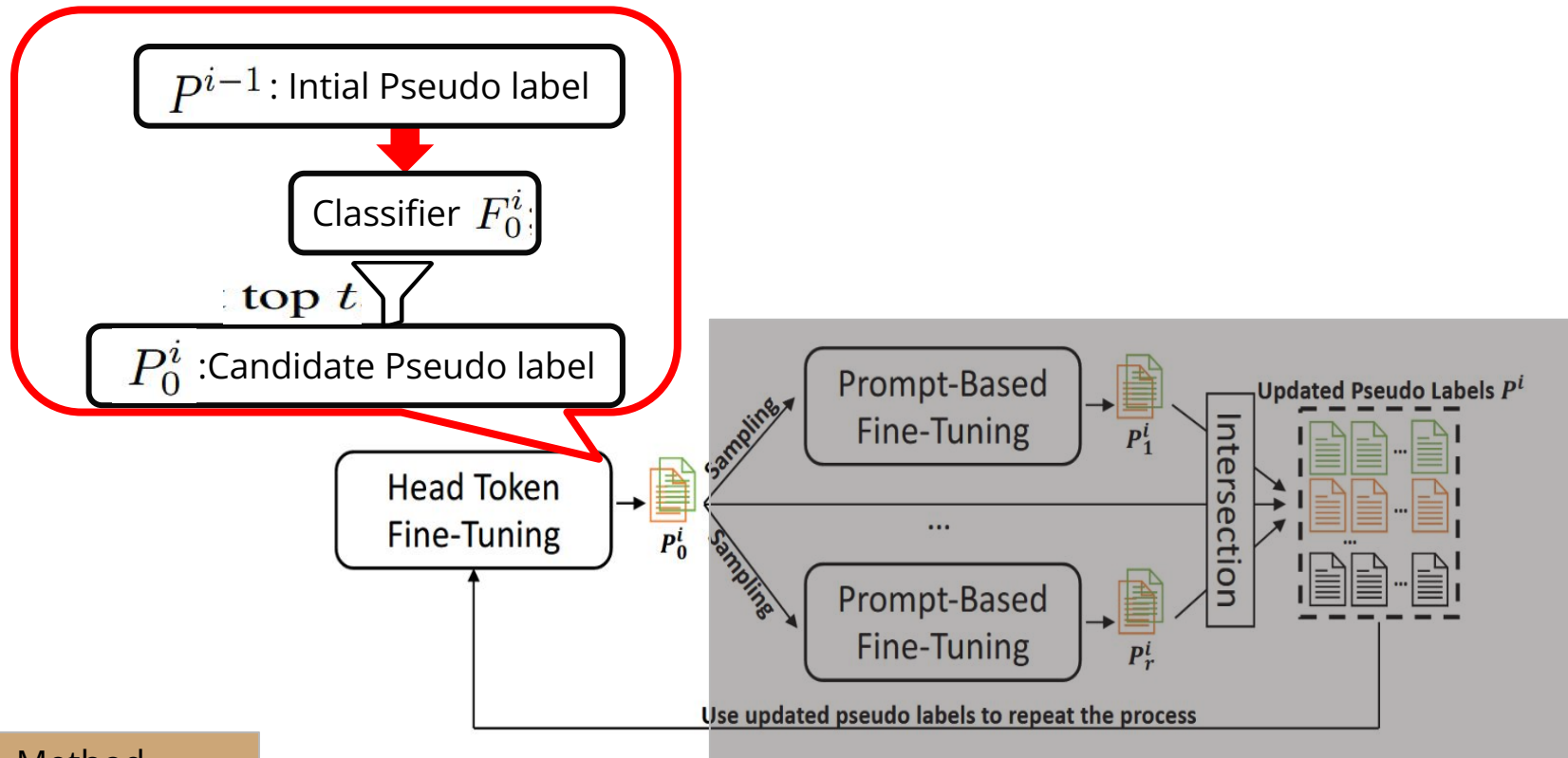
1. **Head token fine-tuning:** Capturing the information of the entire document

