



method

Prototype-Guided Pseudo Labeling for Semi-Supervised Text Classification

task

Advisor : Jia-Ling, Koh

Speaker : Ting-I, Weng

Source : ACL'23

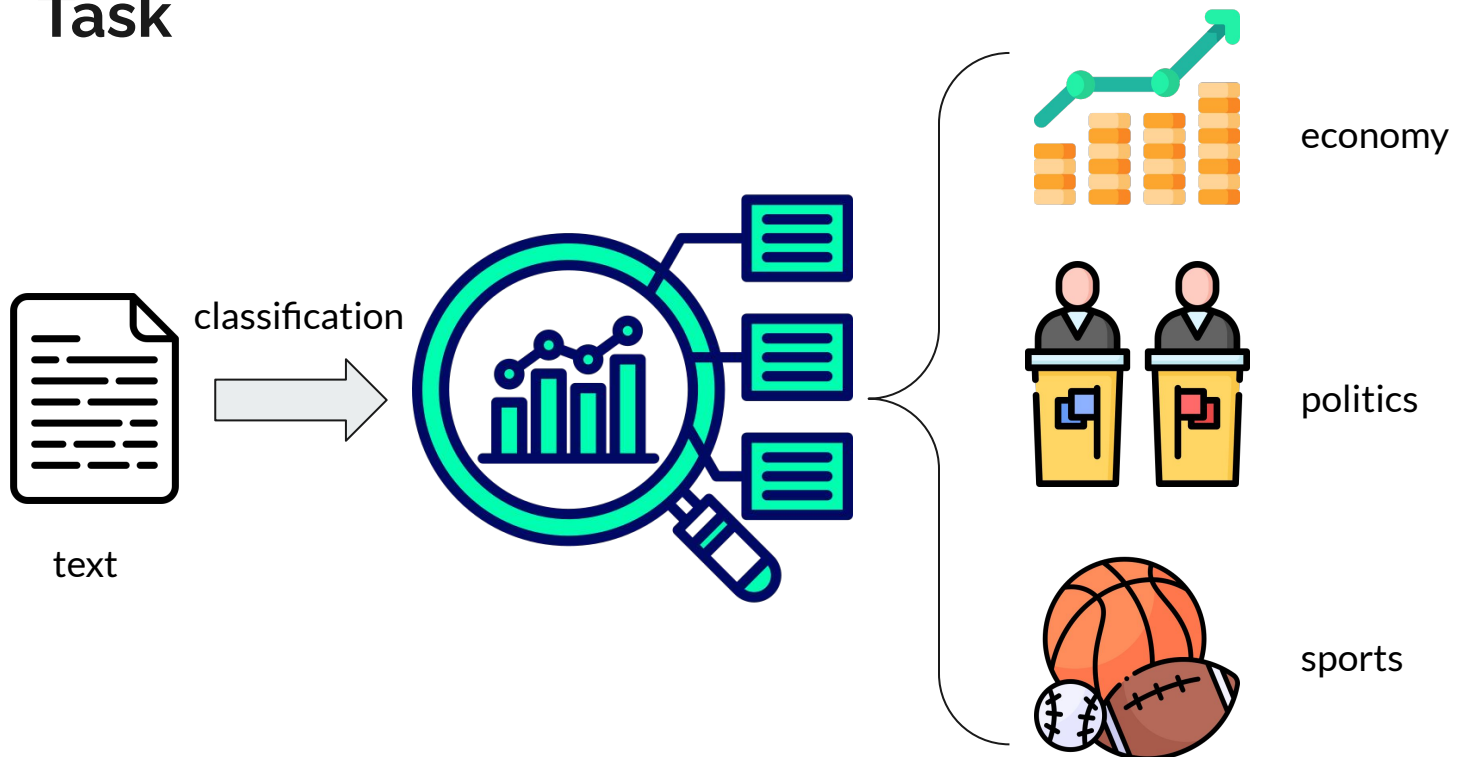
Date : 2023/09/12



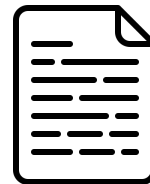
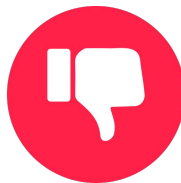
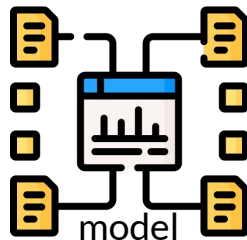
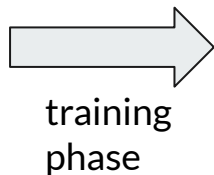
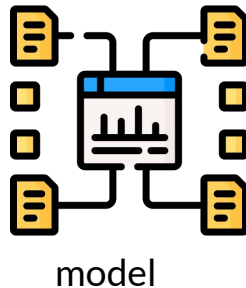
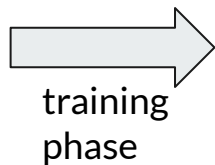
Outline

