
Probabilistic Contrastive Learning Recovers the Correct Aleatoric Uncertainty of Ambiguous Inputs

Michael Kirchhof¹ Enkelejda Kasneci² Seong Joon Oh¹

Abstract

Contrastively trained encoders have recently been proven to invert the data-generating process: they encode each input, e.g., an image, into the true latent vector that generated the image (Zimmermann et al., 2021). However, real-world observations often have inherent ambiguities. For instance, images may be blurred or only show a 2D view of a 3D object, so multiple latents could have generated them. This makes the true posterior for the latent vector probabilistic with heteroscedastic uncertainty. In this setup, we extend the common InfoNCE objective and encoders to predict latent distributions instead of points. We prove that these distributions recover the correct posteriors of the data-generating process, including its level of aleatoric uncertainty, up to a rotation of the latent space. In addition to providing calibrated uncertainty estimates, these posteriors allow the computation of credible intervals in image retrieval. They comprise images with the same latent as a given query, subject to its uncertainty.

1. Introduction

Contrastive learning (Chen et al., 2020) trains encoders to output embeddings that are close to one another for semantically similar inputs and far apart for unsimilar inputs. This general notion of similarity allows transferring the pre-trained encoders to downstream tasks (Wang et al., 2022; Islam et al., 2021; Khosla et al., 2020).

Recently, Zimmermann et al. (2021) corroborated this intuition by a theoretical result: under weak assumptions, the embeddings of an encoder trained under an InfoNCE (Oord et al., 2018) loss are exactly equal to the true latent vectors, up to a rotation of the spherical latent space. This comes from a nonlinear Independent Component Analysis (ICA) perspective (Comon & Jutten, 2010). It assumes an

unknown nonlinear generative process that transforms true latents into our observations. Contrastively trained encoders *invert* this nonlinear function and recover the original latent space, up to a rotation.

This holds for the class of generative processes that are deterministic and injective, so that each image could have been generated by only one latent vector. This is often violated in practice. In Figure 1, the lower image of an animal is in low-resolution, so it is impossible to tell which exact species, i.e., which latent variables, underlie the image. In fact, most scenarios in the wild involve some form of such aleatoric uncertainty, including 3D-to-2D projections (Chen et al., 2021), partially covered objects (Kraus & Dietmayer, 2019), or images with a too low resolution or a bad crop (Li et al., 2021). It also manifests itself outside the image domain, such as in the ambiguity of natural language descriptions (Chun et al., 2022) or measurement noise in general (Meech & Stanley-Marbell, 2021).

This work generalizes the previous theoretical result to this more challenging setting. We do not assume that generative process is an injective and deterministic function, but allow it to be a conditional distribution. We propose Monte-Carlo InfoNCE (MCInfoNCE), a probabilistic analog of InfoNCE for probabilistic encoders that predict distributions over the possible latents, called probabilistic embeddings (Oh et al., 2019; Shi & Jain, 2019). We prove that MCInfoNCE attains its global minimum when the encoder recovers the true posteriors of the generative process, up to a rotation of the latent space; both in terms of both the mean (which latent is most likely to have generated the image) and the variance (the level of aleatoric uncertainty of the individual image). Our work thus generalizes the previous theoretical result in nonlinear ICA to a broader class of generative processes, and provides a theoretical foundation for the recent literature on probabilistic embeddings.

We show empirically that an encoder trained with MCInfoNCE learns the correct posteriors in a controlled experiment with known posteriors. We find that it even provides sensible embeddings when the distribution family or the encoder dimensionality is misspecified and when the generative process may be injective, making it robust in practice. We then show that these predicted uncertainties are consis-

¹University of Tübingen, Germany ²TUM University, Munich, Germany. Correspondence to: Michael Kirchhof <michael dot kirchhof at uni dash tuebingen dot de>.

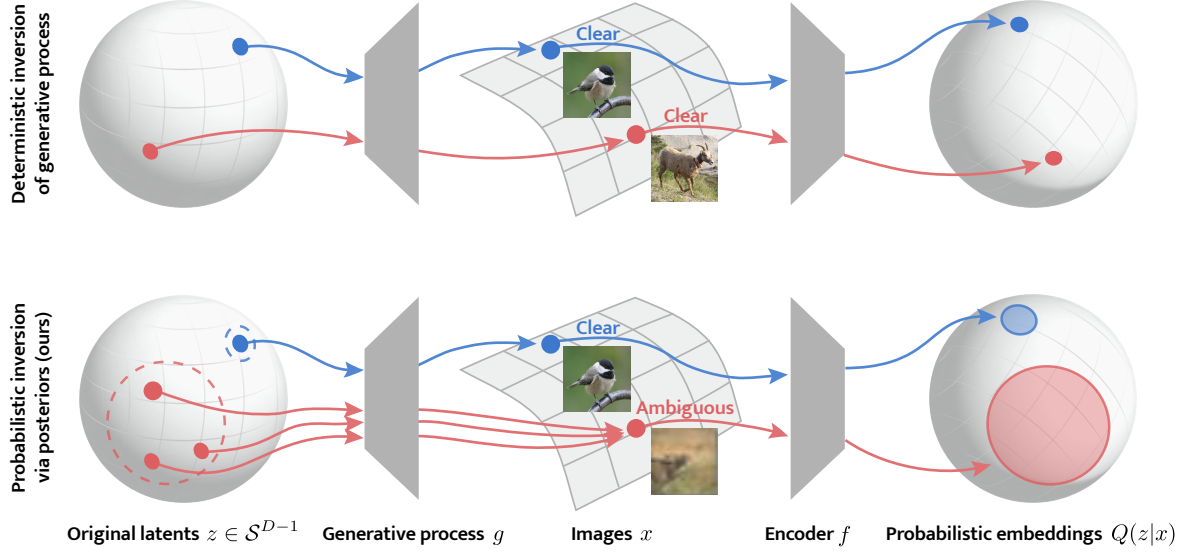


Figure 1. Deterministic encoders embed images to points in the latent space. This recovers the latent vectors that generated them, up to a rotation (top). However, if an image is ambiguous there are multiple possible latents that could have generated it (bottom). An encoder trained with MCInfoNCE correctly recovers this posterior of the generative process, up to a rotation, from contrastive supervision.

tent with human annotator disagreements reported in the recent CIFAR-10H dataset (Peterson et al., 2019), providing a way to handle uncertainty for high-dimensional inputs. We also demonstrate that knowing the true posteriors enables new applications, such as computing credible intervals for image retrieval tasks. They visualize how uncertain we are about a query image by showing other images that represent the region of latents the query is in with a given probability.

In summary, (1) We extend nonlinear ICA to non-injective non-deterministic generative processes to model realistic input ambiguities. (2) We propose MCInfoNCE for training encoders that predict probabilistic embeddings. (3) We show theoretically and empirically that the predicted posteriors are correct and reflect the true amount of aleatoric uncertainty.

2. Related Works

Our work serves as a bridge between the theoretical understanding of contrastive learning via nonlinear ICA, probabilistic embeddings, and recent discussions on the aleatoric uncertainty inherent in vision problems. Below, we discuss how our work extends and connects recent work in these three fields. Extended literature reviews can be found in Kendall & Gal (2017) and Karpukhin et al. (2022).

Nonlinear ICA. From a nonlinear Independent Component Analysis (ICA) perspective (Hyvärinen & Oja, 2000; Comon & Jutten, 2010), images x are modeled as being generated from ground-truth latent components z via an unknown nonlinear generative process. The goal is to invert it to recover the original latents z , which are useful for

downstream tasks. This formalization allows for theoretical proofs of which (contrastive) losses achieve this. Building on Wang & Isola (2020), Zimmermann et al. (2021) recently proved that optimizing a contrastive InfoNCE loss (Oord et al., 2018) recovers z up to a rotation of the latent space, as visualized in Figure 1. This requires certain assumptions about the generative process. A recent strain of literature seeks to reduce these assumptions (Leemann et al., 2022; Roth et al., 2022) to allow modeling broader classes of generative processes, bringing the theoretical results closer to practice. Our work no longer requires the injectivity assumption of Zimmermann et al. (2021) and at the same time we allow stochasticity. This is made possible by modeling the generative process as a conditional distribution $P(x|z)$ instead of a function, which generalizes the class of generative processes. In the vein of Zimmermann et al. (2021), we prove that our contrastive MCInfoNCE loss recovers the correct posterior distribution $P(z|x)$ of the original latents, up to a rotation of the latent space.

Aleatoric Uncertainty. The above generalization allows us to model scenarios in which we encounter aleatoric uncertainty, i.e., the input has reduced information such that z is only recoverable only up to some uncertainty. A prominent practical example is face recognition, where images may be blurred or in low-resolution (Shi & Jain, 2019; Schlett et al., 2022). Other problems with ambiguous inputs include 3D reconstruction from 2D data (Chen et al., 2021), partially occluded traffic participants (Kraus & Dietmayer, 2019), or noisy physical sensors (Meech & Stanley-Marbell, 2021). Such problems with aleatoric uncertainty can be detected by label noise: CIFAR-10H (Peterson et al., 2019) comprises

multiple labels for each image in the CIFAR-10 test-set, and shows that the more ambiguous an image is, the more annotator labels disagree. This finding occurs in several other recent classification datasets (Schmarje et al., 2022; Mehrtens et al., 2023), but also in more complex tasks such as multimodal visual question answering (VQA). Chun et al. (2022) show that there are many possible textual answers to the same visual prompt because language is more ambiguous than vision; i.e., language has more aleatoric uncertainty. Our MCInfoNCE loss explicitly accounts for these uncertainties and learns the correct level of aleatoric uncertainty, which we demonstrate on high-dimensional image inputs.

Probabilistic Embeddings. An emerging approach to modeling this uncertainty is to have encoders predict distributions over the latent space instead of point estimates. There are three main lines of work to learn these probabilistic embeddings. The first idea is to compute a match probability between point estimates, but to integrate it over the predicted distributions. This idea was pioneered via Hedged Instance Embeddings (HIB) (Oh et al., 2019) and has since been successfully extended, e.g., to the above multimodal VQA problem (Chun et al., 2021; Neculai et al., 2022). A second line of works turns existing losses into probabilistic ones by integrating the whole loss over the predicted probabilistic embeddings (Scott et al., 2021). Our MCInfoNCE extension of InfoNCE demonstrates that this blueprint strategy can inherit the properties of the original losses, like Zimmermann et al.’s identifiability theorem. The third line of works provides distribution-to-distribution distances to replace point-to-point distances in losses. The most popular approach is the expected likelihood kernel (ELK) (Jebara & Kondor, 2003; Shi & Jain, 2019). It has recently shown success even in high dimensional embedding spaces (Kirchhof et al., 2022; Karpukhin et al., 2022). Yet, there is no answer to whether and in what sense the predicted probabilistic embeddings, and in particular their variances, are *correct*. Our work answers this question through its proof and a controlled experiment where the true posteriors are recovered. The experiments on CIFAR-10H further ground this theoretical correctness in the human perception of uncertainty. We also show novel practical applications of probabilistic embeddings, such as retrieving credible intervals on which latents the image might show.

3. Probabilistic Generative Processes

In this section, we extend the generative processes commonly used in nonlinear ICA to non-injective, randomized ones. This allows modeling real-world image distributions better and serves as a framework for the upcoming proof.