1. **Prompt-based finetuning:** Focusing more on the context surrounding the

Two fine-tuning strategies for pre-trained language model

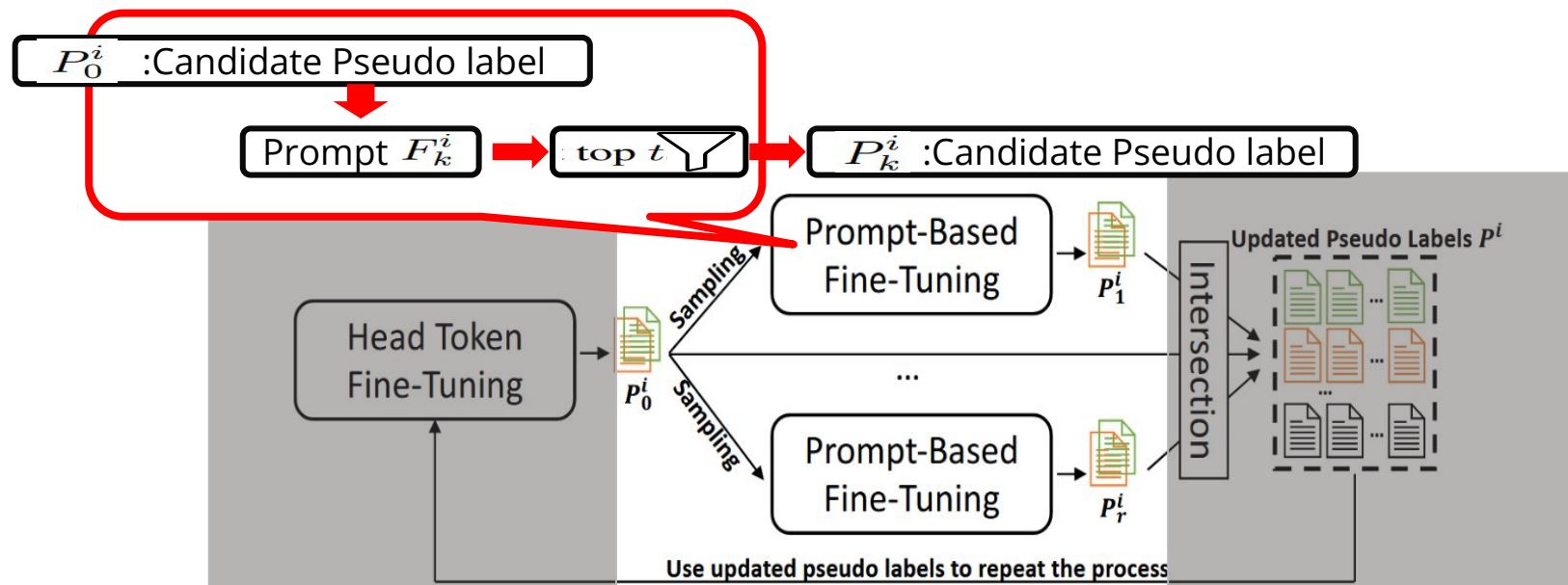


# Noise-Robust Training with Iterative Ensemble



# Noise-Robust Training with Iterative Ensemble

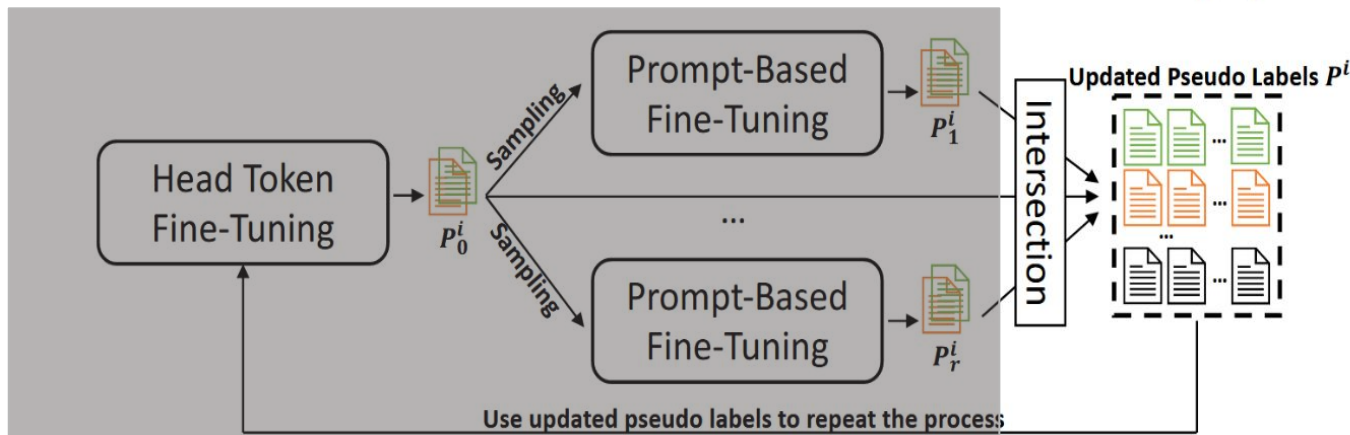
Prompt base only requires a small amount of data to achieve competitive performance with head token fine-tuning



