



# D-CALM: A **Dynamic** approach **Clustering-based Active Learning** **Approach for Mitigating Bias**

Advisor : Jia-Ling, Koh

task

Speaker : Ting-I, Weng

Source : ACL'23

Date : 2023/08/14



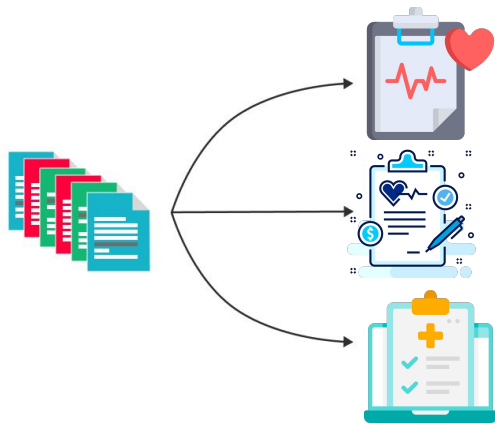
# Outline

- Introduction
- Method
- Experiment
- Conclusion

- Background
- Active Learning
- Challenge
  - bias
    - bias induction
  - effective batch selection
- D-CALM

# Background

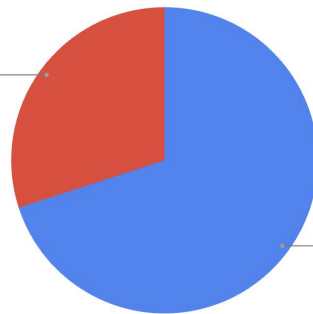
classification task



e.g. Medical text classification

ideal dataset

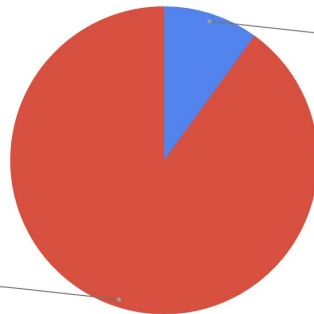
unlabel  
30.0%



label  
70.0%

real world dataset

unlabel  
90.0%



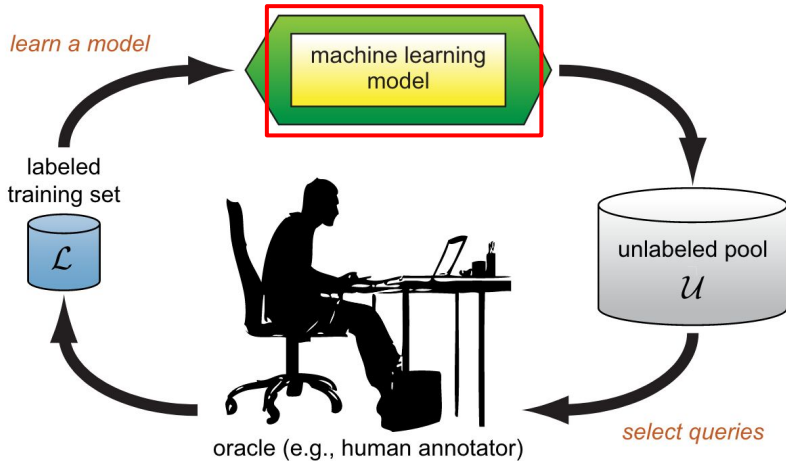
label  
10.0%

labeling



# Active Learning

Active Learning : actively select the most valuable samples for labeling



Goal:

- reduce labeling costs
- find more diverse data



Weakness:



- rely on the judgment of the classifier



Issue:

- whether the data found is helpful for model training

# Challenge - bias

- classifiers may not perform well for underrepresented classes in the data
  - **category** has less data  = 5
  - people of color have a high error rate  = 100%
- Hope to find more **unlabeled data** about people of color using active learning

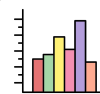


 unlabeled persons of color

training data








reason

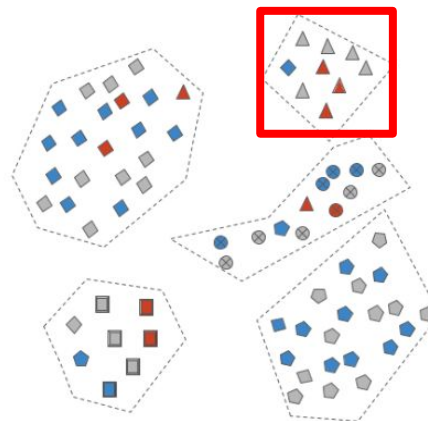


imbalance



didn't learn diverse semantics

 HS targeting **persons of color**  HS targeting **women**  
 HS targeting **Jews**  HS targeting **disability**  HS targeting **migrants**



   Classifier prediction wrong    Classifier prediction correct  
   True label not sampled yet

# Challenge - bias induction



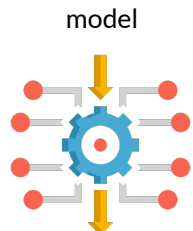
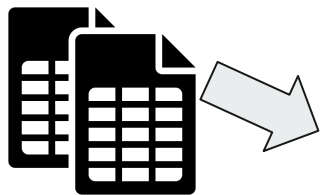
reason