- Introduction
- Method
- Experiment
- Conclusion

Task



Problem - few label data



unlabeled data



annotate



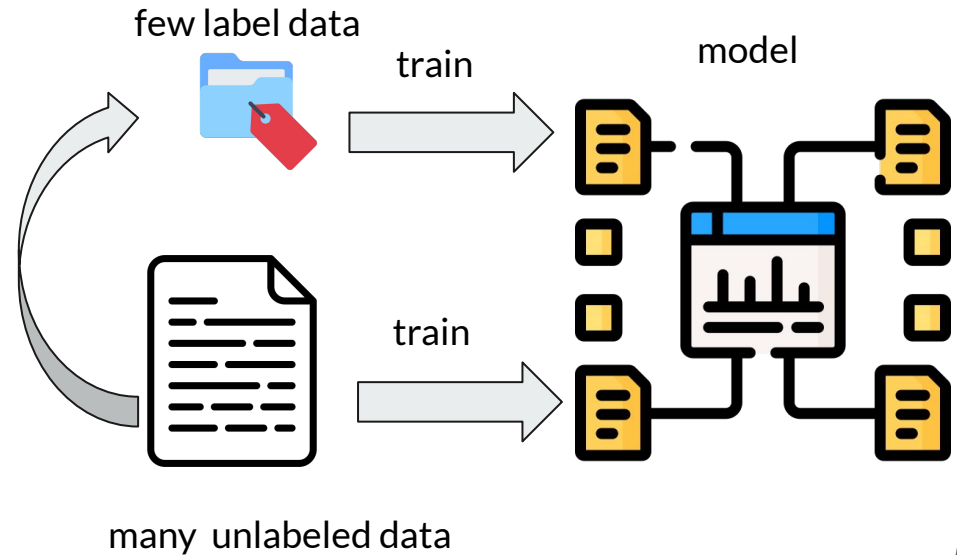
cost expensive

Semi-supervised learning

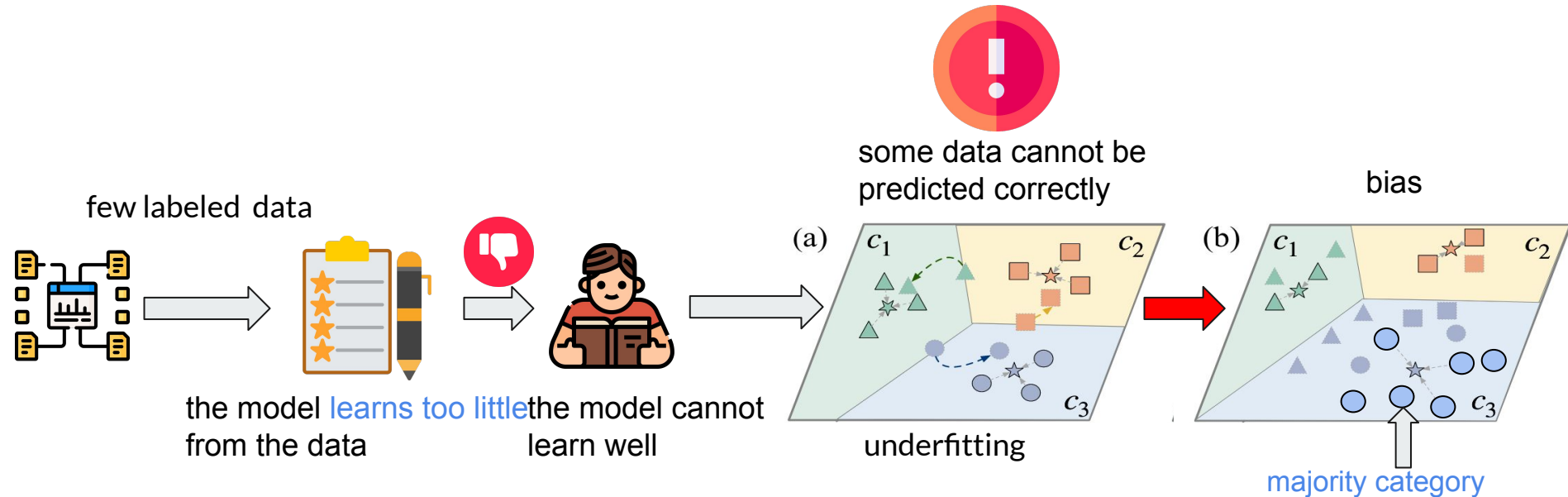


Solve the problem of **too little labeled data**

Hope to use unlabeled data to assist training

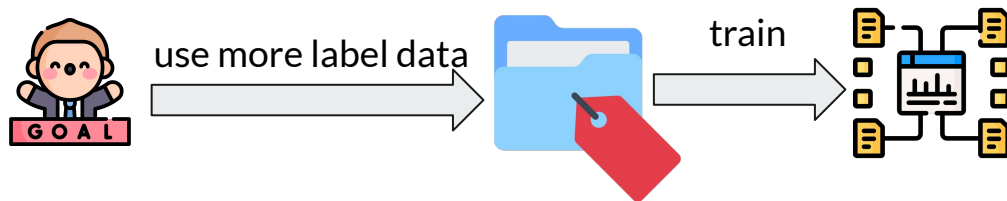


Problem - Semi-supervised learning

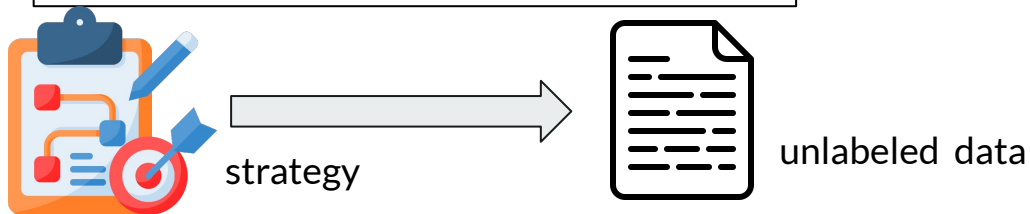


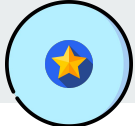
Solution

underfitting

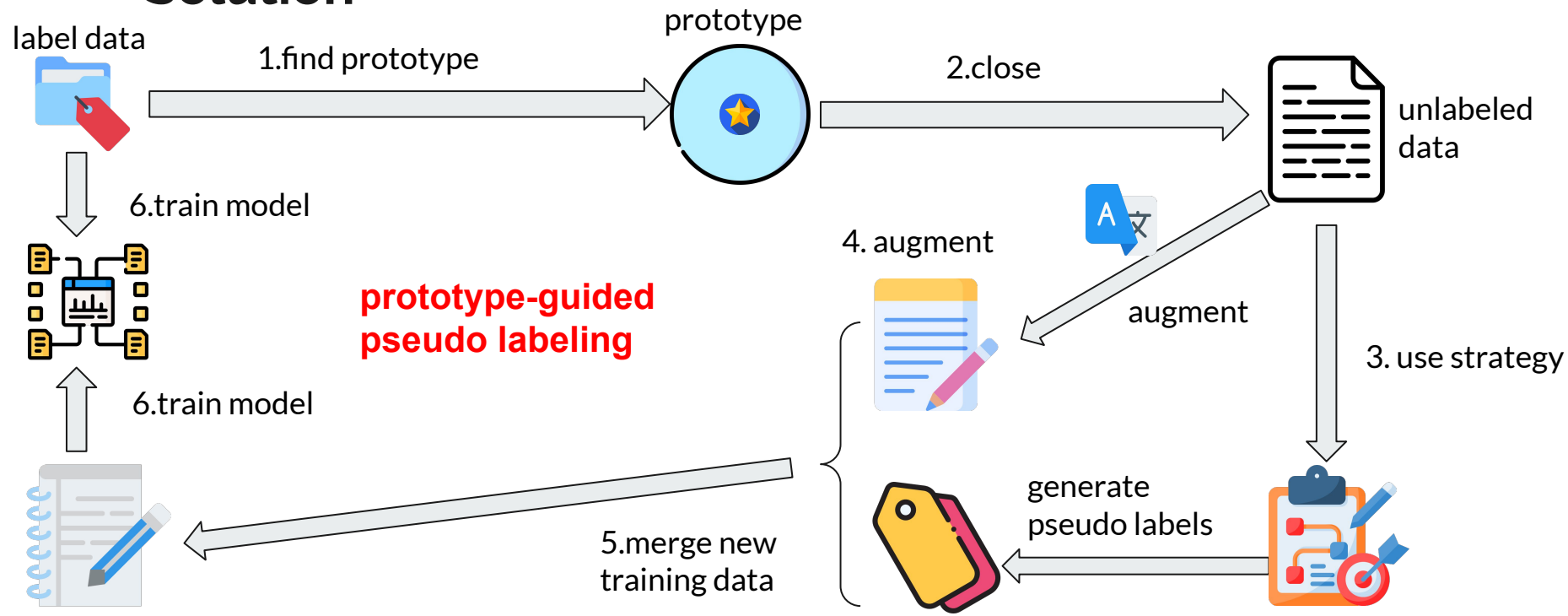


Hope to find **representative** unlabeled data



prototype  \rightarrow the most representative points in each category

Solution





Outline

- Introduction
- **Method**
- Experiment
- Conclusion

Unsupervised Data Augmentation

normal data Augmentation

reason : few label data



hope to get more **labeled data** to train the model





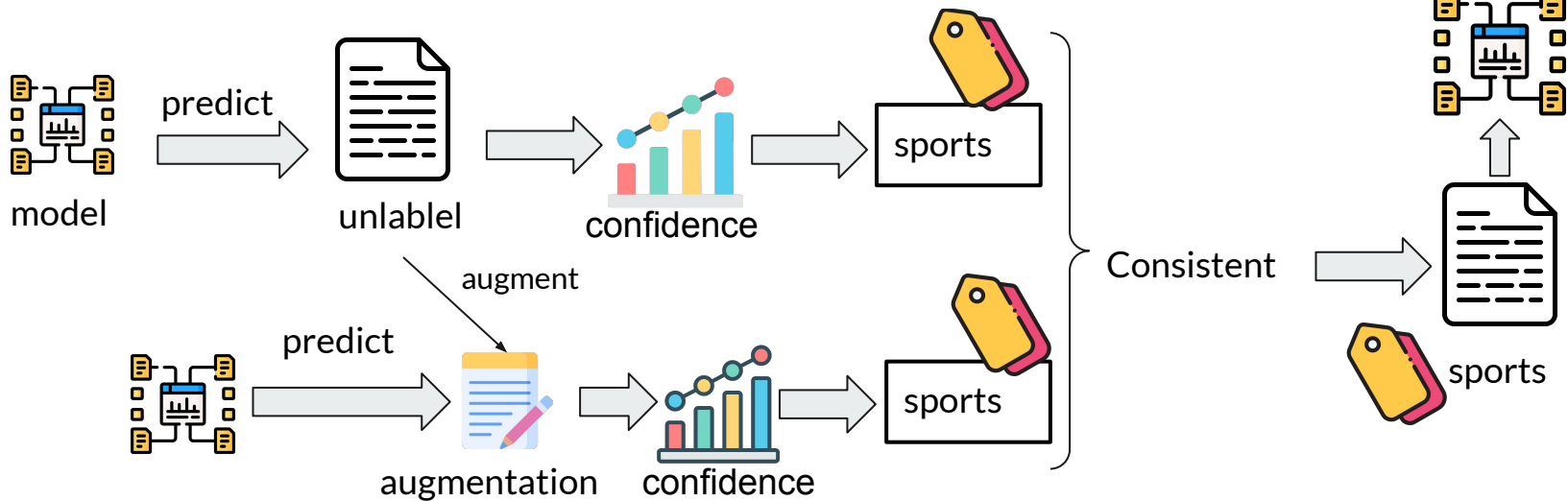
at the beginning,
model trained with label data



