

238 **A Identifiability of latents drawn from a vMF around cluster vectors**

239 In this section, we formally state and prove our core theoretical result. We start off by defining
 240 and discussing a useful notion, then introduce our assumptions on the data generating process. We
 241 proceed with the main statement and finish with the proof.

242 **A.1 Affine Generator Systems**

243 **Definition 1** (Affine Generator System). *A system of vectors $\{\mathbf{v}_c \in \mathbb{R}^d | c \in \mathcal{C}\}$ is called an affine
 244 generator system if the affine hull defined by them is \mathbb{R}^d . More precisely, any vector in \mathbb{R}^d is an
 245 affine linear combination of the vectors in the system. Put into symbols: for any $\mathbf{v} \in \mathbb{R}^d$ there exist
 246 coefficients $\alpha_c \in \mathbb{R}$, such that*

$$\mathbf{v} = \sum_{c \in \mathcal{C}} \alpha_c \mathbf{v}_c \quad \text{and} \quad \sum_{c \in \mathcal{C}} \alpha_c = 1. \quad (3)$$

247 **Lemma 1** (Properties of affine generator systems). *The following hold for any affine generator
 248 system $\{\mathbf{v}_c \in \mathbb{R}^d | c \in \mathcal{C}\}$:*

- 249 1. *for any $a \in \mathcal{C}$ the system $\{\mathbf{v}_c - \mathbf{v}_a | c \in \mathcal{C}\}$ is now a generator system of \mathbb{R}^d ;*
- 250 2. *the invertible linear image of an affine generator system is also an affine generator system.*

251 **A.2 Assumptions and main result**

252 **Assumptions 1C** (DGP with vMF samples around cluster vectors). *Assume the following DGP:*

- 253 (i) *There exists a finite set of classes \mathcal{C} , represented by a set of unit-norm d -dimensional cluster-*
 254 *vectors $\{\mathbf{v}_c | c \in \mathcal{C}\} \subseteq \mathbb{S}^{d-1}$ such that they form an affine generator system of \mathbb{R}^d .*
- 255 (ii) *There is a finite set of instace labels \mathcal{I} and a well-defined, surjective class function $\mathcal{C} : \mathcal{I} \rightarrow \mathcal{C}$*
 256 *(every label belongs to exactly one class and every class is in use).*
- 257 (iii) *Our data sample is labelled with an instance label chosen uniformly, i.e., $I \in \text{Uni}(\mathcal{I})$ and,*
 258 *hence, belongs to class $C = \mathcal{C}(I)$.*
- 259 (iv) *The latent $\mathbf{z} \in \mathbb{S}^{d-1}$ of our data sample with label I is drawn from a vMF distribution around*
 260 *the cluster vector \mathbf{v}_C , where $C = \mathcal{C}(I)$:*

$$\mathbf{z} \sim p(\mathbf{z}|C) \propto e^{\alpha \langle \mathbf{v}_C, \mathbf{z} \rangle}. \quad (4)$$

- 261 (v) *The data sample \mathbf{x} is generated by passing the latent \mathbf{z} through a continuous and injective*
 262 *generator function $\mathbf{g} : \mathbb{S}^{d-1} \rightarrow \mathbb{R}^D$, i.e., $\mathbf{x} = \mathbf{g}(\mathbf{z})$.*

263 Assume that, using the DIET objective (6), we train a continuous encoder $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^d$ on \mathbf{x} and a
 264 linear classification head \mathbf{W} on top of \mathbf{f} . The rows of \mathbf{W} are $\{\mathbf{w}_i^\top | i \in \mathcal{I}\}$. In other words, \mathbf{W}
 265 computes similarities (scalar products) between its rows and the embeddings:

$$\mathbf{W} : \mathbf{f}(\mathbf{x}) \mapsto [\langle \mathbf{w}_i, \mathbf{f}(\mathbf{x}) \rangle | i \in \mathcal{I}]. \quad (5)$$

266 In DIET, we optimize the following objective amongst all possible continuous encoders \mathbf{f} , linear
 267 classifiers \mathbf{W} , and $\beta > 0$:

$$\mathcal{L}(\mathbf{f}, \mathbf{W}, \beta) = \mathbb{E}_{(\mathbf{x}, I)} \left[-\ln \frac{e^{\beta \langle \mathbf{w}_I, \mathbf{f}(\mathbf{x}) \rangle}}{\sum_{j \in \mathcal{I}} e^{\beta \langle \mathbf{w}_j, \mathbf{f}(\mathbf{x}) \rangle}} \right] \quad (6)$$