- rely on the judgment of the classifier
- high probability of selecting unrepresentative samples

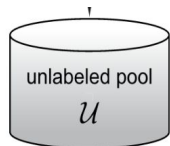
- bias induction : active learning can make bias worse
  - dataset bias
  - semantic bias

women, Jews...

persons of color



model



unlabeled pool  
 $U$

find uncertain data



Reality: the model finds more data here

Ideal: find more data



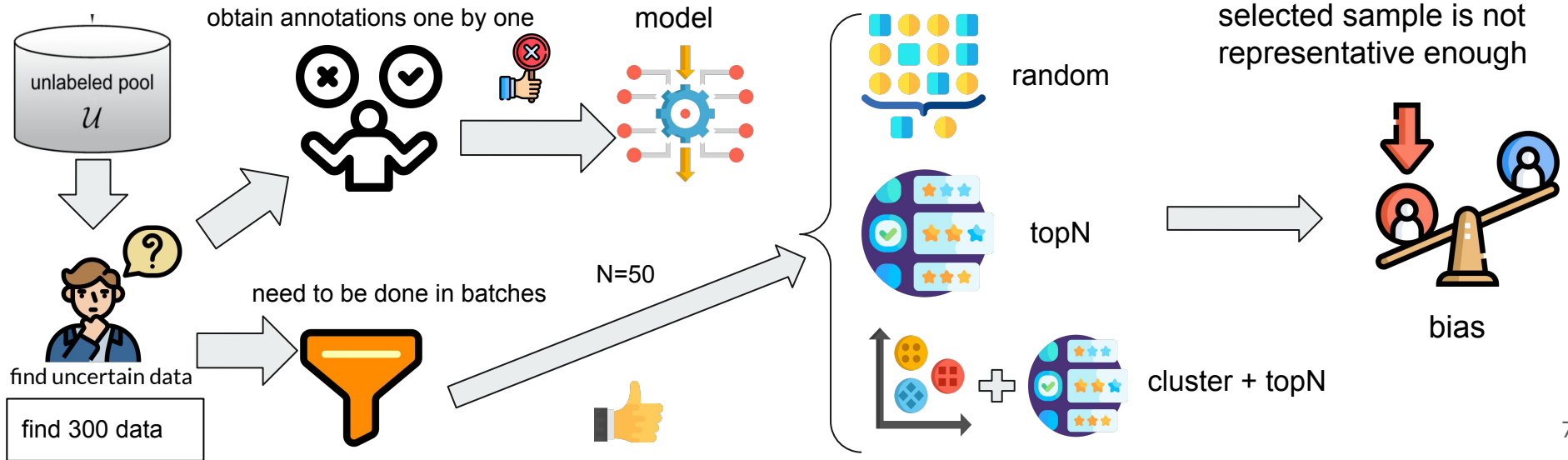
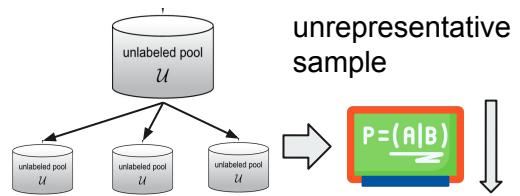
human annotator



# Challenge - effective batch selection

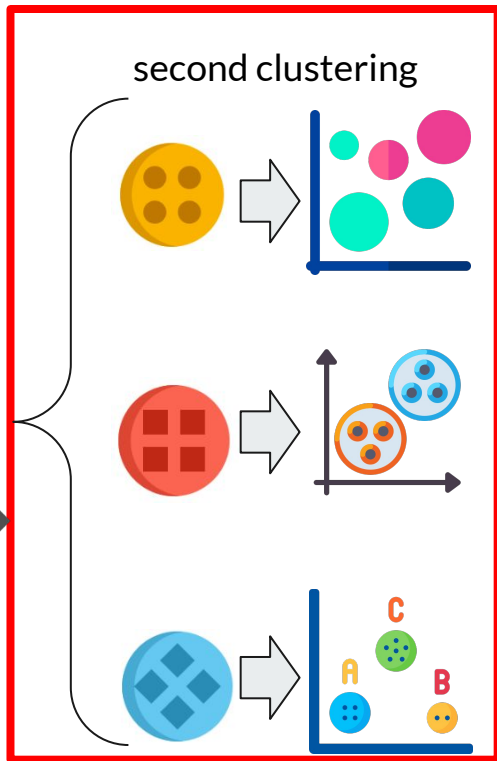
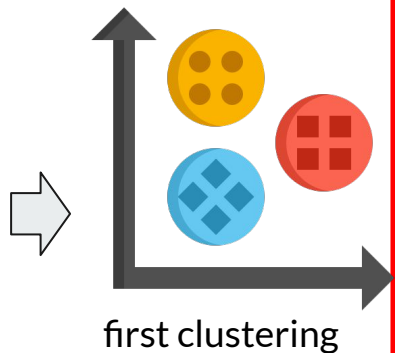
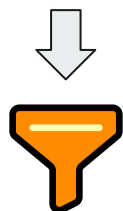
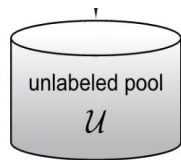
when the batch size is larger, the bias becomes greater

goal :