# Noise-Robust Training with Iterative Ensemble

Only those most confident ones into the pseudo label pool to alleviate the error accumulation problem.

$$\boxed{P_k^i : \text{Candidate Pseudo label}} \xrightarrow{\text{Intersection}} \boxed{\text{Intersection}} \xrightarrow{\text{Intersection}} \mathcal{P}^i = \bigcap_{k=0}^i \mathcal{P}_k^i. \quad (4)$$



# Experiment

# DataSet

- Topic
  - Ag\_News(New topic with 4 class)
  - 20\_News (New topic with 20 class)
  - NYT-Topics (New York Times context: imbalanced with 9 class)
  - NYT-Fine (New York Times context: imbalanced & fine-grained with 9 class)
- Semantic(with 2 class)
  - Yelp(Review:Semantic analysis )
  - IMDB(Movie Review: semantic analysis )
  - Amazon(Amazon Review:semantic analysis )



# Compared Methods

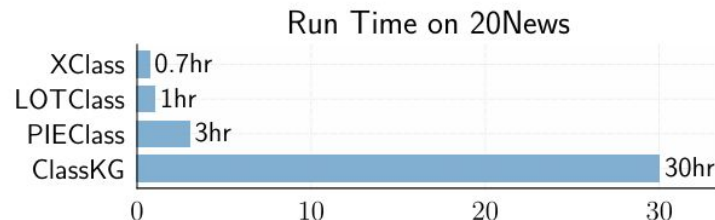
- Weakly method compare
  - WeSTClass
  - ConWea
  - LOTClass
  - XClass
  - ClassKG
- Pre-train model compare
  - RoBERTa (0-shot):Head Token
  - ELECTRA (0-shot):Head Token
  - Fully- Supervised BERT baseline

# Compared Methods

Although ClassKG achieves the better results ClassKG uses more time

Methods	AGNews	20News	NYT-Topics	NYT-Fine	Yelp	IMDB	Amazon
<b>WeSTClass</b>	0.823/0.821	0.713/0.699	0.683/0.570	0.739/0.651	0.816/0.816	0.774/-	0.753/-
<b>ConWea</b>	0.746/0.742	0.757/0.733	<u>0.817/0.715</u>	0.762/0.698	0.714/0.712	-/-	-/-
<b>LOTClass</b>	0.869/0.868	0.738/0.725	0.671/0.436	0.150/0.202	0.878/0.877	0.865/-	0.916/-
<b>XClass</b>	0.857/0.857	0.786/0.778	0.790/0.686	0.857/0.674	0.900/0.900	-/-	-/-
<b>ClassKG<sup>†</sup></b>	0.881/0.881	<b>0.811/0.820</b>	0.721/0.658	0.889/0.705	0.918/0.918	0.888/0.888	<u>0.926/-</u>
<b>PIEClass</b>							
<b>ELECTRA+ELECTRA</b>	<u>0.884/0.884</u>	<b>0.816/0.817</b>	<b>0.832/0.763</b>	<b>0.910/0.776</b>	<b>0.957/0.957</b>	<b>0.931/0.931</b>	<b>0.937/0.937</b>
<b>Fully-Supervised</b>	0.940/0.940	0.965/0.964	0.943/0.899	0.980/0.966	0.957/0.957	0.945/-	0.972/-

