
Probabilistic Variational Contrastive Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

Deterministic embeddings learned by contrastive learning (CL) methods such as SimCLR and SupCon achieve state-of-the-art performance but lack a principled mechanism for uncertainty quantification. We propose *Variational Contrastive Learning* (VCL), a decoder-free framework that maximizes the evidence lower bound (ELBO) by interpreting the InfoNCE loss as a surrogate reconstruction term and adding a KL divergence regularizer to a uniform prior on the unit hypersphere. We model the approximate posterior $q_\theta(\mathbf{z}|\mathbf{x})$ as a projected normal distribution, enabling the sampling of probabilistic embeddings. Our two instantiations—VSimCLR and VSupCon—replace deterministic embeddings with samples from $q_\theta(\mathbf{z}|\mathbf{x})$ and incorporate a normalized KL term into the loss. Experiments on multiple benchmarks demonstrate that VCL mitigates dimensional collapse, enhances mutual information with class labels, and matches or outperforms deterministic baselines in classification accuracy, all the while providing meaningful uncertainty estimates through the posterior model. VCL thus equips contrastive learning with a probabilistic foundation, serving as a new basis for contrastive approaches.

1 Introduction

Deep representation learning seeks to map each input \mathbf{x} into a compact embedding \mathbf{z} that preserves semantic similarity and facilitates downstream tasks such as classification or retrieval [5]. Contrastive learning methods, including SimCLR [8] and SupCon [30], have advanced the state of the art by pulling together positive pairs and pushing apart negatives in the embedding space. However, these approaches rely on deterministic point estimates for each sample, which do not express uncertainty or capture multiple plausible representations.

To address this limitation, we introduce a probabilistic *Variational Contrastive Learning* (VCL) approach, which extends deterministic embeddings to *probabilistic embeddings* by maximizing the evidence lower bound (ELBO) within the contrastive learning framework. Unlike variational autoencoders (VAEs) [31], which employ a decoder to reconstruct inputs from latent variables, VCL omits explicit decoders. Instead, we show that the InfoNCE loss can serve as a surrogate for the ELBO reconstruction term, yielding a principled probabilistic formulation of contrastive learning. Our VCL framework offers several new perspectives on learned embeddings.

Variational Contrastive Learning framework thus provides uncertainty-aware embeddings, a new basis of CL with theoretical insights via the ELBO. Our contributions are summarized as follows:

- We introduce *Variational Contrastive Learning* (VCL), a decoder-free ELBO maximization framework that reinterprets the InfoNCE loss as a surrogate reconstruction term and incorporates a KL divergence regularizer to a uniform prior on the unit hypersphere.
- We propose a probabilistic embedding model using a projected normal posterior that enables sampling, uncertainty quantification, and efficient KL computation on the hypersphere.

- We derive a connection between the optimal InfoNCE critic and the ELBO, showing that minimizing InfoNCE asymptotically maximizes the ELBO reconstruction term.
- We demonstrate that VCL mitigates both dimensional collapse in self-supervised contrastive learning via the KL regularizer, while preserving embedding structure. We show that VCL methods preserve or improve mutual information with labels, match or exceed classification accuracy of deterministic baselines, and provide meaningful implication of distributional embeddings.

2 Preliminaries

Let $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ be a dataset of input $\mathbf{x} \in \mathcal{X}$ and label pairs drawn i.i.d. from the joint distribution $p(\mathbf{x}, \mathbf{y})$. An *encoder* $f_\theta: \mathcal{X} \rightarrow \mathbb{R}^d$, parameterized by θ , maps each input \mathbf{x} to a d -dimensional vector, which we then normalize to unit length: $\mathbf{z} = \frac{f_\theta(\mathbf{x})}{\|f_\theta(\mathbf{x})\|_2}$. Throughout this section, we define the temperature-scaled cosine similarity between embeddings \mathbf{z}_i and \mathbf{z}_j as

$$s(\mathbf{z}_i, \mathbf{z}_j) = \frac{\mathbf{z}_i^\top \mathbf{z}_j}{\tau}, \quad (1)$$

where $\tau > 0$ is the temperature hyperparameter. For any two probability distributions q and p , we denote the Kullback–Leibler (KL) divergence by $D(q \| p) = \mathbb{E}_{\mathbf{z} \sim q} \left[\log \frac{q(\mathbf{z})}{p(\mathbf{z})} \right]$.

2.1 Self-Supervised Contrastive Learning

Self-supervised contrastive learning (SSCL) learns representations from *unlabeled* data by pulling together embeddings of semantically related views (positives) and pushing apart those of unrelated views (negatives). For an anchor \mathbf{x} , let \mathbf{x}'_i denote a positive view sampled from $p(\mathbf{x}'_i | \mathbf{x})$, and let $\{\mathbf{x}'_j\}_{j \neq i}$ be $N - 1$ negative views drawn i.i.d. from the marginal $p(\mathbf{x}')$. The InfoNCE loss [43] for anchor \mathbf{x} is then

$$I_{\text{NCE}}(\mathbf{x}; \mathbf{x}') = -\mathbb{E}_{\substack{\mathbf{x} \sim p(\mathbf{x}) \\ \mathbf{x}'_i \sim p(\mathbf{x}'_i | \mathbf{x}) \\ \{\mathbf{x}'_j\}_{j \neq i} \sim p(\mathbf{x}')}} \left[\log \frac{\exp(s(\mathbf{z}, \mathbf{z}'_i))}{\sum_{j=1}^N \exp(s(\mathbf{z}, \mathbf{z}'_j))} \right], \quad (2)$$

where $\mathbf{z} = f_\theta(\mathbf{x}) / \|f_\theta(\mathbf{x})\|_2$ and $s(\cdot, \cdot)$ is the temperature-scaled cosine similarity.

In practice, following SimCLR [8], we generate positives by applying two random augmentations $t', t'' \sim \mathcal{T}$ to each sample \mathbf{x}_i , yielding $(\mathbf{x}'_i, \mathbf{x}''_i) = (t'(\mathbf{x}_i), t''(\mathbf{x}_i))$.¹ All other $2N - 2$ augmented samples in the mini-batch serve as negatives. Let \mathcal{B} be the set of all $2N$ embeddings in the batch; then InfoNCE can be computed as

$$I_{\text{NCE}} = -\frac{1}{2N} \sum_{\mathbf{z} \in \mathcal{B}} \log \frac{\exp(s(\mathbf{z}, \mathbf{z}_p))}{\sum_{\mathbf{z}_n \in \mathcal{B} \setminus \{\mathbf{z}\}} \exp(s(\mathbf{z}, \mathbf{z}_n))}, \quad (3)$$

where \mathbf{z}_p denotes the positive embedding corresponding to \mathbf{z} . Since InfoNCE lower-bounds the mutual information $I(\mathbf{x}; \mathbf{x}')$ via $I(\mathbf{x}; \mathbf{x}') \geq \log N - I_{\text{NCE}}(\mathbf{x}; \mathbf{x}')$, we can see that minimizing I_{NCE} encourages encoders to preserve the semantic information of \mathbf{x} [48].

2.2 Variational Inference and the Evidence Lower Bound (ELBO)

In variational inference [6, 31], we treat the data distribution $p(\mathbf{x})$ as the marginal of a joint distribution over observed data \mathbf{x} and latent variables \mathbf{z} , i.e., $p(\mathbf{x}) = \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z}$. The latent variable \mathbf{z} captures meaningful structure in \mathbf{x} , serving both as a hidden cause and as a compressed representation for downstream tasks. In representation learning, we interpret \mathbf{z} as the embedding of \mathbf{x} .

The log-evidence can be written with respect to any approximate posterior $q_\phi(\mathbf{z} | \mathbf{x})$ as

$$\log p(\mathbf{x}) = \log \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} \left[\frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z} | \mathbf{x})} \right]. \quad (4)$$

¹Although we adopt the SimCLR augmentation scheme, our method applies to any contrastive framework.

72 Rather than optimizing (4) directly, variational methods maximize the *evidence lower bound* (ELBO)
 73 obtained as a result of applying Jensen’s inequality:

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] - D(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})) = \mathcal{L}^{\text{ELBO}}(\phi), \quad (5)$$

74 where $p(\mathbf{z})$ is a fixed prior (commonly $\mathcal{N}(\mathbf{0}, I_d)$). The ELBO decomposes into a *reconstruction* term
 75 $\mathbb{E}_q[\log p(\mathbf{x}|\mathbf{z})]$ and a *regularizer* $D(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))$. Maximizing $\mathcal{L}^{\text{ELBO}}$ thus balances (i) accurate
 76 reconstruction, (ii) posterior-to-prior regularization, and (iii) posterior accuracy. By

$$\log p(\mathbf{x}) = \mathcal{L}^{\text{ELBO}}(\phi) + D(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x})), \quad (6)$$

77 for fixed $\log p(\mathbf{x})$, maximizing the ELBO minimizes the KL divergence between the approximate
 78 and true posteriors [6].

79 The ELBO provides a tractable surrogate for marginal likelihood that can be optimized by standard
 80 gradient methods. It will serve as the theoretical backbone of our Variational Contrastive Learning
 81 framework, offering both a probabilistic interpretation and explicit control over latent uncertainty.

82 **Relation to contrastive objectives.** Although the ELBO stems from latent-variable modeling,
 83 its two components align naturally with contrastive objectives: the KL divergence term enforces
 84 *uniformity* in the embedding space, while the reconstruction term promotes *alignment* between
 85 embeddings and observations. In Section 3, we leverage this connection by adopting distributional
 86 embeddings in the contrastive framework and incorporating a KL-based regularizer on the posterior.

87 3 Variational Contrastive Learning (VCL)

88 Unlike existing variational contrastive learning methods—which primarily focus on generative models
 89 with explicit decoders [7, 59]—our approach performs *decoder-free* ELBO maximization, making
 90 VCL a truly contrastive learning framework.

91 3.1 Decoder-Free ELBO Maximization

92 Here we describe how to optimize two terms in ELBO (5) within a purely contrastive learning setup.

93 **Reconstruction term.** The reconstruction term $\mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})]$ requires the true conditional
 94 $p(\mathbf{x}|\mathbf{z})$, which is generally intractable. Instead, we approximate it via the embedding conditional

$$p(\mathbf{z}'|\mathbf{z}) = \frac{p(\mathbf{z}, \mathbf{z}')}{\int p(\mathbf{z}, \mathbf{z}') d\mathbf{z}'}, \quad (7)$$

95 where $\mathbf{z}' \sim q_\theta(\cdot | \mathbf{x})$ captures semantics of \mathbf{x} . Thus,

$$\begin{aligned} \mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] &\approx \mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})q_\theta(\mathbf{z}'|\mathbf{x})}[\log p(\mathbf{z}'|\mathbf{z})] \\ &= \mathbb{E}\left[\log \frac{p(\mathbf{z}, \mathbf{z}')}{\int p(\mathbf{z}, \mathbf{z}') d\mathbf{z}'}\right] \approx \mathbb{E}\left[\log \frac{e^{\psi(\mathbf{z}, \mathbf{z}')}}{\sum_j e^{\psi(\mathbf{z}, \mathbf{z}'_j)}}\right], \end{aligned} \quad (8)$$

96 where we approximate $p(\mathbf{z}, \mathbf{z}') \approx e^{\psi(\mathbf{z}, \mathbf{z}')}$ via a critic ψ . Details on parameterizing $p(\mathbf{z}' | \mathbf{z})$
 97 appear in Section 3.2. The following lemma supports the approximation $\mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] \approx$
 98 $\mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})q_\theta(\mathbf{z}'|\mathbf{x})}[\log p(\mathbf{z}'|\mathbf{z})]$. A further discussion on the approximation in (8) and a tightness
 99 condition is in Appendix D.1.