# D-CALM

Focus on categories with high error rates



Pick a representative sample from each group



pick unlabeled data



pick unlabeled data




pick unlabeled data

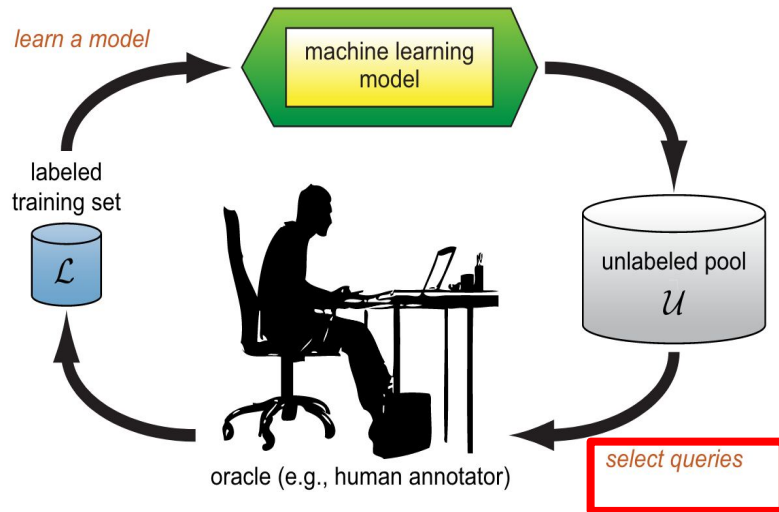




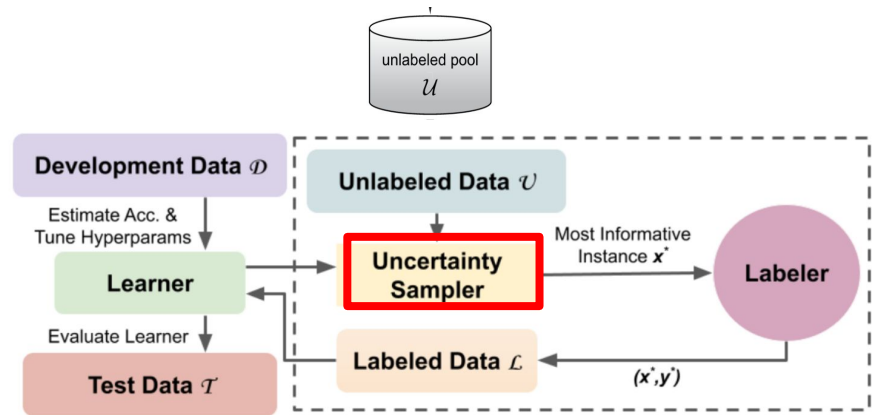
# Outline

- Introduction
  - **Method**
  - Experiment
  - Conclusion
- 
- Active learning - Query Strategy
  - Query-Strategy
    - Least Confident
    - Smallest margin
    - Entropy
  - Active Learning v.s. Clustering-based Active Learning
  - Dynamic Cluster AL v.s. Clustering-based AL

# Active learning - Query Strategy



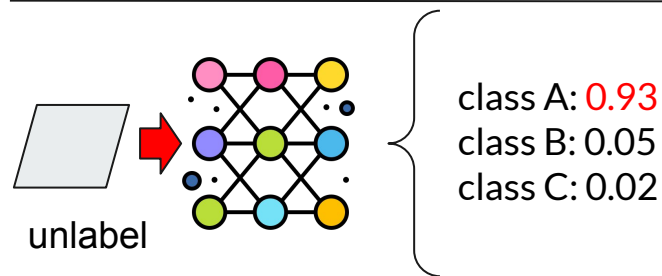
- how to pick ?
- How many ways?



## Query-Strategy - Least Confident

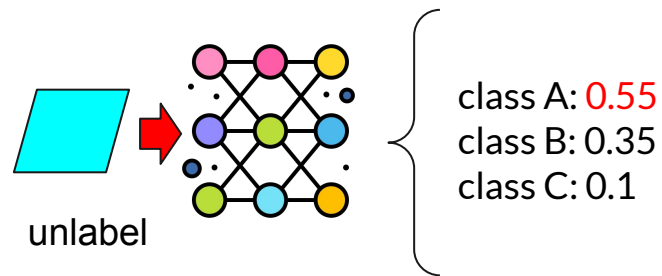
$$\hat{y} = \operatorname{argmax}_y P_\theta(y|x)$$