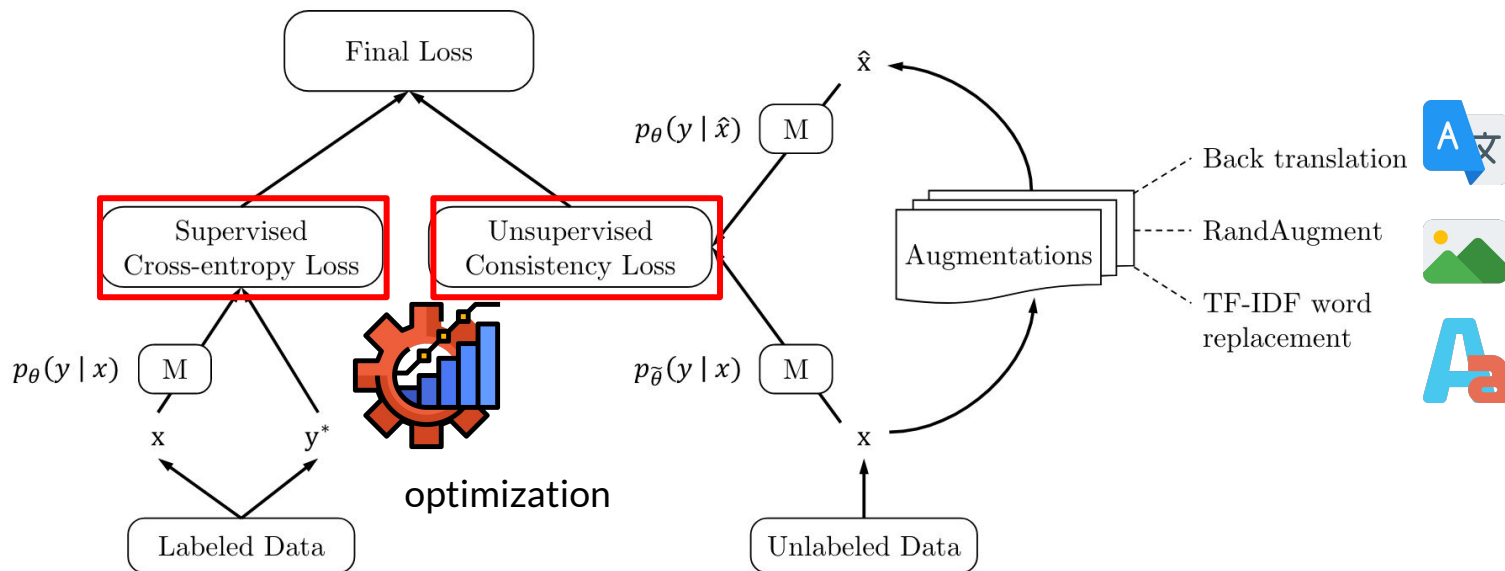
data diversity
Increase model robustness

Unsupervised Data Augmentation

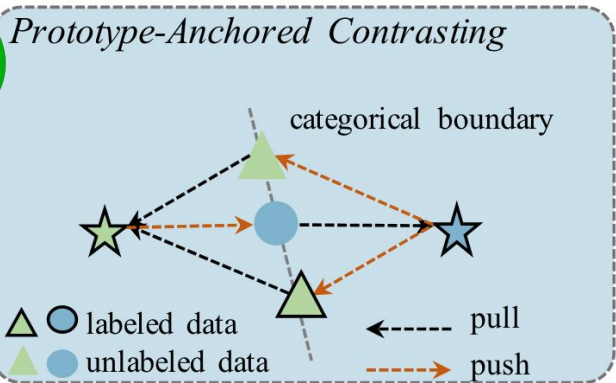
reason : label data has **limited effect** after augmentation