100 **Lemma 3.1.** *Let \mathbf{x} and \mathbf{z} be conditionally independent given \mathbf{z}' . Then, the reconstruction term in*
 101 *Section 3.1 is bounded as*

$$\mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})q(\mathbf{z}'|\mathbf{x})}[\log p(\mathbf{z}'|\mathbf{z})] + \text{const.}, \quad (9)$$

102 where const. is independent of \mathbf{z} .

103 *Proof.* The proof of Proposition 3.1 is in Appendix B.1. □

Noting that the right-hand side of (8) is (up to sign) the InfoNCE surrogate, setting $\psi(\cdot, \cdot) = s(\cdot, \cdot)$ in (8) where $s(\cdot, \cdot)$ is defined in (1) yields

$$\mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] \approx -I_{\text{NCE}}(\mathbf{x}; \mathbf{x}'). \quad (10)$$

Hence, minimizing the InfoNCE loss maximizes the reconstruction term without explicit decoders.

In contrast to VAE embeddings—which often rely on pixel-level reconstruction through expressive decoder [53]—VCL preserves semantics via contrastive objectives. The next proposition (proved in Appendix B.2) provides a theoretical connection between InfoNCE and the reconstruction term.

Proposition 3.2. *Assume that: 1) the critic ψ in InfoNCE is optimal; 2) $p(\mathbf{z}) < \infty$, $\forall \mathbf{z}$; and 3) $0 < \epsilon \leq p(\mathbf{z}|\mathbf{z}') \leq g_+(\mathbf{z})$, $\forall \mathbf{z}, \mathbf{z}'$ with a absolutely integrable $g : \mathcal{Z} \rightarrow (0, \infty)$. Then, as the number of negatives $N \rightarrow \infty$,*

$$-I_{\text{NCE}}(\mathbf{x}; \mathbf{x}') + \log N \rightarrow \mathbb{E} [\log p(\mathbf{z}'|\mathbf{x})] - D(q_\theta(\mathbf{z}'|\mathbf{x}) \| p(\mathbf{z}')) - H(q_\theta(\mathbf{z}'|\mathbf{x})), \quad (11)$$

where the expectation is over $q_\theta(\mathbf{z}|\mathbf{x})q_\theta(\mathbf{z}'|\mathbf{x})$, and $H(\cdot)$ denotes entropy.

Regularization. Maximizing the ELBO requires choosing a prior $p(\mathbf{z})$ and an approximate posterior $q_\theta(\mathbf{z} | \mathbf{x})$. Although both are often taken as Gaussian distributions [31], this choice conflicts with the geometry of contrastive embeddings, which often lie on the unit hypersphere due to the normalization [62]. Instead, we adopt non-Gaussian priors and posteriors—one key distinction from standard VAE approaches.

Motivated by the uniformity property [62] on the unit sphere $\mathcal{S}^{d-1} = \{\mathbf{z} \in \mathbb{R}^d : \|\mathbf{z}\|_2 = 1\}$, we set the prior $p(\mathbf{z})$ to be the uniform distribution over \mathcal{S}^{d-1} . For the approximate posterior, we use the *projected normal* distribution [22], which admits efficient KL-divergence computation while enforcing $\mathbf{z} \in \mathcal{S}^{d-1}$. A random variable $\mathbf{z} \sim \mathcal{PN}(\mu, K)$ is obtained by sampling

$$\mathbf{z} = \frac{\mathbf{u}}{\|\mathbf{u}\|_2} \quad \text{with} \quad \mathbf{u} \sim \mathcal{N}(\mu, K). \quad (12)$$

In particular, $\mathcal{PN}(0, I_d)$ reduces to the uniform distribution on \mathcal{S}^{d-1} , i.e., $\mathcal{PN}(0, I_d) \stackrel{d}{=} \text{Unif}(\mathcal{S}^{d-1})$.

With $q_\theta(\mathbf{z}|\mathbf{x}) = \mathcal{PN}(\mu, K)$, the regularization term becomes

$$D(q_\theta(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z})) = D(\mathcal{PN}(\mu, K) \| \text{Unif}(\mathcal{S}^{d-1})). \quad (13)$$

Since a closed-form KL divergence between projected normals and the uniform sphere is intractable, we instead minimize the Gaussian KL as an upper bound—by the data processing inequality [47]:

$$D(\mathcal{N}(\mu, K) \| \mathcal{N}(0, I_d)) \geq D(\mathcal{PN}(\mu, K) \| \text{Unif}(\mathcal{S}^{d-1})). \quad (14)$$

In Appendix D.2, we analyze the tightness of the gap in (14) and show that the Gaussian KL divergence closely approximates the projected-normal KL divergence; the two exhibit nearly identical behavior throughout VCL training.

For $K = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$, the Gaussian KL admits the closed form

$$D(\mu, K) = \frac{1}{2} \sum_{i=1}^d (\sigma_i^2 + \mu_i^2 - 1 - \log \sigma_i^2). \quad (15)$$

The KL divergence term $D(\mu, K)$ grows linearly with the embedding dimension d , which can destabilize training when d is large. To address this, we normalize the KL term by d , i.e., $\tilde{D}(\mu, K) = \frac{1}{d} D(\mu, K)$, so that its magnitude remains comparable to the InfoNCE loss.

Final objective for maximizing ELBO. By combining (10) and (15), we obtain the following (approximate) lower bound on the ELBO:

$$\mathcal{L}^{\text{ELBO}}(\theta) \geq -I_{\text{NCE}}(\mathbf{x}; \mathbf{x}') - D(\mu_{\mathbf{x}}, K_{\mathbf{x}}), \quad (16)$$

where $\mu_{\mathbf{x}}$ and $K_{\mathbf{x}} = \text{diag}(\sigma_{\mathbf{x},1}, \dots, \sigma_{\mathbf{x},d})$ are the parameters of $q_\theta(\mathbf{z} | \mathbf{x})$. Because this bound is asymmetric in $(\mathbf{x}, \mathbf{x}')$, we symmetrize it to define our final VCL objective:

$$\mathcal{L}^{\text{VCL}} = \frac{1}{2} \left(I_{\text{NCE}}(\mathbf{x}; \mathbf{x}') + I_{\text{NCE}}(\mathbf{x}'; \mathbf{x}) + D(\mu_{\mathbf{x}}, K_{\mathbf{x}}) + D(\mu_{\mathbf{x}'}, K_{\mathbf{x}'}) \right). \quad (17)$$

Minimizing \mathcal{L}^{VCL} therefore maximizes the ELBO. Next, we introduce Variational SimCLR (VSimCLR), which is specifically designed to optimize this objective efficiently.

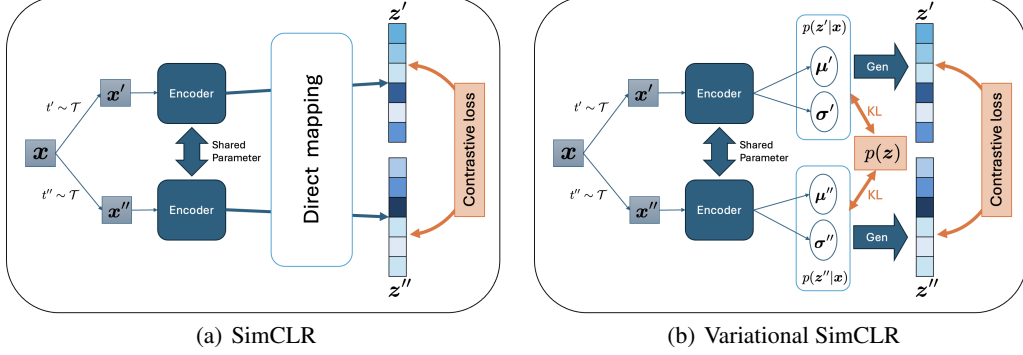


Figure 1: Graphical illustration of SimCLR and Variational SimCLR (VSimCLR).

3.2 Variational SimCLR (VSimCLR)

We propose *Variational SimCLR* (VSimCLR), whose architecture is illustrated in Figure 1(b). VSimCLR minimizes \mathcal{L}^{VCL} in (17), thereby implicitly maximizing the ELBO and bringing the approximate posterior closer to the true posterior by (6). Compared to SimCLR, VSimCLR differs in three key aspects: (i) the encoder outputs the parameters of a variational posterior rather than deterministic embeddings; (ii) embeddings are sampled from this posterior; and (iii) a KL divergence term between the approximate posterior and the prior is included in the loss.

Specifically, during training, each input x is first augmented twice to obtain x' and x'' , as in SimCLR. The encoder then maps x' and x'' to posterior parameters (μ', σ') and (μ'', σ'') , respectively. We then sample

$$z' = \mu' + \text{diag}(\sigma') \epsilon_1, \quad \text{and} \quad z'' = \mu'' + \text{diag}(\sigma'') \epsilon_2, \quad (18)$$

where $\epsilon_1, \epsilon_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, I_d)$. After normalizing z' and z'' to unit length, we compute the InfoNCE loss over the normalized embeddings in the batch and add the KL divergence

$$\frac{1}{d} D(\mathcal{N}(\mu, \text{diag}(\sigma^2)) \parallel \mathcal{N}(\mathbf{0}, I_d)) \quad (19)$$

for each sample. Minimizing this combined objective effectively minimizes \mathcal{L}^{VCL} in (17) and thus maximizes the ELBO. Figure 1 highlights these differences: VSimCLR replaces deterministic embeddings with the projected-normal posterior $\mathcal{PN}(\mu, \text{diag}(\sigma^2))$ and regularizes it via KL divergence to the standard normal.²

4 Experiments

We evaluate VCL with SimCLR and SupCon across five aspects: (i) embedding visualization, (ii) dimensional collapse, (iii) mutual information between embeddings and labels, (iv) classification accuracy, and (v) implications of distributional embeddings. Implementation and training details are provided in Appendix E.1.

4.1 Embedding Visualization

Figure 2 presents t-SNE [57] and UMAP [40] projections of the embeddings learned by SimCLR and VSimCLR on the CIFAR-10 test set. Although VSimCLR incorporates an additional KL-regularizer, it preserves the characteristic cluster structure induced by contrastive learning. This confirms that our distributional embeddings retain the semantic information learned by contrastive methods.

4.2 Dimensional Collapse

²An analogy with SupCon, namely VSupCon, is provided in Appendix C.