$$x_{LC}^* = \operatorname{argmax}_x 1 - P_\theta(\hat{y}|x)$$



$$\hat{y} = 0.93$$

$$x_{LC}^* = 1 - 0.93 = 0.07$$

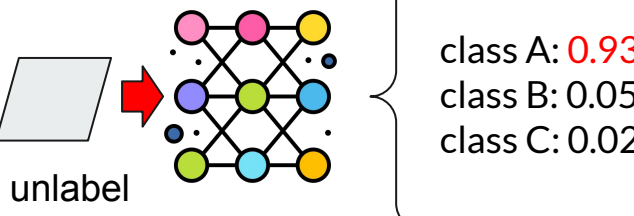
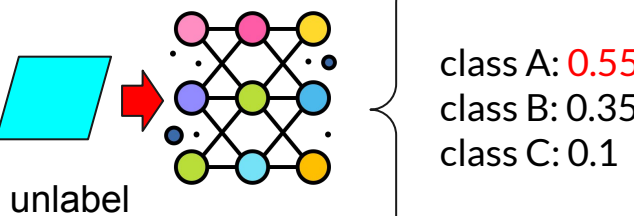


$$\hat{y} = 0.55$$

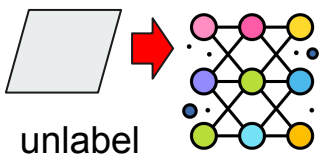
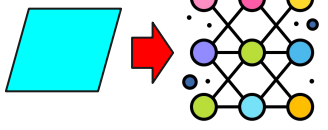
$$x_{LC}^* = 1 - 0.55 = 0.45$$

Confident  $\downarrow$  uncertain  $\uparrow$

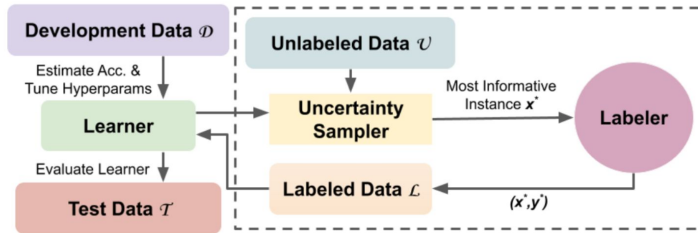
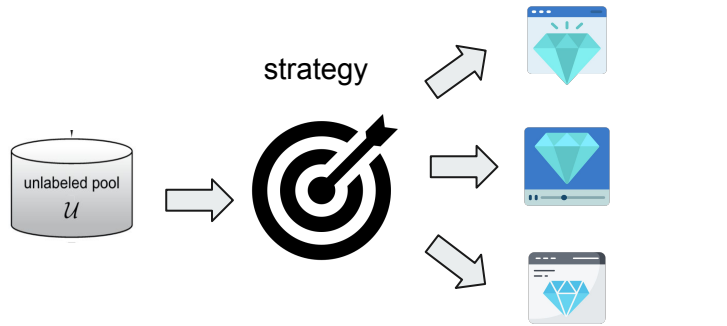
## Query-Strategy - Smallest margin

	$\hat{y} = \operatorname{argmax}_y P_\theta(y x)$	$x_{MS}^* = \operatorname{argmin}_x P_\theta(\hat{y}_1 x) - P_\theta(\hat{y}_2 x)$
 <p>class A: <b>0.93</b> class B: 0.05 class C: 0.02</p>	$\hat{y}_1 = 0.93$ $\hat{y}_2 = 0.05$	$x_{MS}^* = 0.93 - 0.05 = 0.88$
 <p>class A: <b>0.55</b> class B: 0.35 class C: 0.1</p>	$\hat{y}_1 = 0.55$ $\hat{y}_2 = 0.35$	<div style="border: 2px solid red; padding: 5px; display: inline-block;"> <math>x_{MS}^* = 0.55 - 0.35 = 0.2</math> </div> smallest margin ↓      uncertain ↑

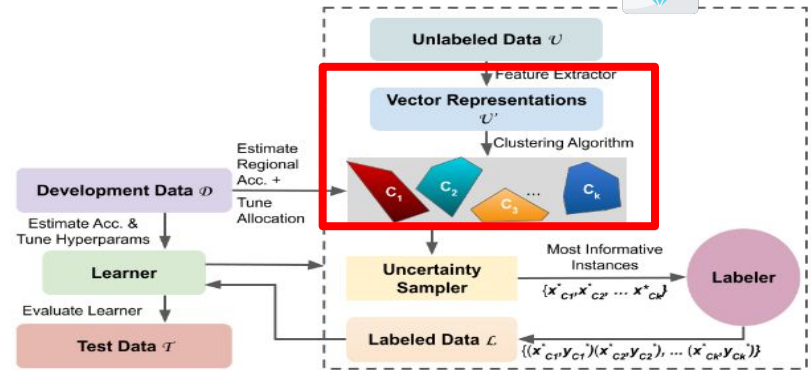
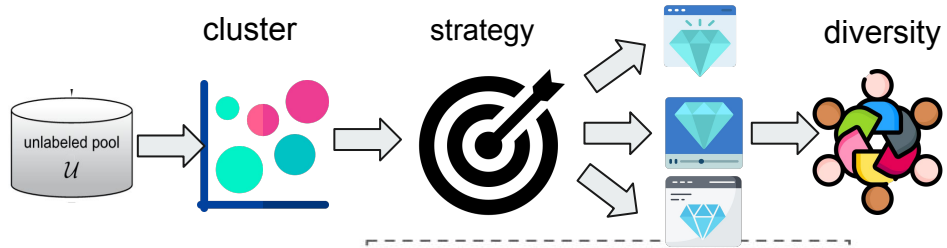
# Query-Strategy - Entropy