Unsupervised Data Augmentation

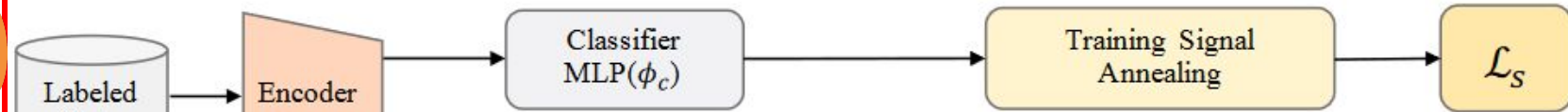


3

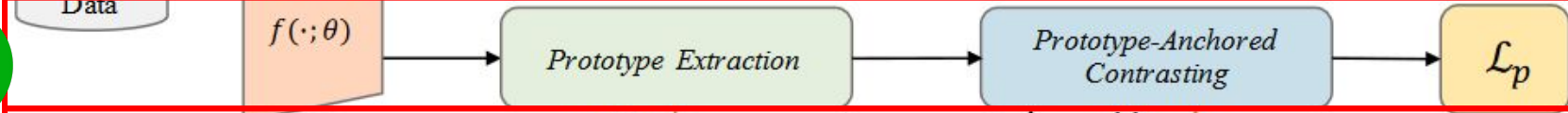


Prototype-Guided Pseudo Labeling

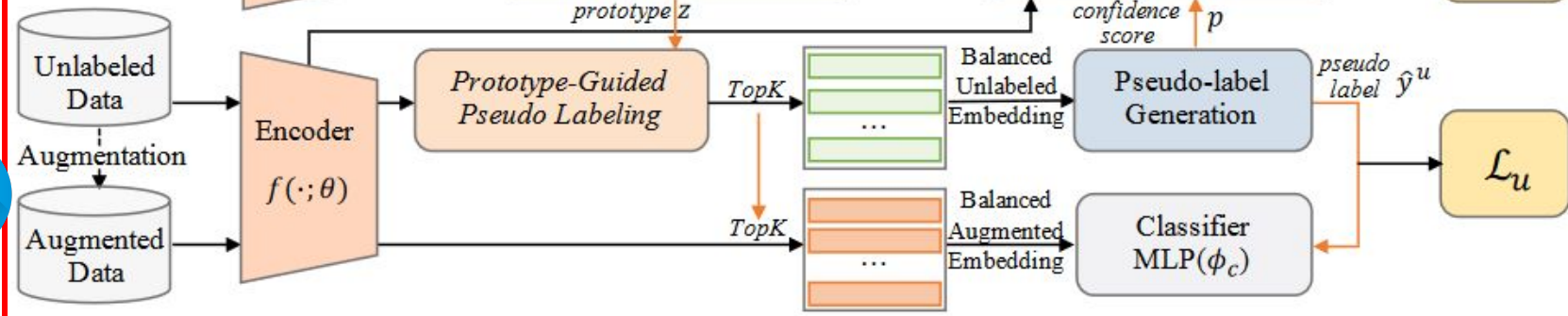
1



3

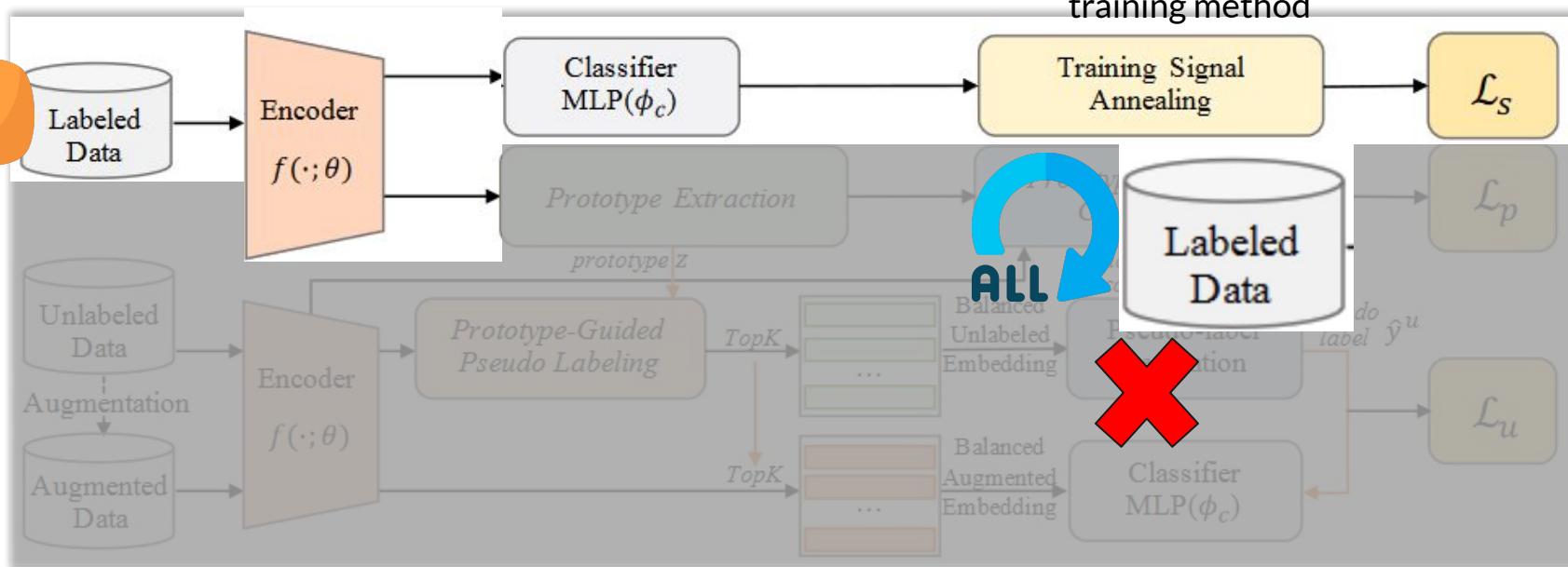


2





Prototype-Guided Pseudo Labeling (Partition 1)

training method

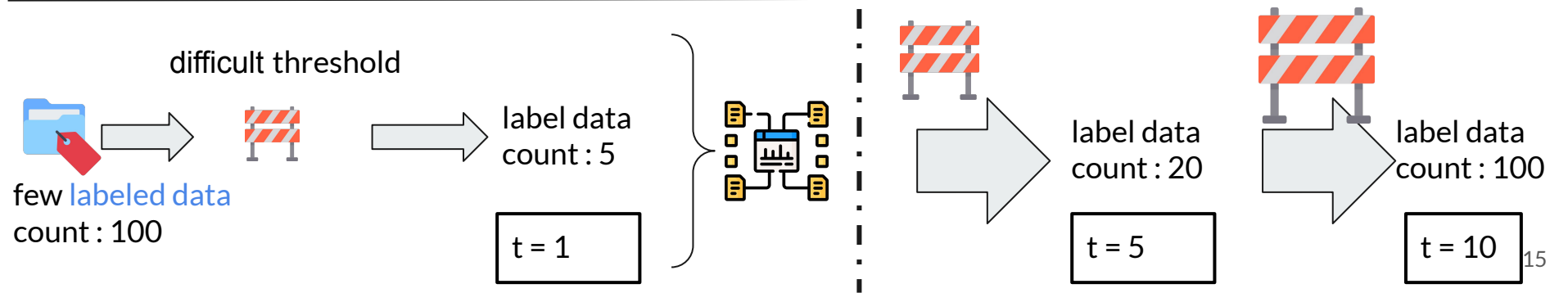
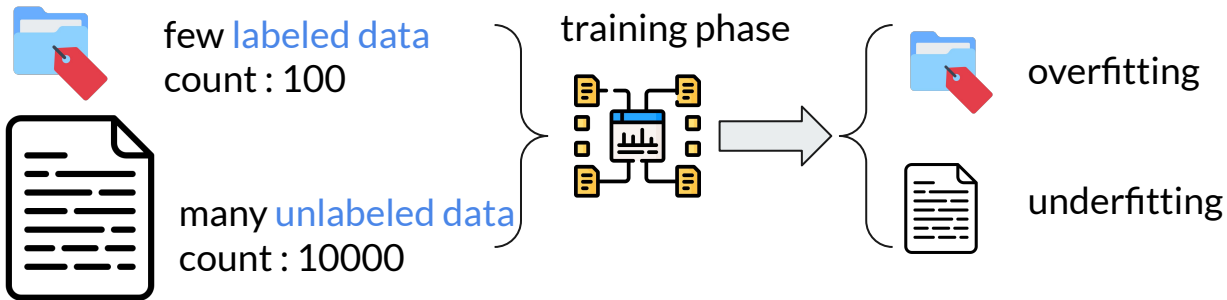


1 Training Signal Annealing

strategy  releases **labeled instances** **gradually** during the training phase

 prevent the model from **overfitting** too quickly

Training signal annealing



Training signal annealing



don't use **all labeled data** at once

e.g. number of categories = 10

In this paper, tau need to $> 1/(\text{number of categories}) \Rightarrow 1/10$, so suppose tau = $5/10 = 0.5$

number of categories	T(total iterations)	t	tau	threshold
10	10	1	0.5	$1/10 \times (1-0.5) + 0.5 = 0.55$
10	10	5	0.5	$5/10 \times (1-0.5) + 0.5 = 0.75$
10	10	10	0.5	$10/10 \times (1-0.5) + 0.5 = 1$

