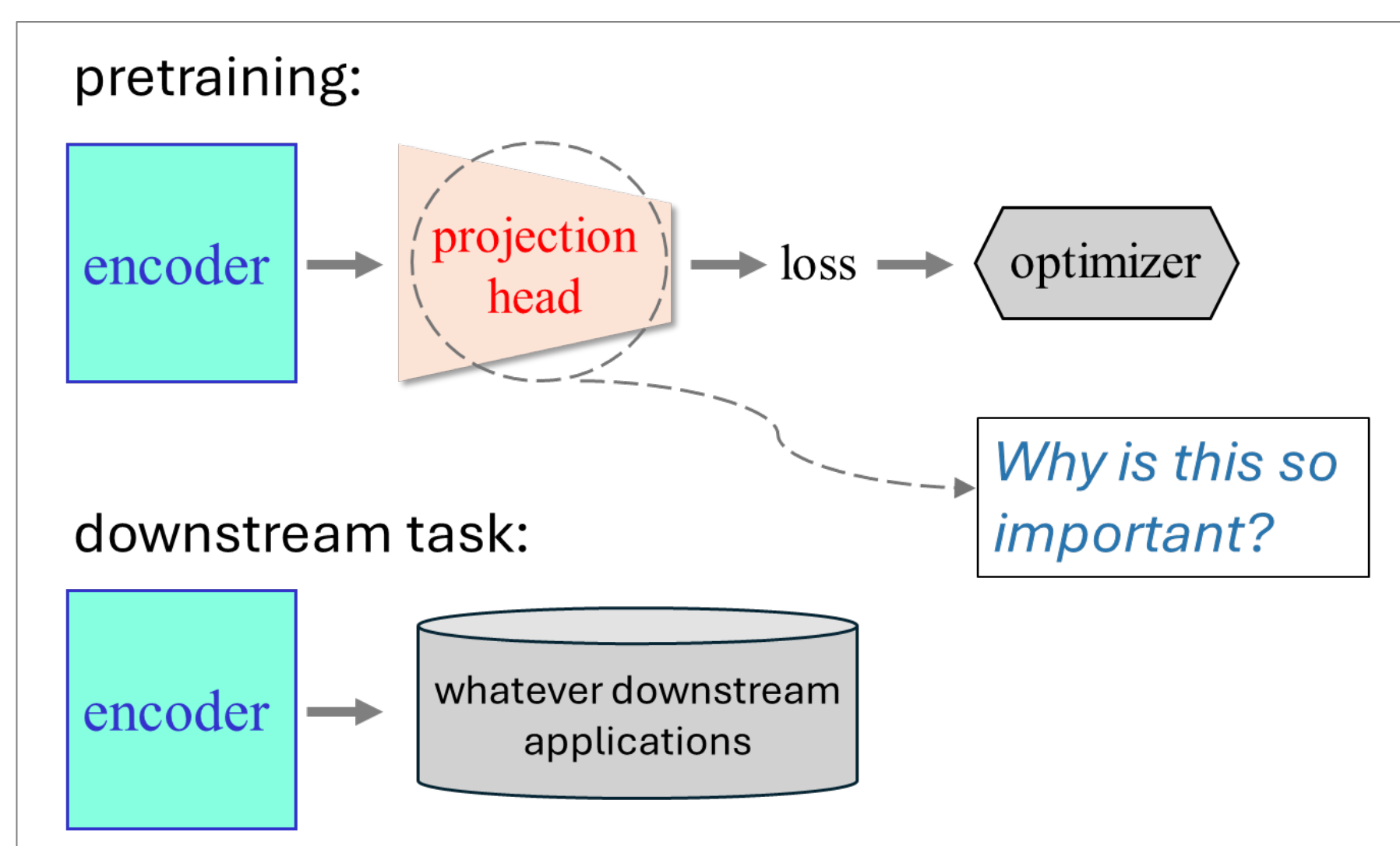# Investigating the Benefits of Projection Head for Representation Learning

Yihao Xue, Eric Gan, Jiayi Ni, Siddharth Joshi, Baharan Mirzasoleiman

## Introduction

An effective technique for obtaining high-quality representations is **adding a projection head** on top of the encoder during pretraining, then **discarding** it and using the pre-projection representations for downstream tasks.



pretraining:

downstream task:

Why is this so important?
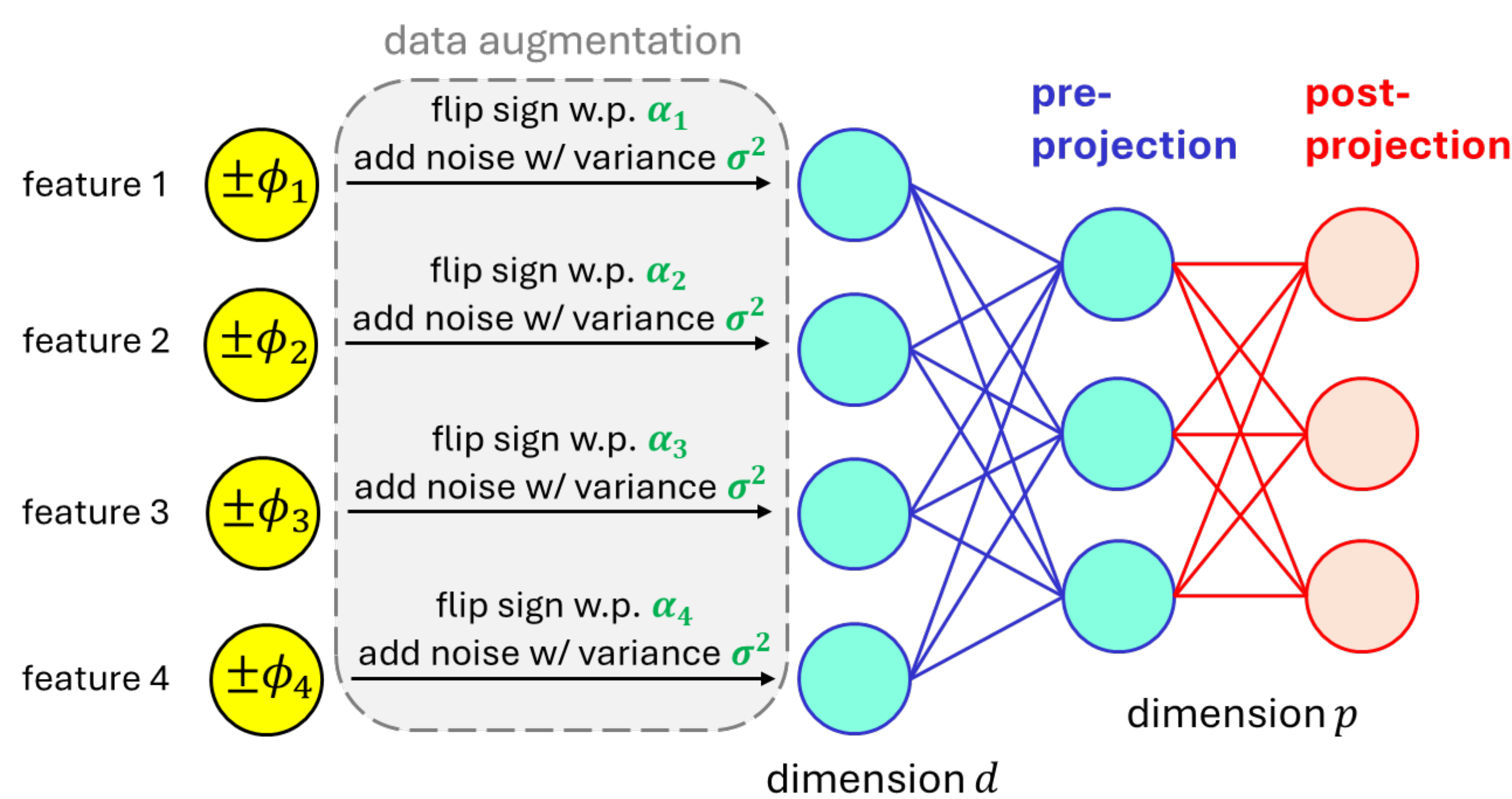
whatever downstream applications

**Our main conclusion:** pre-projection representations represent input features more equally or capture a broader range of features compared to post-projection representations.