268 **Theorem 1C** (Identifiability of latents drawn from a vMF around cluster vectors). *Let $(\mathbf{f}, \mathbf{W}, \beta)$*
 269 *globally minimize the DIET objective (6) under the following additional constraints:*

- 270 C1. *both the embeddings $\mathbf{f}(\mathbf{x})$ and \mathbf{w}_i 's are unit-normalized. Then:*
 - 271 (a) $\mathbf{h} = \mathbf{f} \circ \mathbf{g}$ *is orthogonal linear, i.e., the latents are identified up to an orthogonal linear*
 272 *transformation;*
 - 273 (b) $\mathbf{w}_i = \mathbf{h}(\mathbf{v}_{\mathcal{C}(i)})$ *for any $i \in \mathcal{I}$, i.e., \mathbf{w}_i 's identify the cluster-vectors \mathbf{v}_c up to the same*
 274 *orthogonal linear transformation;*
 - 275 (c) $\beta = \alpha$, *the temperature of the vMF distribution is also identified.*
- 276 C2. *the embeddings $\mathbf{f}(\mathbf{x})$ are unit-normalized, the \mathbf{w}_i 's are unnormalized. Then:*
 - 277 (a) $\mathbf{h} = \mathbf{f} \circ \mathbf{g}$ *is orthogonal linear;*
 - 278 (b) $\mathbf{w}_i = \frac{\alpha}{\beta} \mathbf{h}(\mathbf{v}_{\mathcal{C}(i)}) + \psi$ *for any $i \in \mathcal{I}$, where ψ is a constant vector independent of i .*

279 C3. the embeddings $\mathbf{f}(\mathbf{x})$ are unnormalized, while the \mathbf{w}_i 's are unit-normalized. If the system
 280 $\{\mathbf{v}_c|c\}$ is **diverse enough in the sense of Assum. 2**, then:

281 (a) $\mathbf{w}_i = \mathcal{O}\mathbf{v}_{\mathcal{C}(i)}$, for any $i \in \mathcal{I}$, where \mathcal{O} is orthogonal linear;

282 (b) $\mathbf{h} = \mathbf{f} \circ \mathbf{g} = \frac{\alpha}{\beta}\mathcal{O}$ with the same orthogonal linear transformation, but scaled with $\frac{\alpha}{\beta}$.

283 C4. neither the embeddings $\mathbf{f}(\mathbf{x})$ nor the rows of \mathbf{W} are unit-normalized. Then:

284 (a) $\mathbf{h} = \mathbf{f} \circ \mathbf{g}$ is linear;

285 (b) \mathbf{w}_i identifies $\mathbf{v}_{\mathcal{C}(i)}$ up to an affine linear transformation.

286 Furthermore, in all cases, the row vectors that belong to samples of the same class are equal, i.e., for
 287 any $i, j \in \mathcal{I}$, $\mathcal{C}(i) = \mathcal{C}(j)$ implies $\mathbf{w}_i = \mathbf{w}_j$.

288 **Remark.** In cases C2 and C4, the cluster vectors are unnormalized and, therefore, can absorb the
 289 temperature parameter β . Thus β can be set to 1 without loss of generality. In case C3, it is \mathbf{f} that
 290 can absorb β .

291 **Assumption 2** (Diverse data). The system $\{\mathbf{v}_c|c \in \mathcal{C}\}$ is said to be diverse enough, if the following
 292 $|\mathcal{C}| \times 2d$ matrix has full column rank of $2d$:

$$\begin{pmatrix} \dots\dots\dots & \dots\dots\dots \\ (\mathbf{v}_c \odot \mathbf{v}_c)^\top & \mathbf{v}_c^\top \\ \dots\dots\dots & \dots\dots\dots \end{pmatrix}, \quad (7)$$

293 where $[\mathbf{x} \odot \mathbf{y}]_i = x_i y_i$ is the elementwise- or Hadamard product.

294 As long as $|\mathcal{C}| \geq 2d$, this property holds almost surely w.r.t. the Lebesgue-measure of \mathbb{S}^{d-1} or any
 295 continuous probability distribution of $\mathbf{v}_c \in \mathbb{S}^{d-1}$.