$$\eta_t = \frac{t}{T} (1 - \tau) + \tau,$$

current iterations t
 threshold η_t
 total iterations T
 initial threshold τ

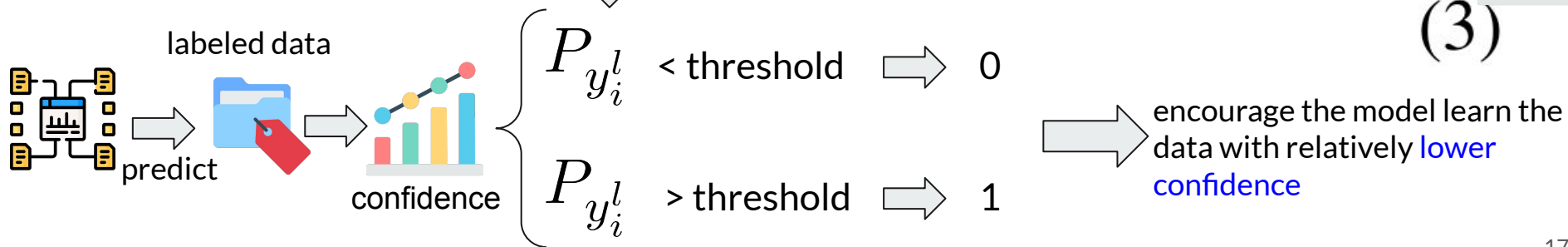


threshold

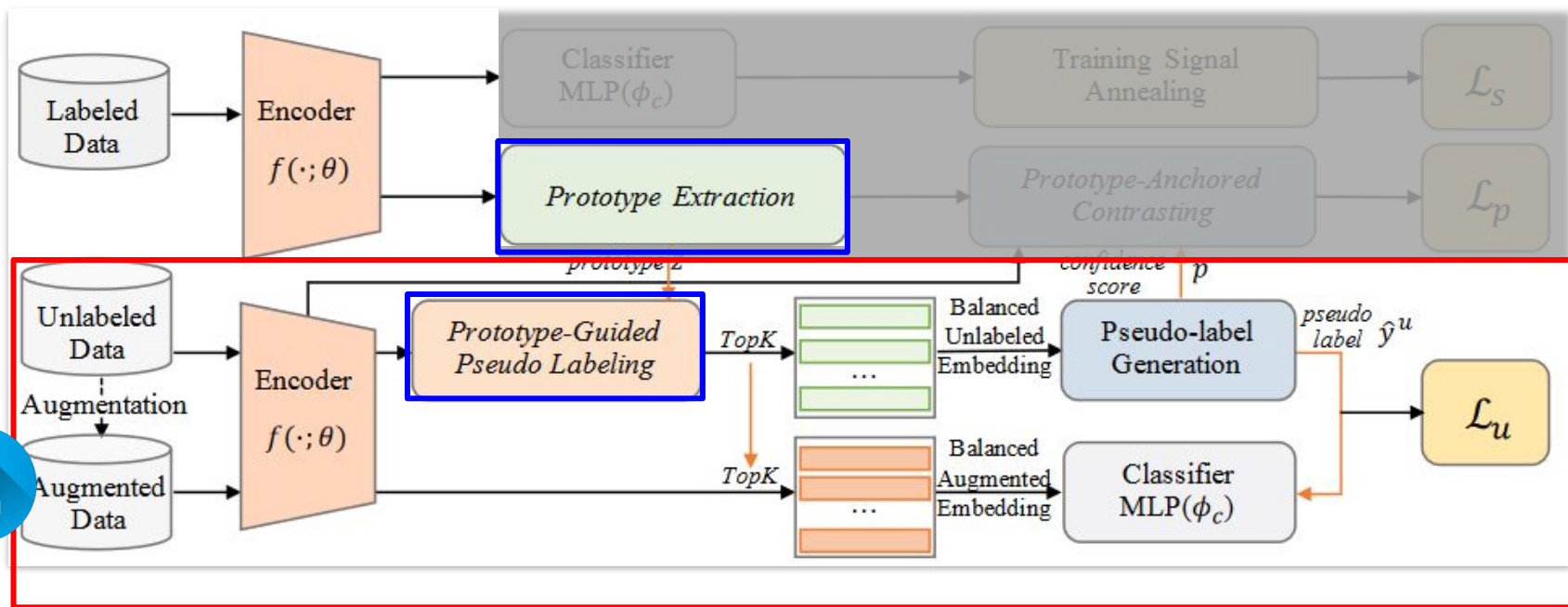
1	\mathcal{L}_s	$P_{y_i^l}$	threshold	$\mathbb{I}(p_{y_i^l} < \eta_t)$	
		labeled data 1	0.4	0.55	1
		labeled data 2	0.8	0.55	0

Annealing Supervised Loss

$$\mathcal{L}_s = -\frac{1}{m} \sum_{i=1}^m \mathbb{I}(p_{y_i^l} < \eta_t) \log \frac{\exp(g(x_i^l, \phi_{y_i^l}, \theta))}{\sum_{c \in \mathcal{C}} \exp(g(x_i^l, \phi_c, \theta))}, \quad (3)$$



Prototype-Guided Pseudo Labeling (Partition 2)



2

2

prototype



the most
representative points
in each category

Pseudo-label
Generation

4

$$\hat{y}_i^u = \arg \max_{c \in \mathcal{C}} \exp(g(x_i^u, \phi_c, \theta))$$

Unlabeled
Data

relevance score

A	0.8
B	0.15
C	0.05

Prototype-Guided Pseudo Labeling

2

Labeled
Data

Encoder
 $f(\cdot; \theta)$

Prototype
Extraction

3

calculate the distance between
prototype and unlabeled data

4

Prototype-Guided
Pseudo Labeling

balance
dataset

Pseudo-label
Generation

pseudo labels

If it is close enough to the prototype

1

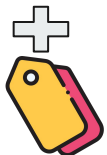
Unlabeled
Data

Encoder
 $f(\cdot; \theta)$



5

Augmented
Data



merge into new training data set



model

\mathcal{L}_u



Find the **most representative data** from each category

Prototype Extraction



labeled data

count



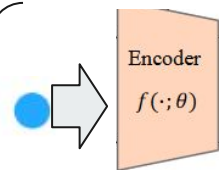
$$n_c = \sum_{y_i^l \in D_l} \mathbb{I}(y_i^l = c)$$

= 6

= 5

$$z_c = \frac{1}{n_c} \sum_{y_i^l = c} f(x_i^l)$$

$$\frac{1}{6} \sum$$

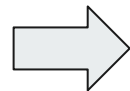
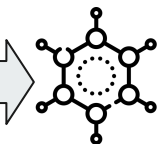
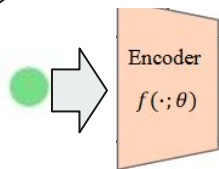


embedding

prototype

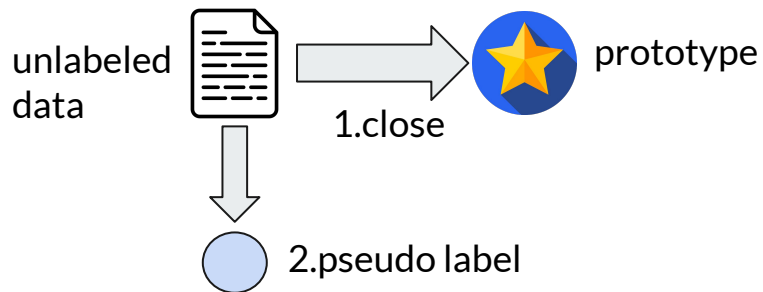


$$\frac{1}{5} \sum$$



why use prototype

- bring **each unlabeled data's embedding** closer to its **associated prototype**