Let us first understand the class of generative processes for which Zimmermann et al. (2021) provide an identifiability proof. They take the nonlinear ICA perspective that there

is a natural generative process g that transforms latent components $z \in \mathcal{Z}$ into the images $x = g(z)$ we observe, as shown in Figure 1. Following the popular cosine-based similarity comparisons (Deng et al., 2019; Teh et al., 2020), \mathcal{Z} is assumed to be a D -dimensional hypersphere $\mathcal{Z} = \mathcal{S}^{D-1}$. We are interested in recovering the latents z that underlie the images x , because they are low-dimensional descriptions useful for downstream tasks. To formalize this problem, they assume that $g : \mathcal{Z} \rightarrow \mathcal{X}$ is an injective (and deterministic) function. Thus, only one latent z can correspond to each image x , and g is invertible. They prove that an encoder f trained with a contrastive InfoNCE loss achieves this inversion and recovers the correct latent z , i.e., $f(x) = f(g(z)) = \hat{z} = Rz$, up to an orthogonal rotation R of the learned embedding space.

However, let us move on to setups where an image x may be motion blurred, low-resolution, or partially obscured. For instance, a 2D projection x of a 3D object z does not show the back part of z , and there are several possible z that could have generated x . In other words, the generative process g is non-injective and the best our encoder can do is to recover the set of possible latents $\{\hat{z} | g(\hat{z}) = x\}$. Further, g may be stochastic. E.g., a random patch of pixels may be occluded, or the image may be zoomed in and show only a random crop of z . The best the encoder can do is to predict a posterior over the possible latents, see Figure 1.

The common denominator of these setups is that g loses information about z and x becomes ambiguous. To subsume them, we can model g as a likelihood $P(x|z)$. This general formulation allows for a large class of operations within g . However, this generality comes at the cost that $P(x|z)$ can be very complicated and difficult to parameterize. We therefore apply a *posterior trick*: instead of explicitly characterizing g by $P(x|z)$ we implicitly characterize it by its posteriors $P(z|x)$. We parameterize $P(z|x)$ by simple von Mises-Fisher distributions $\text{vMF}(z; \mu(x), \kappa(x))$:

$$P(z|x) = C(\kappa(x))e^{\kappa(x)\mu(x)^\top z}. \quad (1)$$

This distribution on \mathcal{S}^{D-1} is unimodal around the location parameter $\mu(x) \in \mathcal{Z}$ with a certain concentration (i.e., an inverse variance) $\kappa(x) \in \mathbb{R}_{>0}$, and a normalizing constant $C(\cdot)$. The functions $\mu : \mathcal{X} \rightarrow \mathcal{S}^{D-1}$ and $\kappa : \mathcal{X} \rightarrow \mathbb{R}_{>0}$ fully parameterize the posterior of each image x . In particular, $\kappa(\cdot)$ represents the aleatoric uncertainty due to information loss, which can be heterogeneous across the images.

The intuition behind modeling the posterior of the generative process as a vMF is that latents of degraded images can usually be located down to sets of semantically similar rather than very dissimilar latents. This is reflected in the unimodality of the vMF and its use of the dot product, which commonly represents how semantically similar two latents are. There may still be images where it is impos-

sible to tell which highly dissimilar latents they show. In these cases, $\kappa(x)$ is low and the posterior spreads broadly across the latent space. At the other end of the spectrum, as $\kappa(x) \rightarrow \infty$, $P(z|x)$ converges to a Dirac distribution. This allows modelling deterministic and injective generative processes as in Zimmermann et al. (2021). This makes the vMF a reasonable and flexible choice for the posterior of generative processes.

4. Probabilistic Contrastive Learning

This section presents our main theoretical result: a probabilistic encoder trained under an MCInfoNCE loss recovers the true posteriors of probabilistic generative processes, up to a rotation, from simple contrastive supervision.

4.1. MCInfoNCE for Probabilistic Contrastive Learning

Let us first formalize the contrastive learning setup. Each training triplet comprises a reference sample x along with a positive (similar) sample x^+ and negative (dissimilar) samples x_1^-, \dots, x_M^- against which it is to be contrasted. As introduced in the previous section, we assume that these samples are generated from corresponding latents $z, z^+, z_1^-, \dots, z_M^-$. Following Zimmermann et al. (2021), the reference z is drawn from the marginal distribution in the latent space, a uniform distribution. The positive sample z^+ is drawn from a close region around z , while negatives z_1^-, \dots, z_M^- are random i.i.d. draws from the marginal:

$$z \sim P(z) = \text{Unif}(z; \mathcal{S}^{D-1}), \quad (2)$$

$$z^+ \sim P(z^+|z) = \text{vMF}(z^+; z, \kappa_{\text{pos}}), \quad (3)$$

$$z_m^- \sim P(z^-|z) =: P(z^-) = \text{Unif}(z^-; \mathcal{S}^{D-1}). \quad (4)$$

The fixed constant $\kappa_{\text{pos}} > 0$ controls how close latents must be to be considered positive to each other¹. The latents $z, z^+, z_1^-, \dots, z_M^-$ are transformed into observations $x, x^+, x_1^-, \dots, x_M^-$ via the generative process $P(x|z)$. This defines $P(x)$, $P(x^+|x)$, and $P(x^-)$.

Our Monte-Carlo InfoNCE (MCInfoNCE) loss is

$$\mathcal{L} := \mathbb{E}_{\substack{x \sim P(x) \\ x^+ \sim P(x^+|x) \\ x_m^- \sim P(x^-), m=1, \dots, M}} (L_f(x, x^+, \{x_m^-\}_{m=1, \dots, M})), \quad \text{with} \quad (5)$$

$$L_f := -\log \mathbb{E}_{\substack{z \sim Q(z|x) \\ z^+ \sim Q(z^+|x^+) \\ z_m^- \sim Q(z_m^-|x_m^-), m=1, \dots, M}} \left(\frac{e^{\kappa_{\text{pos}} z^{\top} z^+}}{\frac{1}{M} e^{\kappa_{\text{pos}} z^{\top} z^+} + \frac{1}{M} \sum_{m=1}^M e^{\kappa_{\text{pos}} z^{\top} z_m^-}} \right). \quad (6)$$

¹ κ_{pos} should not be confused with $\kappa(x)$, which controls the heteroscedastic uncertainty of the generative process.

This is a latent-variable model generalization of the widely used InfoNCE family (Oord et al., 2018), and, in the limit of $M \rightarrow \infty$, SimCLR (Chen et al., 2020). Instead of outputting a point embedding, the encoder f we train outputs probabilistic embeddings $Q(z|x) := \text{vMF}(z; \hat{\mu}(x), \hat{\kappa}(x))$ by predicting $f(x) = (\hat{\mu}(x), \hat{\kappa}(x))$. The InfoNCE fraction within \mathcal{L} is evaluated over these posteriors. In practice, we backpropagate through K Monte-Carlo samples via a reparametrization trick for vMFs (Davidson et al., 2018; Ulrich, 1984):

$$L_f \approx -\log \left(\frac{1}{K} \sum_{k=1}^K \frac{e^{\kappa_{\text{pos}} z_k^{\top} z_k^+}}{\frac{1}{M} e^{\kappa_{\text{pos}} z_k^{\top} z_k^+} + \frac{1}{M} \sum_{m=1}^M e^{\kappa_{\text{pos}} z_k^{\top} z_{m,k}^-}} \right). \quad (7)$$

The only training data for MCInfoNCE are contrastive examples, without any additional supervision on the true aleatoric uncertainty $\kappa(x)$ or the generative latents z .

4.2. Provably Learning the Correct Posteriors

We prove below that the optimizer of this loss learns the *correct* latent posteriors. More precisely, it predicts the correct location $\hat{\mu}(x) = R \cdot \mu(x)$, up to a constant orthogonal rotation R of the latent space, and the correct level of ambiguity $\hat{\kappa}(x) = \kappa(x)$ for each observation x . To prove this, we first show that MCInfoNCE is a cross-entropy over the latent-variable model (Proposition 4.1). This means that the loss matches a predicted marginal value to that of the generative processes. We then show that this marginal is a function that depends only on $(\mu(\cdot)^{\top} \mu(\cdot), \kappa(\cdot))$, resp. $(\hat{\mu}(\cdot)^{\top} \hat{\mu}(\cdot), \hat{\kappa}(\cdot))$ (Proposition 4.2). Due to monotonicity, the predicted function value can only match that of the generative process if their arguments $(\mu(\cdot)^{\top} \mu(\cdot), \kappa(\cdot))$ and $(\hat{\mu}(\cdot)^{\top} \hat{\mu}(\cdot), \hat{\kappa}(\cdot))$ are equal (Proposition 4.3). In summary, the posteriors must be equal, up to a rotation of the latent space (Theorem 4.4).

First, we generalize Zimmermann et al. (2021) and Wang & Isola (2020) to probabilistic generative processes.

Proposition 4.1 (\mathcal{L} is minimized iff marginals match). *Let the latent marginal distributions $P(z) = \int P(z|x) dP(x)$ and $\int Q(z|x) dP(x)$ be uniform. $\lim_{M \rightarrow \infty} \mathcal{L}$ attains its minimum when $\forall x, x^+ \in \{x \in \mathcal{X} | P(x) > 0\}$*

$$\iint Q(z|x) Q(z^+|x^+) P(z^+|z) dz^+ dz = \iint P(z|x) P(z^+|x^+) P(z^+|z) dz^+ dz.$$

This shows that an encoder that minimizes MCInfoNCE has the same marginal as the ground-truth latent model. Next, we need to prove that the equality of these marginals implies that the predicted posteriors Q within them must be equal to

P , up to the rotations mentioned. To this end, we examine the marginals as functions of $\mu(x)$ and $\kappa(x)$.

Proposition 4.2 (The marginal is a function). *Let $P(z|x)$ and $P(z^+|z)$ be vMF distributions as defined in Section 4.1. Given $x, x^+ \in \mathcal{X}$, we can rewrite*

$$\iint P(z|x)P(z^+|x^+)P(z^+|z)dz^+dz \quad (8)$$

$$=: h_{\kappa_{\text{pos}}}(\mu(x)^\top \mu(x^+), \kappa(x), \kappa(x^+)), \quad (9)$$

i.e., as a function $h_{\kappa_{\text{pos}}}$ that depends only on $\mu(x)^\top \mu(x^+)$, $\kappa(x)$, and $\kappa(x^+)$. The same function can be used for $\hat{\mu}(x)^\top \hat{\mu}(x^+)$, $\hat{\kappa}(x)$, $\hat{\kappa}(x^+)$:

$$\iint Q(z|x)Q(z^+|x^+)P(z^+|z)dz^+dz \quad (10)$$

$$= h_{\kappa_{\text{pos}}}(\hat{\mu}(x)^\top \hat{\mu}(x^+), \hat{\kappa}(x), \hat{\kappa}(x^+)). \quad (11)$$

This shows that both sides of the equality in Proposition 4.1 can be expressed by the same function, just with different arguments. All that remains to be shown is that h_{pos} can only attain the same values on both sides of the equality if the arguments match, i.e., $\hat{\kappa}(x) = \kappa(x)$ and $\hat{\mu}(x)^\top \hat{\mu}(x^+) = \mu(x)^\top \mu(x^+)$.

