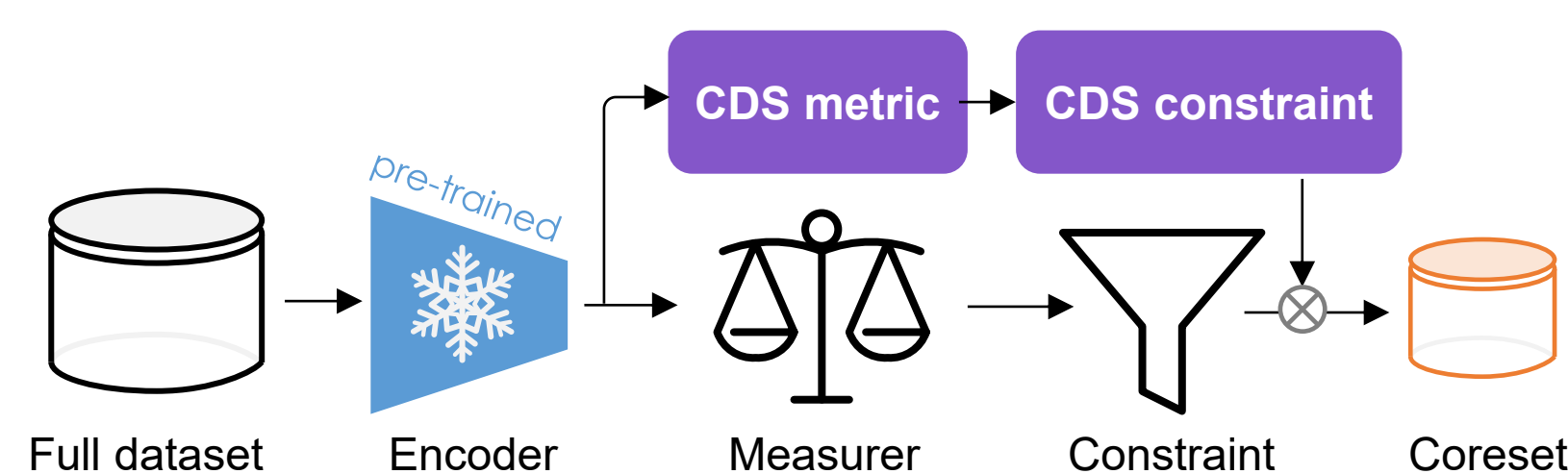


Contributing Dimension Structure of Deep Feature for Coreset Selection

Zhijing Wan¹, Zhixiang Wang^{2,3}, Yuran Wang¹, Zheng Wang¹,
Hongyuan Zhu⁴, and Shin'ichi Satoh^{3,2}

¹Wuhan University, ²The University of Tokyo, ³National Institute of Informatics, ⁴A*STAR