296 **Proof. Step 1: Deriving an equation characterizing the global optimizers of the objective.**

297 **Rewriting the objective in terms of latents:** we plug the expression $\mathbf{x} = \mathbf{g}(\mathbf{z})$ into the optimization
 298 objective (6) to express the dependence in terms of the latents \mathbf{z} :

$$\mathcal{L}(\mathbf{f}, \mathbf{W}, \beta) = \mathbb{E}_{(\mathbf{z}, I)} \left[-\ln \frac{e^{\beta\langle \mathbf{w}_I, \mathbf{f} \circ \mathbf{g}(\mathbf{z}) \rangle}}{\sum_{j \in \mathcal{I}} e^{\beta\langle \mathbf{w}_j, \mathbf{f} \circ \mathbf{g}(\mathbf{z}) \rangle}} \right] = \mathcal{L}_{\mathbf{z}}(\mathbf{f} \circ \mathbf{g}, \mathbf{W}, \beta), \quad (8)$$

299 where the optimization is still over \mathbf{f} (and not $\mathbf{h} = \mathbf{f} \circ \mathbf{g}$).

300 We note that the generator \mathbf{g} is, by assumption, continuously invertible on the compact set \mathbb{S}^{d-1} .
 301 Therefore, its image $\mathbf{g}(\mathbb{S}^{d-1})$ is compact, too, and its inverse \mathbf{g}^{-1} is also continuous. By Tietze's
 302 extension theorem [Wikipedia, 2024b], \mathbf{g}^{-1} can be continuously extended to a function $\mathbf{F} : \mathbb{R}^D \rightarrow$
 303 \mathbb{S}^{d-1} . Therefore, any continuous function $\mathbf{h} : \mathbb{S}^{d-1} \rightarrow \mathbb{R}^d$ can take the role of $\mathbf{f} \circ \mathbf{g}$ by substituting
 304 $\mathbf{f} = \mathbf{h} \circ \mathbf{F}$ continuous, since now $\mathbf{f} \circ \mathbf{g} = \mathbf{h} \circ (\mathbf{F} \circ \mathbf{g}) = \mathbf{h} \circ \text{id}_{\mathbb{S}^{d-1}} = \mathbf{h}$.

305 Hence, minimizing $\mathcal{L}_{\mathbf{z}}(\mathbf{f} \circ \mathbf{g}, \mathbf{W}, \beta)$ (and by extension $\mathcal{L}(\mathbf{f}, \mathbf{W}, \beta)$) for continuous \mathbf{f} equates to
 306 minimizing $\mathcal{L}_{\mathbf{z}}(\mathbf{h}, \mathbf{W}, \beta)$ for continuous \mathbf{h} :

$$\mathcal{L}_{\mathbf{z}}(\mathbf{h}, \mathbf{W}, \beta) = \mathbb{E}_{(\mathbf{z}, I)} \left[-\ln \frac{e^{\beta\langle \mathbf{w}_I, \mathbf{h}(\mathbf{z}) \rangle}}{\sum_{j \in \mathcal{I}} e^{\beta\langle \mathbf{w}_j, \mathbf{h}(\mathbf{z}) \rangle}} \right]. \quad (9)$$

307 **Expressing the condition for global optimality of the objective:** We rewrite the objective (9) by
 308 1) using the indicator variable $\delta_{I=i}$ of the event $\{I = i\}$ and 2) applying the law of total expectation:

$$\mathcal{L}_{\mathbf{z}}(\mathbf{h}, \mathbf{W}, \beta) = \mathbb{E}_{(\mathbf{z}, I)} \left[-\sum_{i \in \mathcal{I}} \delta_{I=i} \ln \frac{e^{\beta\langle \mathbf{w}_i, \mathbf{h}(\mathbf{z}) \rangle}}{\sum_{j \in \mathcal{I}} e^{\beta\langle \mathbf{w}_j, \mathbf{h}(\mathbf{z}) \rangle}} \right] \quad (10)$$

$$= \mathbb{E}_{\mathbf{z}} \left[\mathbb{E}_I \left[-\sum_{i \in \mathcal{I}} \delta_{I=i} \ln \frac{e^{\beta\langle \mathbf{w}_i, \mathbf{h}(\mathbf{z}) \rangle}}{\sum_{j \in \mathcal{I}} e^{\beta\langle \mathbf{w}_j, \mathbf{h}(\mathbf{z}) \rangle}} \mid \mathbf{z} \right] \right]. \quad (11)$$