Proposition 4.3 (Arguments of h_{pos} must be equal). *Define h_{pos} as in Proposition 4.2. Let $\mathcal{X}' \subseteq \mathcal{X}$, $\mu, \hat{\mu} : \mathcal{X}' \rightarrow \mathcal{Z}$, $\kappa, \hat{\kappa} : \mathcal{X}' \rightarrow \mathbb{R}_{>0}$, $\kappa_{\text{pos}} > 0$. If $h_{\text{pos}}(\hat{\mu}(x)^\top \hat{\mu}(x^+), \hat{\kappa}(x), \hat{\kappa}(x^+)) = h_{\text{pos}}(\mu(x)^\top \mu(x^+), \kappa(x), \kappa(x^+)) \forall x, x^+ \in \mathcal{X}'$, then*

$$\hat{\mu}(x)^\top \hat{\mu}(x^+) = \mu(x)^\top \mu(x^+) \text{ and} \quad (12)$$

$$\hat{\kappa}(x) = \kappa(x) \quad \forall x, x^+ \in \mathcal{X}'. \quad (13)$$

We can now combine these ingredients to derive our main result: if an encoder minimizes the probabilistic InfoNCE loss, then it must have identified the correct posteriors, up to a constant orthogonal rotation of the latent space.

Theorem 4.4 (\mathcal{L} identifies the correct posteriors). *Let $\mathcal{Z} = \mathcal{S}^{D-1}$ and $P(z) = \int P(z|x)dP(x)$ and $\int Q(z|x)dP(x)$ be the uniform distribution over \mathcal{Z} . Let g be a probabilistic generative process defined in Formulas 2, 3, and 4 with known² κ_{pos} . Let g have vMF posteriors $P(z|x) = \text{vMF}(z; \mu(x), \kappa(x))$ with $\mu : \mathcal{X} \rightarrow \mathcal{S}^{D-1}$ and $\kappa : \mathcal{X} \rightarrow \mathbb{R}_{>0}$. Let an encoder $f(x)$ parametrize vMF distributions $\text{vMF}(z; \hat{\mu}(x), \hat{\kappa}(x))$. Then $f^* = \arg \min_f \lim_{M \rightarrow \infty} \mathcal{L}$ has the correct posteriors up to a rotation of \mathcal{Z} , i.e., $\hat{\mu}(x) = R\mu(x)$ and $\hat{\kappa}(x) = \kappa(x)$, where R is an orthogonal rotation matrix, $\forall x \in \{x \in \mathcal{X} | P(x) > 0\}$.*

This proof generalizes the recent results of Zimmermann et al. (2021) to a broader family of generative processes. MCInfoNCE recovers not only the correct (mean) embeddings $\mu(x)$ under heterogeneous noise and non-injectivity, but also the level of aleatoric uncertainty $\kappa(x)$.

²In practice, κ_{pos} is a tuneable temperature hyperparameter.

5. Experiments

5.1. MCInfoNCE Learns the Correct Posteriors

In this section, we experimentally confirm the theoretical result that *probabilistic embeddings learned under a MCInfoNCE loss recover the correct posteriors up to a rotation*. We also test its robustness to violated assumptions. Code for reproduction is available under https://github.com/mkirchhof/Probabilistic_Contrastive_Learning.

Setup. To test whether MCInfoNCE recovers the correct posteriors, we need a controlled experiment where the true posteriors of the generative process are known. Previous nonlinear ICA experiments randomly initialize a multi-layer perceptron (MLP) as the nonlinear data-generating process and train a second one to invert it (Hyvarinen & Morioka, 2017; Zimmermann et al., 2021). In our probabilistic setup we randomly initialize two MLPs to parameterize $\mu(x)$ and $\kappa(x)$ of the vMF posteriors of the generative process. The MLP for $\mu(x)$ outputs normalized vectors of dimension $D = 10$ and the MLP for $\kappa(x)$ outputs a scalar $\tilde{\kappa}(x)$ wrapped in an exponential Softplus function $\kappa(x) = 1 + \exp(\tilde{\kappa}(x))$ to ensure the strict positivity of $\kappa(x)$ (Li et al., 2021; Shi & Jain, 2019). We sample contrastive training data $(x, x^+, (x_m^-)_{m=1, \dots, M})$ from the generative process parameterized by $\mu(x)$ and $\kappa(x)$ via rejection sampling, as explained in the supplementary. On this data, we train two MLPs to predict $\hat{\mu}(x)$ and $\hat{\kappa}(x)$. All hyperparameters of the generative process and MLP architectures follow the deterministic counterpart of this experiment in Zimmermann et al. (2021) and are reported in the supplementary.

Metrics. To quantify if the predicted posteriors are correct up to a rotation, i.e., $\hat{\kappa}(x) = \kappa(x)$ and $\hat{\mu}(x) = R\mu(x)$ where R is an orthogonal rotation matrix, we compare $\hat{\kappa}(x)$ to $\kappa(x)$ on 10^4 samples of x and compare $\hat{\mu}(x_1)^\top \hat{\mu}(x_2)$ to $\mu(x_1)^\top \mu(x_2)$ on all pairs (x_1, x_2) of the 10^4 samples. We use the root mean square error (RMSE) to test for exact correctness and Spearman’s rank correlation (Rank Corr.) to test for correct ordering. The latter is sufficient in practical scenarios that are invariant to the scale of the latents, such as retrieval based on embedding distances $\hat{\mu}(x_1)^\top \hat{\mu}(x_2)$ or abstention from prediction based on a threshold of the predicted certainty $\hat{\kappa}(x)$. All results report mean and standard error over five randomly generated generative processes.

Results. Table 1 shows that MCInfoNCE recovers the correct posteriors of ambiguous inputs up to a high rank correlation of 0.99 for $\hat{\mu}(x)$ and 0.82 for $\hat{\kappa}(x)$. Figure 2 visualizes this in a simplified case with $D = 2$ latent dimensions. The learned latent space equals the true latent space up to a rotation. However, we can see in Table 1 that $\hat{\kappa}(x)$ tends to get overestimated especially for high values of $\kappa(x) \in [64, 128]$, leading to an RMSE of 125.02. However, the ranking is

Table 1. MCInfoNCE recovers the generative processes’ true posteriors for various degrees of ambiguity and even in the limit of an injective generative process. Mean \pm std. err. for five seeds.

Generative Process Ambiguity	Predicted Location $\hat{\mu}(x)$		Predicted Certainty $\hat{\kappa}(x)$	
	RMSE \downarrow	Rank Corr. \uparrow	RMSE \downarrow	Rank Corr. \uparrow
Ambiguous ($\kappa(x) \in [16, 32]$)	0.04 ± 0.00	0.99 ± 0.00	6.15 ± 0.61	0.82 ± 0.04
Clear ($\kappa(x) \in [64, 128]$)	0.05 ± 0.00	0.98 ± 0.00	125.02 ± 10.64	0.64 ± 0.04
Injective ($\kappa(x) = \infty$)	0.05 ± 0.01	0.98 ± 0.00	$\hat{\kappa}(x) \rightarrow \infty$	

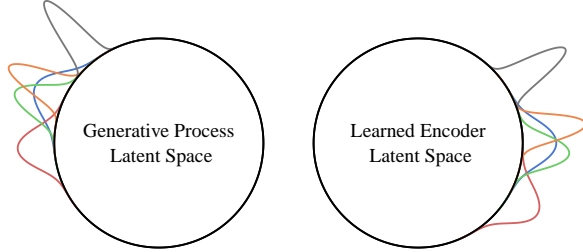


Figure 2. Five posteriors of the generative process and the encoder trained in a run with a 2D latent space. The encoder correctly predicts the posteriors of the generative process, up to a rotation: Rank corr. between $\hat{\mu}(x)$ and the true $\mu(x)$ is 1.00 ± 0.00 (RMSE 0.05 ± 0.00) and that of $\hat{\kappa}(x)$ is 0.82 ± 0.05 (RMSE 2.89 ± 0.56).

still largely preserved with a rank correlation of 0.64. This is because the MC approximation in Formula 7 is a biased estimator of the loss in Formula 6. This is also known as marginal likelihood estimation problem (Perrakis et al., 2014; Burda et al., 2015). This bias decreases with the number of MC samples, as shown in Figure 3. In the standard setup with $\kappa(x) \in [16, 32]$, it is largely mitigated with 16 samples (RMSE = 4.55), or already with 4 samples if only the relative ordering of the samples matters in practice (rank corr. = 0.77). This coincides with the range of number of MC samples used by other probabilistic embedding losses: Oh et al. (2019) use 10 and Kirchoff et al. (2022) use 5. In summary, MCInfoNCE behaves as theoretically expected and fulfills our main theoretical hypothesis.

Violated Assumptions. We test MCInfoNCE in setups where its assumptions are violated. First, we change the posterior distributions of the generative process to Gaussian and Laplace distributions while the encoder is still predicting vMFs. We compare the predicted vMF concentrations $\hat{\kappa}(x)$ to the Gaussian’s inverse variances $1/\sigma^2(x)$ and the Laplace’s inverse diversity $1/b(x)$, so that a higher value in each case means a more concentrated distribution. Table 2 shows that the vMF posteriors model Gaussians almost as well as vMFs (rank corr. 0.78 vs 0.82), since Gaussians with normalized outputs are closely approximated by vMFs (Mardia et al., 2000). For Laplace, the encoder predicts vMF distributions with high concentrations, resulting in an RMSE of 10^7 . This is because the Laplace distribution is more concentrated around its mode than the vMF the encoder uses. Second, we over- and underparameterize the

Table 2. MCInfoNCE predicts sensible vMF posteriors if the true generative posteriors are non-vMF. Mean \pm std. err. for five seeds.

Posterior	Predicted Location $\hat{\mu}(x)$		Predicted Certainty $\hat{\kappa}(x)$	
	RMSE \downarrow	Rank Corr. \uparrow	RMSE \downarrow	Rank Corr. \uparrow
vMF	0.04 ± 0.00	0.99 ± 0.00	6.15 ± 0.61	0.82 ± 0.04
Gaussian	0.04 ± 0.00	0.99 ± 0.00	12.76 ± 0.66	0.78 ± 0.02
Laplace	0.05 ± 0.01	0.98 ± 0.00	$10^7 \pm 10^7$	0.54 ± 0.14

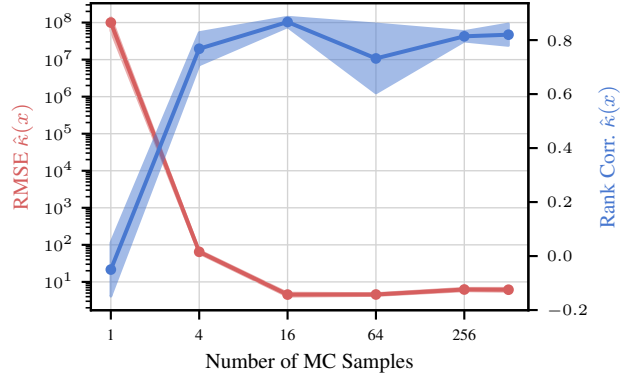


Figure 3. The marginal likelihood approximation bias diminishes with sufficient MC samples. Mean \pm std. err. for five seeds.

latent space dimension of the encoder compared to that of the generative process ($D = 10$). Figure 4 shows that encoder dimensions between 8 and 32 all yield $\hat{\kappa}$ predictions with a rank correlation of at least 0.6. Third, we test the behaviour of MCInfoNCE when the generative process is injective and deterministic, i.e., when all posteriors are Dirac distributions. This is a limiting case of the vMF posteriors the encoder uses for its posteriors. Table 1 shows that the predicted vMFs converge to infinite concentrations $\hat{\kappa}(x)$, recovering the Diracs. In summary, these results indicate that MCInfoNCE serves as a robust starting point when characteristics of the generative process such as its (non-) injectivity, posterior family, or dimension are unknown.

