

On the Surrogate Gap between Contrastive and Supervised Losses



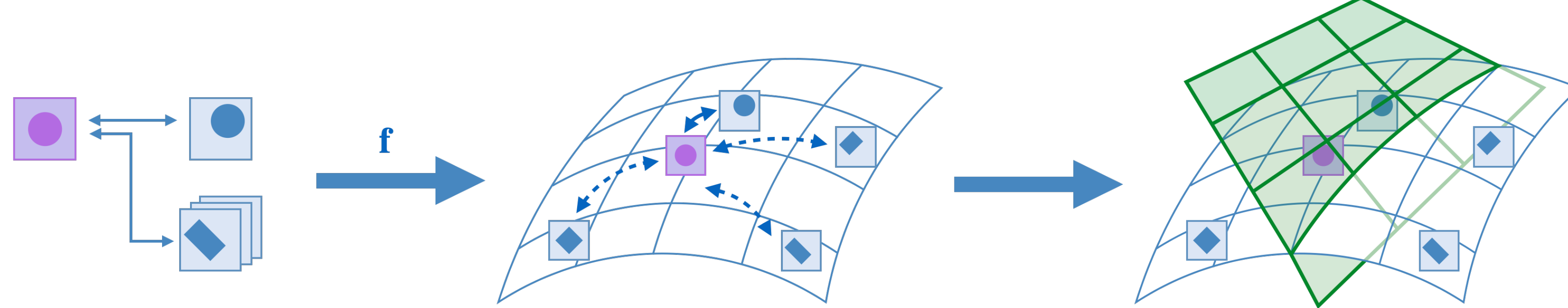
Han Bao^{*1,2} Yoshihiro Nagano^{*1,2} Kento Nozawa^{*1,2}

¹The University of Tokyo, Japan / ²RIKEN AIP, Japan / ^{*}Equal contribution

arXiv: 2110.02501

Introduction

Contrastive Unsupervised Representation Learning (CURL)