	$\log P_{\theta}(y_i x)$	$P_{\theta}(y_i x)\log P_{\theta}(y_i x)$	$x_E^* = \underset{x}{\operatorname{argmax}} - \sum_i P_{\theta}(y_i x)\log P_{\theta}(y_i x)$
 <p>unlabel</p> <p>class A: <b>0.93</b> class B: 0.05 class C: 0.02</p>	<p>class A: -0.104 class B: -4.321 class C: -5.6438</p>	<p>class A: -0.09672 class B: -0.21605 class C: -0.11287</p>	<p><math>-(0.09672+0.21605+0.11287) = -0.4256</math> <math>x_E^* = -(-0.4256) = 0.4256</math></p>
 <p>unlabel</p> <p>class A: <b>0.55</b> class B: 0.35 class C: 0.1</p>	<p>class A: -0.8624 class B: -1.5145 class C: -3.3219</p>	<p>class A: -0.47432 class B: -0.53007 class C: -0.33219</p>	<p><math>-(0.47432+0.53007+0.33219) = -1.33658</math> <math>x_E^* = -(-1.33658) = 1.33658</math></p> <p>entropy ↑      uncertain ↑</p>

# Active Learning v.s. Clustering-based Active Learning

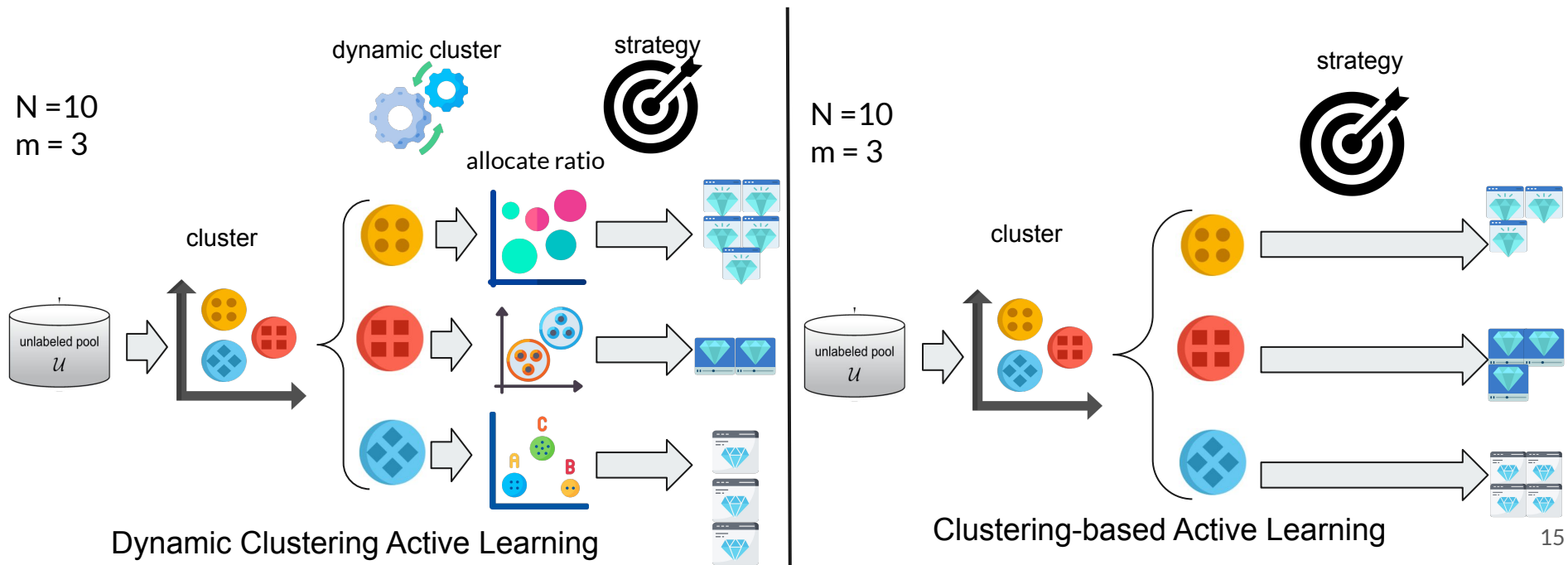