Micro-F1/Macro-F1



# Compared Methods

Methods	AGNews	20News	NYT-Topics	NYT-Fine	Yelp	IMDB	Amazon
RoBERTa (0-shot)	0.581/0.529	0.507/0.445 <sup>‡</sup>	0.544/0.382	-/- <sup>‡</sup>	0.812/0.808	0.784/0.780	0.788/0.783
ELECTRA (0-shot)	0.810/0.806	0.558/0.529	0.739/0.613	0.765/0.619	0.820/0.820	0.803/0.802	0.802/0.801
<b>PIEClass</b>							
ELECTRA+BERT	0.884/0.884	0.789/0.791	0.807/0.710	0.898/0.732	0.919/0.919	0.905/0.905	0.858/0.858
RoBERTa+RoBERTa	<b>0.895/0.895</b>	0.755/0.760 <sup>‡</sup>	0.760/0.694	-/- <sup>‡</sup>	0.920/0.920	0.906/0.906	0.912/0.912
ELECTRA+ELECTRA	0.884/0.884	<b>0.816/0.817</b>	<b>0.832/0.763</b>	<b>0.910/0.776</b>	<b>0.957/0.957</b>	<b>0.931/0.931</b>	<b>0.937/0.937</b>
Fully-Supervised	0.940/0.940	0.965/0.964	0.943/0.899	0.980/0.966	0.957/0.957	0.945/-	0.972/-

Micro-F1/Macro-F1

# Ablation Study

- **Two-Stage:** Directly trains classifier using pseudo labels from zero-shot prompting
- **Single-View ST:** Standard self-training method(only using zero-shot pseudo label)
- **Co-Training:** W/O Regularize in step Intersection

# Ablation Study

- The single-view and two-stage method is not stable.
- Co-training ensures the consistency of model predictions, yielding great results.

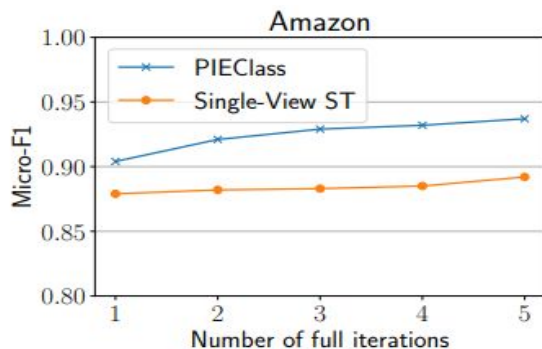
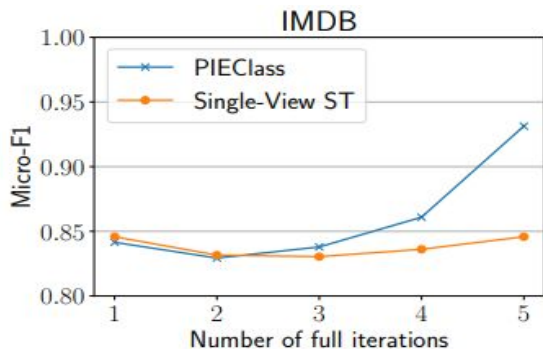
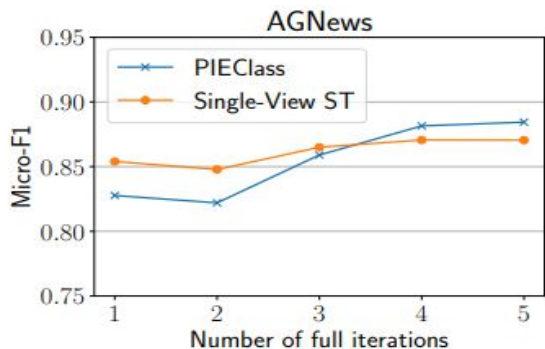
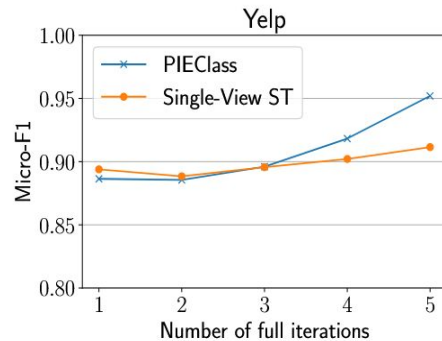
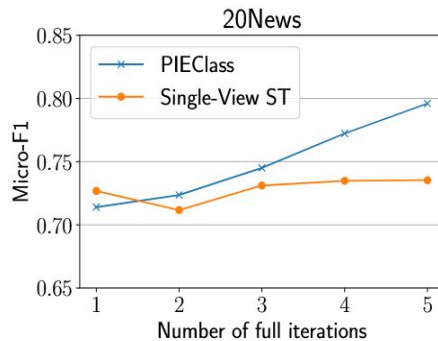
Methods	AGNews	20News	NYT-Topics	NYT-Fine	Yelp	IMDB	Amazon
<b>Two-Stage</b>	0.847/0.847	0.739/0.733	0.776/0.664	0.838/0.678	0.913/0.913	0.870/0.870	0.836/0.835
<b>Single-View ST</b>	0.871/0.871	0.736/0.737	0.757/0.668	0.853/0.695	0.912/0.912	0.846/0.846	0.892/0.892
<b>Co-Training</b>	0.877/0.877	0.795/0.791	0.818/0.715	0.877/0.744	0.948/0.948	0.925/0.925	0.930/0.930
<b>PIEClass</b>	<b>0.884/0.884</b>	<b>0.816/0.817</b>	<b>0.832/0.763</b>	<b>0.910/0.776</b>	<b>0.957/0.957</b>	<b>0.931/0.931</b>	<b>0.937/0.937</b>