Further losses. Recent literature has proposed other losses to predict probabilistic embeddings. While they have given sensible empirical results, we can use our experimental setup to investigate whether their predictions *exactly* match the true posteriors. We reimplement Hedged Instance Embeddings (HIB) (Oh et al., 2019) and Expected Likelihood Kernels (ELK) (Kirchoff et al., 2022) and modify them to our contrastive setup, as detailed in the supplementary. All losses are hyperparameter-tuned via grid search. Table 3 shows that all losses recover $\mu(x)$ with a rank correlation ≥ 0.82 despite the considerable noise in our experimental setup. We find that, besides MCInfoNCE, ELK also recovers $\kappa(x)$ well (rank corr. = 0.92). This is the first confirmation that ELK predicts correct posteriors in a controlled setup and opens space for future theoretical investigations.

Table 3. Besides MCInfoNCE, ELK also gives correct probabilistic embeddings. Mean \pm std. err. for five seeds.

Loss	Predicted Location $\hat{\mu}(x)$		Predicted Certainty $\hat{\kappa}(x)$	
	RMSE \downarrow	Rank Corr. \uparrow	RMSE \downarrow	Rank Corr. \uparrow
HIB	0.18 ± 0.02	0.82 ± 0.03	$10^{14} \pm 10^{14}$	-0.02 ± 0.09
ELK	0.02 ± 0.00	1.00 ± 0.00	21.70 ± 0.31	0.92 ± 0.00
MCInfoNCE	0.04 ± 0.00	0.99 ± 0.00	6.15 ± 0.61	0.82 ± 0.04

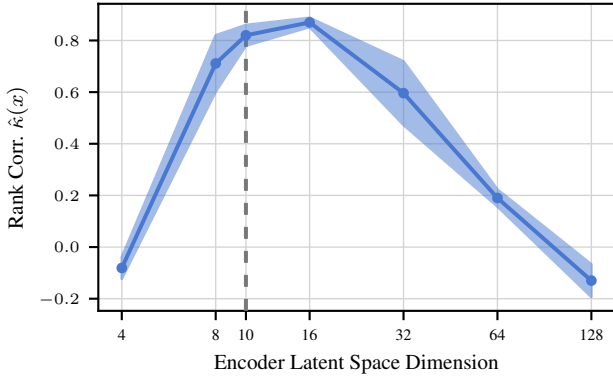


Figure 4. MCInfoNCE learns good $\hat{\kappa}(x)$ even when the encoder latent space dimension mismatches the true generative dimensionality ($D = 10$). Mean \pm std. err. for five seeds.

5.2. Posteriors Reflect Aleatoric Uncertainty in Practice

After confirming that the predicted posteriors are correct, this section shows that they resemble the aleatoric uncertainty in image data. We also show that this enables novel applications such as credible intervals for image retrieval.

Measuring Aleatoric Uncertainty. In the upcoming experiment, we do not have access to any ground-truth $\kappa(x)$ against which to compare $\hat{\kappa}(x)$. Instead, we need to compare it to various indicators of the amount of aleatoric uncertainty in each image. We use three different indicators that capture human uncertainty, information loss, and performance decrease with respect to the amount of aleatoric uncertainty. First, if an image is ambiguous, we can expect human annotators to disagree about the latent that it shows. We therefore conduct our experiment on CIFAR-10H (Peterson et al., 2019). It comprises fifty class annotations for each image. This gives a soft-label distribution whose entropy reflects the ambiguity of the image. We compute the rank correlation between $\hat{\kappa}(x)$ and this annotator entropy to measure how well $\hat{\kappa}(x)$ reflects human-perceived input ambiguity. Second, we can induce controlled information loss by cropping images. Specifically, we crop test images to percentages $\text{crop_size} \sim \text{Unif}([0.25, 1])$ of their original width and height. The aleatoric uncertainty increases the more the image is cropped. We thus report the rank correlation between $\hat{\kappa}(x)$ and the crop size as a second met-

Table 4. Predicted certainties $\hat{\kappa}(x)$ of MCInfoNCE correlate with human annotator disagreement and information reduction via cropping images smaller. Rank correlation on unseen test data.

Loss	Annotator Entropy \uparrow	Crop Size \uparrow
HIB	0.28	0.72
ELK	0.10	0.55
MCInfoNCE	0.33	0.70

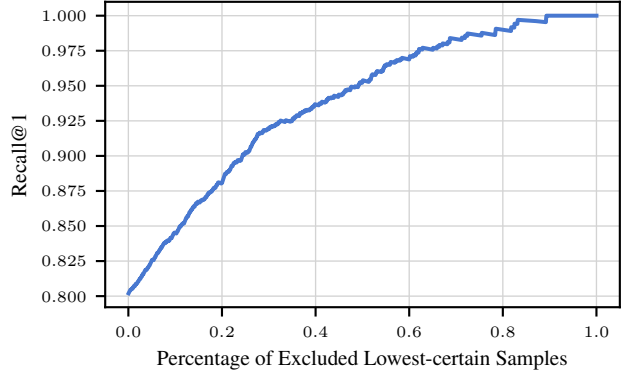


Figure 5. Rejecting images with predicted certainty values $\hat{\kappa}(x)$ below a threshold improves the performance on the remaining data. This shows that $\hat{\kappa}(x)$ is predictive of performance.

ric. Third, ambiguous images inevitably lead to decreased performance. To investigate whether $\hat{\kappa}(x)$ is indicative of performance, we calculate the Recall@1 (Jegou et al., 2010) on the $p\%$ images with the highest $\hat{\kappa}(x)$. If $\hat{\kappa}(x)$ correctly reflects aleatoric uncertainty, removing ambiguous images should improve performance, so the Recall@1 should increase monotonically with p . This metric also illustrates the popular use case of abstaining from prediction on images that are uncertain and likely to be misclassified.

Architecture and Training. We need to translate the CIFAR-10H classification task into a contrastive task. To do this, we consider images to be positive if they are in the same class and negative otherwise. We create training examples $(x, x^+, x_1^-, \dots, x_M^-)$ by drawing class labels for each image from its soft class distribution, selecting a random image x , an image x^+ with the same class label, and M images x_m^- with different class labels. On this contrastive data, we train a ResNet-18 (He et al., 2016) pre-trained on CIFAR-10 (Phan, 2021). $\hat{\mu}(x)$ is the normalized output of a final linear layer. Following recent common practice for probabilistic embeddings (Kirchhof et al., 2022; Scott et al., 2021; Li et al., 2021), we parameterize $\hat{\kappa}(x)$ by the norm of the output. We randomly split CIFAR-10H into train, validation, and test sets. We train the model for 175 epochs on the train data and select the best epoch by calculating the rank correlation with the crop size on the validation data. We choose this metric over annotator entropy because it can

Table 5. $\hat{\kappa}(x)$ can be learned by MCInfoNCE from both soft and hard labels. Rank correlation on unseen test data.

Labels	Annotator Entropy \uparrow	Crop Size \uparrow
CIFAR-10H Soft Labels	0.33	0.70
CIFAR-10H Hard Labels	0.26	0.71
CIFAR-10 Hard Labels	0.32	0.72

be computed synthetically on any dataset without additional supervision. All details on generating the contrastive data and the hyperparameter search are in the supplementary.

Results. Table 4 shows that $\hat{\kappa}(x)$ learned via MCInfoNCE has a high rank correlation of 0.72 with the amount of information lost due to cropping, i.e., images with less information return more uncertain posteriors. The correlation with the human annotator entropy is lower (0.33), but positive. HIB achieves a similar performance, while ELK shows lower correlations with both ground-truths (0.55 and 0.1, resp.). Figure 5 provides the performance decrease metric. Up to noise, the Recall@1 increases monotonically as images with the lowest $\hat{\kappa}(x)$ are rejected. This means that $\hat{\kappa}(x)$ is a good predictor of performance. As an additional qualitative metric and sanity check, the supplementary provides the images with the lowest and highest $\hat{\kappa}(x)$ of each class. MCInfoNCE learns from labeling noise in this experiment, since the image class was drawn anew from its soft label distribution each time the image was used. In practice, we may have only one annotation per image, so that labeling noise occurs across examples rather than on each individual image. To this end, we further train on hard labels. These are either the most likely class of each soft label distribution on CIFAR-10H or the classical class labels on the CIFAR-10. Table 5 shows that MCInfoNCE can learn under both of these circumstances with a performance roughly equal to that when soft labels are available.

Credible Intervals for Image Retrieval. Since we are estimating posteriors $Q(z|x)$, we can also introduce credible intervals (Lee, 1989) from Bayesian statistics to our image representation task. Such an interval $\text{CI}_p(x) \subset \mathcal{Z}$ contains the true generative latent z of x with a user-defined probability $p \in [0, 1]$, i.e., $P(z \in \text{CI}_p(x)) = p$ for $x \sim P(x|z)$. Note that these intervals depend on both the aleatoric uncertainty of the input and the quality of the model’s posterior estimates. Credible intervals help understand the degree to which our model can identify the latent that x shows. We can visualize these latents by searching for images whose posteriors’ modal value $\hat{\mu}(x)$ fall within CI_p . Figure 6 shows such intervals on our MCInfoNCE model for CIFAR-10H. A clear image (top) has a sharp posterior and thus a small credible interval containing only one image from the same class. A more ambiguous query image, like the third one, leads to more possible matches. It tells us that the model places the image in the region of cats, but that it could



Figure 6. We can an image’s posterior to define the credible interval that its latents lie in with a given probability. Clear query images (top) have small credible intervals containing images of the same class as the query. More ambiguous queries (bottom) return larger credible intervals with images from multiple possible classes.

also be a dog. Highly ambiguous queries, like the last one, lead to wide credible intervals that span multiple possible classes. Credible intervals augment image retrieval with uncertainty-awareness by providing the number of possible similar images to display, subject to ambiguity of the query.

6. Conclusion

This work presented MCInfoNCE, a probabilistic contrastive loss that predicts posteriors instead of points. We proved that it learns the generative processes’ true posteriors. This provides a theoretical grounding for the recent probabilistic embeddings literature and connects it to a probabilistic extension of nonlinear ICA. In practice, the posteriors allow predicting the level of aleatoric uncertainty in ambiguous inputs as well as estimating credible intervals with flexible sizes depending on a query’s ambiguity in image retrieval. These are only two usages that correct posteriors enable and further usages are a promising area for future research. Aleatoric uncertainty is not only faced in computer vision and retrieval. We hope that the blueprint way of enhancing InfoNCE into MCInfoNCE inspires applications in further tasks with intrinsic ambiguities in their inputs.

Acknowledgements

Kay Choi has helped designing Figure 1. This work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy – EXC number 2064/1 – Project number 390727645. The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting Michael Kirchhof.