309 Using the properties that $\mathbb{E}[A f(B)|B] = \mathbb{E}[A|B]f(B)$ and that $\mathbb{E}[\delta_{I=i}] = \mathbb{P}(I = i)$, we conclude
 310 that:

$$\mathcal{L}_{\mathbf{z}}(\mathbf{h}, \mathbf{W}, \beta) = \mathbb{E}_{\mathbf{z}} \left[- \sum_{i \in \mathcal{I}} \mathbb{E}_I \left[\delta_{I=i} \ln \frac{e^{\beta \langle \mathbf{w}_i, \mathbf{h}(\mathbf{z}) \rangle}}{\sum_{j \in \mathcal{I}} e^{\beta \langle \mathbf{w}_j, \mathbf{h}(\mathbf{z}) \rangle}} \mid \mathbf{z} \right] \right] \quad (12)$$

$$= \mathbb{E}_{\mathbf{z}} \left[- \sum_{i \in \mathcal{I}} \mathbb{E}_I \left[\delta_{I=i} \mid \mathbf{z} \right] \ln \frac{e^{\beta \langle \mathbf{w}_i, \mathbf{h}(\mathbf{z}) \rangle}}{\sum_{j \in \mathcal{I}} e^{\beta \langle \mathbf{w}_j, \mathbf{h}(\mathbf{z}) \rangle}} \right] \quad (13)$$

$$= \mathbb{E}_{\mathbf{z}} \left[- \sum_{i \in \mathcal{I}} \mathbb{P}(I = i | \mathbf{z}) \ln \frac{e^{\beta \langle \mathbf{w}_i, \mathbf{h}(\mathbf{z}) \rangle}}{\sum_{j \in \mathcal{I}} e^{\beta \langle \mathbf{w}_j, \mathbf{h}(\mathbf{z}) \rangle}} \right]. \quad (14)$$

311 By Gibbs' inequality [Wikipedia, 2024a], the cross-entropy inside the expectation is globally mini-
 312 mized if and only if

$$\frac{e^{\beta \langle \mathbf{w}_i, \mathbf{h}(\mathbf{z}) \rangle}}{\sum_{j \in \mathcal{I}} e^{\beta \langle \mathbf{w}_j, \mathbf{h}(\mathbf{z}) \rangle}} = \mathbb{P}(I = i | \mathbf{z}), \quad \text{for any } i \in \mathcal{I}. \quad (15)$$

313 Moreover, the entire expectation is globally minimized if and only if the above equality (15) holds
 314 almost everywhere for $\mathbf{z} \in \mathbb{S}^{d-1}$.

315 Using that instance label I is uniformly distributed, or $\mathbb{P}(I = j) = \mathbb{P}(I = i)$, the likelihood of the
 316 sample being in class i can be expressed via Bayes' theorem as:

$$\mathbb{P}(I = i | \mathbf{z}) = \frac{p(\mathbf{z} | I = i) \mathbb{P}(I = i)}{\sum_{j \in \mathcal{I}} p(\mathbf{z} | I = j) \mathbb{P}(I = j)} = \frac{p(\mathbf{z} | I = i)}{\sum_{j \in \mathcal{I}} p(\mathbf{z} | I = j)}. \quad (16)$$

317 Substituting (16) into (15) yields that for any $i \in \mathcal{I}$ and almost everywhere w.r.t. $\mathbf{z} \in \mathbb{S}^{d-1}$:

$$\frac{e^{\beta \langle \mathbf{w}_i, \mathbf{h}(\mathbf{z}) \rangle}}{\sum_{j \in \mathcal{I}} e^{\beta \langle \mathbf{w}_j, \mathbf{h}(\mathbf{z}) \rangle}} = \frac{p(\mathbf{z} | I = i)}{\sum_{j \in \mathcal{I}} p(\mathbf{z} | I = j)}. \quad (17)$$

318 We now divide the equation (17) for the probability of a sample having label i with that of having
 319 label k and take the logarithm. This yields that $\mathcal{L}_{\mathbf{z}}(\mathbf{h}, \mathbf{W}, \beta)$ is globally minimized if and only if

$$\beta \langle \mathbf{w}_i - \mathbf{w}_k, \mathbf{h}(\mathbf{z}) \rangle = \ln \frac{p(\mathbf{z} | I = i)}{p(\mathbf{z} | I = k)} \quad (18)$$

320 holds for any $i, k \in \mathcal{I}$ and almost everywhere w.r.t. $\mathbf{z} \in \mathbb{S}^{d-1}$.