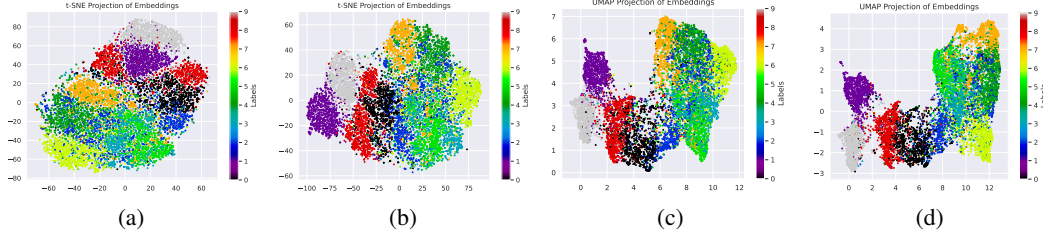


Figure 2: Embedding visualization for SimCLR and VSimCLR on CIFAR-10 test set. (a) t-SNE of SimCLR; (b) t-SNE of VSimCLR; (c) UMAP of SimCLR; (d) UMAP of VSimCLR. VSimCLR preserves the characteristic cluster structure of contrastive learning while introducing probabilistic embeddings regularized by (15).

Table 1: Classification accuracy on various datasets. We report top-1 and top-5 accuracies of SimCLR, VSimCLR, SupCon, and VSupCon across the datasets.

METHOD	CIFAR-10		CIFAR-100		TINY-IMAGENET		STL10		CALTECH256	
	TOP1	TOP5	TOP1	TOP5	TOP1	TOP5	TOP1	TOP5	TOP1	TOP5
SIMCLR	78.42	98.52	49.56	78.84	38.95	66.89	60.44	95.80	43.14	66.15
VSIMCLR	81.48	98.95	54.58	82.87	37.70	66.06	60.11	92.00	48.50	69.99
SUPCON	93.60	99.71	70.79	89.11	57.60	77.16	75.88	98.51	87.06	91.64
VSUPCON	93.85	99.68	71.66	89.42	48.30	72.84	75.76	96.99	83.06	91.29

Contrastive learning methods such as SimCLR often suffer from *dimensional collapse*, where embeddings concentrate in a low-dimensional subspace, underutilizing the full capacity of the representation space [29]. To quantify this effect, let $\{z_i\}_{i=1}^N$ be the test-set embeddings and their covariance matrix $C = \frac{1}{N} \sum_{i=1}^N (z_i - \bar{z})(z_i - \bar{z})^\top$, with $\bar{z} = \frac{1}{N} \sum_{i=1}^N z_i$. Figure 3 shows the singular values of C for SimCLR and VSimCLR. VSimCLR produces a substantially flatter spectrum, indicating a higher effective rank and thus mitigating dimensional collapse. Remarkably, on CIFAR-100, VSimCLR nearly doubles the number of dominant components compared to SimCLR. These results demonstrate that VSimCLR not only preserves semantic clustering but also leverages the embedding space more fully, and can be combined with other collapse-mitigation strategies for further gains. Additional experiments on Caltech-256 and Tiny-ImageNet (Figure 8, Appendix E.2) exhibit similar behavior.

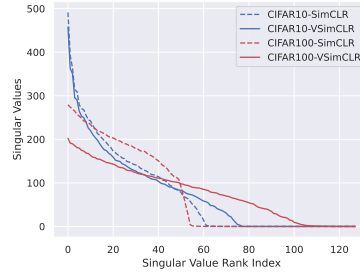


Figure 3: Singular-value spectrum.

4.3 Mutual Information Comparison

Figure 4 reports the estimated mutual information $I(z; c)$ between the learned embeddings z and their true class labels c of CIFAR-10. We compute this using the Mixed KSG estimator [13], which is well-suited for mixed or multimodal distributions.

Both VSimCLR and VSupCon achieve mutual information on par with—or slightly exceeding—their non-variational counterparts. In particular, during the first 200 epochs, VSimCLR exhibits lower mutual information than SimCLR, reflecting the added optimization challenge of the KL regularizer. After this initial phase, VSimCLR surpasses SimCLR and maintains higher mutual information for the remainder of training. These results indicate that VSimCLR ultimately preserves—or even improves—information between embeddings and labels, while also producing rich distributional representations.

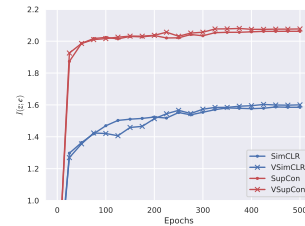


Figure 4: Estimate of $I(z; c)$.

4.4 Classification

For classification, we use the posterior mean μ_x as the embedding and train a linear classifier.

Table 1 reports Top-1 and Top-5 accuracies on CIFAR-10, CIFAR-100, Tiny-ImageNet, STL-10, and Caltech-256. VSimCLR outperforms SimCLR on CIFAR-10 (78.42 \rightarrow 81.48) and CIFAR-100 (49.56 \rightarrow 54.58) in Top-1 accuracy, with similar gains in Top-5. On Caltech-256, VSimCLR also improves Top-1 accuracy substantially. Performance on Tiny-ImageNet and STL-10 remains comparable, with slight decreases (within experimental variance) likely due to the KL regularizer.

SupCon provides supervised baselines, and VSupCon further improves Top-1 accuracy on CIFAR-10 (93.60 \rightarrow 93.85) and CIFAR-100 (70.79 \rightarrow 71.66). Modest declines on Tiny-ImageNet, STL-10, and Caltech-256 reflect the trade-off of adding the KL term on datasets with higher complexity or fewer samples.

Although VCL is not explicitly designed to boost classification accuracy, VSimCLR consistently match or exceed their deterministic counterparts. This demonstrates that distributional embeddings preserve the alignment and uniformity properties [62], while providing meaningful uncertainty proxy.

4.5 Implications of Distributional Embeddings

We illustrate the interpretability of distributional embeddings using examples from CIFAR-10. Figure 9 displays sample images alongside the log-determinant $\log \det(K)$ of their posterior covariance K learned by VSimCLR. Top-row images are common class members and exhibit larger $\log \det(K)$ —indicating broader posterior dispersion—whereas bottom-row images are atypical or uncommon with smaller $\log \det(K)$, reflecting more concentrated posteriors.³

We quantify the relationship between posterior covariance and uncertainty using CIFAR-10H [46] and CIFAR-10C [19]. Figure 5 plots $\log \det(K)$ against the entropy of the CIFAR-10H soft labels [24, 25]; the negative slope of the linear fit (red line) indicates that images with lower $\log \det(K)$ —i.e., more concentrated posteriors—tend to have higher label entropy and thus greater ambiguity. Next, using CIFAR-10C, we examine how posterior covariance varies with corruption severity, which correlates with label uncertainty. Figures 6 and 11 show that $\log \det(K)$ decreases as corruption strength increases, implying that lower posterior dispersion corresponds to higher uncertainty, consistent with Figure 5.

These results demonstrate that the dispersion of the learned posterior correlates with semantic uncertainty, highlighting the practical interpretability of VCL’s distributional embeddings. As an example application of posterior covariance, we consider CIFAR-100 under a label-scarce setting in which only a small number of labels per class are available to train a linear classifier. Table 2 reports accuracies for SimCLR, VSimCLR, and VSimCLR+wt, with classifiers trained using cross-entropy (CE). Here, “+wt” denotes a weighted CE in which sample weights are proportional to posterior covariance to downweight ambiguous examples. Specifically, we use

$$\mathcal{L}_{wCE} = \sum_{i=1}^N w_i \log \phi_{c_i}(z_i), \text{ with } w_i \propto \log \det(K) \text{ (after normalization)}, \quad (20)$$

where $\phi_{c_i}(z_i)$ is the estimated probability of the true class. Table 2 shows that VCL variants improve over SimCLR and SupCon, with smaller gains for SupCon since it already leverages labels during pretraining. Moreover, weighting by posterior covariance further improves performance, supporting

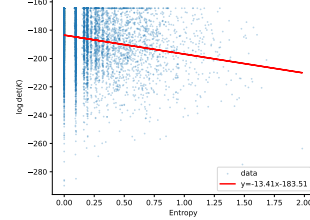


Figure 5: Posterior dispersion versus label ambiguity. Each point plots $\log \det(K)$ against the entropy of human-annotated class probabilities from CIFAR-10H, with a first-order linear fit (red line).

³ $\log \det K$ quantifies the *dispersion of the posterior in embedding space*, which reflects *typicality* rather than label uncertainty. Larger values correspond to more “typical” samples with many latent realizations consistent with the data manifold, whereas smaller values indicate more “unique” or outlier samples with tightly concentrated posteriors. A generative analogy may help understanding: if an outlier image had an extremely large posterior variance, then samples drawn from the prior would reproduce that outlier far too often—contradicting its rarity. Hence, larger variance corresponds to “typical” not “uncertain” inputs.

Table 2: Classification accuracy on CIFAR-100 with label scarcity. We use ResNet-18 back-bone and same augmentations for all experiments. We sample the labelled subset once and report the mean accuracy of five runs with (standard error).

METHODS	1 LABELS / CLASS	3 LABELS / CLASS	5 LABELS / CLASS	10 LABELS / CLASS	20 LABELS / CLASS
SIMCLR	12.22 (0.12)	21.37 (0.15)	26.37 (0.01)	33.09 (0.11)	38.00 (0.06)
VSimCLR	15.57 (0.09)	25.70 (0.19)	30.89 (0.11)	37.40 (0.08)	42.13 (0.10)
VSimCLR+WT	15.97 (0.08)	26.07 (0.20)	31.12 (0.06)	37.48 (0.08)	42.36 (0.03)
SUPCON	71.55 (0.04)	71.56 (0.05)	71.64 (0.02)	71.65 (0.03)	72.07 (0.05)
VSupCON	71.77 (0.12)	71.79 (0.10)	71.96 (0.09)	72.07 (0.05)	72.16 (0.04)
VSupCON+WT	71.87 (0.02)	71.78 (0.07)	71.94 (0.07)	72.07 (0.07)	72.16 (0.06)

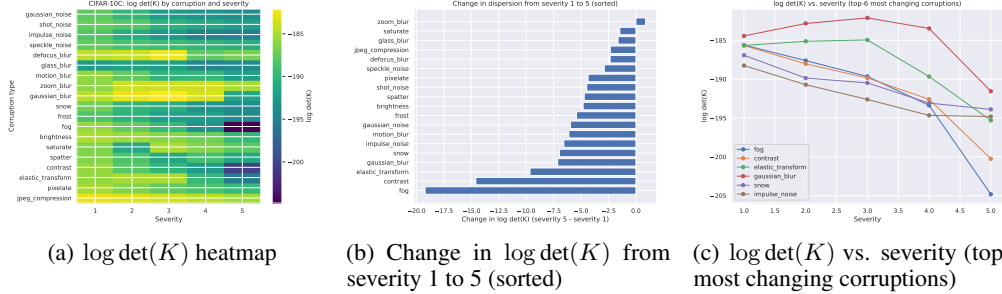


Figure 6: $\log \det(K)$ of VSimCLR embeddings on CIFAR-10C under different corruption types and severities. “Severity” denotes the corruption level. The observed negative correlation between $\log \det(K)$ and severity is consistent with our finding that more uncertain samples exhibit smaller posterior covariance dispersion. Exact $\log \det(K)$ values are in Table 7.

distributional embeddings as a confidence proxy. Additional experiments and discussion on posterior distributions and label uncertainty are provided in Appendix E.5.