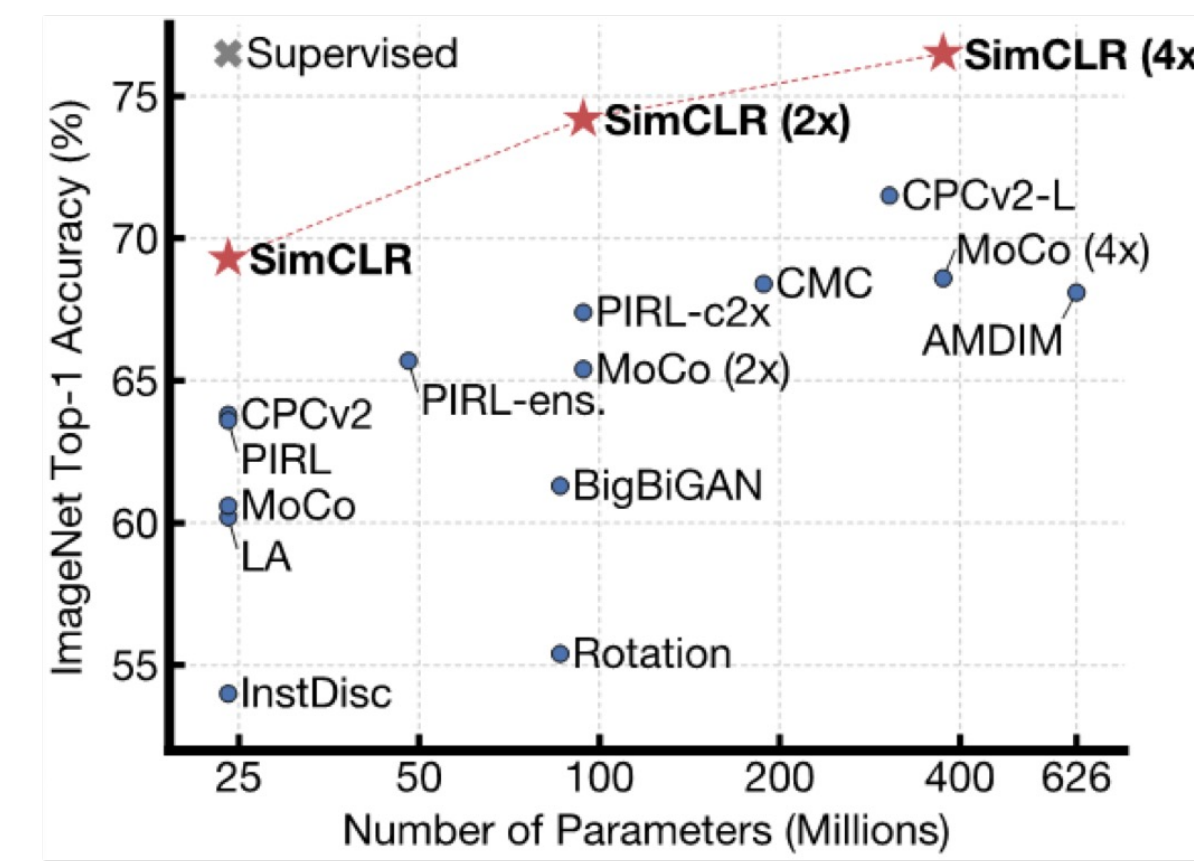


Learning Scheme

CURL aims to learn a representation function f by making **semantically similar (positive) pair closer** while **randomly drawn (negative) pair further**.

Inference

By performing fine-tuning linear classifier on top of the learned f , we can get good empirical performance for the downstream task.



[Chen et al., 2020]