References

- Burda, Y., Grosse, R., and Salakhutdinov, R. Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*, 2015.
- Chen, H., Huang, Y., Tian, W., Gao, Z., and Xiong, L. Monorun: Monocular 3d object detection by reconstruction and uncertainty propagation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10379–10388, 2021.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, 2020.
- Chun, S., Oh, S. J., De Rezende, R. S., Kalantidis, Y., and Larlus, D. Probabilistic embeddings for cross-modal retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- Chun, S., Kim, W., Park, S., Chang, M. C., and Oh, S. J. ECCV caption: Correcting false negatives by collecting machine-and-human-verified image-caption associations for ms-coco. In *European Conference on Computer Vision (ECCV)*, 2022.
- Comon, P. and Jutten, C. *Handbook of Blind Source Separation: Independent component analysis and applications*. 2010.
- Davidson, T. R., Falorsi, L., De Cao, N., Kipf, T., and Tomczak, J. M. Hyperspherical variational auto-encoders. *34th Conference on Uncertainty in Artificial Intelligence (UAI)*, 2018.
- Deng, J., Guo, J., Xue, N., and Zafeiriou, S. ArcFace: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, pp. 4690–4699, 2019.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, pp. 770–778, 2016.
- Hyvarinen, A. and Morioka, H. Nonlinear ica of temporally dependent stationary sources. In *Artificial Intelligence and Statistics (AISTATS)*, 2017.
- Hyvärinen, A. and Oja, E. Independent component analysis: algorithms and applications. *Neural Networks*, 13(4-5): 411–430, 2000.
- Islam, A., Chen, C.-F. R., Panda, R., Karlinsky, L., Radke, R., and Feris, R. A broad study on the transferability of visual representations with contrastive learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 8845–8855, 2021.
- Jebara, T. and Kondor, R. Bhattacharyya and expected likelihood kernels. In *Learning Theory and Kernel Machines*. 2003.
- Jegou, H., Douze, M., and Schmid, C. Product quantization for nearest neighbor search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(1):117–128, 2010.
- Karpukhin, I., Dereka, S., and Kolesnikov, S. Probabilistic embeddings revisited. *arXiv preprint arXiv:2202.06768*, 2022.
- Kendall, A. and Gal, Y. What uncertainties do we need in bayesian deep learning for computer vision? *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017.
- Khosla, P., Teterwak, P., Wang, C., Sarna, A., Tian, Y., Isola, P., Maschinot, A., Liu, C., and Krishnan, D. Supervised contrastive learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:18661–18673, 2020.
- Kirchhof, M., Roth, K., Akata, Z., and Kasneci, E. A non-isotropic probabilistic take on proxy-based deep metric learning. In *European Conference on Computer Vision (ECCV)*, 2022.
- Kraus, F. and Dietmayer, K. Uncertainty estimation in one-stage object detection. In *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pp. 53–60, 2019.
- Lee, P. M. *Bayesian statistics*. Oxford University Press London, 1989.
- Leemann, T., Kirchhof, M., Rong, Y., Kasneci, E., and Kasneci, G. Disentangling embedding spaces with minimal distributional assumptions. *arXiv preprint arXiv:2206.13872*, 2022.
- Li, S., Xu, J., Xu, X., Shen, P., Li, S., and Hooi, B. Spherical confidence learning for face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

- Mardia, K. V., Jupp, P. E., and Mardia, K. *Directional statistics*, volume 2. Wiley Online Library, 2000.
- Meech, J. T. and Stanley-Marbell, P. An algorithm for sensor data uncertainty quantification. *IEEE Sensors Letters*, 6(1):1–4, 2021.
- Mehrtens, H. A., Kurz, A., Bucher, T.-C., and Brinker, T. J. Benchmarking common uncertainty estimation methods with histopathological images under domain shift and label noise. *arXiv preprint arXiv:2301.01054*, 2023.
- Neculai, A., Chen, Y., and Akata, Z. Probabilistic compositional embeddings for multimodal image retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition MULA Workshop (CVPR MULA)*, pp. 4547–4557, 2022.
- Oh, S. J., Gallagher, A. C., Murphy, K. P., Schroff, F., Pan, J., and Roth, J. Modeling uncertainty with hedged instance embeddings. In *International Conference on Learning Representations (ICLR)*, 2019.
- Oord, A. v. d., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Perrakis, K., Ntzoufras, I., and Tsionas, E. G. On the use of marginal posteriors in marginal likelihood estimation via importance sampling. *Computational Statistics & Data Analysis*, 77:54–69, 2014.
- Peterson, J. C., Battleday, R. M., Griffiths, T. L., and Rusakovsky, O. Human uncertainty makes classification more robust. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (CVPR)*, pp. 9617–9626, 2019.
- Phan, H. Pytorch CIFAR-10 v3.0.1, 2021. URL <https://doi.org/10.5281/zenodo.4431043>.
- Romanazzi, M. Discriminant analysis with high dimensional von mises-fisher distributions. In *8th Annual International Conference on Statistics*, 2014.
- Roth, K., Ibrahim, M., Akata, Z., Vincent, P., and Bouchacourt, D. Disentanglement of correlated factors via hausdorff factorized support. *arXiv preprint arXiv:2210.07347*, 2022.
- Schlett, T., Rathgeb, C., Henniger, O., Galbally, J., Fierrez, J., and Busch, C. Face image quality assessment: A literature survey. *ACM Computing Surveys (CSUR)*, 54(10s):1–49, 2022.
- Schmarje, L., Grossmann, V., Zelenka, C., Dippel, S., Kiko, R., Oszust, M., Pastell, M., Stracke, J., Valros, A., Volkman, N., et al. Is one annotation enough? A data-centric image classification benchmark for noisy and ambiguous label estimation. *arXiv preprint arXiv:2207.06214*, 2022.
- Scott, T. R., Gallagher, A. C., and Mozer, M. C. von Mises-Fisher loss: An exploration of embedding geometries for supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- Shi, Y. and Jain, A. K. Probabilistic face embeddings. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- Teh, E. W., DeVries, T., and Taylor, G. W. ProxyNCA++: Revisiting and revitalizing proxy neighborhood component analysis. In *European Conference on Computer Vision (ECCV)*, pp. 448–464, 2020.
- Ulrich, G. Computer generation of distributions on the m-sphere. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 33(2):158–163, 1984.
- Wang, T. and Isola, P. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, 2020.
- Wang, Y., Tang, S., Zhu, F., Bai, L., Zhao, R., Qi, D., and Ouyang, W. Revisiting the transferability of supervised pretraining: an mlp perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9183–9193, 2022.
- Zimmermann, R. S., Sharma, Y., Schneider, S., Bethge, M., and Brendel, W. Contrastive learning inverts the data generating process. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, 2021.

A. Proofs

A.1. Proof of Proposition 4.1

Proposition 4.1 (\mathcal{L} is minimized iff marginals match) Let the latent marginal distributions $P(z) = \int P(z|x)dP(x)$ and $\int Q(z|x)dP(x)$ be uniform. $\lim_{M \rightarrow \infty} \mathcal{L}$ attains its minimum when $\forall x, x^+ \in \{x \in \mathcal{X} | P(x) > 0\}$

$$\begin{aligned} \iint Q(z|x)Q(z^+|x^+)P(z^+|z)dz^+dz = \\ \iint P(z|x)P(z^+|x^+)P(z^+|z)dz^+dz . \end{aligned}$$

Proof. All of the above densities are integrable, so we can write the loss function \mathcal{L} in the form of Riemann integrals.

$$\lim_{M \rightarrow \infty} \mathcal{L} = - \lim_{M \rightarrow \infty} \int P(x)P(x^+|x) \int \prod_{m=1}^M P(x_m^-) \log \int Q(z|x)Q(z^+|x^+) \quad (14)$$

$$\prod_{m=1}^M Q(z_m^-|x_m^-) \frac{e^{\kappa_{\text{pos}} z^\top z^+}}{\frac{1}{M} e^{\kappa_{\text{pos}} z^\top z^+} + \frac{1}{M} \sum_{m=1}^M e^{\kappa_{\text{pos}} z^\top z_m^-}} dz_1^- \dots z_M^- dz^+ dz dx_1^- \dots dx_M^- dx^+ dx \quad (15)$$

We know that $\kappa_{\text{pos}} < \infty$, $\kappa(x) < \infty \forall x \in \mathcal{X}$, the normalization constants $C(\kappa) < \infty \forall \kappa < \infty$, and the dot products are bounded. This implies that all densities inside these integrals as well as the exponentials in the fraction are bounded. Thus, the whole term inside the outmost integral is bounded. Due to the dominated convergence theorem we can pull the limit into the integral.

$$= - \int P(x)P(x^+|x) \lim_{M \rightarrow \infty} \int \prod_{m=1}^M P(x_m^-) \log \int Q(z|x)Q(z^+|x^+) \quad (16)$$

$$\prod_{m=1}^M Q(z_m^-|x_m^-) \frac{e^{\kappa_{\text{pos}} z^\top z^+}}{\frac{1}{M} e^{\kappa_{\text{pos}} z^\top z^+} + \frac{1}{M} \sum_{m=1}^M e^{\kappa_{\text{pos}} z^\top z_m^-}} dz_1^- \dots z_M^- dz^+ dz dx_1^- \dots dx_M^- dx^+ dx \quad (17)$$

The strong law of large numbers and the fact that $\int Q(z^-|x^-)P(x^-)dx^- = P(z)$ imply

$$= - \int P(x)P(x^+|x) \lim_{M \rightarrow \infty} \log \int Q(z|x)Q(z^+|x^+) \frac{e^{\kappa_{\text{pos}} z^\top z^+}}{\frac{1}{M} e^{\kappa_{\text{pos}} z^\top z^+} + \mathbb{E}_{z^- \sim P(z)}(e^{\kappa_{\text{pos}} z^\top z^-})} dz^+ dz dx^+ dx . \quad (18)$$

Both densities and the fraction inside the inner integral are positive and bounded, so the integral is, too. In this range, i.e., $(0, \infty)$, the logarithm is continuous, so the continuous mapping theorem gives

$$= - \int P(x)P(x^+|x) \log \lim_{M \rightarrow \infty} \int Q(z|x)Q(z^+|x^+) \frac{e^{\kappa_{\text{pos}} z^\top z^+}}{\frac{1}{M} e^{\kappa_{\text{pos}} z^\top z^+} + \mathbb{E}_{z^- \sim P(z)}(e^{\kappa_{\text{pos}} z^\top z^-})} dz^+ dz dx^+ dx . \quad (19)$$

With the arguments from above, the inside of the inner integral is bounded, so we can again apply the dominated convergence theorem.

$$= - \int P(x)P(x^+|x) \log \int Q(z|x)Q(z^+|x^+) \lim_{M \rightarrow \infty} \frac{e^{\kappa_{\text{pos}} z^\top z^+}}{\frac{1}{M} e^{\kappa_{\text{pos}} z^\top z^+} + \mathbb{E}_{z^- \sim P(z)}(e^{\kappa_{\text{pos}} z^\top z^-})} dz^+ dz dx^+ dx \quad (20)$$

$$= - \int P(x)P(x^+|x) \log \int Q(z|x)Q(z^+|x^+) \frac{e^{\kappa_{\text{pos}} z^\top z^+}}{\mathbb{E}_{z^- \sim P(z)}(e^{\kappa_{\text{pos}} z^\top z^-})} dz^+ dz dx^+ dx \quad (21)$$

Since $P(z) = \text{Unif}(S^{D-1}) = \frac{1}{\|S^{D-1}\|}$, which we define as $\frac{1}{S}$ in shorthand, we get

$$= - \int P(x)P(x^+|x) \log S \int Q(z|x)Q(z^+|x^+) \frac{e^{\kappa_{\text{pos}} z^\top z^+}}{\int_{S^{D-1}} e^{\kappa_{\text{pos}} z^\top z^-} dz^-} dz^+ dz dx^+ dx \quad (22)$$

$$= - \int P(x)P(x^+|x) \log S \int Q(z|x)Q(z^+|x^+)P(z^+|z) dz^+ dz dx^+ dx . \quad (23)$$