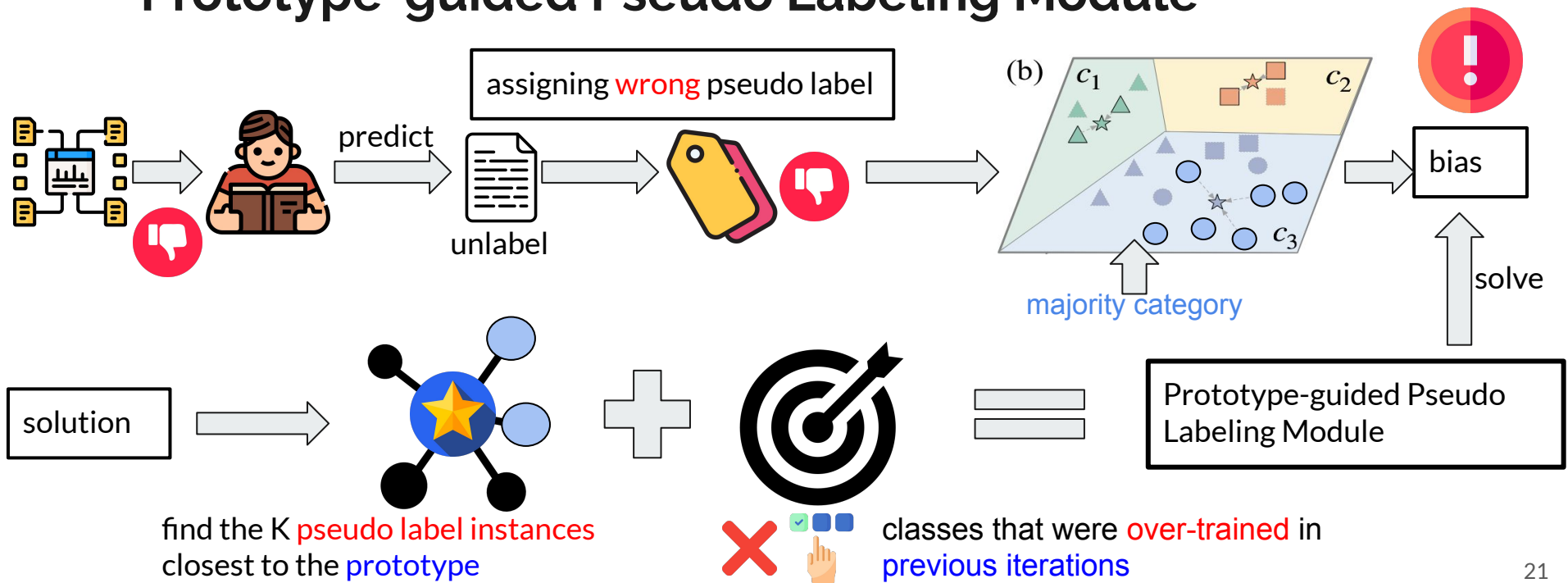


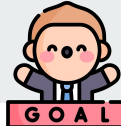


over-trained →

class	<t	t	final
A	50	5	50
B	20	4	22
C	5	10	15

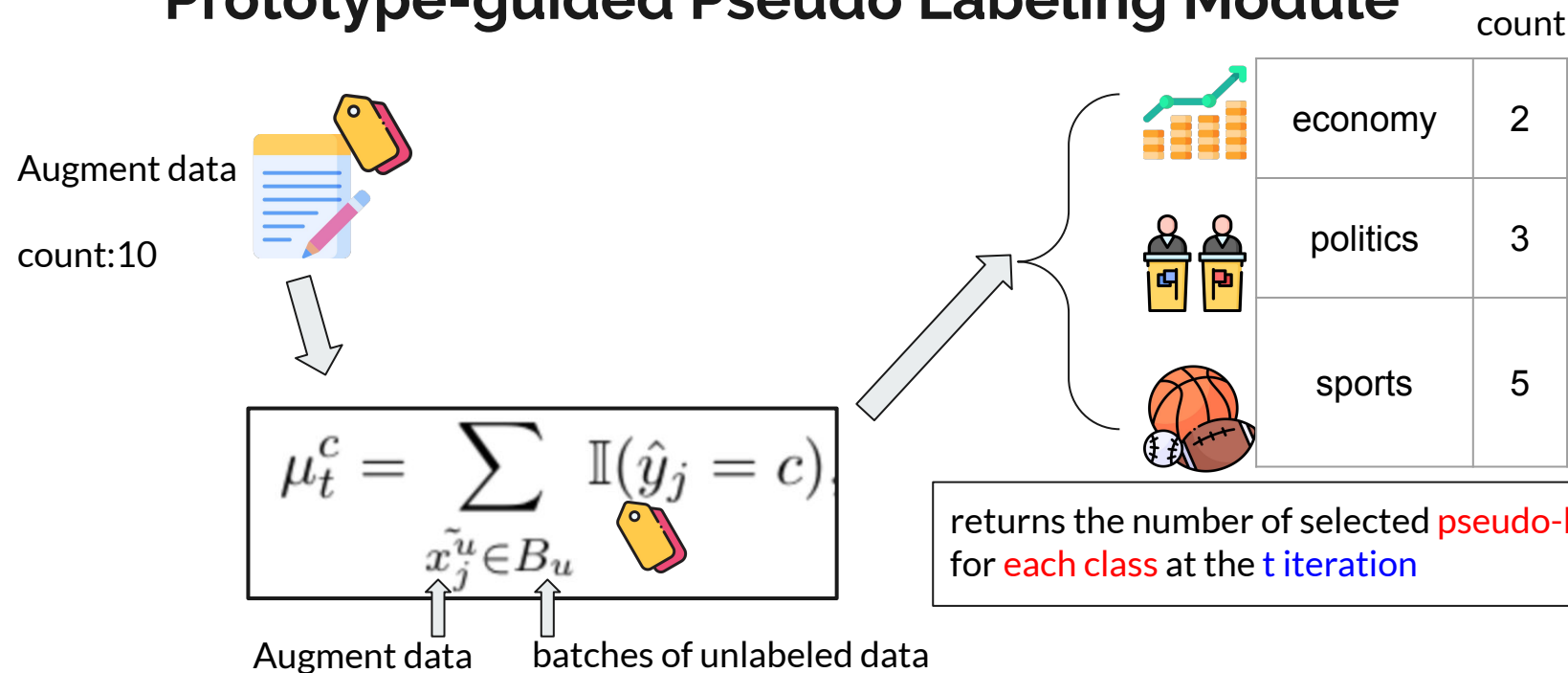
Prototype-guided Pseudo Labeling Module





calculate how many **pseudo labels** each category should receive **in this iteration**

Prototype-guided Pseudo Labeling Module





Prototype-guided Pseudo Labeling Module

$$\mu_{<t}^c$$



Get the **total number of pseudo label** for each category from the **previous iteration**

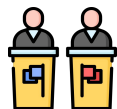
$$\gamma_t = \arg \min_{c \in \mathcal{C}} \mu_{<t}^c$$



The **previous iteration** finds the **smallest number of pseudo-labels in each category**

Assign the number of pseudo labels to each category

$$k_c = \begin{cases} \mu_t^c & \text{if } \mu_{<t}^c - \gamma_t = 0 \\ \mu_t^c - (\mu_{<t}^c - \gamma_t) & \text{if } 0 \leq \mu_{<t}^c - \gamma_t < \mu_t^c \\ 0 & \text{if } \mu_{<t}^c - \gamma_t \geq \mu_t^c \end{cases}$$



economy

 μ_t^c

2

politics

3

sports

5

 $\gamma_t = 2$ (sports)

$$\gamma_t = \arg \min_{c \in \mathcal{C}} \mu_{<t}^c$$

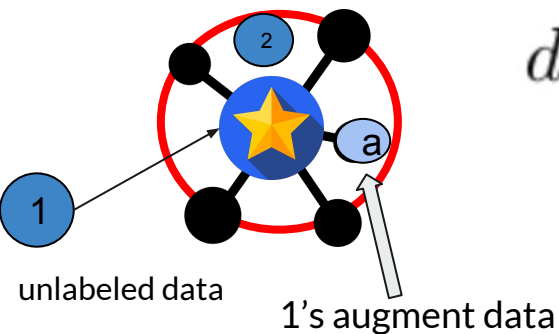
Prototype-guided Pseudo Labeling Module

 $\mu_{<t}^c$ μ_t^c

	Previous t times	pseudo label currently found	$\mu_{<t}^c - \gamma_t$	condition	Kc
economy	6	2	6-2 = 4	$\mu_{<t}^c - \gamma_t \geq \mu_t^c$	0
politics	4	3	4-2=2	$0 \leq \mu_{<t}^c - \gamma_t < \mu_t^c$	$3-(4-2) = 1$ $\mu_t^c - (\mu_{<t}^c - \gamma_t)$
sports	2 γ_t	5	2-2=0	$\mu_{<t}^c - \gamma_t = 0$	5 μ_t^c

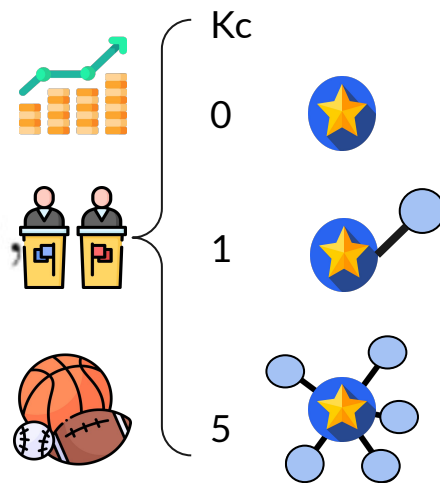


Prototype-guided Pseudo Labeling Module



$$d_c = \text{TopK}(d(f(\tilde{x}_j^u), z_c), k_c)$$

↑
augment data



Choose the **length** closest to the **prototype**

Selective Unsupervised Loss

$$\mathcal{L}_u = -\frac{1}{n} \sum_{i=1}^n \mathbb{I}(x_i^u) \log \frac{\exp(g(\tilde{x}_i^u, \phi_{y_i^u}, \theta))}{\sum_{c \in \mathcal{C}} \exp(g(\tilde{x}_i^u, \phi_c, \theta))}$$