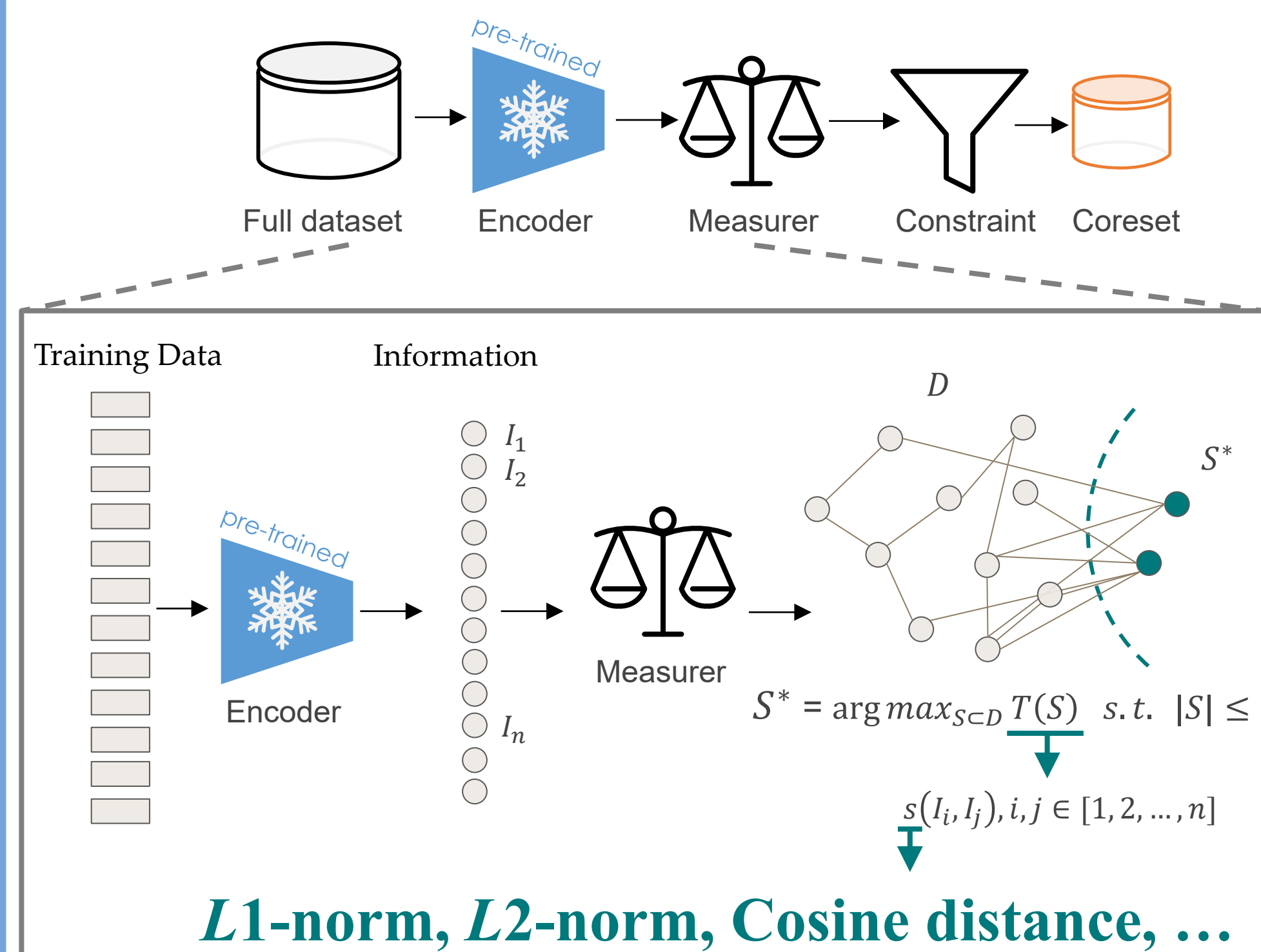
0. Highlights



- Propose CDS metric and constraint to successfully improve the current coreset selection pipeline;
- Propose the CDS metric to introduce the information of the Contributing Dimension Structure (CDS)
- Propose CDS constraint to enrich the diversity of CDS in the coreset