Active Learning

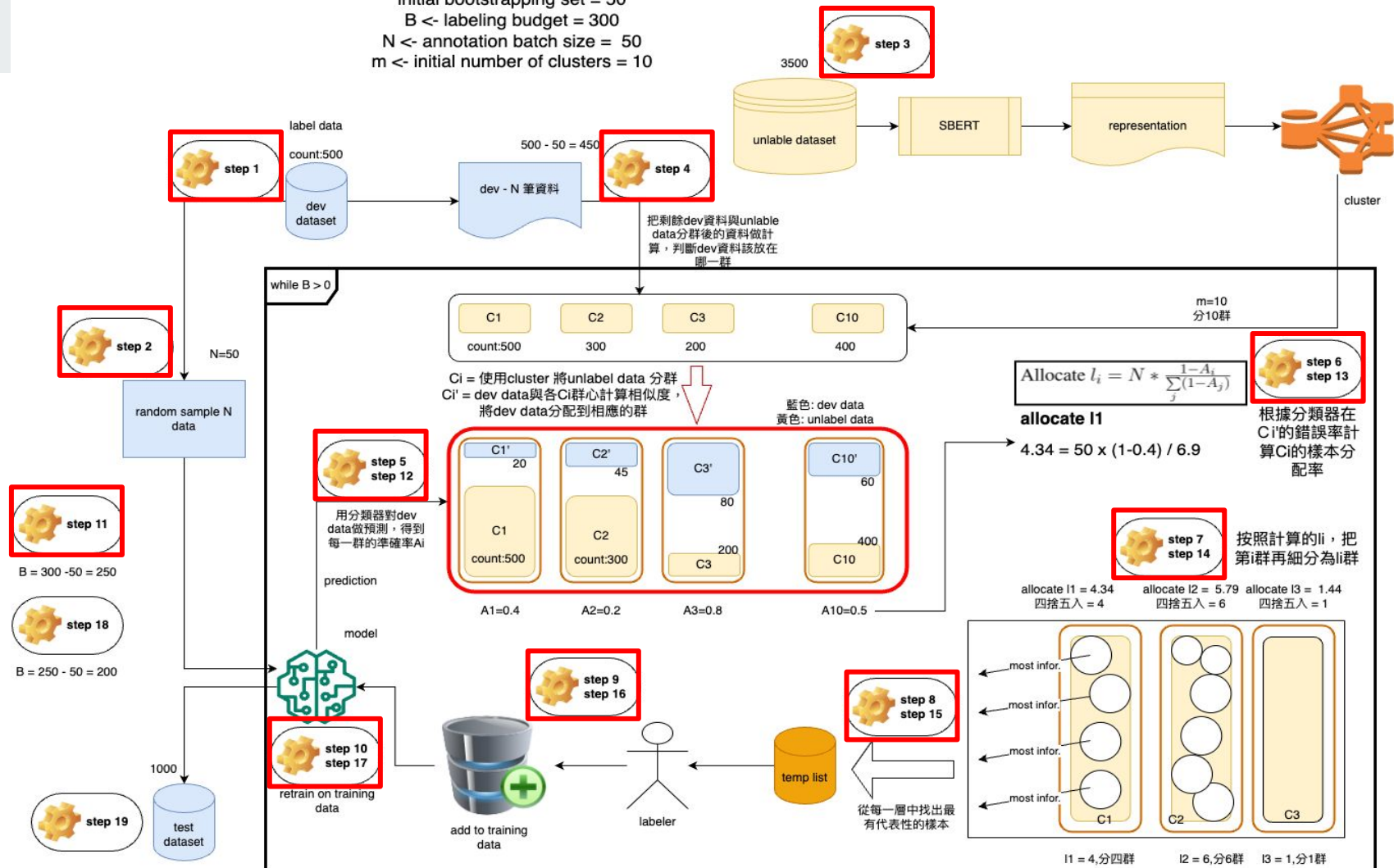


Clustering-based Active Learning

# Dynamic Cluster AL v.s. Clustering-based AL



initial bootstrapping set = 50  
 B <- labeling budget = 300  
 N <- annotation batch size = 50  
 m <- initial number of clusters = 10







# Outline

- Introduction
  - Method
  - **Experiment**
  - Conclusion
- Dataset
  - D-CALM v.s. Baseline with BERT
  - D-CALM v.s. Baseline with SVM



## Dataset

count : 70% : 10% : 20%

<b>Dataset</b>	<b>classes</b>	<b>Pool</b>	<b>Dev</b>	<b>Test</b>
BOOK32	32	14K	2K	4K
CONAN	8	3.5K	0.5K	1K
CARER	6	16K	2K	4K
CoLA	2	8.5K	0.5K	0.5K
Hatespeech	3	17.2K	2.4K	4.9K
MRDA	5	14K	2K	4K
Q-Type	6	4.9K	0.5K	0.5K
Subjectivity	2	7K	1K	2K

- book title
- hatespeech
- emotion detection
- acceptable sentence
- tweets for hatespeech
- dialog act
- question
- snippets from IMDB reviews