## Self-supervised Contrastive Learning

### Pretraining loss
max sim. b/w augmentations of the same examples
min sim. b/w augmentations of different examples

### Simple data & linear model



data augmentation

feature 1  $\pm\phi_1$
flip sign w.p. $\alpha_1$
add noise w/ variance $\sigma^2$

feature 2  $\pm\phi_2$
flip sign w.p. $\alpha_2$
add noise w/ variance $\sigma^2$

feature 3  $\pm\phi_3$
flip sign w.p. $\alpha_3$
add noise w/ variance $\sigma^2$

feature 4  $\pm\phi_4$
flip sign w.p. $\alpha_4$
add noise w/ variance $\sigma^2$

pre-projection

post-projection

dimension $d$

dimension $p$

## Self-supervised Contrastive Learning (Cont'd)

**Theorem:** Define $\beta_i = \frac{(1-\alpha_i)^2 \phi_i^2}{\phi_i^2 + \sigma^2}$, $\gamma_i = \sqrt{\frac{(1-\alpha_i)\phi_i}{\phi_i^2 + \sigma^2}}$

selection

feature $i$ is weighted by
$$\begin{cases} 0, & \text{if } \beta_i \text{ is not among the } p \text{ largest } p \ \beta's \\ \boldsymbol{\gamma_i} \ \text{pre} - \text{projection and} \ \boldsymbol{\gamma_i^2} \ \text{post} - \text{projection} \end{cases}$$

weighting

### Key insights

- The model selects and weights features based on the interplay between feature strengths ($\phi$), noise ($\sigma$), and data augmentation ($\alpha$)
- Features are weighted **more equally** pre-projection than post-projection

### When is it beneficial to use pre-projection representations?

Assume that **feature $i^*$** is the only one useful for the downstream task. Ideally, pretraining should assign a large weight to it relative to other features. If this doesn't occur, i.e., pretraining assigns it a small weight (small $\gamma_{i^*}$), it would be better to use the pre-projection representations.

Here are some concrete scenarios considering the interaction between $\phi, \sigma, \alpha$
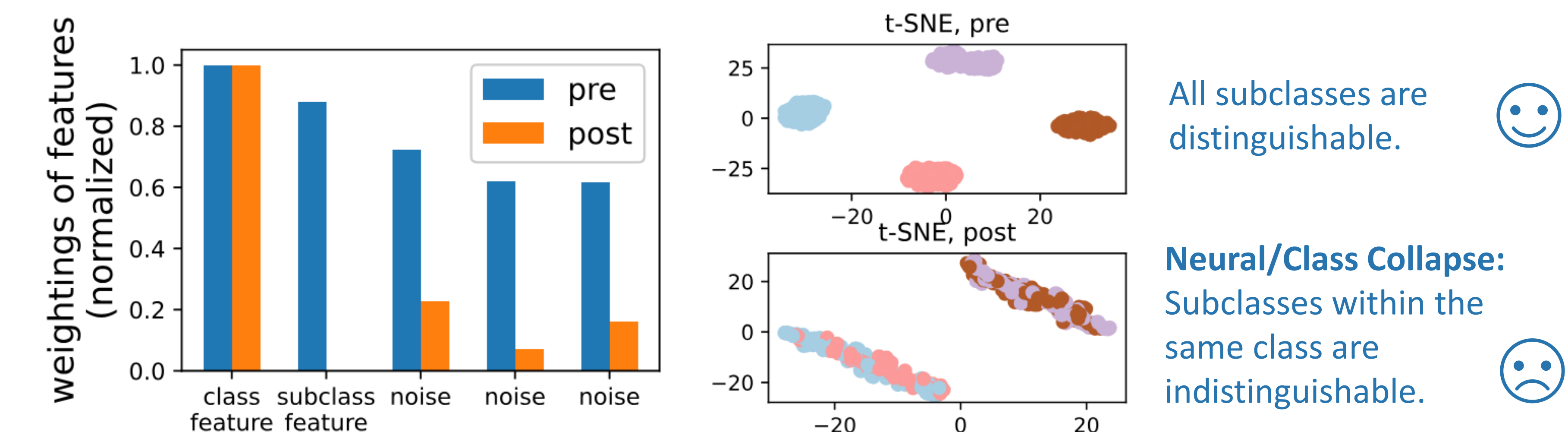
**Corollary (Informal):**
1. Data augmentation disrupts feature $i^*$ too much (large $\alpha_{i^*} \to$ small $\gamma_{i^*}$)
2. Feature $i^*$ is too weak in pretraining data (very small $\phi_{i^*} \to$ small $\gamma_{i^*}$)
3. Feature $i^*$ is too strong in pretraining data (very large $\phi_{i^*} \to$ small $\gamma_{i^*}$)