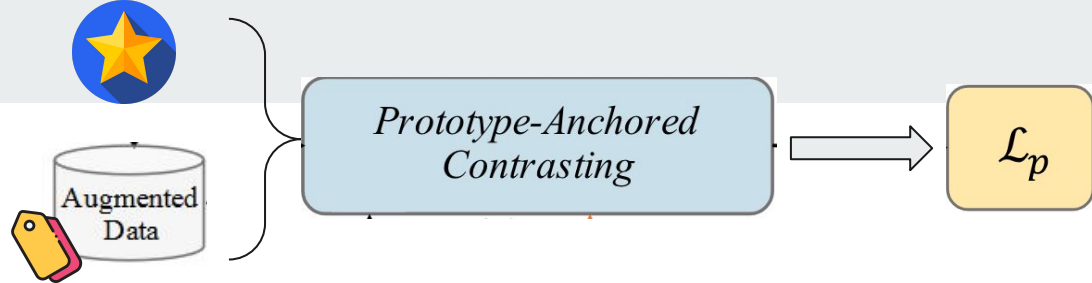
entropy

if **pseudo label** = sports
& **true label** = sports

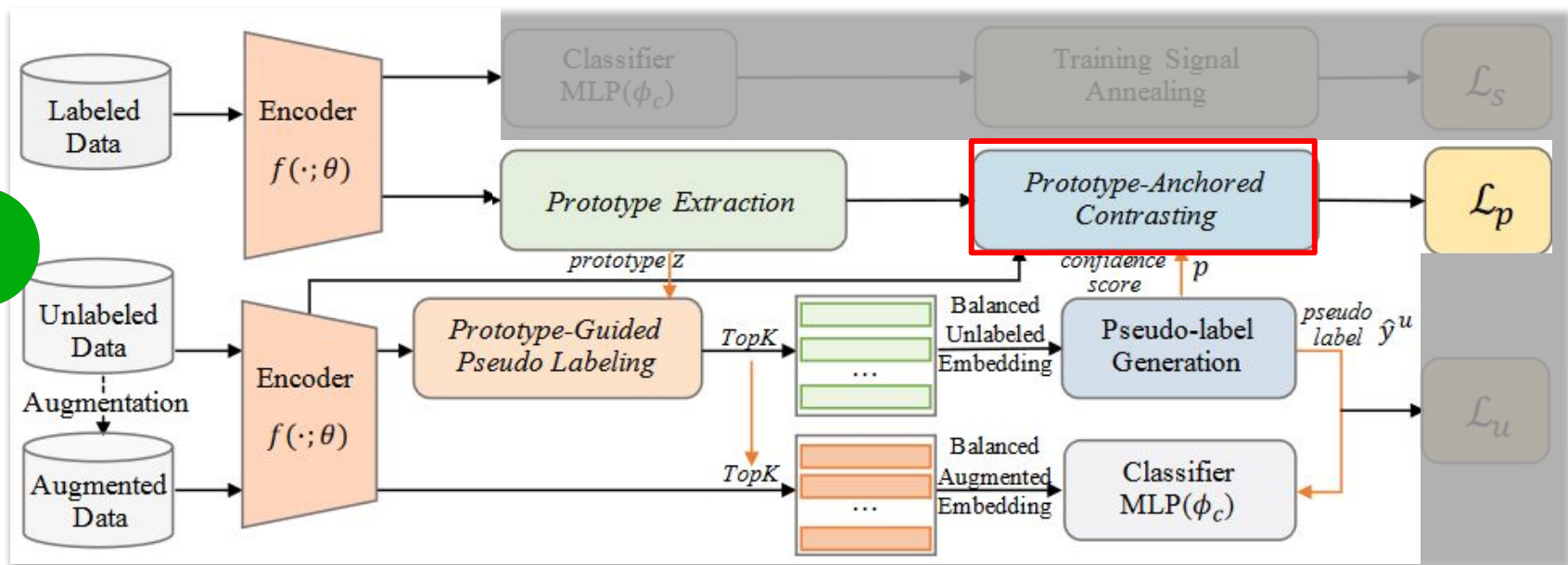
$$\mathbb{I}(x_i^u) = \mathbb{I}(\hat{y}_i^u = c) \wedge \mathbb{I}(d(f(x_j^u), z_c) \leq d_c)$$

1 or 0

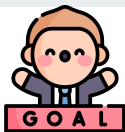
The distance between **unlabeled data** and **prototype** is smaller than the distance between **augmentation** and **prototype**.



Prototype-Guided Pseudo Labeling (Partition 3)



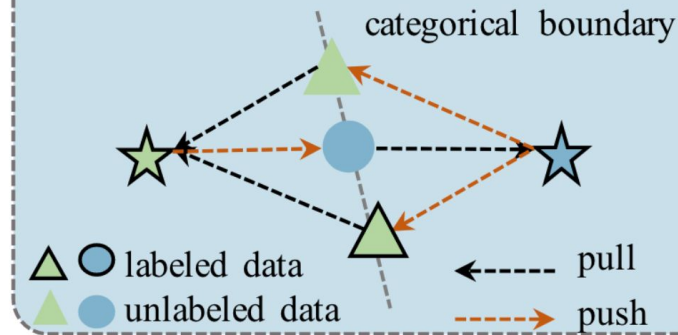
3

 \mathcal{L}_p 

make **each instance** closer to the **prototype of the category**

Prototype-Anchored Contrasting (PAC)

Prototype-Anchored Contrasting



$$\mathcal{L}_p = -\frac{1}{m} \sum_i \sum_c \mathbb{I}(y_i^l = c) \log \frac{\exp(-d(f(\tilde{x}_i^l), z_c))}{\sum_{k \in \mathcal{C}} \exp(-d(f(\tilde{x}_i^l), z_k))}$$

loss of
label data

$$-\frac{\lambda}{n} \sum_j \sum_c p_{\hat{y}_j^u} \mathbb{I}(\hat{y}_j^u = c) \log \frac{\exp(-d(f(\tilde{x}_j^u), z_c))}{\sum_{k \in \mathcal{C}} \exp(-d(f(\tilde{x}_j^u), z_k))}$$