Let us turn our attention to $P(x^+|x)$. By marginalization, factorization, and the conditional independencies of the data-generating process, we get

$$P(x^+|x) \quad (24)$$

$$= \int P(x^+, z^+, z|x) dz^+ dz \quad (25)$$

$$= \int P(x^+|z^+, z, x)P(z^+|z, x)P(z|x) dz^+ dz \quad (26)$$

$$= \int P(x^+|z^+)P(z^+|z)P(z|x) dz^+ dz . \quad (27)$$

After a multiplication with 1, Bayes Theorem, and using $P(z) = \frac{1}{S}$, we get

$$= \int \frac{P(x^+|z^+)P(z^+)P(x^+)}{P(z^+)P(x^+)} P(z^+|z)P(z|x) dz^+ dz \quad (28)$$

$$= \int P(z|x)P(z^+|x^+)P(z^+|z) \frac{P(x^+)}{P(z^+)} dz^+ dz \quad (29)$$

$$= P(x^+) S \int P(z|x)P(z^+|x^+)P(z^+|z) dz^+ dz . \quad (30)$$

We can insert this into Formula 23.

$$- \int P(x)P(x^+) S \int P(z|x)P(z^+|x^+)P(z^+|z) dz^+ dz \quad (31)$$

$$\log S \int Q(z|x)Q(z^+|x^+)P(z^+|z) dz^+ dz dx^+ dx \quad (32)$$

$$= \mathbb{E}_{\substack{x \sim P(x) \\ x^+ \sim P(x^+)}} \left(S \int P(z|x)P(z^+|x^+)P(z^+|z) dz^+ dz \log S \int Q(z|x)Q(z^+|x^+)P(z^+|z) dz^+ dz \right) . \quad (33)$$

Note that both terms are conditional on x, x^+ and the expected value is taken over both of these. I.e., \mathcal{L} in the limit is a (non-normalized) cross-entropy between $\int P(z|x)P(z^+|x^+)P(z^+|z) dz^+ dz$ and $\int Q(z|x)Q(z^+|x^+)P(z^+|z) dz^+ dz$. The loss is minimized iff the two terms match for all values in the outmost expected value, i.e., $\forall x, x^+ \in \{x \in \mathcal{X} | P(x) > 0\}$. \square

A.2. Proof of Proposition 4.2

Proposition 4.2 (The marginal is a function) Let $P(z|x)$ and $P(z^+|z)$ be vMF distributions as defined in Section 4.1. Given $x, x^+ \in \mathcal{X}$, we can rewrite

$$\iint P(z|x)P(z^+|x^+)P(z^+|z) dz^+ dz \quad (34)$$

$$=: h_{\kappa_{\text{pos}}}(\mu(x)^\top \mu(x^+), \kappa(x), \kappa(x^+)), \quad (35)$$

i.e., as a function $h_{\kappa_{\text{pos}}}$ that depends only on $\mu(x)^\top \mu(x^+)$, $\kappa(x)$, and $\kappa(x^+)$. The same function can be used for $\hat{\mu}(x)^\top \hat{\mu}(x^+), \hat{\kappa}(x), \hat{\kappa}(x^+)$:

$$\iint Q(z|x)Q(z^+|x^+)P(z^+|z) dz^+ dz \quad (36)$$

$$= h_{\kappa_{\text{pos}}}(\hat{\mu}(x)^\top \hat{\mu}(x^+), \hat{\kappa}(x), \hat{\kappa}(x^+)). \quad (37)$$

Proof. Let us first insert the vMF densities.

$$\iint P(z|x)P(z^+|x^+)P(z^+|z)dz^+dz \quad (38)$$

$$=C(\kappa(x^+))C(\kappa_{\text{pos}}) \iint C(\kappa(x)) \exp[\kappa(x)\mu(x)^\top z + \kappa(x^+)\mu(x^+)^\top z^+ + \kappa_{\text{pos}}z^\top z^+]dz^+dz \quad (39)$$

$$=C(\kappa(x^+))C(\kappa_{\text{pos}}) \int C(\kappa(x)) \exp(\kappa(x)\mu(x)^\top z) \int \exp[(\kappa(x^+)\mu(x^+) + \kappa_{\text{pos}}z)^\top z^+]dz^+dz \quad (40)$$

The term inside the inner integral can be rewritten into an unnormalized vMF density if we specify $\mu^* := \frac{\kappa(x^+)\mu(x^+) + \kappa_{\text{pos}}z}{\|\kappa(x^+)\mu(x^+) + \kappa_{\text{pos}}z\|}$ and $\kappa^* := \|\kappa(x^+)\mu(x^+) + \kappa_{\text{pos}}z\|$. The integral over this density is 1.

$$=C(\kappa(x^+))C(\kappa_{\text{pos}}) \int C(\kappa(x)) \exp(\kappa(x)\mu(x)^\top z) \frac{1}{C(\kappa^*)} \int C(\kappa^*) \exp[\kappa^* \mu^{*\top} z^+]dz^+dz \quad (41)$$

$$=C(\kappa(x^+))C(\kappa_{\text{pos}}) \int C(\kappa(x)) \exp(\kappa(x)\mu(x)^\top z) \frac{1}{C(\kappa^*)} dz \quad (42)$$

$$=C(\kappa_{\text{pos}}) \mathbb{E}_{z \sim \text{vMF}(\mu(x), \kappa(x))} \left(\frac{C(\kappa(x^+))}{C(\sqrt{\kappa(x^+)^2 + \kappa_{\text{pos}}^2 + 2\kappa(x^+)\kappa_{\text{pos}}\mu(x^+)^\top z})} \right) \quad (43)$$

$$=: h_{\kappa_{\text{pos}}}(\mu(x)^\top \mu(x^+), \kappa(x), \kappa(x^+)) \quad (44)$$

In the last step, the expected value is over $\mu(x^+)^\top z, z \sim \text{vMF}(\mu(x), \kappa(x))$. This depends only on the distance $\mu(x)^\top \mu(x^+)$ instead of the full location parameters $\mu(x)$ and $\mu(x^+)$ because the vMF is rotationally symmetric and we can perform a suitable Householder rotation, see also [Romanazzi \(2014\)](#). \square

A.3. Proof of Proposition 4.3

Proposition 4.3 (Arguments of h_{pos} must be equal) Define h_{pos} as in Proposition 4.2. Let $\mathcal{X}' \subseteq \mathcal{X}$, $\mu, \hat{\mu} : \mathcal{X}' \rightarrow \mathcal{Z}$, $\kappa, \hat{\kappa} : \mathcal{X}' \rightarrow \mathbb{R}_{>0}$, $\kappa_{\text{pos}} > 0$. If $h_{\text{pos}}(\hat{\mu}(x)^\top \hat{\mu}(x^+), \hat{\kappa}(x), \hat{\kappa}(x^+)) = h_{\text{pos}}(\mu(x)^\top \mu(x^+), \kappa(x), \kappa(x^+)) \forall x, x^+ \in \mathcal{X}'$, then

$$\hat{\mu}(x)^\top \hat{\mu}(x^+) = \mu(x)^\top \mu(x^+) \text{ and} \quad (45)$$

$$\hat{\kappa}(x) = \kappa(x) \quad \forall x, x^+ \in \mathcal{X}'. \quad (46)$$

Proof. (a) The normalization constant of the vMF $C(\kappa) = \frac{\kappa^{D/2-1}}{(2\pi)^{D/2} I_{D/2-1}(\kappa)}$, where I_o is the modified Bessel function of the first kind and order o , is strictly monotonically decreasing and convex ([Kirchhof et al., 2022](#)).

(b) Consider arbitrary $x = x^+, x \in \mathcal{X}'$. In this case, $\mu(x)^\top \mu(x^+) = \hat{\mu}(x)^\top \hat{\mu}(x^+) = 1$, and both sides of the equality simplify

$$\iint Q(z|x)Q(z^+|x^+)P(z^+|z)dz^+dz = \iint P(z|x)P(z^+|x^+)P(z^+|z)dz^+dz \quad (47)$$

$$\iff h_{\kappa_{\text{pos}}}(1, \kappa(x), \kappa(x)) = h_{\kappa_{\text{pos}}}(1, \hat{\kappa}(x), \hat{\kappa}(x)) \quad (48)$$

$$\iff \tilde{h}_{\kappa_{\text{pos}}}(\kappa(x)) = \tilde{h}_{\kappa_{\text{pos}}}(\hat{\kappa}(x)) \quad (49)$$

with $\tilde{h}_{\kappa_{\text{pos}}}(\kappa) := h_{\kappa_{\text{pos}}}(1, \kappa, \kappa)$. Due to (a), the denominator in Formula 43 grows strictly faster than the numerator. So \tilde{h} is strictly monotonically increasing. Thus, $\tilde{h}_{\kappa_{\text{pos}}}(\kappa(x)) = \tilde{h}_{\kappa_{\text{pos}}}(\hat{\kappa}(x))$ only if $\kappa(x) = \hat{\kappa}(x)$.

(c) Let $x, x^+ \in \mathcal{X}'$ be arbitrary. From (b) we know $\hat{\kappa}(x) = \kappa(x)$, so we can simplify

$$h_{\kappa_{\text{pos}}}(\mu(x)^\top \mu(x^+), \kappa(x), \kappa(x^+)) = h_{\kappa_{\text{pos}}}(\hat{\mu}(x)^\top \hat{\mu}(x^+), \hat{\kappa}(x), \hat{\kappa}(x^+)) \quad (50)$$

$$\iff h_{\kappa_{\text{pos}}, \kappa(x), \kappa(x^+)}^*(\mu(x)^\top \mu(x^+)) = h_{\kappa_{\text{pos}}, \kappa(x), \kappa(x^+)}^*(\hat{\mu}(x)^\top \hat{\mu}(x^+)) \quad (51)$$

with $h_{\kappa_{\text{pos}}, \kappa(x), \kappa(x^+)}^*(\cdot) := h_{\kappa_{\text{pos}}}(\cdot, \kappa(x), \kappa(x^+))$. In other words, both sides of the equality are the same function $h_{\kappa_{\text{pos}}, \kappa(x), \kappa(x^+)}^*$ with only one free variable. Due to (a), the denominator in Formula 43 strictly decreases with increasing $\mu(x)^\top \mu(x^+)$ if $\kappa(x^+) > 0$ and $\kappa_{\text{pos}} > 0$. So, $h_{\kappa_{\text{pos}}, \kappa(x), \kappa(x^+)}^*$ is strictly monotonically increasing and $h_{\kappa_{\text{pos}}, \kappa(x), \kappa(x^+)}^*(\mu(x)^\top \mu(x^+)) = h_{\kappa_{\text{pos}}, \kappa(x), \kappa(x^+)}^*(\hat{\mu}(x)^\top \hat{\mu}(x^+))$ implies $\mu(x)^\top \mu(x^+) = \hat{\mu}(x)^\top \hat{\mu}(x^+)$. \square

A.4. Proof of Theorem 4.4

Theorem 4.4 (\mathcal{L} identifies the correct posteriors) Let $\mathcal{Z} = \mathcal{S}^{D-1}$ and $P(z) = \int P(z|x)dP(x)$ and $\int Q(z|x)dP(x)$ be the uniform distribution over \mathcal{Z} . Let g be a probabilistic generative process defined in Formulas 2, 3, and 4 with known κ_{pos} . Let g have vMF posteriors $P(z|x) = \text{vMF}(z; \mu(x), \kappa(x))$ with $\mu : \mathcal{X} \rightarrow \mathcal{S}^{D-1}$ and $\kappa : \mathcal{X} \rightarrow \mathbb{R}_{>0}$. Let an encoder $f(x)$ parametrize vMF distributions $\text{vMF}(z; \hat{\mu}(x), \hat{\kappa}(x))$. Then $f^* = \arg \min_f \lim_{M \rightarrow \infty} \mathcal{L}$ has the correct posteriors up to a rotation of \mathcal{Z} , i.e., $\hat{\mu}(x) = R\mu(x)$ and $\hat{\kappa}(x) = \kappa(x)$, where R is an orthogonal rotation matrix, $\forall x \in \{x \in \mathcal{X} | P(x) > 0\}$.