### What about non-linear models?

**Theorem (Informal):** Non-linear models allow pre-projection representations to capture features that are **entirely absent** post-projection

### Experiments, toy example



$\alpha_i$ (disruption by augmentation)



$\phi_i$ (strength of features)

## Supervised Contrastive Learning and Supervised Learning

Similar conclusions hold for both linear and non-linear models. Notably, the conclusions concerning non-linear models suggest that using pre-projection representations can mitigate **Neural/Class Collapse**, thereby enhancing *coarse-to-fine transferability*.
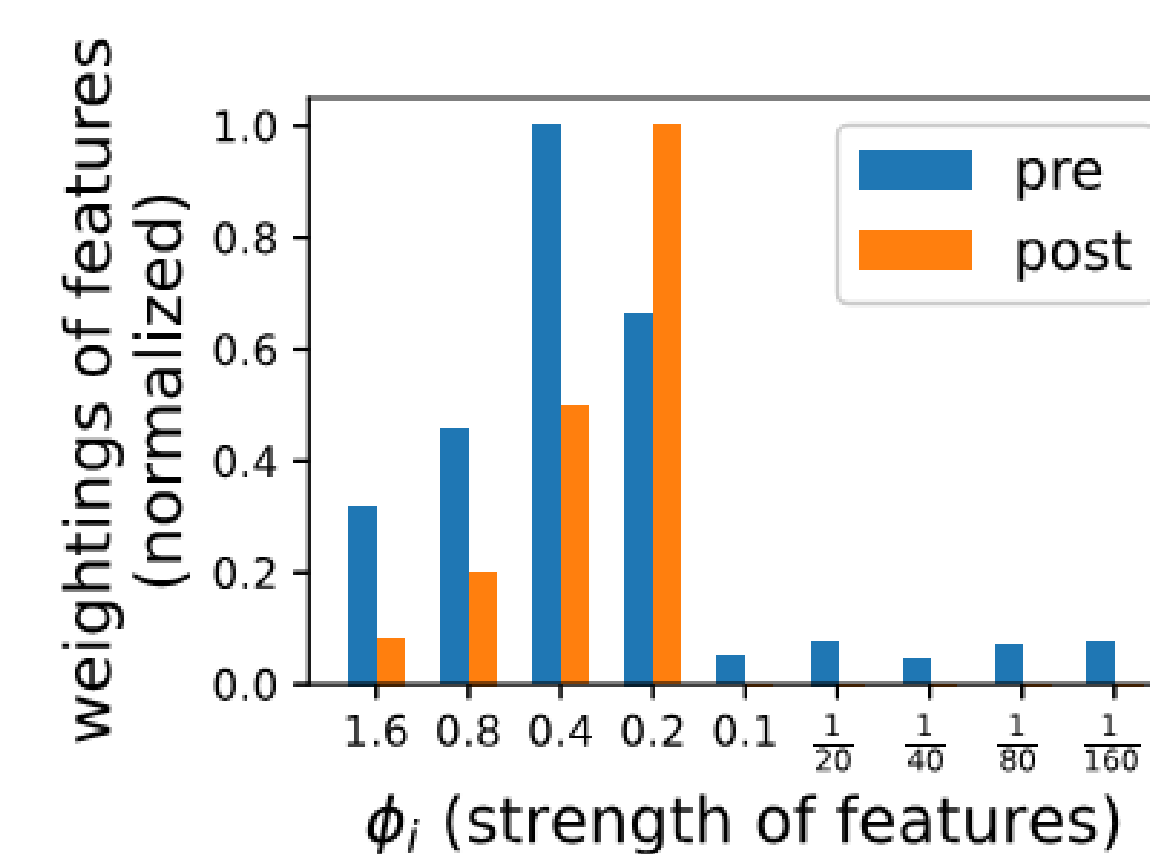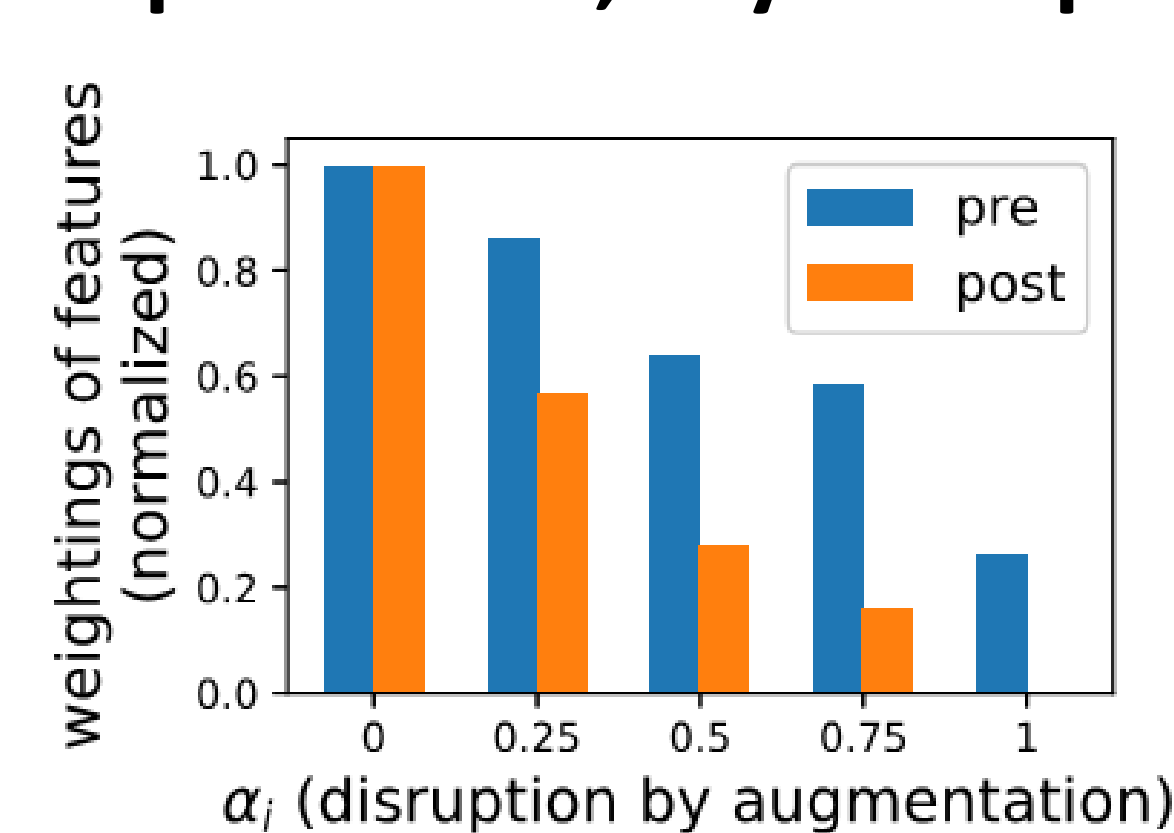