loss of
augment
data



predict



unlabeled



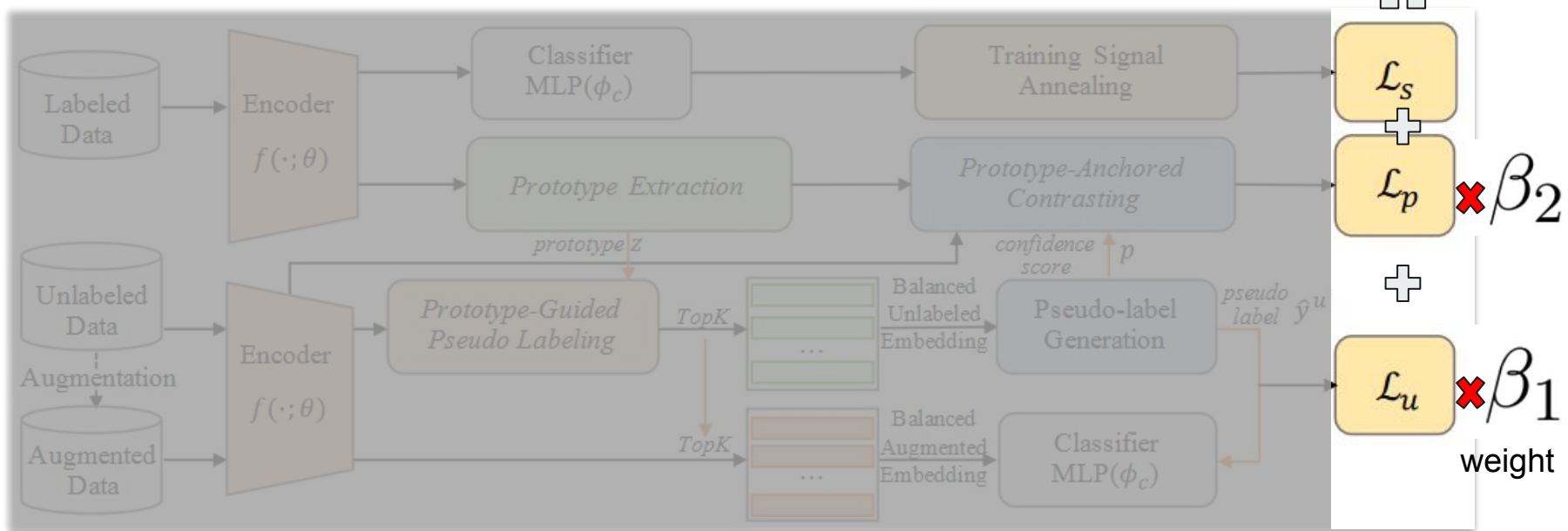
confidence

Use the confidence of **unlabeled data** as weight



balance

Prototype-Guided Pseudo Labeling





Outline

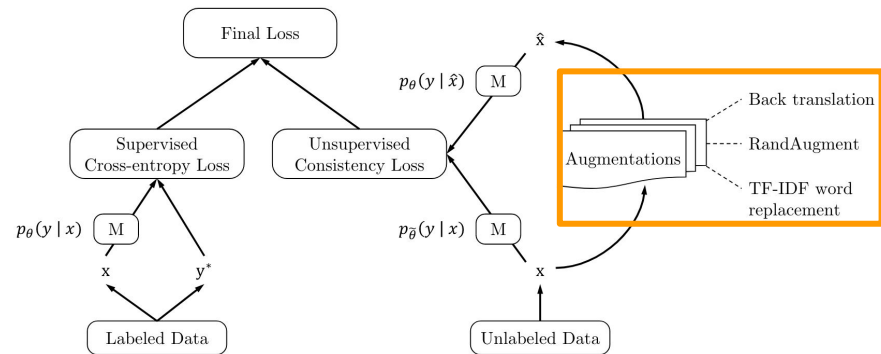
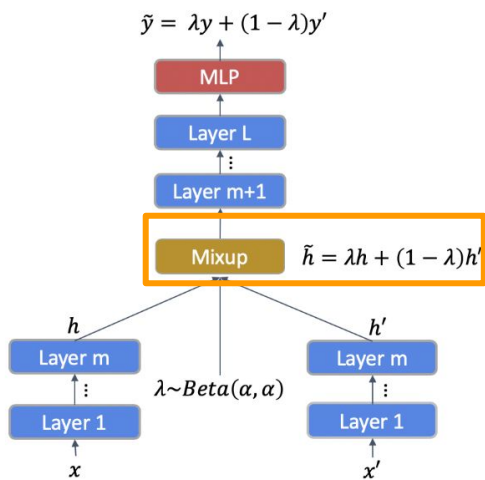
- Introduction
- Method
- **Experiment**
- Conclusion

Dataset

Dataset	Classification Type	Class	Train	Unlabeled	Dev	Test
AG News	News Topic	4	200	5000	2000	1900
DBpedia	Wikipedia Topic	14	200	5000	2000	5000
Yahoo! Answer	QA Topic	10	200	5000	2000	6000
IMDB	Movie Review Sentiment	2	200	5000	2000	12500

Dataset	AG News	IMDB	Yahoo! Answer	DBpedia
description	collection of news	<ul style="list-style-type: none"> ● movie reviews ● Binary emotion classification 	topic classification	extract structured content from the information created in Wikipedia
category	<ul style="list-style-type: none"> ● World ● Sports ● Business ● Sci/Tech 	<ul style="list-style-type: none"> ● positive ● negative 	<ul style="list-style-type: none"> ● Society & Culture ● Science & Mathematics ● Health ● Education & Reference ● Sports ● Business & Finance ● 	14 class <ul style="list-style-type: none"> ● company ● education ●

Baseline



Model		learning method
BERT	trained with only labeled data	supervised
UDA	make different augment to unlabeled data	semi-supervised
Mixtext	TMix: A new data augmentation method that interpolates two input vectors in hidden space to generate a new vector	semi-supervised

- metric : accuracy

Experiment Results

the number of
labeled data used

Model	AG News			IMDB			Yahoo! Answer			DBpedia		
	10	30	200	10	30	200	10	30	200	10	30	200
supervised Bert	81.0	84.3	87.2	70.6	73.3	86.1	60.1	64.1	69.3	96.6	98.2	98.6
semi-supervised UDA	86.4	86.4	88.3	86.4	86.4	88.7	64.3	68.3	70.2	97.8	98.3	98.8
semi-supervised Mixtext	87.3	87.4	88.2	74.2	85.3	89.1	67.7	68.5	70.6	98.5	98.8	98.9
this paper PGPL	87.8	88.5	89.2	88.9	90.2	90.3	67.4	69.1	70.7	98.7	99.0	99.0

↑
unsupervised > supervised

↑
the gap between unsupervised and supervised narrows

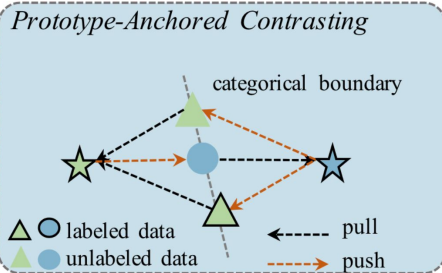
PGPL is stable in all aspects except Yahoo

- metric : accuracy

Evaluation Results with other pre-trained models

Dataset	Supervised (all labels)		Supervised (10 labels)		Semi-supervised (10 labels)	
	BERT	RoBERTa	BERT	RoBERTa	PGPL(BERT)	PGPL(RoBERTa)
AG News	91.2	92.4	80.2	80.7	87.8	88.4
IMDB	90.4	93.5	70.9	71.2	88.9	91.2
Yahoo!Answer	73.7	74.2	60.1	61.0	67.4	67.8
DbPedia	99.1	99.1	96.6	96.1	98.7	98.8
Average	88.6	89.8	76.9	77.3	85.7	86.6

this paper



Ablation Study

1. PGP and PAC can independently improve model performance
2. TSA helps too, besides DBpedia

Training Signal
Annealing

Data	AG News	IMDB
PGPL	88.3	89.7
w/o PGP	86.5	88.2
w/o PAC	86.2	88.9
w/o TSA	87.2	87.9

Data	Yahoo!Answer	DBpedia
PGPL	68.3	98.4
w/o PGP	65.7	98.2
w/o PAC	67.4	98.2
w/o TSA	68.0	98.6

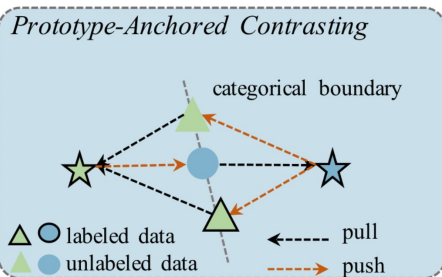


Outline

- Introduction
- Method
- Experiment
- **Conclusion**

Prototype-Guided Pseudo Labeling

PGP



PAC

Conclusion

1. A **semi-supervised model PGPL** that **combines PAC and PGP strategies** is proposed for **semi-supervised text classification tasks**.
2. After constructing the **prototype**, use **PAC to group text embeddings** belonging to the same category together to alleviate the **problem of underfitting**.
3. **PGP selects reliable pseudo-labeled data nearby prototypes** to address the **training bias** from the imbalanced data