321 **Plugging in the vMF distribution:** Plugging the assumed conditional distribution from (4) into
 322 (18) yields the equivalent expression:

$$\beta \langle \mathbf{w}_i - \mathbf{w}_k, \mathbf{h}(\mathbf{z}) \rangle = \alpha \langle \mathbf{v}_{\mathcal{C}(i)} - \mathbf{v}_{\mathcal{C}(k)}, \mathbf{z} \rangle \quad (19)$$

323 holds for any $i, k \in \mathcal{I}$ and almost everywhere w.r.t. $\mathbf{z} \in \mathbb{S}^{d-1}$. Since \mathbf{h} is continuous, the equation
 324 holds almost everywhere w.r.t. \mathbf{z} if and only if it holds for all $\mathbf{z} \in \mathbb{S}^{d-1}$.

325 Observe that if $\mathbf{h} = id|_{\mathbb{S}^{d-1}}$, $\mathbf{w}_i = \mathbf{v}_{\mathcal{C}(i)}$ for any $i \in \mathcal{I}$, and $\beta = \alpha$, then the equation is satisfied.
 326 Thus, we can conclude that the global minimum of the cross-entropy loss is achieved.

327 **Step 2: Solving the equation for \mathbf{h} , \mathbf{W} and proving identifiability.**

328 We now find all solutions to prove the identifiability of the latent variables and that of the cluster
 329 vectors. Denote $\tilde{\mathbf{w}}_i = \frac{\beta}{\alpha} \mathbf{w}_i$ to simplify the above equation to:

$$\langle \tilde{\mathbf{w}}_i - \tilde{\mathbf{w}}_k, \mathbf{h}(\mathbf{z}) \rangle = \langle \mathbf{v}_{\mathcal{C}(i)} - \mathbf{v}_{\mathcal{C}(k)}, \mathbf{z} \rangle. \quad (20)$$

330 **\mathbf{h} is injective and has full-dimensional image:** We prove that \mathbf{h} is injective. Assume that
 331 $\mathbf{h}(\mathbf{z}_1) = \mathbf{h}(\mathbf{z}_2)$ for some $\mathbf{z}_1, \mathbf{z}_2 \in \mathbb{S}^{d-1}$. Plugging \mathbf{z}_1 and \mathbf{z}_2 into (20) and subtracting the two
 332 equations yields:

$$0 = \langle \tilde{\mathbf{w}}_i - \tilde{\mathbf{w}}_k, \mathbf{h}(\mathbf{z}_1) - \mathbf{h}(\mathbf{z}_2) \rangle = \langle \mathbf{v}_{\mathcal{C}(i)} - \mathbf{v}_{\mathcal{C}(k)}, \mathbf{z}_1 - \mathbf{z}_2 \rangle, \quad (21)$$

333 for any i, k . However, as the cluster vectors $\{\mathbf{v}_c | c\}$ form an affine generator system, the vectors
 334 $\{\mathbf{v}_{\mathcal{C}(i)} - \mathbf{v}_{\mathcal{C}(k)} | i, k\}$ form a generator system of \mathbb{R}^d (see Lem. 1). Therefore, $\langle \mathbf{y}, \mathbf{z}_1 - \mathbf{z}_2 \rangle = 0$, for
 335 any $\mathbf{y} \in \mathbb{R}^d$, which holds if and only if $\mathbf{z}_1 = \mathbf{z}_2$. Hence, \mathbf{h} is injective.

336 By the Borsuk-Ulam theorem, for any continuous map from \mathbb{S}^{d-1} to a space of dimensionality at
 337 most $d - 1$ there exists some pair of antipodal points that are mapped to the same point. Consequently,
 338 no such function can be injective at the same time. Since $h : \mathbb{S}^{d-1} \rightarrow \mathbb{R}^d$ is injective, the linear span
 339 of its image must be \mathbb{R}^d .

340 **Collapse of w_i 's:** We prove that $\tilde{w}_i = \tilde{w}_k$ if $\mathcal{C}(i) = \mathcal{C}(k)$, i.e., samples from the same cluster will
 341 have equal rows of \mathbf{W} associated with them.