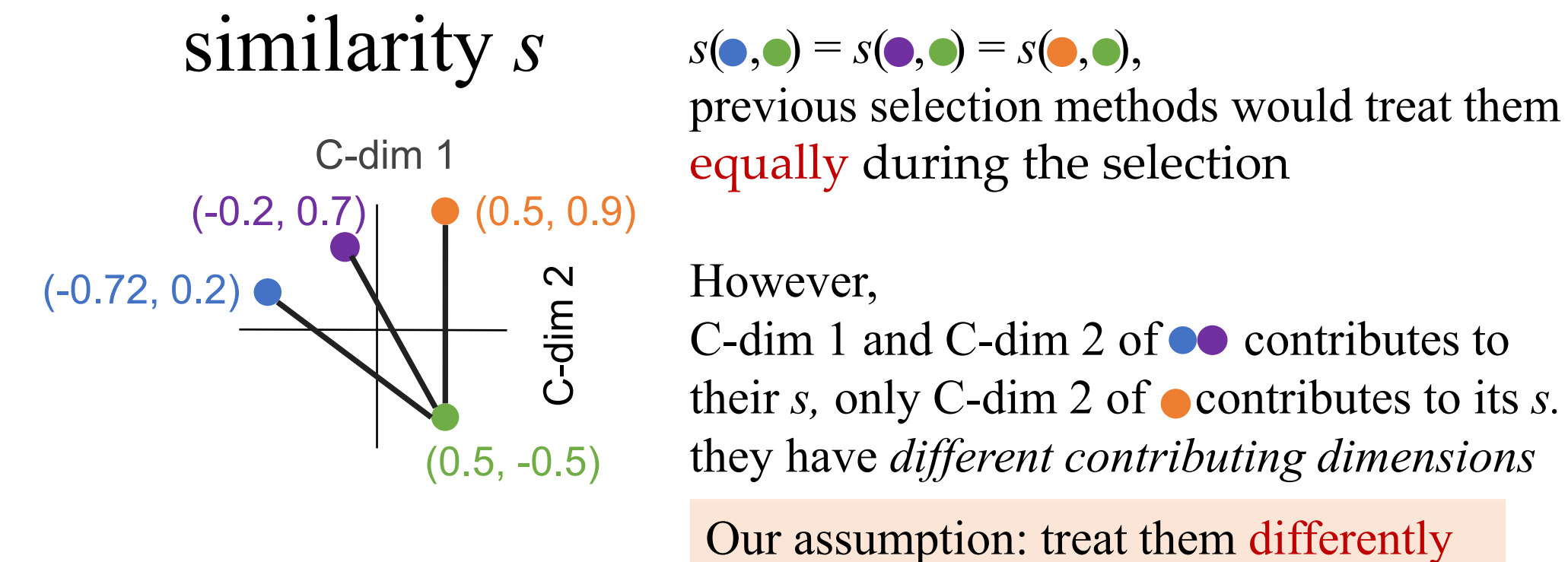
1. Introduction

A. Coreset selection



B. Problems* when using similarity metrics: Similarity metrics ignore

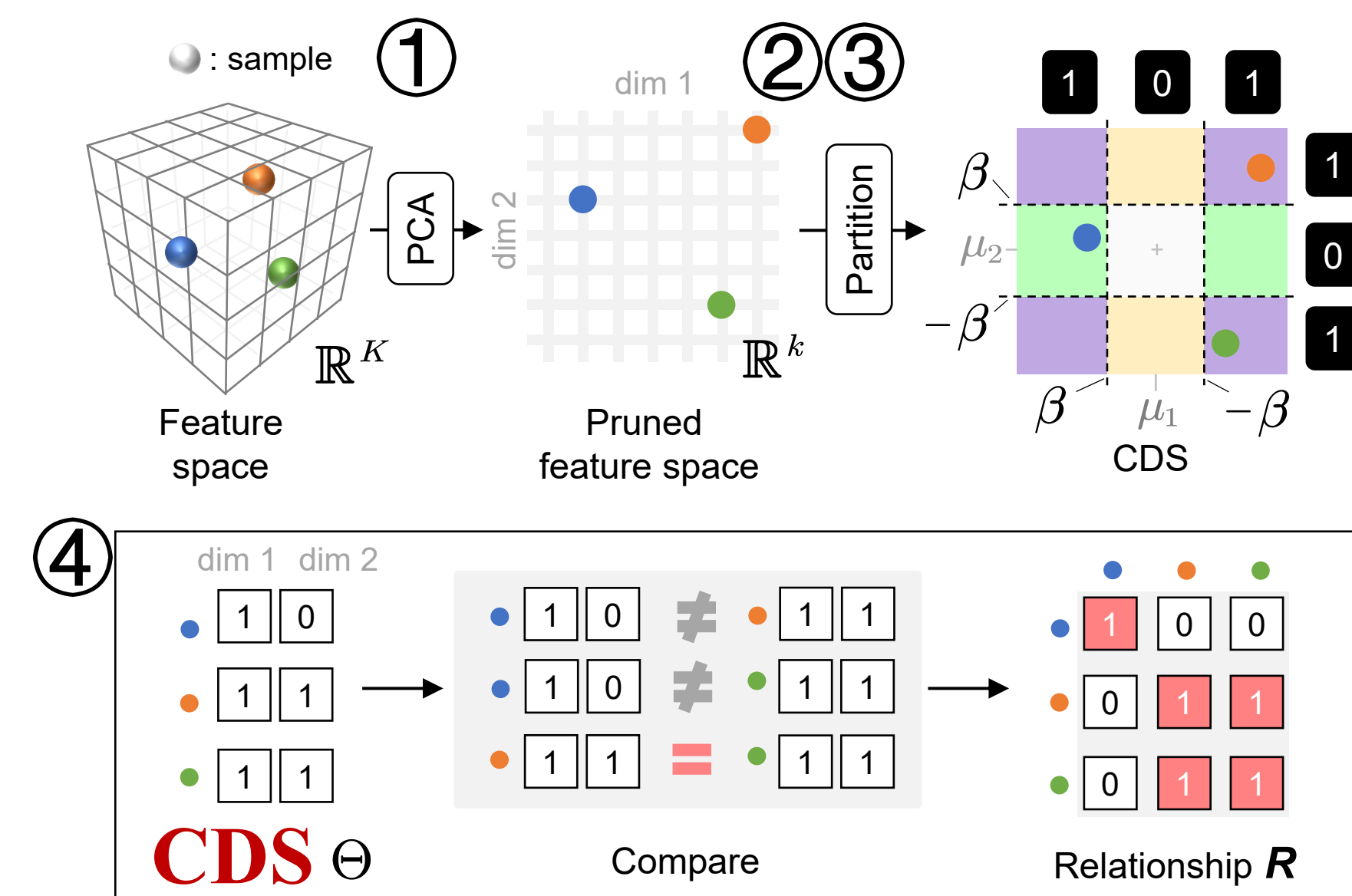
- the redundancy in the feature dimensions;
- disparities among the dimensions that significantly contribute to the final similarity s



*Studied specifically with feature-based selection methods using L2-norm

2. Methodology

A. CDS Metric of Deep Feature



① Dimension Reduction

② Deviation from the Mean

$\sigma = [|f_i^0 - \mu_0|, \dots, |f_i^{k-1} - \mu_{k-1}|] \in \mathbb{R}^k$,
Where $i \in \{0, 1, \dots, N_c - 1\}$.

③ Partition

$\theta(x_i) = [\mathbb{I}(|f_i^0 - \mu_0|), \dots, \mathbb{I}(|f_i^{k-1} - \mu_{k-1}|)]$