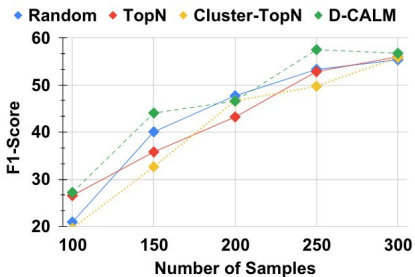
unlabel : label : unlabel

- metric : F1-score

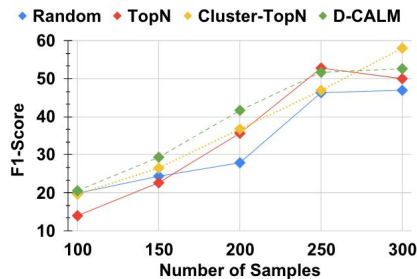
- Strategy : entropy

# D-CALM v.s. Baseline with BERT

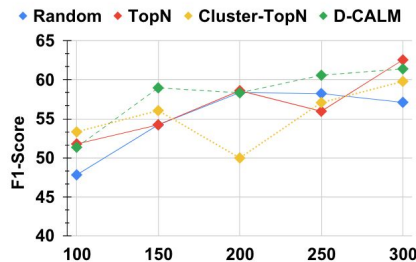
BOOK32



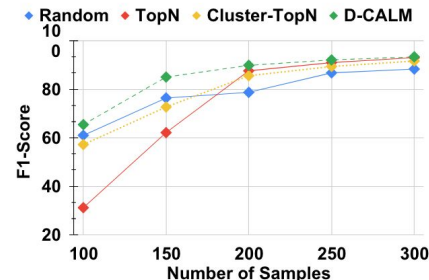
CARER



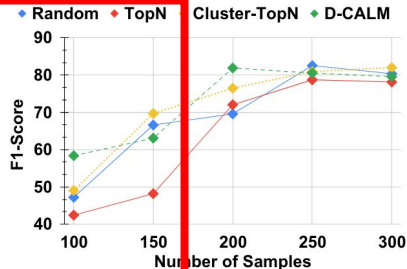
CoLA



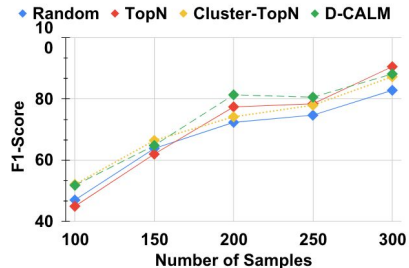
CONAN



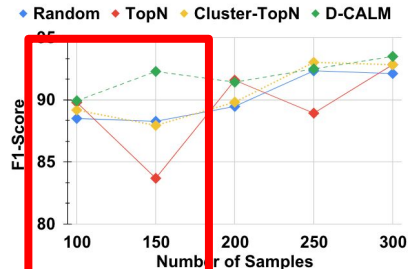
MRDA



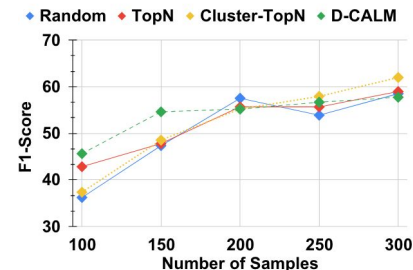
Q-TYPE



Subjectivity



Hatespeech



# D-CALM v.s. Baseline with SVM

BOOK32

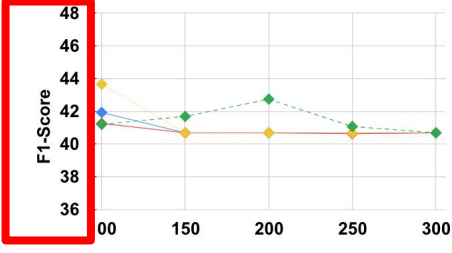
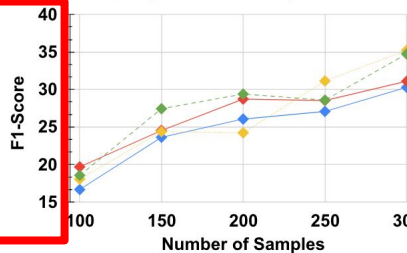
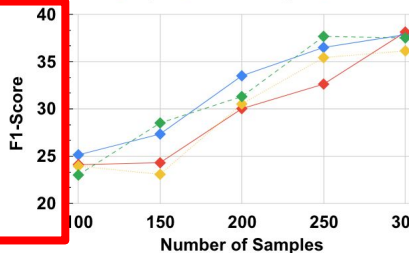
CARER