Proof. If f^* optimizes \mathcal{L} , then by Proposition 4.1 $\forall x, x^+ \in \{x \in \mathcal{X} | P(x) > 0\}$ we have

$$\iint Q(z|x)Q(z^+|x^+)P(z^+|z)dz^+dz = \iint P(z|x)P(z^+|x^+)P(z^+|z)dz^+dz. \quad (52)$$

Then by Proposition 4.3 with $\mathcal{X}' := \{x \in \mathcal{X} | P(x) > 0\}$ we get $\hat{\kappa}(x) = \kappa(x)$ and $\mu(x)^\top \mu(x^+) = \hat{\mu}(x)^\top \hat{\mu}(x^+)$. With the extended Mazur-Ulam Theorem (Zimmermann et al., 2021), the latter implies $\hat{\mu}(x) = R\mu(x)$ with an orthogonal rotation matrix $R \in \mathbb{R}^{D \times D}$. \square

B. Controlled Experiment

B.1. Network Architectures

We use MLPs to parametrize the generative processes' posteriors $\mu(x)$ and $\kappa(x)$ as well as the encoder $\hat{\mu}(x)$ and $\hat{\kappa}(x)$.

For $\mu(x)$ and $\hat{\mu}(x)$ we follow Zimmermann et al. (2021). The MLP for $\mu(x)$ has three linear layers with 10 dimensions and leaky ReLU activations. To prevent collapsed initializations we take 1000 exemplary samples for $\mu(x)$ and re-initiate it if the smallest cosine similarity $x_1^\top x_2$ between any pair x_1, x_2 of them is bigger than 0.5. $\hat{\mu}(x)$ has six hidden linear layers with leaky ReLU activations plus an input and an output layer with the input and output dimensions $[D \rightarrow 10 \cdot D, 10 \cdot D \rightarrow 50 \cdot D, 50 \cdot D \rightarrow 50 \cdot D, 50 \cdot D \rightarrow 50 \cdot D, 50 \cdot D \rightarrow 50 \cdot D, 50 \cdot D \rightarrow 10 \cdot D, 10 \cdot D \rightarrow D]$. The outputs of both networks are normalized to an L_2 norm of 1 to ensure they are on the unit sphere.

The MLPs for $\kappa(x)$ and $\hat{\kappa}(x)$ have the same architecture as $\mu(x)$ and $\hat{\mu}(x)$, but $\kappa(x)$ has one less hidden layer than $\mu(x)$. The last layer of both networks outputs only a scalar instead of a D -dimensional vector. It is postprocessed by $\tilde{\kappa}(x) = 1 + \exp(\kappa(x))$ to ensure their strict positivity. Before training, $\hat{\kappa}(x)$ is normalized to output the same range of values as $\kappa(x)$ to improve training stability.

B.2. Generating Contrastive Training Data

The generative process in Section 4.1 first draws latents z and then generates observations x to create contrastive training data. However, we want to control our generative processes' posteriors. Thus, we need to first sample x and then $z \sim P(z|x)$. A method to sample backwards like this while still obtaining samples as if they were from the forward generative process is rejection sampling. We first draw random candidates (x, x^+) from $\mathcal{X} = [0, 1]^D$, then draw (z, z^+) from their corresponding posteriors. To ensure that they form a valid positive example as per the distributions in Formulas 2 and 3, we accept or reject them with a probability proportional to

$$\frac{C(\kappa_{\text{pos}})e^{\kappa_{\text{pos}}z^\top z^+}}{C(\kappa_{\text{pos}})e^{\kappa_{\text{pos}}z^\top z^+} + C(0)}. \quad (53)$$

This is the probability that z and z^+ are positive to one another. The proposal distribution's density for rejection sampling is dropped here due to the uniform priors. Negative examples $(x_m^-)_{m=1, \dots, M}$ are drawn randomly from \mathcal{X} due to Formula 4.

B.3. Experiment Parameters

Following Zimmermann et al. (2021), all experiments used $\kappa_{\text{pos}} = 20$ and the above network architectures. The learning rate was 0.0001 and was decreased after each 25% of training progress by a factor of 0.1. Performance was measured at the end of the training without early stopping on 10000 sampled x points. All experiments were implemented in Python 3.8.11, PyTorch 1.9.0 on NVIDIA-RTX 2080TI GPUs with 12GB VRAM. Table 6 below summarizes the remaining parameters used by all ablations of the controlled experiment.

Probabilistic Contrastive Learning Recovers the Correct Aleatoric Uncertainty

Experiment	Gen. D	Enc. D'	Posterior	$\min(\kappa(x))$	$\max(\kappa(x))$	Batchsize	Number of Batches	Number MC Samples	Comment
Ambiguous ($\kappa(x) \in [16, 32]$)	10	10	vMF	16	32	512	100000	512	Also used for HIB, ELK, InfoNCE
Clear ($\kappa(x) \in [64, 128]$)	10	10	vMF	64	128	512	100000	512	
Injective ($\kappa(x) = \infty$)	10	10	vMF/Dirac	∞	∞	512	100000	512	
$D = 2$	2	2	vMF	16	32	512	8192	512	
Gaussian	10	10	Gaussian	16	32	512	100000	512	$\sigma^2 = 1/\kappa(x)$
Laplace	10	10	Laplace	16	32	512	100000	512	$b = 1/\kappa(x)$
MC Samples	10	10	vMF	16	32	512	100000	x	$x \in \{1, 4, 16, 64, 256, 512\}$
Encoder Dim	10	x	vMF	16	32	512	100000	512	for $x \in \{4, 8, 10, 16, 32\}$
— " —						512		256	for $x = 64$
— " —						256		256	for $x = 128$
High Dim	x	x	vMF	16	32	512	100000	512	$x \in \{10, 16\}$
— " —						256		256	for $x \in \{32, 40, 48, 56, 64\}$

Table 6. Parameters of the generative process and loss in the controlled experiments. x denotes variable parameters. Batchsize and number of MC samples were reduced in high dimensions to not exceed the available VRAM.

B.4. Contrastive Hedged Instance Embeddings

HIB (Oh et al., 2019) is formulated similarly to MCInfoNCE in that it also draws samples of a posterior and computes a probability score with them. HIB originally uses Gaussians and compares L_2 distances between samples. We adapt this to vMFs and cosine distances to align it with the spherical formulation of the latent space. The reformulated HIB loss is

$$\mathcal{L}_{\text{HIB}} := \mathbb{E}_{\substack{x \sim P(x) \\ x^+ \sim P(x^+|x) \\ x_m^- \sim P(x^-), m=1, \dots, M}} \left(-\log \mathbb{E}_{\substack{z \sim Q(z|x) \\ z^+ \sim Q(z^+|x^+)}} (s(a \cdot z^\top z^+ + b)) - \frac{1}{M} \sum_{m=1}^M \log \mathbb{E}_{\substack{z \sim Q(z|x) \\ z^+ \sim Q(z^+|x_m^-)}} (1 - s(a \cdot z^\top z_m^- + b)) \right), \quad (54)$$

where $s(\cdot)$ is the Sigmoid function and a and b are tuneable hyperparameters. We excluded the KL regularizer originally proposed by Oh et al. since none of the other losses receive prior information on $\kappa(x)$.

B.5. Contrastive Expected Likelihood Kernel

The ELK is commonly used inside a classification cross-entropy loss (Kirchhof et al., 2022). Its key characteristic is that it replaces the point-to-point distance, e.g., cosine distance, by the expected likelihood distance. An analytical solution to compare two vMFs is provided in the supplementary of Kirchhof et al.. We can plug this distance $d_{\text{EL-vMF}}(\hat{\mu}(x_1), \hat{\kappa}(x_1), \hat{\mu}(x_2), \hat{\kappa}(x_2))$ into InfoNCE and transform it into a similarity by multiplying it with -1 to obtain our contrastive ELK loss:

$$\mathcal{L}_{\text{ELK}} := \mathbb{E}_{\substack{x \sim P(x) \\ x^+ \sim P(x^+|x) \\ x_m^- \sim P(x^-), m=1, \dots, M}} \left(-\log \frac{e^{-\kappa_{\text{pos}} d_{\text{EL-vMF}}(\hat{\mu}(x), \hat{\kappa}(x), \hat{\mu}(x^+), \hat{\kappa}(x^+))}}{\frac{1}{M} e^{-\kappa_{\text{pos}} d_{\text{EL-vMF}}(\hat{\mu}(x), \hat{\kappa}(x), \hat{\mu}(x^+), \hat{\kappa}(x^+))} + \frac{1}{M} \sum_{m=1}^M e^{-\kappa_{\text{pos}} d_{\text{EL-vMF}}(\hat{\mu}(x), \hat{\kappa}(x), \hat{\mu}(x_m^-), \hat{\kappa}(x_m^-))}} \right). \quad (55)$$

B.6. Hyperparameter Tuning

All losses were tuned on the "Standard" experiment setup via grid search. The seed for the generative process was exclusive and not used in the five seeds of the final results. Table 7 below gives the hyperparameters along with the chosen best setup according to the rank correlation between $\kappa(x)$ and $\hat{\kappa}(x)$.

There are two interesting results in this tuning. First, the true generative κ_{pos} was indeed the best choice. All methods performed worse when they learned it themselves (starting from the true value) or when given a different value (not shown here). Second, MCInfoNCE performs best with a high number of negative samples. This corroborates the theoretical study of its limiting behaviour as $M \rightarrow \infty$.

Phasewise training is the empirical strategy of first learning $\hat{\mu}(x)$ during the first half of epochs, then fixing it and learning $\hat{\kappa}(x)$ (Shi & Jain, 2019; Li et al., 2021). MCInfoNCE showed an improved performance with this strategy. This is likely because the training signal of $\kappa(x)$ is far lower in the loss than that of $\mu(x)$. During the training phase of $\hat{\mu}(x)$, it turned out beneficial to use negatives from the same batch, i.e., $M = 0$.

	HIB	ELK	MCInfoNCE
Number of negatives M	{0, 1, 32}	{0, 1, 32}	{0, 1, 32}
κ_{pos} learnable	{yes, no }	{yes, no }	{yes, no }
Phasewise training	{yes, no }	{yes, no }	{ yes , no}
a	{0.5, 1 , 2, 4}		
b	{-8, -4, -2, -1, 0 , 1, 2, 4, 8}		

Table 7. Possible hyperparameters and best-performing hyperparameters (**bold**). $M = 0$ corresponds to not sampling negatives, but using one sample from the same batch as a negative. HIB’s additional hyperparameters were tuned after the first three parameters to reduce the number of grid-search evaluations.