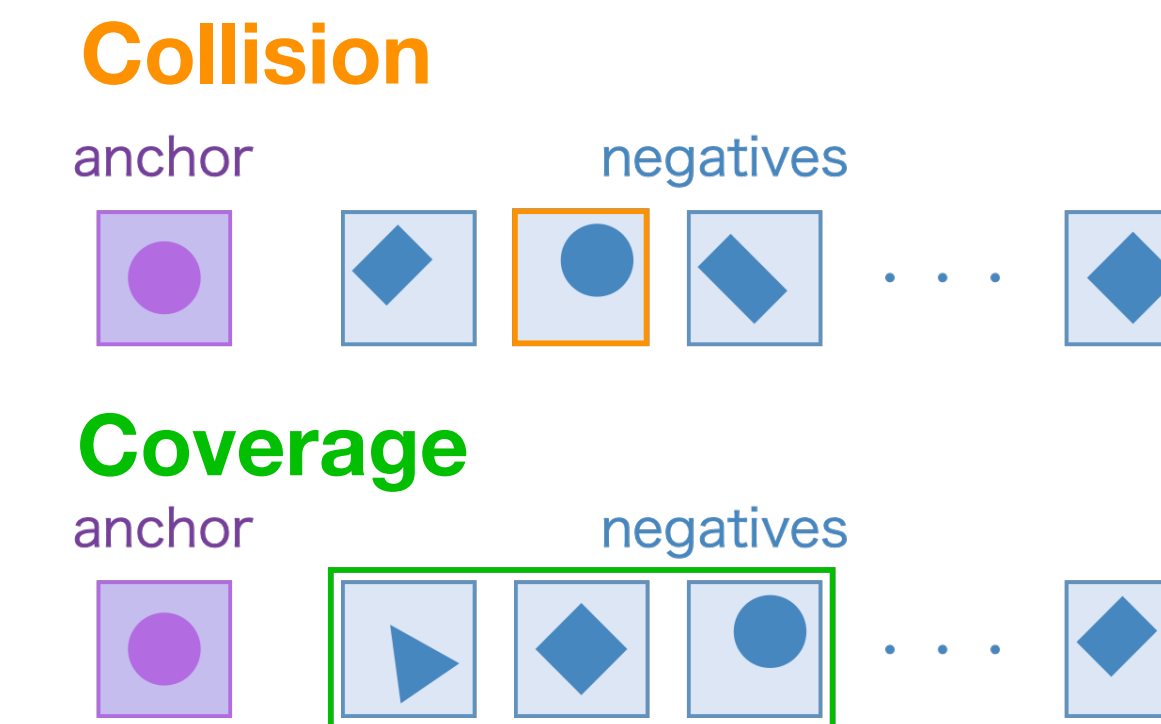
Analysis

Existing Work: Collision-Coverage Formulation

Collision-Coverage Formulation

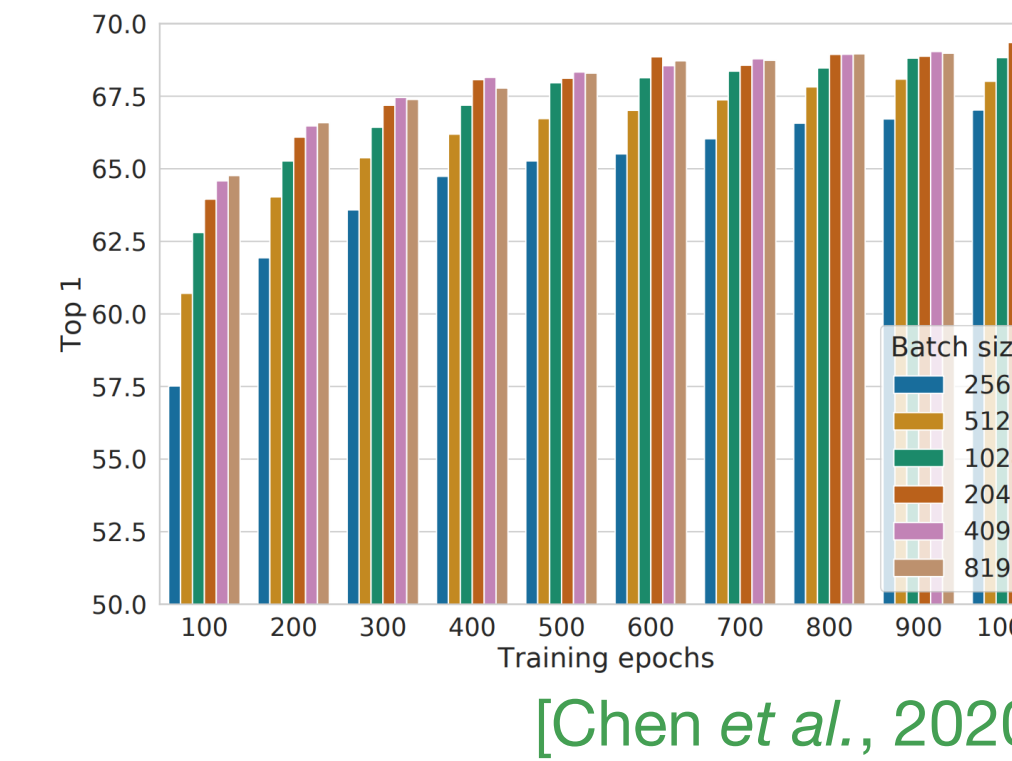
Rewrite the contrastive loss using the conditions under which the label **collision/coverage** occurs.

- ❖ **Collision**: randomly drawn negative class collides with the anchor class.
- ❖ **Coverage**: negative classes covers entire label space of the downstream classification.



Issue: Disagreement with Experiment

- ❖ Theory predicts the downstream performance degrades with increase in K because of the **label collision**, while larger K helps performance in practice.
- ❖ Upper bound becomes **exponentially loose** in K .



[Chen et al., 2020]

	UPPER BOUND	REFERENCE
$R_{\mu\text{-supv}}(\mathbf{f}) \leq$	$\frac{1}{(1-\tau_K)v_{K+1}} \{R_{\text{cont}}(\mathbf{f}) - \mathbb{E} \ln(\text{Col} + 1)\}$	Arora et al. (2019)
	$\frac{1}{v_{K+1}} \{2R_{\text{cont}}(\mathbf{f}) - \mathbb{E} \ln(\text{Col} + 1)\}$	Nozawa & Sato (2021)
	$\frac{2}{1-\tau_K} \left[\frac{2(C-1)H_{C-1}}{K} \right] \{R_{\text{cont}}(\mathbf{f}) - \mathbb{E} \ln(\text{Col} + 1)\}$	Ash et al. (2022)

τ_K : collision prob.
 v_{K+1} : coverage prob.