This counterintuitive finding—that typical (i.e., common) samples exhibit larger posterior dispersion—parallels the observation in concurrent work by Guth et al. [16], albeit under different settings: (i) Quantity: we analyze latent-space posterior covariance via $\log \det K$, whereas they study input-space marginal density $p(x)$; (ii) Observation: typical samples have larger $\log \det K$, while they have lower marginal density. Although the quantities are measured in different spaces, both results indicate that typical samples are not the highest-density points. In our case, typical images yield larger posterior dispersion and atypical images smaller dispersion; since dispersion is inversely related to peak density, our result aligns with Guth et al.’s observation. Hence, in both settings, “typical” \neq “highest-density.”

5 Conclusion

We have introduced *Variational Contrastive Learning* (VCL), a decoder-free ELBO-maximization framework that endows contrastive learning with principled probabilistic embeddings. By interpreting InfoNCE as a surrogate reconstruction term and regularizing with a KL divergence to a uniform prior on the unit sphere, VCL enables distributional encodings without explicit decoders. We instantiated VCL in two variants—VSimCLR and VSupCon—by replacing deterministic embeddings with samples from $q_\theta(z | x)$ and adding a normalized KL term.

Theoretical and empirical results show that VCL preserves the properties of contrastive embeddings, mitigates dimensional collapse, maintains or improves mutual information with labels, and matches or exceeds deterministic baselines in classification accuracy, while also providing meaningful posterior uncertainty estimates. We further analyzed the implications of probabilistic embeddings—spanning label uncertainty, typicality, and OOD behavior—through posterior-covariance dispersion. We also observed a counterintuitive but consistent pattern, echoed in concurrent diffusion-model work [16]: lower posterior-covariance dispersion is associated with higher sample uniqueness (i.e., more atypical or outlier examples), whereas typical samples exhibit larger posterior covariance dispersion.

References

- [1] L. Aitchison and S. Ganey. Infonce is variational inference in a recognition parameterised model. *arXiv preprint arXiv:2107.02495*, 2021.
- [2] C. A. Barbano, B. Dufumier, E. Tartaglione, M. Grangetto, and P. Gori. Unbiased supervised contrastive learning. *arXiv preprint arXiv:2211.05568*, 2022.
- [3] A. Bardes, J. Ponce, and Y. LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- [4] P. L. Bartlett, D. J. Foster, and M. J. Telgarsky. Spectrally-normalized margin bounds for neural networks. *Advances in neural information processing systems*, 30, 2017.
- [5] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- [6] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- [7] L. Chen, P. Wang, X. Han, and L. Xu. Multi-relational variational contrastive learning for next poi recommendation. In *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5, 2025.
- [8] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PmLR, 2020.
- [9] S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with application to face verification. In *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR’05)*, volume 1, pages 539–546. IEEE, 2005.
- [10] C.-Y. Chuang, J. Robinson, Y.-C. Lin, A. Torralba, and S. Jegelka. Debiased contrastive learning. *Advances in neural information processing systems*, 33:8765–8775, 2020.
- [11] A. Coates, A. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223. JMLR Workshop and Conference Proceedings, 2011.
- [12] A. Ermolov, A. Siarohin, E. Sangineto, and N. Sebe. Whitening for self-supervised representation learning. In *International conference on machine learning*, pages 3015–3024. PMLR, 2021.
- [13] W. Gao, S. Kannan, S. Oh, and P. Viswanath. Estimating mutual information for discrete-continuous mixtures. *Advances in neural information processing systems*, 30, 2017.
- [14] R. Girdhar, A. El-Nouby, Z. Liu, M. Singh, K. V. Alwala, A. Joulin, and I. Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15180–15190, 2023.
- [15] G. Griffin, A. Holub, P. Perona, et al. Caltech-256 object category dataset. Technical report, Technical Report 7694, California Institute of Technology Pasadena, 2007.
- [16] F. Guth, Z. Kadkhodaie, and E. P. Simoncelli. Learning normalized image densities via dual score matching. *arXiv preprint arXiv:2506.05310*, 2025.
- [17] J. He, J. Du, and W. Ma. Preventing dimensional collapse in self-supervised learning via orthogonality regularization. *arXiv preprint arXiv:2411.00392*, 2024.
- [18] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [19] D. Hendrycks and T. Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019.

- [20] D. Hendrycks, N. Mu, E. D. Cubuk, B. Zoph, J. Gilmer, and B. Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. In *Proc. International Conference on Learning Representations (ICLR)*, 2020.
- [21] A. Hermans, L. Beyer, and B. Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017.
- [22] D. Hernandez-Stumpfhauser, F. J. Breidt, and M. J. van der Woerd. The general projected normal distribution of arbitrary dimension: Modeling and bayesian inference. 2017.
- [23] N. M. Hieu, A. Ledent, Y. Lei, and C. Y. Ku. Generalization analysis for deep contrastive representation learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 17186–17194, 2025.
- [24] T. Ishida, I. Yamane, N. Charoenphakdee, G. Niu, and M. Sugiyama. Is the performance of my deep network too good to be true? a direct approach to estimating the bayes error in binary classification. In *The Eleventh International Conference on Learning Representations*, 2023.
- [25] M. Jeong, M. Cardone, and A. Dytso. Demystifying the optimal performance of multi-class classification. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 31638–31664. Curran Associates, Inc., 2023.
- [26] M. Jeong, M. Cardone, and A. Dytso. Data-driven estimation of the false positive rate of the bayes binary classifier via soft labels. In *2024 IEEE International Symposium on Information Theory (ISIT)*, pages 368–373, 2024.
- [27] M. Jeong and A. Hero. Generalizing supervised contrastive learning: A projection perspective. *arXiv preprint arXiv:2506.09810*, 2025.
- [28] M. Jeong, M. Namgung, Z. M. Kim, D. Kang, Y.-Y. Chiang, and A. Hero. Anchors aweigh! sail for optimal unified multi-modal representations. *arXiv preprint arXiv:2410.02086*, 2024.
- [29] L. Jing, P. Vincent, Y. LeCun, and Y. Tian. Understanding dimensional collapse in contrastive self-supervised learning. *arXiv preprint arXiv:2110.09348*, 2021.
- [30] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673, 2020.
- [31] D. P. Kingma, M. Welling, et al. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019.
- [32] M. Kirchhof, E. Kasneci, and S. J. Oh. Probabilistic contrastive learning recovers the correct aleatoric uncertainty of ambiguous inputs. In *International Conference on Machine Learning*, pages 17085–17104. PMLR, 2023.
- [33] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [34] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, volume 25, pages 1097–1105, 2012.
- [35] Y. Le and X. Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- [36] Y. Lei, T. Yang, Y. Ying, and D.-X. Zhou. Generalization analysis for contrastive representation learning. In *International Conference on Machine Learning*, pages 19200–19227. PMLR, 2023.
- [37] R. Linsker. Self-organization in a perceptual network. *Computer*, 21(3):105–117, 1988.
- [38] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.

- [39] Z. Ma and M. Collins. Noise contrastive estimation and negative sampling for conditional models: Consistency and statistical efficiency. In E. Riloff, D. Chiang, J. Hockenmaier, and J. Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3698–3707, Brussels, Belgium, Oct.-Nov. 2018. Association for Computational Linguistics.
- [40] L. McInnes, J. Healy, and J. Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- [41] M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of Machine Learning*. Adaptive Computation and Machine Learning. MIT Press, Cambridge, MA, 2012.
- [42] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, A. Y. Ng, et al. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, page 4. Granada, 2011.
- [43] A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [44] J. Park, J. Lee, I.-J. Kim, and K. Sohn. Probabilistic representations for video contrastive learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14711–14721, 2022.
- [45] L. Perez and J. Wang. The effectiveness of data augmentation in image classification using deep learning. *arXiv preprint arXiv:1712.04621*, 2017.
- [46] J. C. Peterson, R. M. Battleday, T. L. Griffiths, and O. Russakovsky. Human uncertainty makes classification more robust. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9617–9626, 2019.
- [47] Y. Polyanskiy and Y. Wu. *Information theory: From coding to learning*. Cambridge university press, 2025.
- [48] B. Poole, S. Ozair, A. Van Den Oord, A. Alemi, and G. Tucker. On variational bounds of mutual information. In *International conference on machine learning*, pages 5171–5180. PMLR, 2019.
- [49] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [50] B. Rhodes and M. U. Gutmann. Variational noise-contrastive estimation. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2741–2750. PMLR, 2019.
- [51] N. Saunshi, O. Plevrakis, S. Arora, M. Khodak, and H. Khandeparkar. A theoretical analysis of contrastive unsupervised representation learning. In *International Conference on Machine Learning*, pages 5628–5637. PMLR, 2019.
- [52] C. Shorten and T. M. Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6(1):60, 2019.
- [53] T. Song, J. Sun, X. Liu, and W. Peng. Scale-vae: Preventing posterior collapse in variational autoencoder. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 14347–14357, 2024.
- [54] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, and P. Isola. What makes for good views for contrastive learning? *Advances in neural information processing systems*, 33:6827–6839, 2020.
- [55] M. Tschannen, J. Djolonga, P. K. Rubenstein, S. Gelly, and M. Lucic. On mutual information maximization for representation learning. *arXiv preprint arXiv:1907.13625*, 2019.
- [56] M. Tschannen, J. Djolonga, P. K. Rubenstein, S. Gelly, and M. Lucic. On mutual information maximization for representation learning. In *International Conference on Learning Representations*, 2020.

- 405 [57] L. Van der Maaten and G. Hinton. Visualizing data using t-sne. *Journal of machine learning*
406 *research*, 9(11), 2008.
- 407 [58] W. I. Walker, H. Soulat, C. Yu, and M. Sahani. Unsupervised representation learning with
408 recognition-parametrised probabilistic models. In *International Conference on Artificial Intelli-*
409 *gence and Statistics*, pages 4209–4230. PMLR, 2023.
- 410 [59] B. Wang, Z. Tian, A. Ye, F. Wen, S. Du, and Y. Gao. Generative variational-contrastive learning
411 for self-supervised point cloud representation. *IEEE Transactions on Pattern Analysis and*
412 *Machine Intelligence*, 46(9):6154–6166, 2024.
- 413 [60] B. Wang, Z. Tian, A. Ye, F. Wen, S. Du, and Y. Gao. Generative variational-contrastive learning
414 for self-supervised point cloud representation. *IEEE Transactions on Pattern Analysis and*
415 *Machine Intelligence*, 2024.
- 416 [61] Q. Wang, S. R. Kulkarni, and S. Verdu. Divergence estimation for multidimensional densities
417 via k -nearest-neighbor distances. *IEEE Transactions on Information Theory*, 55(5):2392–2405,
418 2009.
- 419 [62] T. Wang and P. Isola. Understanding contrastive representation learning through alignment
420 and uniformity on the hypersphere. In *International conference on machine learning*, pages
421 9929–9939. PMLR, 2020.
- 422 [63] X. Wang, R. Zhang, C. Shen, T. Kong, and L. Li. Dense contrastive learning for self-supervised
423 visual pre-training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
424 *recognition*, pages 3024–3033, 2021.
- 425 [64] D. Ward, M. Beaumont, and M. Fasiolo. SoftCVI: Contrastive variational inference with self-
426 generated soft labels. In *The Thirteenth International Conference on Learning Representations*,
427 2025.
- 428 [65] S. Xie and J. H. Giraldo. Variational graph contrastive learning. *arXiv preprint*
429 *arXiv:2411.07150*, 2024.
- 430 [66] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow twins: Self-supervised learning via
431 redundancy reduction. In *International conference on machine learning*, pages 12310–12320.
432 PMLR, 2021.
- 433 [67] R. S. Zimmermann, Y. Sharma, S. Schneider, M. Bethge, and W. Brendel. Contrastive learning
434 inverts the data generating process. In *International conference on machine learning*, pages
435 12979–12990. PMLR, 2021.