B.7. Ablation with High Latent Space Dimension

We use the latent space dimension $D = 10$ for most experiments following Zimmermann et al. (2021). Below in Figure 7, we increase the latent space dimension of the generative process and encoder up to 64. We notice considerable performance drops for $D \geq 40$. Other losses than MCInfoNCE also suffer this. Hence, it is likely because of our experimental setup: We use uniformly distributed negatives instead of sophisticated negative mining and the rejection sampling has lower success probabilities in high dimensions, making it harder to generate valid contrastive examples.

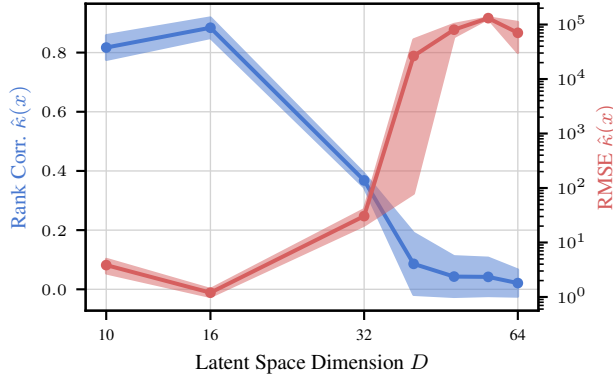


Figure 7. The metrics worsen if the generative process has a latent space of dimension $D \geq 40$. This is likely not due to MCInfoNCE, but a limitation of the contrastive setup of our controlled experiment. Mean \pm std. err. for five seeds.

C. CIFAR-10H Experiment

C.1. Contrastive Learning on CIFAR

To test whether the predicted certainty $\hat{\kappa}(x)$ aligns with human-judged aleatoric uncertainty, we require a dataset that provides a ground-truth. CIFAR-10H (Peterson et al., 2019) provides 50 annotations for each test-set image of CIFAR-10. We use the entropy of the probability distribution over these annotations as a measure of aleatoric uncertainty in each image, and compare its negative to the predicted certainty $\hat{\kappa}(x)$ via rank correlation. Since the annotations were only collected for the 10000 images of the test set of CIFAR-10, we randomly split them into 6000 images for training, 2000 for validation, and 2000 for test. To prevent confusions with the CIFAR-10 train and test set, we refer to these as the CIFAR-10H train, validation, and test sets. The image indices that belong to each set are provided in our code repository.

This leaves us with the task of redefining the CIFAR classification task into a contrastive learning problem. To this end, we simply assume that images are positive to one another if they belong to the same class and negative if they do not. CIFAR-10H, however, has soft class distributions for each image instead of a crisp class. Thus, we first draw a class c from the class distribution $P(C|x)$ of a reference image x from the train set. We then draw a positive image x^+ from a multinomial distribution over all train images weighed by their probabilities of that class $P(C = c|x^+)$. Negative images x^- are selected the same way, but weighed by the probability of *not* being class c , i.e., $1 - P(C = c|x^-)$. This provides the contrastive data generator required for training.

Since the human annotation data might be noisy in how well it captures the aleatoric uncertainty, we complement it with a synthetical way to introduce aleatoric uncertainty. In a second test dataset, we copy the CIFAR-10H test images, but perform a random crop and rescale that reduces the image to a proportion $\text{crop_size} \sim \text{Unif}([0.25, 1])$ of its original width and length. This directly reduces the information available in the image and therefore increases its aleatoric uncertainty, without introducing artifacts that might let the image go out-of-distribution. We calculate the rank correlation of the reduction in size crop_size and the (negative) predicted certainty $-\hat{\kappa}(x)$ as an alternative way to evaluate whether $\hat{\kappa}(x)$ reflects loss in information in the input, and therefore aleatoric uncertainty.

C.2. Hyperparameters

We use a ResNet-18 (He et al., 2016) pretrained on the CIFAR-10 train dataset (Phan, 2021) and replace the classification layer by a linear layer with the input and output dimensions $[512, D]$. We then train the linear layer and the ResNet backbone under each loss for 8192 batches of batchsize 128, which corresponds to roughly 175 epochs on the 6000 CIFAR-10H train images. We use the CIFAR-10H validation set to select the best model, evaluated after each 16 batches. The criterion is the rank correlation between $\hat{\kappa}(x)$ and the crop size in the synthetically deteriorated CIFAR-10H validation set. We chose this metric rather than the human annotator disagreement since it can be generated on arbitrary datasets without new annotations. All losses use 128 MC samples and, according to the results in Appendix B.6, a fixed κ_{pos} . We use the same Adam optimizer with a learning rate of 0.0001, learning rate scheduling, and (optional) phase-wise training as in B.6. The remaining hyperparameters were tuned via grid search. The best choices are highlighted in Table 8.

Loss	HIB	ELK	MCInfoNCE	MCInfoNCE	MCInfoNCE
Train Dataset / Label Type	CIFAR-10H soft	CIFAR-10H soft	CIFAR-10H soft	CIFAR-10H hard	CIFAR-10 hard
Latent Dim D	{ 8 , 16}	{ 8 , 16}	{ 8 , 16}	{8, 16 }	{ 8 , 16}
Number of negatives M	{ 0 , 1, 32}	{0, 1 , 32}	{0, 1, 32 }	{ 0 , 1, 32}	{ 0 , 1, 32}
κ_{pos}	{16, 32 , 64}	{16, 32 , 64}	{ 16 , 32, 64}	{ 16 , 32, 64}	{16, 32, 64 }
Phasewise training	{ yes , no}	{yes, no }	{ yes , no}	{yes, no }	{yes, no}
a	{0.5, 1, 2 , 4}				
b	{-2, -1, 0, 1 , 2}				

Table 8. Possible hyperparameters and best-performing hyperparameters (**bold**). $M = 0$ corresponds to not sampling negatives, but using one sample from the same batch as a negative. HIB’s additional hyperparameters were tuned after the first four parameters to reduce the number of grid-search evaluations.

C.3. Certainty Intervals

Since we have a (estimated) posterior distribution $P(z|x)$, we can give a credible interval $\text{CI}_p \subseteq \mathcal{Z}$ that the latent z of x falls into with a probability $p \in [0, 1]$, i.e., $P(z \in \text{CI}_p) = p$. We center this interval around the mode of the posterior vMF, such that it is a highest posterior density interval (HPDI). Due to the rotational symmetry of the vMF, for a given $\kappa(x)$ and credible level p , this interval has the form $\text{CI}_p = \{z \in \mathcal{Z} | z^\top \mu(x) \leq t\}$, i.e., all latents z closer to the mode $\mu(x)$ than a certain threshold $t \in [-1, 1]$ measured by cosine similarity. This threshold is the (approximated) $(1 - p)$ quantile of the vMF.

To visualize this latent interval, we define the credible images interval (CII). This is a pre-image of the corresponding CI and gives all images whose mode is within the CI, i.e., $\text{CII}_p := \{x \in \mathcal{X} | \mu(x) \in \text{CI}_p\}$. This can either be visualized via a GAN conditional on $z \in \text{CI}_p$ or by images from the dataset with $\mu(x) \in \text{CII}_p$. We note that this does not reflect the aleatoric uncertainty of those images. We leave this extension for future work.

C.4. Qualitative Evaluation of Aleatoric Uncertainty

Besides the quantitative metrics reported in the main text, we can also take a qualitative look at whether $\hat{\kappa}(x)$ represents aleatoric uncertainty in the inputs. Figure 8 visualizes the five images with the lowest and highest $\hat{\kappa}(x)$ in each class in the CIFAR-10H test set, i.e., on unseen data. It can be seen that images with a low $\hat{\kappa}(x)$ tend to hide characteristic parts of the object via bad crops, being too far away from the object, or an uncommon perspective. Images with a high $\hat{\kappa}(x)$ show characteristic features clearly, making it less ambiguous to tell what they show. In other words, they indeed have a lower aleatoric uncertainty.

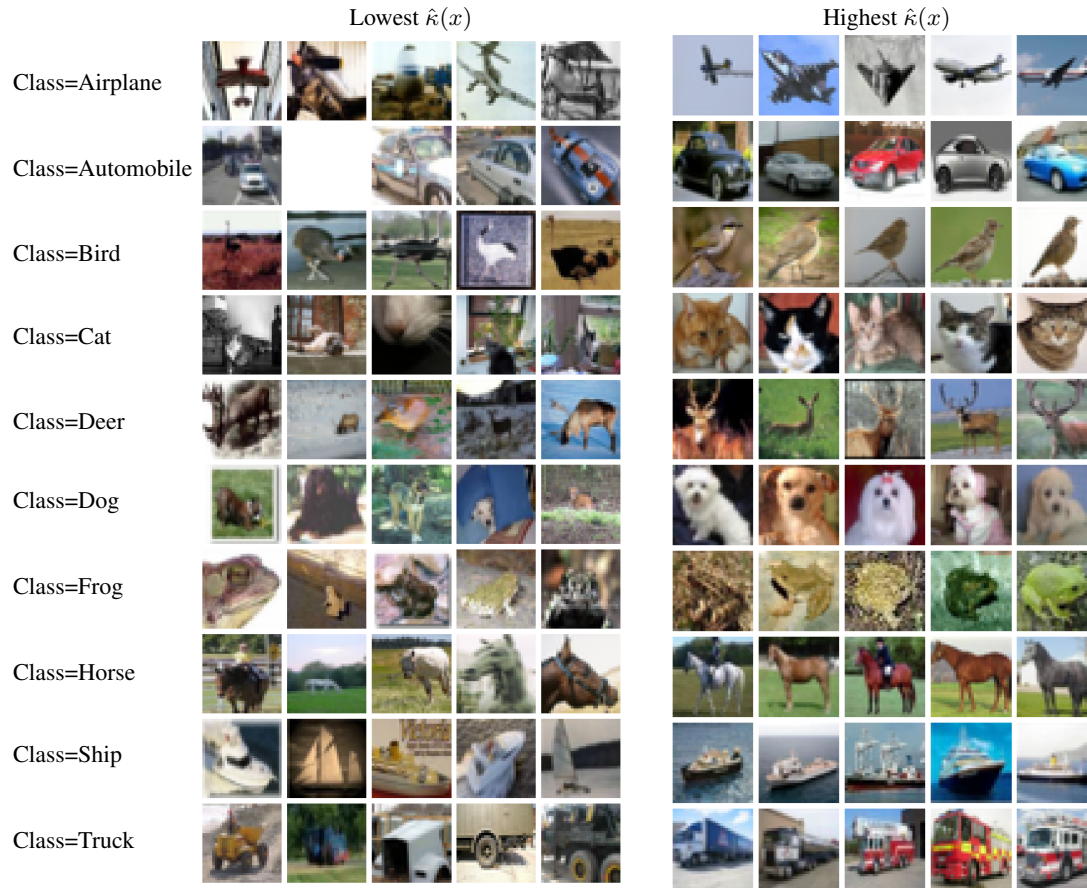


Figure 8. Images for which MCInfoNCE predicts the highest aleatoric uncertainty, i.e., lowest $\hat{\kappa}(x)$, (left) per class qualitatively look more ambiguous than those with the highest predicted $\hat{\kappa}(x)$ (right).