class feature, subclass feature, noise, noise, noise

t-SNE, pre

All subclasses are distinguishable.

t-SNE, post

**Neural/Class Collapse:** Subclasses within the same class are indistinguishable.

## Replacing the Projection Head with a Fixed Reweighting Head



$1/\kappa^0$
$1/\kappa^1$
$1/\kappa^2$

| Scenario | Dataset | Alg. | Performance Measure | Performance | | |
|---|---|---|---|---|---|---|
| | | | | vanilla | proj | reweight |
| synthetic | M-on-C | SSCL | digit clf. acc. | 77.0 | 97.3 | 97.3 |
| coarse-to-fine | CIFAR100 | SCL | fine-grained clf. acc. | 21.8 | 36.0 | 30.2 |
| coarse-to-fine | CIFAR100 | SL | fine-grained clf. acc. | 31.44 | 33.7 | 32.2 |
| distribution shift | UrbanCars | SL | few-shot adaption acc. | 82.2 | 86.1 | 87.0 |

The fixed reweighting head can achieve improvements that are comparable to those of the projection head across many tasks and algorithms.

### Experiments, semi-synthetic



Keep the digit for one
Aug.
Drop the digit for the other with prob. $p_{drop}$
$s = 0.4$  $s = 1$

(a) An illustration
(b) Effect of data aug.
(c) Effect of strength
(d) Effect of wd

probability of dropping digits

strength of mnist digits

weight decay

post-projection
pre-projection