CoLA

Random TopN Cluster-TopN D-CALM

Random TopN Cluster-TopN D-CALM

Random TopN Cluster-TopN D-CALM



learner : SVM

BOOK32

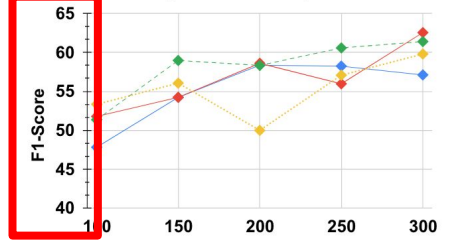
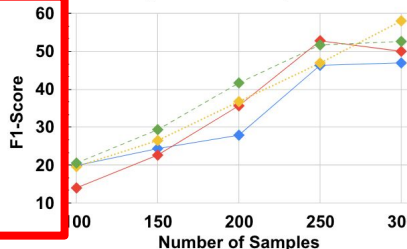
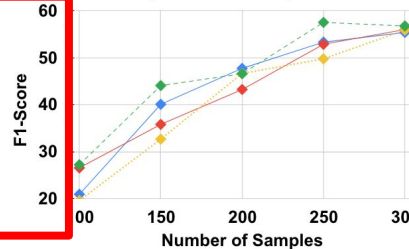
CARER

CoLA

Random TopN Cluster-TopN D-CALM

Random TopN Cluster-TopN D-CALM

Random TopN Cluster-TopN D-CALM

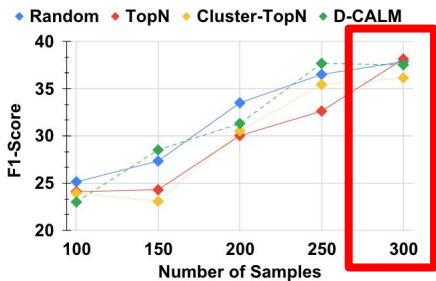


learner: BERT

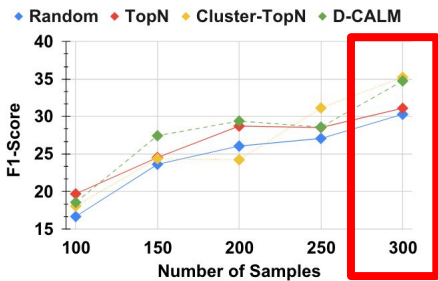
- Strategy : entropy

# D-CALM v.s. Baseline with SVM

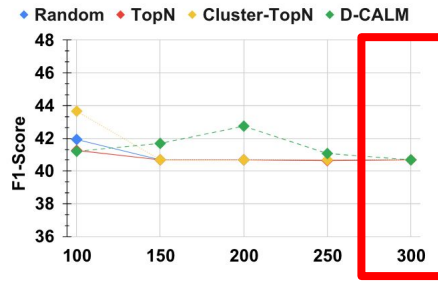
BOOK32



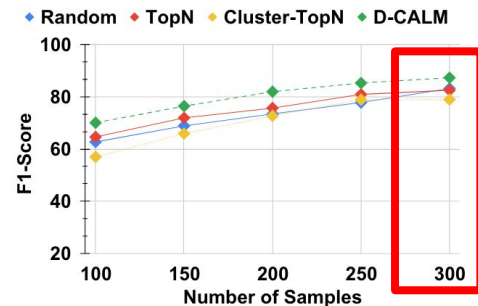
CARER



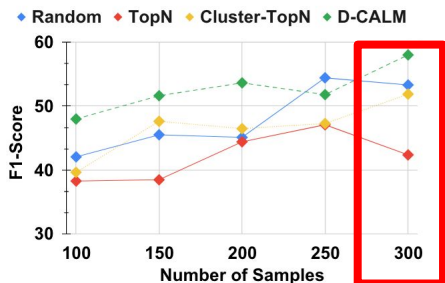
CoLA



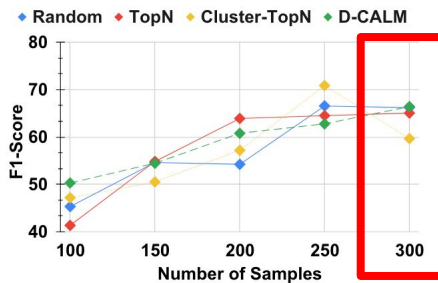
CONAN



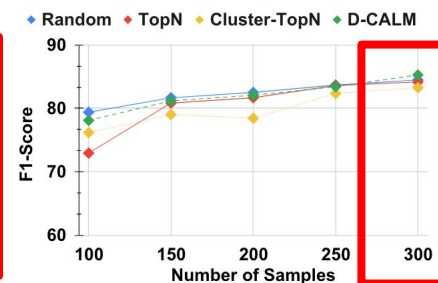
MRDA



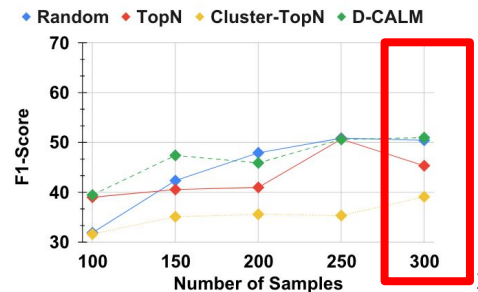
Q-TYPE



Subjectivity



Hatespeech





# Outline

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



## Conclusion

1. The model trained by DCALM can **second clustering** the unlabeled data through the **error rate**, reduce the bias against underrepresented groups in the unlabeled data.
2. DCALM improvements are model-independent