A Related work

A.1 Contrastive learning

Self-supervised contrastive learning methods [8, 54] train an encoder $f: \mathcal{X} \rightarrow \mathcal{S}^{d_z-1}$ by drawing semantically related views (positives) together in the embedding space while pushing unrelated views (negatives) apart. In the standard setup, each example is treated as its own category, and only its augmented copies count as positives. A variety of contrastive objectives—such as InfoNCE [43], Debaised Contrastive Loss [10], Unbiased Contrastive Loss [2], triplet-based losses [9, 21], and others—have been used to learn robust representations for tasks ranging from dense prediction in computer vision [63] to multimodal alignment [49, 14, 28]. InfoNCE [43] in particular has been shown to lower-bound mutual information [48], and subsequent work has revealed that its empirical success hinges on a balance of *alignment* and *uniformity* in the learned embeddings [55, 62]. In the supervised setting, SupCon [30] extends this idea by using class labels to define positive pairs among same-class samples, often surpassing cross-entropy training in downstream performance. ProjNCE, a generalization of SupCon [27], modifies SupCon loss so that it becomes a proper mutual information lower bound.

A.2 Probabilistic contrastive learning

A growing body of work has begun to integrate probabilistic latent-variable modeling with contrastive objectives. In the video domain, Park et al. represent each video clip as a Gaussian and combine them into a mixture model, learning these distributions via a stochastic contrastive loss that captures clip-level uncertainty and obviates complex augmentation schemes [44]. For 3D point clouds, Wang et al. propose a Generative Variational-Contrastive framework that models latent features as Gaussians, enforces distributional consistency across positive pairs by combining the variational autoencoder and contrastive learning [60]. In graph representation learning, Xie and Giraldo introduce Subgraph Gaussian Embedding Contrast, which maps subgraphs into a structured Gaussian space and employs optimal-transport distances for robust contrastive objectives, yielding improved classification and link-prediction performance [65].

On the theoretical front, Zimmermann et al. prove that contrastive objectives invert the data-generating process under mild conditions, uncovering a deep connection to nonlinear independent component analysis [67]. With a more generalized setting, Kirchhof et al. extend the InfoNCE loss so that the encoder predicts a full posterior distribution rather than a point, and prove that these distributions asymptotically recover the true aleatoric uncertainty of the data-generating process [32].

A.3 Variational Inference and Contrastive Learning

The most closely related line of work frames contrastive learning within a latent-variable inference paradigm via Recognition-Parametrised Models (RPMs) [1, 58]. Aitchison and Ganev introduce RPMs as a class of Bayesian models whose (unnormalized) likelihood is defined implicitly through a recognition network [1]. They show that, under RPMs, the ELBO decomposes into mutual information minus a KL term (up to a constant), and that for a suitable choice of prior the infinite-sample InfoNCE objective coincides with this ELBO. Walker et al. consider RPMs by assuming conditional independence of observations given latent variables, and develop an EM algorithm that achieves exact maximum-likelihood learning for discrete latents along with principled posterior inference [58].

Other works recast variational inference itself as a contrastive estimation task. Rhodes and Gutmann’s Variational Noise-Contrastive Estimation (VNCE) derives a variational lower bound to the standard NCE objective, enabling joint learning of model parameters and latent posteriors in unnormalized models [50]. More recently, Ward et al. propose SoftCVI, which treats VI as a classification problem: they generate “soft” pseudo-labels from the unnormalized posterior and optimize a contrastive-style objective that yields zero-variance gradients at the optimum [64].

A.4 Dimensional collapse

In contrastive self-supervised learning, several approaches have been proposed to prevent dimensional collapse by regularizing either the embedding projector or the second-order statistics of the

representations. Jing *et al.* [29] first demonstrated that, despite the repulsive effect of negative samples, embeddings can still collapse to a low-dimensional subspace due to a combination of strong augmentations and implicit low-rank bias in weight updates. They introduced DirectCLR, which fixes a low-rank diagonal projector during training; this projector enforces the embeddings to occupy a predetermined subspace and was shown empirically to outperform SimCLR’s learned linear projector.

Following this, several works have designed novel loss functions that explicitly regularize the covariance or cross-correlation of the embedding vectors. Ermolov *et al.* [12] apply a whitening MSE loss so that positive pairs match under mean-square error while enforcing identity covariance. Barlow Twins [66] minimize the deviation of the normalized cross-correlation matrix from the identity, effectively performing “soft whitening” to reduce redundancy. VICReg [3] further augments this idea by combining variance, invariance, and covariance regularizers to avoid collapse without using negative samples; notably, VICReg allows its two branches to use different architectures or even modalities, enabling joint embedding across data types. More recently, He *et al.* [17] showed that orthogonal regularization of encoder weight matrices preserves representation diversity and prevents collapse.

B Proofs

B.1 Proof of Lemma 3.1

Proof. With any auxiliary probability function $r(\mathbf{z}'|\mathbf{x})$ and Jensen’s inequality, we have

$$\begin{aligned}\mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] &\geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})r(\mathbf{z}'|\mathbf{x})}\left[\log \frac{p(\mathbf{z}'|\mathbf{x})p(\mathbf{x}|\mathbf{z}')}{r(\mathbf{z}'|\mathbf{x})}\right] \\ &\stackrel{(a)}{=} \mathbb{E}_{q(\mathbf{z}|\mathbf{x})r(\mathbf{z}'|\mathbf{x})}[\log p(\mathbf{z}'|\mathbf{z})] + \mathbb{E}_{r(\mathbf{z}'|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z}')] + H(r(\mathbf{z}'|\mathbf{x})) \\ &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})q(\mathbf{z}'|\mathbf{x})}[\log p(\mathbf{z}'|\mathbf{z})] + \text{const.},\end{aligned}\tag{21}$$

where (a) follows by choosing $r(\mathbf{z}'|\mathbf{x}) = q(\mathbf{z}'|\mathbf{x})$. This proves Lemma 3.1. \square

B.2 Proof of Proposition 3.2

Proof. Optimal critic [39] for InfoNCE satisfies that

$$\psi^*(\mathbf{x}, \mathbf{z}) \propto \log \frac{p(\mathbf{x}|\mathbf{z})}{p(\mathbf{x})} + \alpha(\mathbf{z}),\tag{22}$$

where $\alpha(\mathbf{z})$ only depends on \mathbf{z} . With the optimal critic, we then have

$$\begin{aligned}I_{\text{NCE}}(\mathbf{x}; \mathbf{x}') &= -\mathbb{E}\left[\log \frac{e^{\psi(\mathbf{z}, \mathbf{z}_i')}}{\sum_{j=1}^N e^{\psi(\mathbf{z}, \mathbf{z}_j')}}\right] \\ &= -\mathbb{E}\left[\log \frac{p(\mathbf{z}|\mathbf{z}_i')}{\sum_{j=1}^N p(\mathbf{z}|\mathbf{z}_j')}\right] \\ &= -\mathbb{E}\left[\log \frac{p(\mathbf{z}|\mathbf{z}_i')}{\frac{1}{N} \sum_{j=1}^N p(\mathbf{z}|\mathbf{z}_j')}\right] + \log N.\end{aligned}\tag{23}$$

Given \mathbf{z} , since $p(\mathbf{z}|\mathbf{z}_j')$, $j \in \{1, 2, \dots, N\}$ are i.i.d. with $\mathbb{E}[p(\mathbf{z}|\mathbf{z}_j')] = p(\mathbf{z}) < \infty$, the strong law of large numbers yields

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{j=1}^N p(\mathbf{z}|\mathbf{z}_j') = p(\mathbf{z}).\tag{24}$$

The continuous mapping theorem then gives

$$\lim_{N \rightarrow \infty} \log \frac{p(\mathbf{z}|\mathbf{z}_i')}{\frac{1}{N} \sum_{j=1}^N p(\mathbf{z}|\mathbf{z}_j')} = \log \frac{p(\mathbf{z}|\mathbf{z}_i')}{p(\mathbf{z})}.\tag{25}$$

511 Rearranging (22) and taking $N \rightarrow \infty$, we obtain

$$\begin{aligned}
\lim_{N \rightarrow \infty} \{I_{\text{NCE}}(\mathbf{x}; \mathbf{x}') + \log N\} &= \lim_{N \rightarrow \infty} \mathbb{E} \left[\log \frac{p(\mathbf{z}|\mathbf{z}'_i)}{\frac{1}{N} \sum_{j=1}^N p(\mathbf{z}|\mathbf{z}'_j)} \right] \\
&\stackrel{(a)}{=} \mathbb{E} \left[\lim_{N \rightarrow \infty} \log \frac{p(\mathbf{z}|\mathbf{z}'_i)}{\frac{1}{N} \sum_{j=1}^N p(\mathbf{z}|\mathbf{z}'_j)} \right] \\
&= \mathbb{E} \left[\log \frac{p(\mathbf{z}|\mathbf{z}'_i)}{p(\mathbf{z})} \right], \tag{26}
\end{aligned}$$

512 where the equality (a) follows by dominated convergence theorem that is verifiable using the fact that

$$\begin{aligned}
\mathbb{E} \left[\log \frac{p(\mathbf{z}|\mathbf{z}'_i)}{\frac{1}{N} \sum_{j=1}^N p(\mathbf{z}|\mathbf{z}'_j)} \right] &= \mathbb{E} \left[\log p(\mathbf{z}|\mathbf{z}'_i) - \log \frac{1}{N} \sum_{j=1}^N p(\mathbf{z}|\mathbf{z}'_j) \right] \\
&\leq \mathbb{E} [\log g(\mathbf{z}) - \log \epsilon] \\
&\leq \log \mathbb{E} [g(\mathbf{z})] - \log \epsilon \\
&< \infty. \tag{27}
\end{aligned}$$