342 Assume that $\mathcal{C}(i) = \mathcal{C}(k)$ and substitute them into (20):

$$\langle \tilde{w}_i - \tilde{w}_k, \mathbf{h}(z) \rangle = 0 \quad \text{for any } z \in \mathbb{S}^{d-1}. \quad (22)$$

343 However, we have just seen that the linear span of the image of \mathbf{h} is \mathbb{R}^d , which implies that $\tilde{w}_i = \tilde{w}_k$.
 344 Consequently, we may abuse our notation by setting $\tilde{w}_c = \tilde{w}_i$ if $\mathcal{C}(i) = c$, which yields a new form
 345 for (20):

$$\langle \tilde{w}_a - \tilde{w}_b, \mathbf{h}(z) \rangle = \langle \mathbf{v}_a - \mathbf{v}_b, z \rangle, \quad (23)$$

346 for any $a, b \in \mathcal{C}$ and any $z \in \mathbb{S}^{d-1}$.

347 **Linear transformation from $\mathbf{v}_a - \mathbf{v}_b$ to $\tilde{w}_a - \tilde{w}_b$:** We now prove the existence of a linear map
 348 \mathcal{A} on \mathbb{R}^d such that $\mathcal{A}(\mathbf{v}_a - \mathbf{v}_b) = \tilde{w}_a - \tilde{w}_b$ for any $a, b \in \mathcal{C}$. For this, we prove that the following
 349 mapping is well-defined:

$$\mathcal{A}: \sum_{a,b \in \mathcal{C}} \lambda_{ab}(\mathbf{v}_a - \mathbf{v}_b) \mapsto \sum_{a,b \in \mathcal{C}} \lambda_{ab}(\tilde{w}_a - \tilde{w}_b). \quad (24)$$

350 Since the system $\{\mathbf{v}_a - \mathbf{v}_b | a, b\}$ is not necessarily linearly independent, we have to prove that
 351 the mapping is independent of the choice of the linear combination. More precisely if for some
 352 coefficients $\lambda_{ab}, \lambda'_{ab}$

$$\sum_{a,b \in \mathcal{C}} \lambda_{ab}(\mathbf{v}_a - \mathbf{v}_b) = \sum_{a,b \in \mathcal{C}} \lambda'_{ab}(\mathbf{v}_a - \mathbf{v}_b) \quad (25)$$

353 holds, then it should be implied that

$$\sum_{a,b \in \mathcal{C}} \lambda_{ab}(\tilde{w}_a - \tilde{w}_b) = \sum_{a,b \in \mathcal{C}} \lambda'_{ab}(\tilde{w}_a - \tilde{w}_b). \quad (26)$$

354 Assume that (25) holds. Then, the difference of the two sides is:

$$0 = \sum_{a,b \in \mathcal{C}} (\lambda_{ab} - \lambda'_{ab})(\mathbf{v}_a - \mathbf{v}_b). \quad (27)$$

355 Taking the scalar product with an arbitrary $z \in \mathbb{S}^{d-1}$ and using the linearity of the scalar product
 356 gives us:

$$0 = \left\langle \sum_{a,b \in \mathcal{C}} (\lambda_{ab} - \lambda'_{ab})(\mathbf{v}_a - \mathbf{v}_b), z \right\rangle = \sum_{a,b \in \mathcal{C}} (\lambda_{ab} - \lambda'_{ab}) \langle \mathbf{v}_a - \mathbf{v}_b, z \rangle. \quad (28)$$

357 Now using (23) yields:

$$0 = \sum_{a,b \in \mathcal{C}} (\lambda_{ab} - \lambda'_{ab}) \langle \tilde{w}_a - \tilde{w}_b, \mathbf{h}(z) \rangle = \left\langle \sum_{a,b \in \mathcal{C}} (\lambda_{ab} - \lambda'_{ab})(\tilde{w}_a - \tilde{w}_b), \mathbf{h}(z) \right\rangle. \quad (29)$$

358 However, the linear span of the image of \mathbf{h} is \mathbb{R}^d , which implies that

$$\sum_{a,b \in \mathcal{C}} (\lambda_{ab} - \lambda'_{ab})(\tilde{w}_a - \tilde{w}_b) = 0, \quad (30)$$

359 equivalent to (26). Therefore, the mapping is well-defined. The linearity of \mathcal{A} follows trivially.