Our Approach: Surrogate Bound

Main Result

Directly transform the contrastive loss to the supervised loss by **linearizing the log-sum-exp functions**.

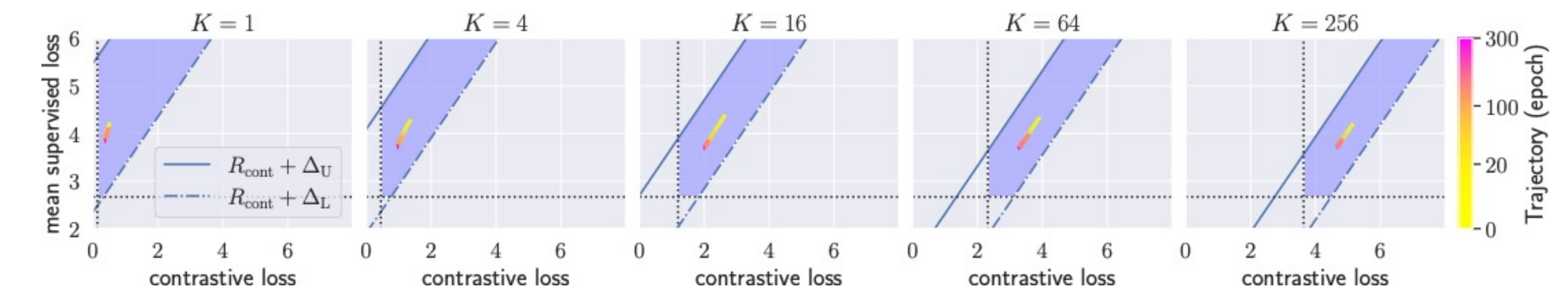
$$R_{\text{cont}}(\mathbf{f}) + \Delta_L \leq R_{\mu\text{-supv}}(\mathbf{f}) \leq R_{\text{cont}}(\mathbf{f}) + \Delta_U$$

$$\Delta_L = \Delta_U = O\left(\ln \frac{1}{K}\right)$$

- ❖ We can interpret the contrastive loss as the surrogate estimator of the mean supervised loss in a sense that these two losses behave similarly.
- ❖ Coefficients of the bounds are constant with respect to C and K .
- ❖ Surrogate gap (intercept) decreases as K increases; agrees with experimental facts.

Experiments

Synthetic Dataset



Setting

- ❖ 2D synthetic dataset circle with $C = 10$.
- ❖ f : 3-layer MLP (# hidden units is 256) with ReLU activation.