513 Rewriting (26) gives

$$\begin{aligned}
&\lim_{N \rightarrow \infty} \{I_{\text{NCE}}(\mathbf{x}; \mathbf{x}') + \log N\} \\
&= \mathbb{E} \left[\log \frac{p(\mathbf{z}|\mathbf{z}'_i)}{p(\mathbf{z})} \right] \\
&= \mathbb{E} \left[\log \frac{p(\mathbf{z}'_i|\mathbf{z})}{p(\mathbf{z}'_i)} \right] \\
&= \mathbb{E}_{q_\theta(\mathbf{z}'_i|\mathbf{x})q_\theta(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{z}'_i|\mathbf{z})] + \mathbb{E}_{q_\theta(\mathbf{z}'_i|\mathbf{x})} [\log p(\mathbf{z}'_i)] \\
&= \mathbb{E}_{q_\theta(\mathbf{z}'_i|\mathbf{x})q_\theta(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{z}'_i|\mathbf{z})] + \mathbb{E}_{q_\theta(\mathbf{z}'_i|\mathbf{x})} \left[\log \frac{p(\mathbf{z}'_i)}{q_\theta(\mathbf{z}'_i|\mathbf{x})} \right] + \mathbb{E}_{q_\theta(\mathbf{z}'_i|\mathbf{x})} [\log q_\theta(\mathbf{z}'_i|\mathbf{x})] \\
&= \mathbb{E}_{q_\theta(\mathbf{z}'_i|\mathbf{x})q_\theta(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{z}'_i|\mathbf{z})] - D(q_\theta(\mathbf{z}'_i|\mathbf{x}) \| p(\mathbf{z}'_i)) - H(q_\theta(\mathbf{z}'_i|\mathbf{x})). \tag{28}
\end{aligned}$$

514 Substituting \mathbf{z}'_i into \mathbf{z}' , this concludes the proof of Proposition 3.2 \square

515 C Variational SupCon

516 C.1 Supervised Contrastive Learning

517 Khosla et al. [30] extend the InfoNCE loss from the self-supervised setting to a supervised context,
 518 calling the resulting method *Supervised Contrastive Learning* (SupCon). When class labels $y_i \in$
 519 $\{1, \dots, C\}$ are available, all samples sharing the same label can serve as positives.

Given a mini-batch $\{(\mathbf{x}_i, y_i)\}_{i=1}^B$, define for each anchor index a

$$\mathcal{A}(a) = \{1, 2, \dots, B\} \setminus \{a\}, \text{ and } \mathcal{P}(a) = \{p \in \mathcal{A}(a) : y_p = y_a\},$$

520 so that $\mathcal{P}(a)$ contains the indices of all positives for anchor a . The SupCon loss for anchor \mathbf{x}_a is then

$$I_{\text{SUP}}(\mathbf{x}_a) = -\frac{1}{|\mathcal{P}(a)|} \sum_{p \in \mathcal{P}(a)} \log \frac{\exp(s(\mathbf{z}_a, \mathbf{z}_p))}{\sum_{j \in \mathcal{A}(a)} \exp(s(\mathbf{z}_a, \mathbf{z}_j))}. \tag{29}$$

521 Averaging over all anchors in the batch yields the full objective:

$$\mathcal{L}^{\text{sup}} = \frac{1}{B} \sum_{a=1}^B I_{\text{SUP}}(\mathbf{x}_a). \tag{30}$$

522 C.2 Variational SupCon (VSupCon)

523 Building on the variational embedding pipeline of VSimCLR, VSupCon simply swaps the unsuper-
 524 vised InfoNCE term for the supervised contrastive loss while retaining the KL regularizer. Concretely,
 525 for each input \mathbf{x} with two augmentations \mathbf{x}' , \mathbf{x}'' , let the encoder output posterior parameters $(\boldsymbol{\mu}', K')$
 526 and $(\boldsymbol{\mu}'', K'')$, and sample normalized embeddings

$$\mathbf{z}' \sim \mathcal{PN}(\boldsymbol{\mu}', K'), \quad \mathbf{z}'' \sim \mathcal{PN}(\boldsymbol{\mu}'', K''). \quad (31)$$

527 Then the VSupCon objective is the symmetrized supervised loss plus the averaged, normalized KL:

$$\mathcal{L}^{\text{VSup}} = \frac{1}{2} \left(\mathcal{L}^{\text{sup}}(\mathbf{z}', \mathbf{z}'') + \mathcal{L}^{\text{sup}}(\mathbf{z}'', \mathbf{z}') \right) + \frac{1}{2d} \left(D(\boldsymbol{\mu}', K') + D(\boldsymbol{\mu}'', K'') \right). \quad (32)$$

528 Minimizing $\mathcal{L}^{\text{VSup}}$ therefore aligns same-class embeddings and regularizes their posterior distribu-
 529 tions toward the uniform prior on the sphere.

530 D Discussion on the approximation in Section 3.1

531 D.1 Discussion on (8)

532 The key step in our decoder-free ELBO maximization is the approximation

$$\mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] \approx \mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})q_\theta(\mathbf{z}'|\mathbf{x})} [\log p(\mathbf{z}'|\mathbf{z})] \quad (33)$$

533 **Lower-bound view.** As shown in Lemma 3.1, this approximation admits a lower bound up to an
 534 additive constant independent of \mathbf{z} :

$$\mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] \geq \mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})q_\theta(\mathbf{z}'|\mathbf{x})} [\log p(\mathbf{z}'|\mathbf{z})] + \text{const}. \quad (34)$$

535 Consequently, maximizing the right-hand side with respect to θ implicitly maximizes the reconstruc-
 536 tion term $\mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})]$, which is the objective of ELBO maximization. Moreover, using (10)
 537 (see Section 3.1), the surrogate is negatively related to InfoNCE:

$$\mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] \approx -I_{\text{NCE}}(\mathbf{x}; \mathbf{x}'), \quad (35)$$

538 so minimizing the InfoNCE loss increases the reconstruction term.

539 **Change-of-variables view.** Another perspective on the reconstruction approximation (8) comes
 540 from a change of variables. Let g be an invertible, differentiable mapping such that $\mathbf{x} = g(\mathbf{z}')$. Then,
 541 by the change-of-variables formula,

$$p(\mathbf{x}|\mathbf{z}) = p(\mathbf{z}'|\mathbf{z}) |\det J_{g^{-1}}(\mathbf{x})| = p(\mathbf{z}'|\mathbf{z}) |\det J_g(\mathbf{z}')|^{-1}, \quad (36)$$

542 where J_g and $J_{g^{-1}}$ denote the Jacobians of g and g^{-1} , respectively, and $\mathbf{z}' = g^{-1}(\mathbf{x})$. Taking
 543 logarithms yields

$$\log p(\mathbf{x}|\mathbf{z}) = \log p(\mathbf{z}'|\mathbf{z}) + \log |\det J_{g^{-1}}(\mathbf{x})| = \log p(\mathbf{z}'|\mathbf{z}) - \log |\det J_g(\mathbf{z}')|, \quad (37)$$

544 where the second term depends only on \mathbf{x} (equivalently, on \mathbf{z}') and is independent of \mathbf{z} .

545 **Sufficient condition (tightness).** If, in addition to invertibility, g is *volume-preserving*, i.e.,
 546 $|\det J_{g^{-1}}(\mathbf{x})| \equiv 1$ (equivalently, $|\det J_g(\mathbf{z}')| \equiv 1$) on the data manifold, then the additive term
 547 in (37) vanishes and we obtain the tight equality $\log p(\mathbf{x}|\mathbf{z}) = \log p(\mathbf{z}'|\mathbf{z})$. More generally,
 548 when $|\det J_{g^{-1}}(\mathbf{x})|$ is approximately constant over the data manifold, the additive term acts as
 549 (approximately) a constant shift independent of \mathbf{z} , yielding a tight surrogate for optimization.

550 This assumption is plausible in practice under the commonly observed *dimension-collapse* phe-
 551 nomenon: the embeddings \mathbf{z}' have effective rank (intrinsic dimension) much smaller than the ambient
 552 embedding dimension yet retain nearly all task-relevant information about the features \mathbf{x} . When the
 553 feature and embedding manifolds have (approximately) the same intrinsic dimension and g behaves
 554 near-isometrically between them, the Jacobian determinant varies weakly, making the surrogate
 555 in (37) tight in practice.

Table 3: Gaussian KL (G-KL) vs. projected normal KL (PN-KL) on synthetic data.

	G-KL	PN-KL	Gap (G-KL–PN-KL)	Ratio (G-KL/PN-KL)
mean	106.86	97.37	9.49	0.91
std	9.56	7.63	-	-

D.2 Gaussian KL Surrogate for Projected-Normal KL

We study the tightness of the bound in (14), repeated here:

$$D(\mathcal{N}(\mu, K) \parallel \mathcal{N}(0, I_d)) \geq D(\mathcal{PN}(\mu, K) \parallel \text{Unif}(\mathcal{S}^{d-1})). \quad (38)$$

Before analyzing tightness, we note several practical benefits of using the Gaussian KL as a surrogate for the projected-normal KL:

- **Closed form.** It is trivial to implement and numerically stable.
- **Aligned optima.** The Gaussian KL and projected-normal KL share the same minimizer (e.g., at $\mu = 0$ and $K = I_d$), so optimizing the surrogate steers the model toward the same optimum.
- **Efficiency.** Unlike Monte Carlo or k -NN estimators needed for the projected-normal KL, the Gaussian KL requires no sampling.

Moreover, the KL term acts only as a regularizer, whereas InfoNCE directly drives semantic similarity; thus modest approximation error in the KL has limited effect on downstream performance.

We assess tightness by comparing the closed-form Gaussian KL with an estimated projected-normal KL using a divergence estimator [61] in two settings: synthetic data and CIFAR-10 under VCL training.

KL gap on synthetic data. We approximate $D(\mathcal{PN}(\mu, K) \parallel \text{Unif}(\mathcal{S}^{d-1}))$ numerically using 10^5 samples in dimension $d = 128$ for random (μ, K) draws, with $\mu \sim \mathcal{N}(0, I_d)$ and

$$K = \frac{1}{d}AA^\top + 0.1 I_d, \quad A_{ij} \sim \mathcal{N}(0, 0.5) \quad \forall i, j. \quad (39)$$

We employ the k -nearest-neighbor divergence estimator [61] with $k = 1$, compute both the Gaussian KL (analytically) and the projected-normal KL (using the estimator) on the same samples, and repeat over 20 random trials to reduce variance.

Table 3 reports the gap between the two KLs on synthetic data: the average absolute gap is approximately 9.49 (about a 10% relative difference). Thus, the Gaussian KL surrogate closely tracks the projected-normal KL while retaining the practical advantages noted above.

KL gap on CIFAR-10. Beyond the synthetic study, we measure the gap during VCL training on CIFAR-10 using the same experimental settings (Appendix E.1); results are shown in Figure 7. After only a few epochs, the Gaussian KL and the projected-normal KL closely track each other. This indicates that minimizing the Gaussian-KL surrogate effectively minimizes the projected-normal KL—the quantity we aim to reduce—while retaining the practical advantages of the surrogate.

E Experiments

E.1 Training Details and Hyperparameters

Datasets and preprocessing. Experiments are conducted on CIFAR-10 [33], CIFAR-10C [19], CIFAR-10H [46], CIFAR-100 [33], STL-10 [11], Tiny-ImageNet [35], and Caltech-256 [15]. We train VCL models on CIFAR-10/100, Tiny-ImageNet, and Caltech-256, Tiny-ImageNet, and STL10. Following SimCLR, we sample two views per image via random resized crop (image size 32×32 and scale $[0.2, 1.0]$), horizontal flip ($p=0.5$), color jitter (brightness/contrast/saturation/hue = 0.4, applied with $p=0.8$), Gaussian blur (kernel size 9), and random grayscale ($p=0.2$). Inputs are normalized with dataset-specific means/standard deviations.