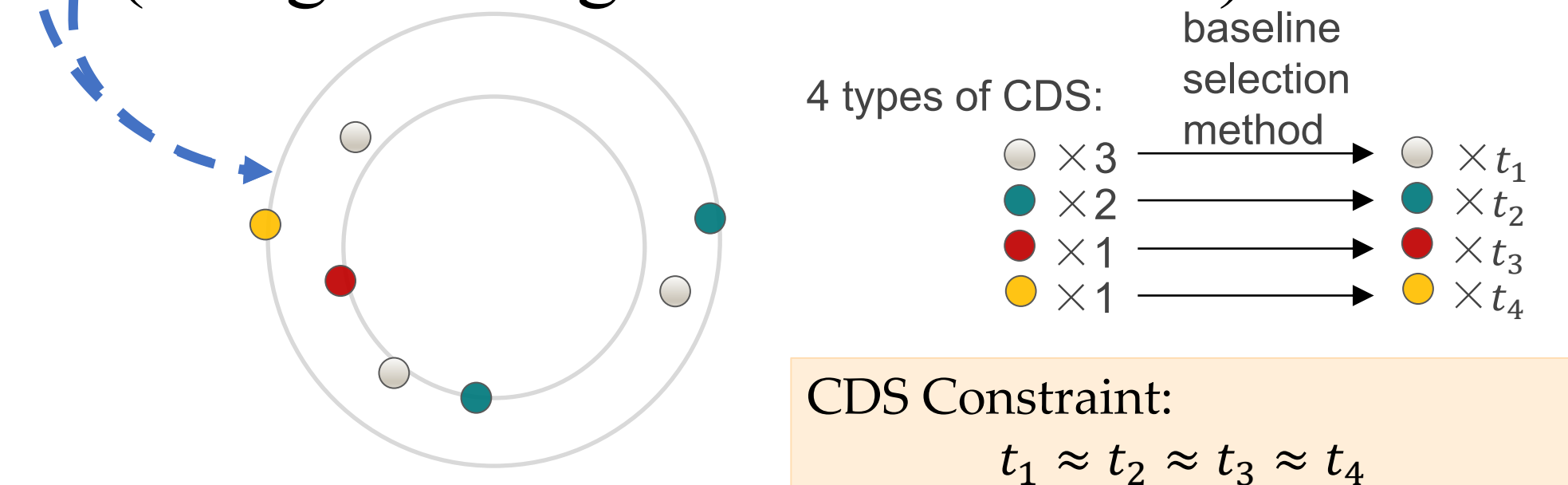
CDS $\mathbb{I}(\Delta f) = \begin{cases} 1, \Delta f > \beta \\ 0, \Delta f \leq \beta \end{cases}$

④ Comparison

C. Coreset Selection with CDS Constraints

Hard version

- 1st stage clustering
- 2nd stage clustering
- data selection (using existing selection method)



+ K-Center Greedy, Least Confidence, Moderate-DS

Soft version

Algorithm 1: Coreset Selection with Soft CDS Constraint (Class Balanced Sampling)

Input: Train set: \mathcal{D} ; data classes: C ; budget: b ; pre-trained model: θ_{pre} ; objective function: T ; constraint function: H ; parameter β .