Micro-F1/Macro-F1

# Ablation Study

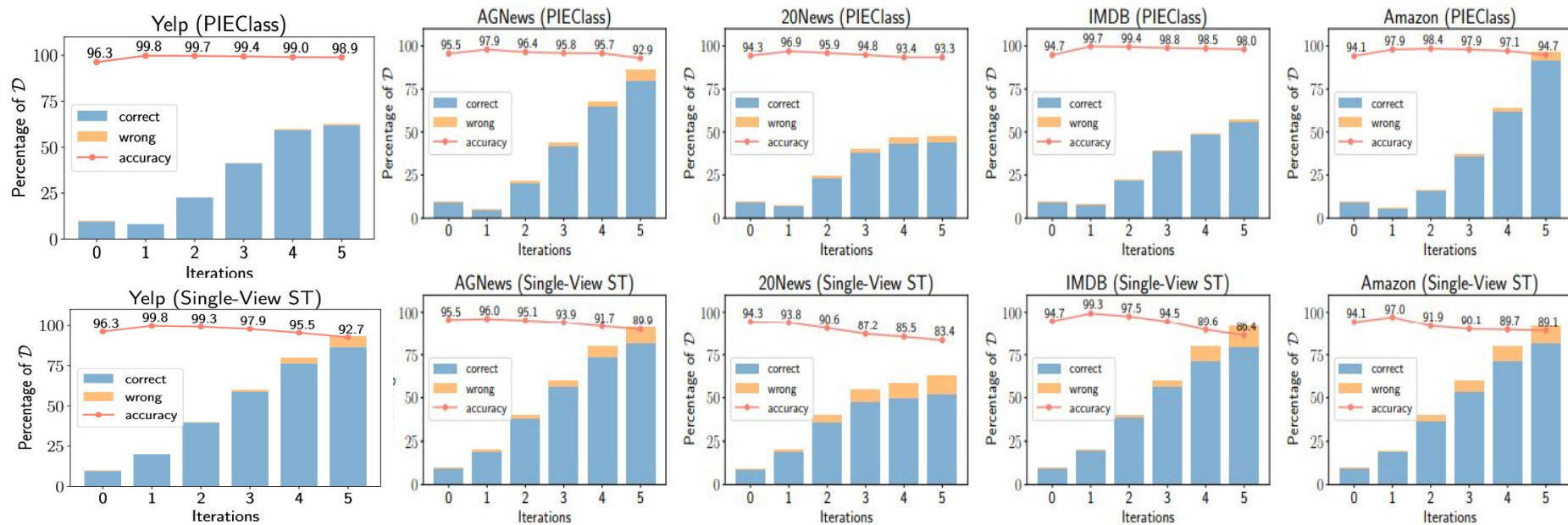
The PIEClass can surpass the bottleneck of traditional self-learning.

Traditional self-learning micor-f1 will be flattened after several iterations.



# Quantities and qualities of the pseudo labels

We can see at the **first servals iteration** the pseudo label qualities in well.



# Conclusion



# Conclusion

1. Using zero-shot PLM prompting to assign pseudo labels based on contextualized text understanding.
2. Implementing a noise-robust iterative ensemble to expand pseudo labels while ensuring their quality.

# Personal Comment

- In this paper, the noise-robust approach is crucial. Fully embracing it could significantly improve model adaptability in noisy environments.