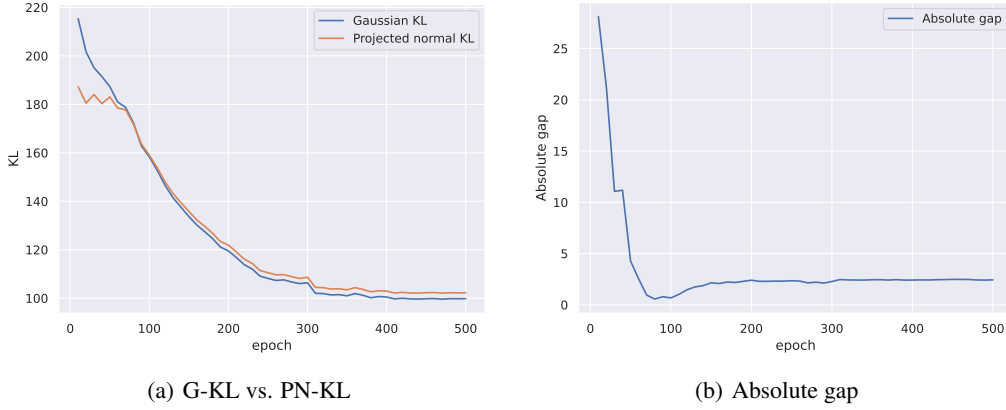


Figure 7: Tracking Gaussian KL (G-KL) and projected normal KL (PN-KL) during VCL training on CIFAR-10. (a) G-KL vs. PN-KL; (b) Absolute gap, $|G-KL - PN-KL|$. This shows that minimizing Gaussian KL leads to minimizing projected normal KL.

Architectures. We use ResNet-18 [18] as encoder and embedding dimension $d = 128$, and employ a linear classifier for downstream evaluations.

Optimization. We use AdamW [38] with base LR 10^{-2} (encoder and head), weight decay 10^{-4} , batch size $B=512$, and $T=500$ epochs for pretraining and $T = 100$ for training linear classifier. Temperature for InfoNCE loss is $\tau=0.07$. We set $m=1$ posterior samples per view for VSimCLR and VSupCon by default (ablation in Table 5). No momentum encoder or queue is used; all negatives are in-batch. For training stability, we clip the posterior log-variance ($\log \sigma^2$) to $[-5, 5]$ to bound variances, and clip gradient global norm at 1.0.

E.2 Additional Results on Dimension Collapse

In addition to the singular spectrum of VCL embeddings on CIFAR-10 and CIFAR-100 in Figure 3, Figure 8 reports results on Caltech-256 and Tiny-ImageNet. In both datasets, VCL mitigates the dimension-collapse phenomenon commonly observed in contrastive learning.

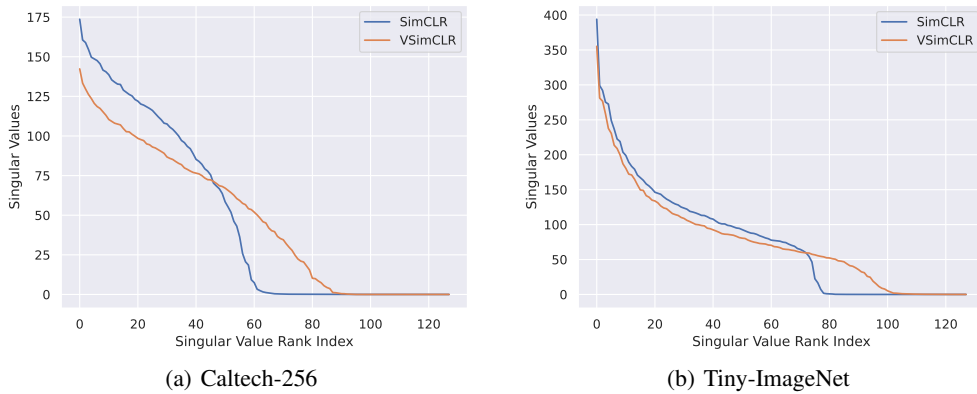


Figure 8: Singular-value spectrum of the embedding covariance on Cartech-256 and Tiny-ImageNet. VSimCLR mitigates dimensional collapse on both datasets.

E.3 Distributional Contrastive Loss

In addition to the contrastive loss on embeddings, it is worthwhile to contrast the posterior distributions within the VCL framework. Specifically, we aim to pull together the posteriors corresponding to

Table 4: Log-determinant of average posterior covariance K for each CIFAR-10 class.

Index	Class	$\log \det(K)$
0	airplane	-182.207
1	automobile	-181.691
2	bird	-183.713
3	cat	-191.317
4	deer	-184.969
5	dog	-185.432
6	frog	-182.125
7	horse	-179.331
8	ship	-185.991
9	truck	-188.179

Table 5: Classification accuracy on STL10 with different number of embedding generation from posterior. We report top-1 and top-5 accuracies of SimCLR, VSimCLR, SupCon, and VSupCon across the datasets with different m and DistNCE (40).

METHOD	STL10	
	TOP1	TOP5
SIMCLR	60.44	95.80
VSIMCLR ($m = 1$)	60.11	92.00
VSIMCLR ($m = 4$)	57.86	88.29
VSIMCLR ($m = 16$)	59.13	92.85
VSIMCLR ($m = 64$)	56.91	86.63
VSIMCLR WITH DISTNCE (40)	36.54	80.25
VSIMCLR (ASYM)	57.38	88.78
SUPCON	75.88	75.88
VSUPCON ($m = 1$)	75.76	96.99
VSUPCON ($m = 4$)	74.35	97.14
VSUPCON ($m = 16$)	76.11	98.39
VSUPCON ($m = 64$)	77.96	98.44

different augmentations of the same input and to push apart posteriors from distinct inputs. To incorporate this into VCL, we introduce the *DistNCE* loss, a contrastive objective over posterior parameters, defined as

$$D_{\text{DistNCE}}(\theta) = -\mathbb{E} \left[\log \frac{\exp(s(\theta, \theta^+))}{\sum_j \exp(s(\theta, \theta_j))} \right], \quad (40)$$

where θ denotes the posterior parameters (μ, K) , θ^+ is the positive-pair parameter for the same input, and $\{\theta_j\}_{j \neq +}$ are negative-pair parameters from other inputs. The expectation is taken over the joint distribution $p(\theta, \theta^+) \prod_{j \neq +} p(\theta_j)$.

Moreover, we increase the number of posterior samples used for the InfoNCE loss. Specifically, we draw m samples $\{\mathbf{z}^{(k)}\}_{k=1}^m$ from each posterior, resulting in an m -fold increase in effective batch size, and compute the InfoNCE loss over this enlarged set of embeddings. The classification results are reported in Table 5.

We also evaluate the performance of the asymmetric lower bound (16) (denoted ASYM) in place of the symmetrized objective (17). These results are also shown in Table 5.

From these experiments, we did not observe any significant differences when applying DistNCE (40), using the asymmetric loss, or sampling multiple embeddings per posterior. Based on these findings, we proceed with the basic VCL variants from the main text for all subsequent experiments.

Table 6: Classification accuracy on STL10 with different number of embedding generation from posterior. We report top-1 and top-5 accuracies of SimCLR, VSimCLR, SupCon, and VSupCon across the datasets with different m and DistNCE (40).

β	TOP-1 ACCURACY	TOP-5 ACCURACY
1	47.90	72.34
0.1	47.24	71.90
0.01	50.35	73.27
0.001	51.34	73.09



Figure 9: Sample images from the CIFAR-10, organized by class (columns) and sorted by their corresponding $\log \det(K)$ (rows). In each column, the top image has the highest $\log \det(K)$, the bottom image the lowest; the overlaid numbers indicate each image’s $\log \det(K)$.

623 E.4 Effect of KL Regularizer on Classification

624 As shown in Table 1, VSupCon exhibits reduced classification accuracy on some datasets, whereas
625 VSimCLR remains stable. We attribute this degradation to two factors:

- 626 1. VSimCLR’s objective coincides with the VCL objective in (17), but VSupCon’s does not,
627 creating a mismatch that can impede proper ELBO maximization.
- 628 2. SupCon optimizes embeddings directly for classification; adding a KL term can conflict
629 with this objective.

630 We therefore hypothesize that weakening the KL regularizer improves VSupCon’s accuracy. To test
631 this, we scale the KL term by $\beta \in \{1, 10^{-1}, 10^{-2}, 10^{-3}\}$,

$$\mathcal{L}^{\text{vsup}}(\beta) = \mathcal{L}^{\text{sup}} + \beta D_{\text{KL}}(q_{\theta}(\mathbf{z} \mid \mathbf{x}) \parallel p(\mathbf{z})), \quad (41)$$

632 and evaluate the resulting embeddings. As expected, smaller β (i.e., a weaker KL effect) yields higher
633 accuracy. Thus, for pure classification tasks, SupCon may not benefit from a VCL variant unless the
634 KL weight is carefully tuned.

635 E.5 Implications of Distributional Embeddings

636 Distributional (probabilistic) embeddings provide useful capabilities, including uncertainty quantifi-
637 cation and probability-based distances between samples and classes. We analyze them along three
638 axes: uncertainty, typicality, and out-of-distribution (OOD) behavior.

639 **Posterior covariance vs. uncertainty.** As shown in Figure 9, different samples exhibit varying
640 degrees of posterior dispersion (e.g., the log-determinant of the covariance, $\log \det(K)$), which can

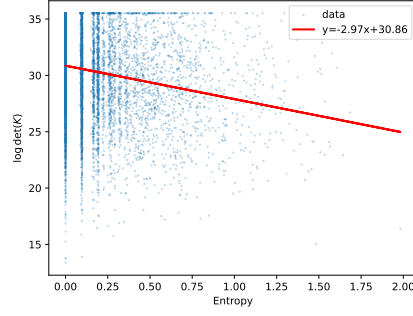


Figure 10: Relationship between posterior dispersion and label ambiguity. Each point plots the trace of K ($\text{tr}(K)$) against the entropy of human-annotated class probabilities from CIFAR-10H [46], with a first-order linear fit (red line). Similar to the result in Figure 5, the dispersion is negatively correlated with label ambiguity.