Output: Coreset $S \subset \mathcal{D}$.

```

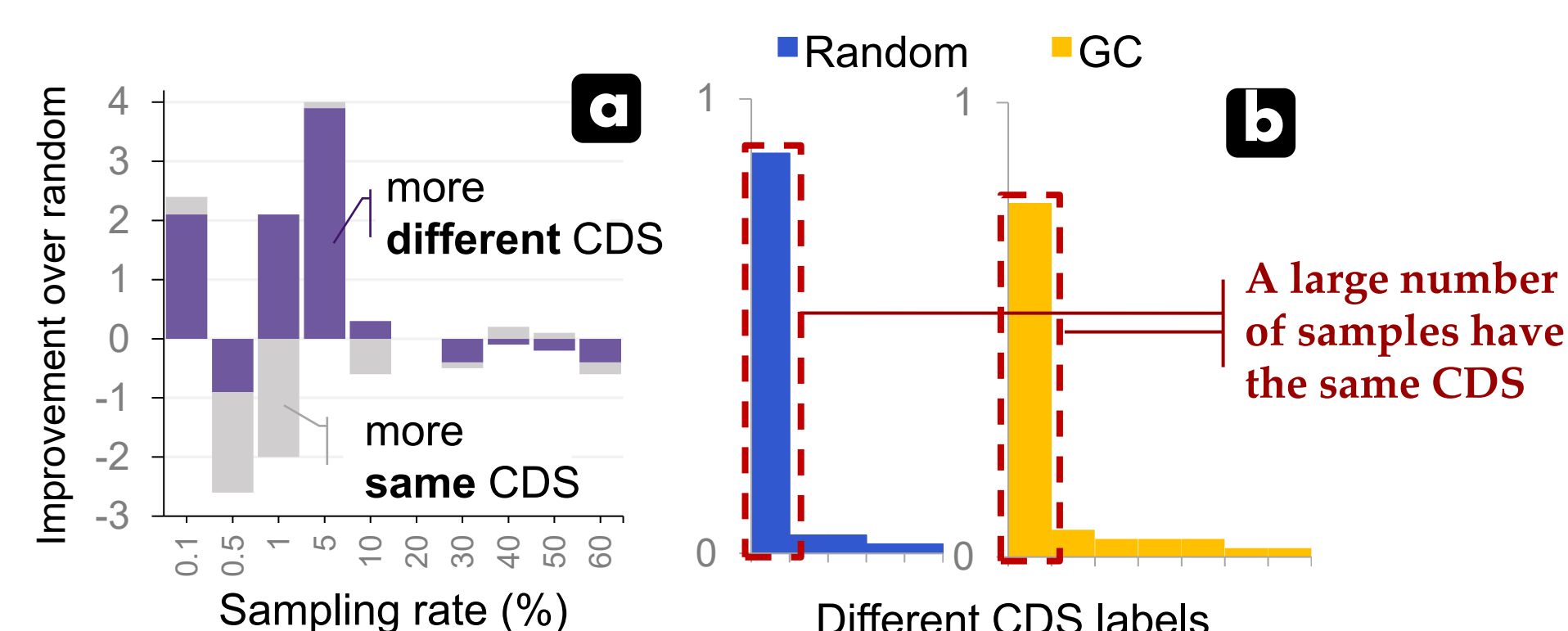
1 Initial  $S \leftarrow S_0$ 
2 for  $i \leftarrow 0$  to  $C - 1$  by 1 do
3    $F \leftarrow M(\mathcal{D}^i; \theta_{pre})$ 
4    $R^i \leftarrow \text{CDS-Metric}(F, \beta)$ 
5   Initial  $S' \leftarrow S_0$ 
6   while  $|S'| < b$  do
7      $V \leftarrow \mathcal{D}^i \setminus S'$ 
8      $e \leftarrow \arg \max_{e \in V} T(e|S') \times H(R^i)$ 
9      $S' \leftarrow S' \cup \{e\}$ 
10  end
11  $S \leftarrow S'$ 
12 end
```

+ CRAIG*

+ GC*

*See paper for more details

B. Empirical findings



- More data with different CDSs need to be sampled into the coreset
- The coresets selected by the existing SOTA methods are sub-optimal

Propose the **CDS Constraints** to improve SOTA selection methods

3. Experiments

A. Class-balanced sampling

Method	Sampling rates				
	0.1%	0.5%	1%	5%	20%
Random	18.9±0.2	29.5±0.4	39.3±1.5	62.4±1.7	74.7±1.9
KCG	18.7±2.9	27.4±1.0	31.6±2.1	53.5±2.9	73.2±1.3
Forgetting	21.8±1.7	29.2±0.7	35.0±1.1	50.7±1.7	66.8±2.5
LC	14.8±2.4	19.6±0.8	20.9±0.4	37.4±1.9	56.0±2.0
CRAIG	21.1±2.4	27.2±1.0	31.5±1.5	45.0±2.9	58.9±3.6
Cal	20.8±2.8	32.0±1.9	39.1±3.2	60.7±0.8	72.2±1.5
Glistar	19.5±2.1	29.7±1.1	33.2±1.1	47.1±2.6	65.7±1.7
GC	22.9±1.4	34.0±1.3	42.0±3.0	66.2±1.0	75.6±1.4
M-DS	21.0±3.0	31.8±1.2	37.7±1.4	63.4±2.2	78.0±1.3
GC+Ours	24.6±1.7	36.4±1.0	43.1±1.8	67.1±0.6	76.9±0.2
Δ	1.7 \uparrow	2.4 \uparrow	1.1 \uparrow	0.9 \uparrow	1.3 \uparrow
M-DS+Ours	22.0±2.0	33.0±1.3	40.7±1.0	64.9±0.8	79.6±0.4
Δ	1.0 \uparrow	1.2 \uparrow	3.0 \uparrow	1.5 \uparrow	0.0 \uparrow