Result

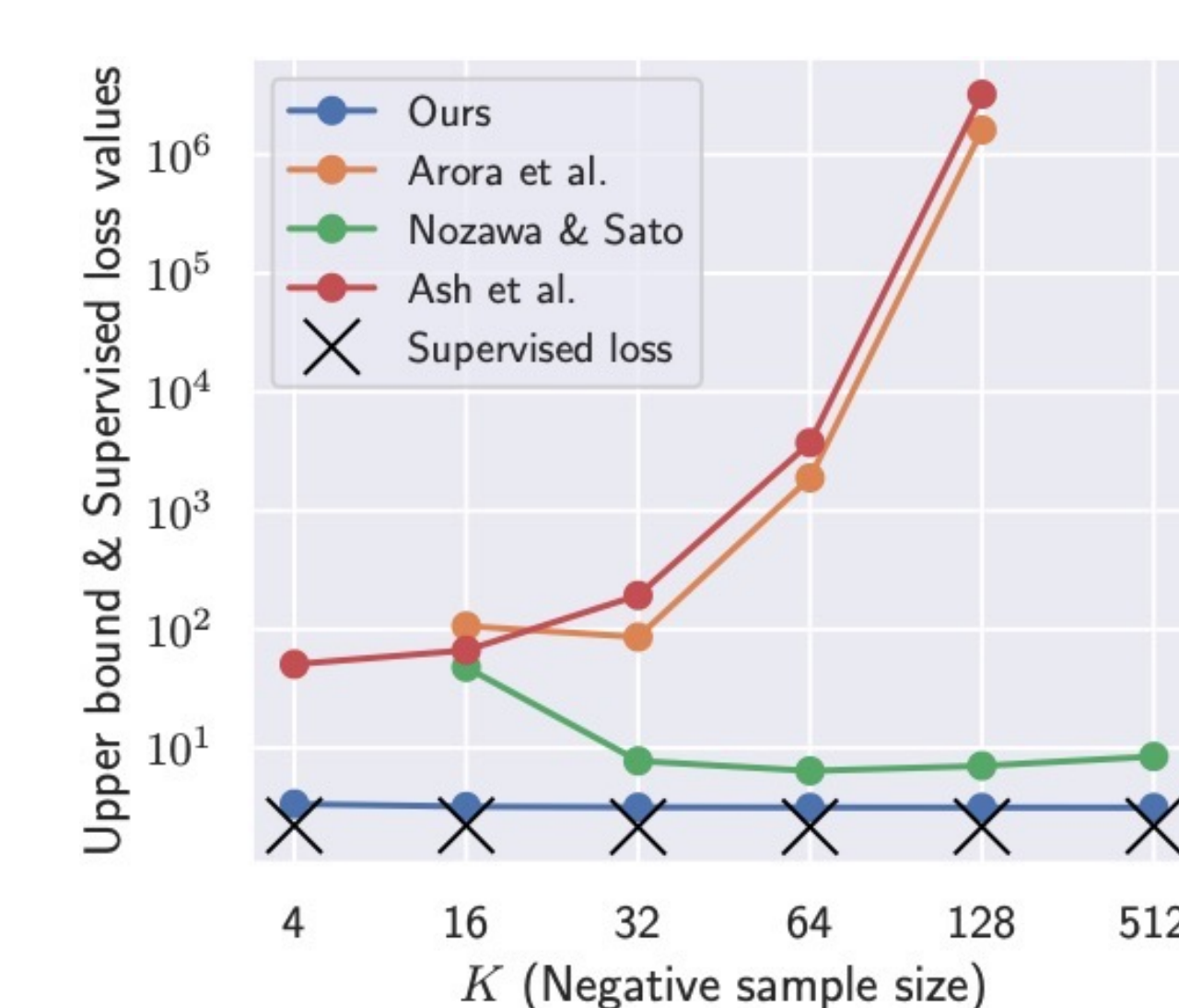
- ❖ Our surrogate bounds capture the learning dynamics well in different negative sample sizes.

Vision & Language Datasets

Setting

- ❖ Dataset: CIFAR-10/100 (vision) & Wiki-3029 (language).
- ❖ f : ResNet-18-based [He et al., 2016] (vision) & fasttext-based [Joulin et al., 2017] (language)

Vision Dataset



Language Dataset

C	mean classifier				linear classifier			
	8	64	256	1024	8	64	256	1024
3029	37.86 (0.19)	40.36 (0.13)	40.85 (0.09)	40.94 (0.13)	41.70 (0.14)	43.00 (0.15)	43.26 (0.14)	43.30 (0.21)
2000	42.09 (0.17)	44.06 (0.16)	44.38 (0.09)	44.24 (0.23)	45.13 (0.29)	46.33 (0.14)	46.46 (0.15)	46.37 (0.16)
1000	48.07 (0.12)	48.77 (0.23)	48.75 (0.05)	48.66 (0.10)	50.86 (0.25)	51.21 (0.33)	51.06 (0.24)	50.94 (0.23)
500	52.66 (0.71)	52.66 (0.69)	53.04 (0.67)	53.72 (0.63)	55.41 (0.39)	55.41 (0.43)	55.52 (0.51)	55.69 (0.47)

- ❖ Existing theories result in exponentially loose prediction of the downstream supervised loss for the test data in the vision dataset.
- ❖ Proposed upper bound agrees with the actual supervised loss well in all range of K .
- ❖ Larger K moderately helps performance as predicted from our theory.