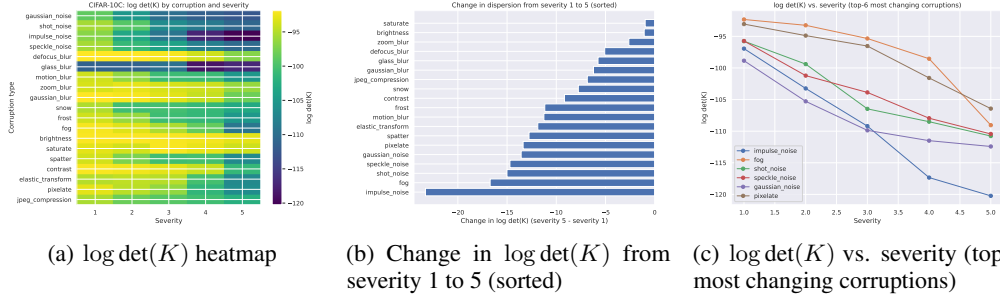


Figure 11: $\log \det(K)$ of VSupCon embeddings on CIFAR-10C [19] under different corruption types and severities. “Severity” denotes the corruption level. The observed negative correlation between $\log \det(K)$ and severity is consistent with our finding that more uncertain samples exhibit smaller posterior covariance dispersion. Exact $\log \det(K)$ values are in Table 8.

serve as an uncertainty measure. To examine how uncertainty and posterior covariance are related, we conduct experiments on two benchmark datasets, CIFAR-10H [46] and CIFAR-10C [19]:

- **CIFAR-10H:** The test set provides soft labels [24, 25, 26] aggregated from multiple annotators. Using these soft labels, we compute the per-sample label entropy as a measure of uncertainty about the underlying class.
- **CIFAR-10C:** The test set provides systematically corrupted images with multiple corruption types and severities (higher severity = stronger corruption), which induces greater label ambiguity and thus higher uncertainty.

Beyond comparing $\log \det(K)$ with label entropy in Figure 5, we also compare the trace of K (denoted $\text{tr}(K)$) against label entropy in Figure 10. In both cases, we observe a *negative* slope under a first-order linear fit. This indicates that VSimCLR assigns **lower** posterior dispersion to inputs with greater label uncertainty. Conversely, inputs that humans classify unambiguously—i.e., prototypical class examples—exhibit posteriors with **larger** dispersion, suggesting their latent representations span a broader region of the class-specific embedding space; ambiguous or outlier inputs yield **smaller** dispersion, reflecting more concentrated latent distributions.

A similar pattern appears in Figures 6 and 11, which relate $\log \det(K)$ to corruption severity on CIFAR-10C. We train VSimCLR and VSupCon on CIFAR-10 and evaluate their embeddings on CIFAR-10C. Because higher severity entails stronger corruption and greater label ambiguity, these figures further support the finding that posterior covariance dispersion is negatively correlated with uncertainty. Tables 7 and 8 report the mean $\log \det(K)$ for each corruption type and severity level.

Table 7: Average $\log \det K$ of VSimCLR embeddings on CIFAR-10C for each corruption type and severity (higher severity = stronger corruption).

Corruption	Severity 1	Severity 2	Severity 3	Severity 4	Severity 5
gaussian_noise	-187.74	-189.85	-192.23	-193.05	-193.70
shot_noise	-187.49	-188.11	-190.18	-190.95	-191.97
impulse_noise	-188.25	-190.71	-192.61	-194.66	-194.82
speckle_noise	-187.59	-188.64	-189.21	-189.93	-190.48
defocus_blur	-184.41	-183.84	-182.67	-187.67	-186.76
glass_blur	-192.35	-191.76	-192.03	-194.36	-193.98
motion_blur	-185.83	-187.53	-189.88	-189.78	-191.94
zoom_blur	-185.95	-183.85	-183.86	-183.75	-185.07
gaussian_blur	-184.43	-182.83	-182.11	-183.47	-191.56
snow	-186.92	-189.86	-190.48	-193.08	-193.89
frost	-188.43	-190.13	-192.08	-192.16	-193.85
fog	-185.61	-187.61	-189.65	-193.37	-204.82
brightness	-184.89	-185.43	-186.17	-187.16	-189.70
saturate	-186.40	-191.14	-185.02	-186.36	-187.87
spatter	-186.32	-188.43	-191.12	-188.88	-191.03
contrast	-185.67	-188.03	-189.84	-192.59	-200.25
elastic_transform	-185.66	-185.12	-184.95	-189.66	-195.31
pixelate	-185.10	-186.44	-187.62	-188.58	-189.46
jpeg_compression	-182.94	-183.30	-183.73	-184.38	-185.28

Table 8: Average $\log \det K$ of VSupCon embeddings on CIFAR-10C for each corruption type and severity (higher severity = stronger corruption).

Corruption	Severity 1	Severity 2	Severity 3	Severity 4	Severity 5
gaussian_noise	-98.85	-105.28	-109.87	-111.50	-112.42
shot_noise	-95.76	-99.39	-106.47	-108.50	-110.77
impulse_noise	-96.94	-103.24	-109.20	-117.34	-120.23
speckle_noise	-95.73	-101.21	-103.87	-107.95	-110.44
defocus_blur	-91.95	-91.90	-92.33	-93.94	-97.03
glass_blur	-111.32	-111.29	-109.63	-118.74	-117.08
motion_blur	-93.95	-96.48	-100.86	-100.96	-105.21
zoom_blur	-93.66	-92.94	-93.67	-94.06	-96.29
gaussian_blur	-91.95	-92.31	-93.14	-94.40	-98.17
snow	-95.28	-100.62	-100.32	-101.30	-103.04
frost	-93.98	-96.23	-100.71	-101.33	-105.15
fog	-92.33	-93.25	-95.34	-98.54	-109.05
brightness	-92.04	-92.06	-92.16	-92.40	-93.11
saturate	-93.05	-93.80	-92.14	-92.82	-94.02
spatter	-93.86	-97.46	-100.59	-100.27	-106.63
contrast	-92.14	-92.54	-93.10	-94.30	-101.31
elastic_transform	-95.01	-94.65	-94.96	-100.26	-106.89
pixelate	-93.06	-94.88	-96.53	-101.58	-106.43
jpeg_compression	-95.47	-98.31	-99.28	-100.59	-102.32

661 This counterintuitive observation—that typical (i.e., common) samples exhibit larger posterior dis-
662 persion—parallels the concurrent findings of Guth et al. [16], albeit under different settings: (i)
663 **Quantity:** we analyze latent-space posterior dispersion via $\log \det K$, whereas they study input-space
664 marginal density $p(x)$; (ii) **Observation:** typical samples have larger $\log \det K$ (ours), while they
665 have lower $p(x)$ (theirs). Although these quantities live in different spaces, both results indicate that
666 typical samples are not the highest-density points. In our case, typical images yield larger dispersion
667 and atypical images smaller dispersion; since dispersion is inversely related to peak density, our result
668 is consistent with Guth et al. Hence, in both settings, “typical” \neq “highest-density.” Consequently,

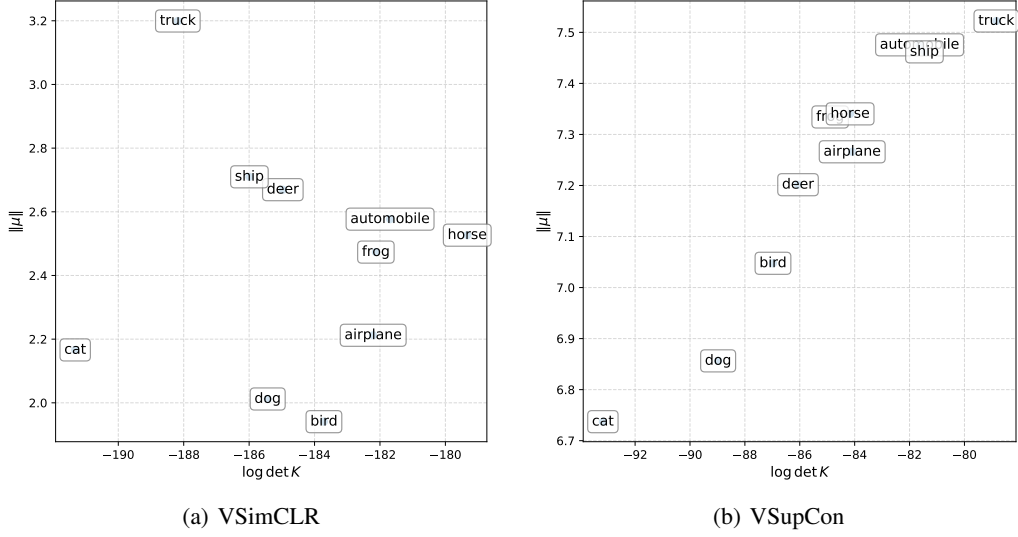


Figure 12: Norm of the posterior mean $\|\mu\|$ versus the log-determinant of the covariance $\log \det(K)$, averaged per class. Both μ and K are computed by averaging over all samples belonging to the same class.

posterior dispersion serves as a useful uncertainty signal; see Table 2 for an application under label scarcity.

Class-wise average posterior parameters. Figure 12 reports class-wise averages of the posterior parameters—the mean norm $\|\mu\|$ and the covariance dispersion $\log \det K$ —for VSimCLR and VSupCon. Classes exhibit distinct dispersion profiles. Despite being trained independently, the two methods yield similar class-wise patterns in both quantities: for example, the *cat* and *dog* classes show comparatively lower $\|\mu\|$ and $\log \det K$, whereas *truck* attains the largest $\|\mu\|$. Table 4 provides detailed per-class $\log \det K$ values.

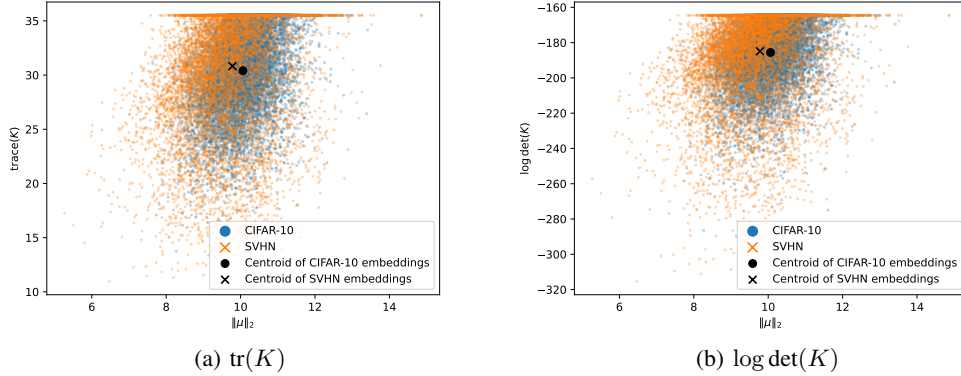


Figure 13: Posterior parameters of CIFAR-10 and SVHN datasets. We use the same encoder of VSimCLR trained with CIFAR-10.

Posterior on in-distribution vs. out-of-distribution. We compare per-sample posterior parameters under VSimCLR for in-distribution (ID; CIFAR-10) versus out-of-distribution (OOD; SVHN [42]) inputs. VSimCLR is trained on the CIFAR-10 training set, after which we extract (μ, K) on the CIFAR-10 and SVHN test sets. Figure 13 plots the pairs $(\|\mu\|, \log \det K)$ for each dataset; black markers denote dataset-wise means. While the mean values $\text{avg}(\|\mu\|)$ and $\text{avg}(\log \det K)$ are similar

682 across CIFAR-10 and SVHN, the SVHN points exhibit substantially greater spread (dispersion)
683 across samples, indicating a broader posterior-parameter distribution for OOD data.