Table 1: Comparison on the class-balanced sampling set- Figure 5: Performance improvement over baselines. We im-
ting. We train randomly initialized ResNet-18 on coresets of prove current methods with our proposed CDS metric and
CIFAR-10 selected by different methods and then test them constraint. We compare the improved versions with respec-
on the test set of CIFAR-10. Green **GC** emphasizes the best baseline on CIFAR-10 (a-c) and TinyImageNet (d-f)
performance at each sampling rate. Δ denotes the improve- under the class-balanced sampling setting. The improved
ment of baseline+Ours over baseline. versions consistently outperform baselines, suggesting that
increasing the diversity of CDS in the coreset can univer-
sally enhance existing coreset selection methods.

B. Ablation and Parameter Studies

	dim reduction	partition	CDS-r	cons- traint	
(v1)	\times	\times	\times	\times	34.0±1.3
(v2)	\times	\times	\times	\checkmark	32.8±0.7
(v3)	\checkmark	\times	\times	\checkmark	34.3±2.5
(v4)	\checkmark	\checkmark	\times	\checkmark	33.5±1.1
full	\checkmark	\checkmark	\checkmark	\checkmark	36.4±1.0

Table 2: Ablation study on 0.5% of the CIFAR-10

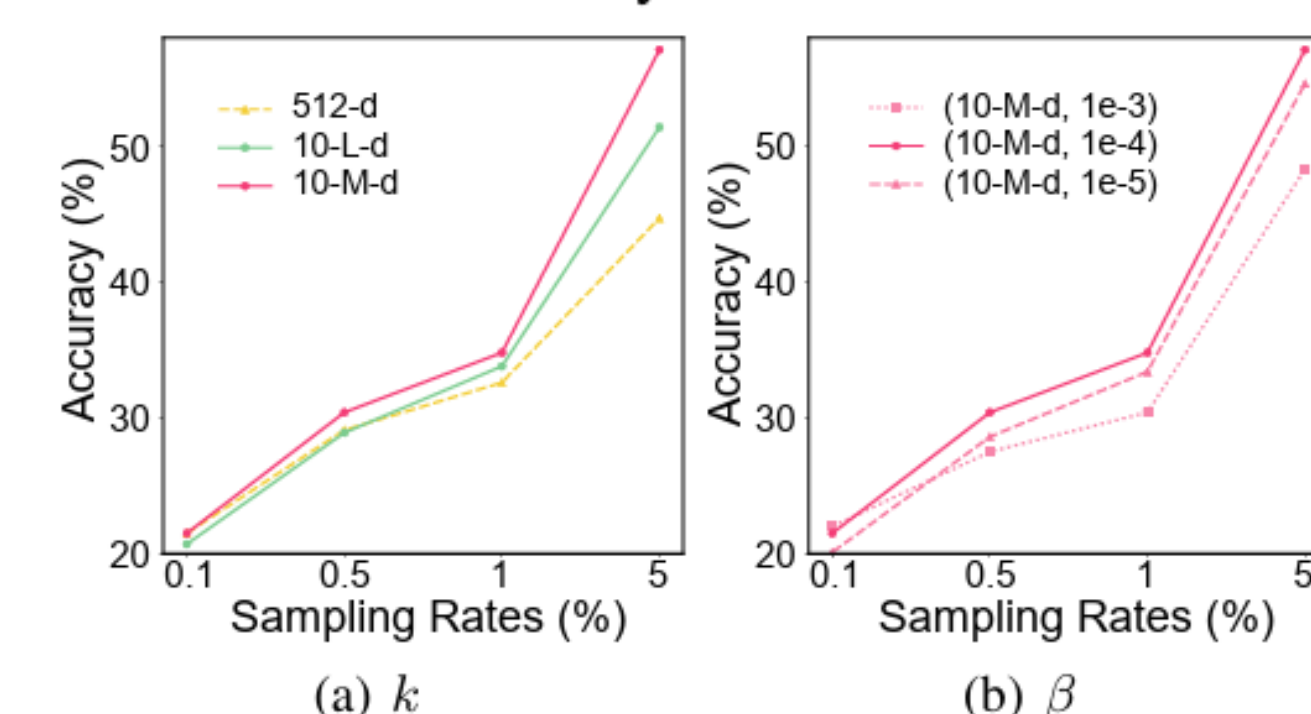
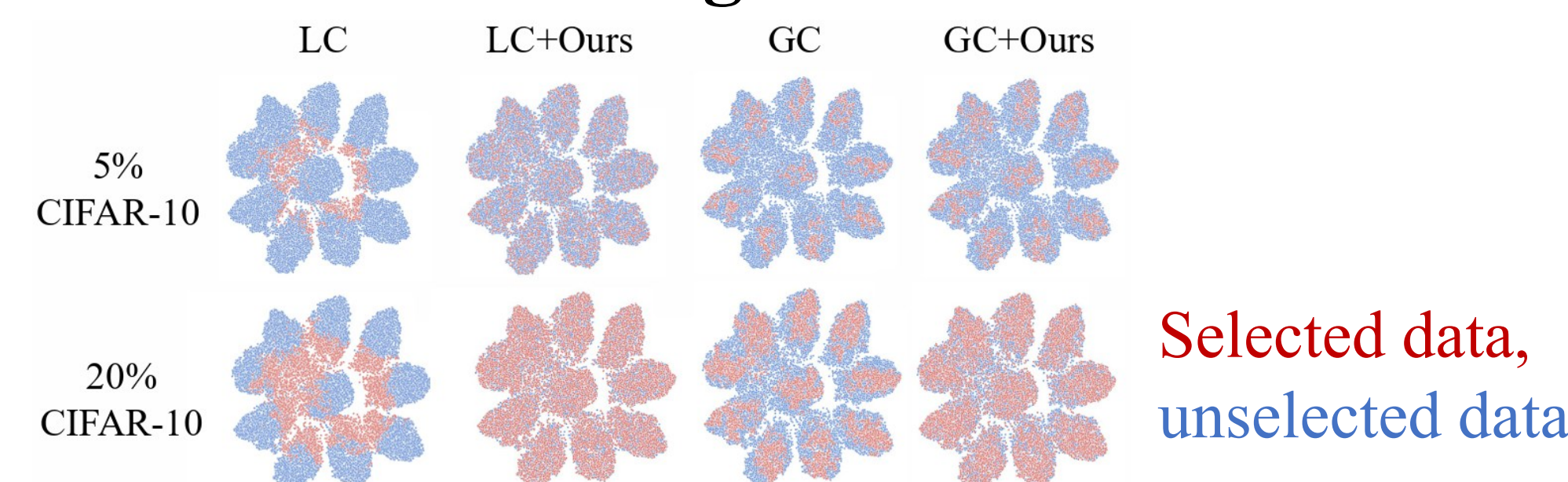


Figure 7: Parameter analysis. It shows that our method achieves the best improvement compared to the baseline method (CRAIG) when $K=10$ -M-D and $\beta=1e-4$.

C. TSNE embeddings of coreset



Contact information: wanzjwhu@whu.edu.cn