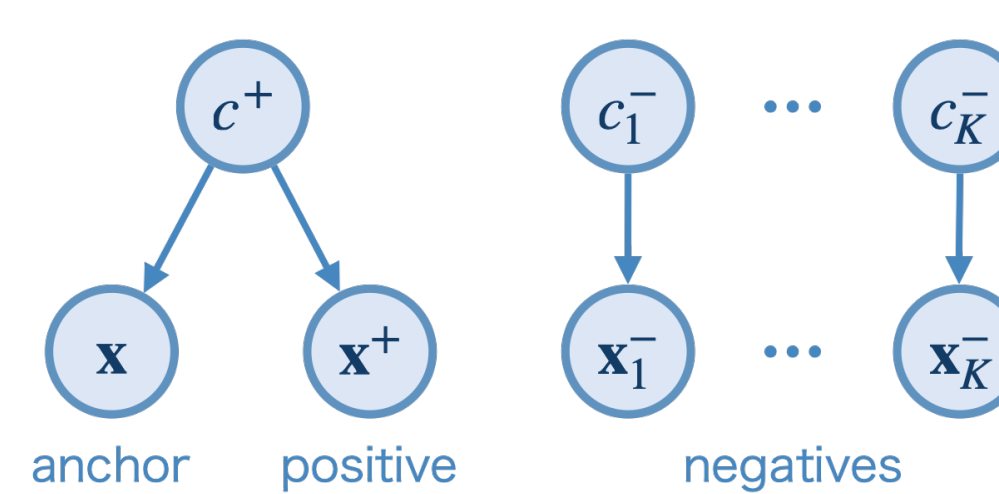
A. Contrastive loss behaves as the surrogate estimator.

Q. What is the underlying mechanism of the success?

Problem Setup Based on [Arora et al., 2019]

Data Generating Process

- ❖ Draw 1 **positive**/ K **negative** classes $c^+, \{c_k^-\}_{k \in [K]} \sim \mathbb{P}(Y)$
- ❖ Draw an anchor/positive sample $\mathbf{x}, \mathbf{x}^+ \sim \mathbb{P}(X|Y = c^+)$
- ❖ Draw K negative samples $\mathbf{x}_k^- \sim \mathbb{P}(X|Y = c_k^-)$



Training Objective

Train the representation function f by minimizing the following objective:

$$R_{\text{cont}}(\mathbf{f}) = \mathbb{E}_{\substack{c^+, \{c_k^-\} \\ \mathbf{x}, \mathbf{x}^+, \{\mathbf{x}_k^-\}}} \left[-\ln \frac{e^{\mathbf{f}(\mathbf{x})^\top \mathbf{f}(\mathbf{x}^+)}}{e^{\mathbf{f}(\mathbf{x})^\top \mathbf{f}(\mathbf{x}^+)} + \sum_{k \in [K]} e^{\mathbf{f}(\mathbf{x})^\top \mathbf{f}(\mathbf{x}_k^-)}} \right]$$

Downstream performance

Evaluate the learned f by the downstream mean supervised loss:

$$R_{\mu\text{-supv}}(\mathbf{f}) = \mathbb{E}_{\mathbf{x}, y \sim \mathbb{P}} \left[-\ln \frac{e^{\mathbf{f}(\mathbf{x})^\top \mu_y}}{\sum_{c \in \mathcal{Y}} e^{\mathbf{f}(\mathbf{x})^\top \mu_c}} \right]$$

Notations

C : # classes $\mu_c = \mathbb{E}_{\mathbf{x}|c}[\mathbf{f}(\mathbf{x})]$ $\mathbf{W}^\mu = [\mu_1 \dots \mu_C]^\top$ $\inf_{\mathbf{W} \in \mathbb{R}^{C \times d}} R_{\text{supv}}(\mathbf{W}\mathbf{f}) \leq R_{\mu\text{-supv}}(\mathbf{f})$