360 **\mathbf{h} is linear:** Equation (23) becomes:

$$\langle \mathcal{A}(\mathbf{v}_a - \mathbf{v}_b), \mathbf{h}(z) \rangle = \langle \mathbf{v}_a - \mathbf{v}_b, z \rangle, \quad (31)$$

361 for any $a, b \in \mathcal{C}$ and any $z \in \mathbb{S}^{d-1}$. Nevertheless, $\{\mathbf{v}_a - \mathbf{v}_b | a, b \in \mathcal{C}\}$ is a generator system of \mathbb{R}^d ,
 362 and, hence, (31) is equivalent to

$$\langle \mathcal{A}\mathbf{y}, \mathbf{h}(z) \rangle = \langle \mathbf{y}, z \rangle, \quad \text{for any } \mathbf{y} \in \mathbb{R}^d \text{ and any } z \in \mathbb{S}^{d-1}. \quad (32)$$

363 This is further equivalent to

$$\langle \mathbf{y}, \mathcal{A}^\top \mathbf{h}(z) \rangle = \langle \mathbf{y}, z \rangle. \quad (33)$$

364 Since \mathbf{y} is arbitrary, we conclude that $\mathcal{A}^\top \mathbf{h}(z) = z$ for any $z \in \mathbb{S}^{d-1}$. Therefore \mathcal{A} is an invertible
 365 transformation and $\mathbf{h} = (\mathcal{A}^\top)^{-1}$ is linear.

366 **Proving Thm. 1C case C4:** We have shown that \mathbf{h} is linear. Furthermore, from (31) it follows, by
 367 fixing b and defining $\boldsymbol{\psi} = \mathcal{A}\mathbf{v}_b - \mathbf{w}_b$, that

$$\tilde{\mathbf{w}}_a = \mathcal{A}\mathbf{v}_a + \boldsymbol{\psi}, \quad \text{for any } a \in \mathcal{C}, \quad (34)$$

368 which proves case C4 of Thm. 1C.

369 **Proving Thm. 1C case C2:** As a special case of the previous one, now we assume that $\mathbf{h}(z)$
 370 is unit-normalized and maps \mathbb{S}^{d-1} to \mathbb{S}^{d-1} . That amounts to $\mathbf{h} = (\mathcal{A}^\top)^{-1}$ being linear, norm-
 371 preserving, and therefore orthogonal. Consequently \mathcal{A} is also orthogonal, $\mathbf{h} = \mathcal{A}$ and (34) simplifies
 372 to $\frac{\beta}{\alpha}\mathbf{w}_a = \tilde{\mathbf{w}}_a = \mathcal{A}\mathbf{v}_a + \boldsymbol{\psi} = \mathbf{h}(\mathbf{v}_a) + \boldsymbol{\psi}$, which proves C2 of Thm. 1C.

373 **Proving Thm. 1C case C1:** We now assume that both \mathbf{h} and \mathbf{w}_i 's are unit-normalized. Conse-
 374 quently, $\mathbf{h} = \mathcal{A}$ is orthogonal linear and $\mathbf{w}_a = \frac{\alpha}{\beta}\mathcal{A}\mathbf{v}_a + \boldsymbol{\psi}$.

375 Therefore, on one hand, the \mathbf{w}_a 's lie on a d -dimensional hypersphere of radius $\frac{\alpha}{\beta}$ and center $\boldsymbol{\psi}$. On
 376 the other hand, by definition, \mathbf{w}_a 's also lie on the unit hypersphere \mathbb{S}^{d-1} .

377 Since the system $\{\mathbf{w}_a | a \in \mathcal{C}\}$ is the bijective affine linear image of the affine generator system
 378 $\{\mathbf{v}_a | a \in \mathcal{C}\}$, $\{\mathbf{w}_a | a \in \mathcal{C}\}$ is also an affine generator system (Lem. 1). Consequently, there could be
 379 at most one hypersphere in \mathbb{R}^d which contains all the \mathbf{w}_a 's. Hence $\frac{\alpha}{\beta} = 1$, $\boldsymbol{\psi} = \mathbf{0}$, and $\mathbf{w}_a = \mathbf{h}(\mathbf{v}_a)$,
 380 which proves C1 of Thm. 1C.

381 **Proving Thm. 1C case C3:** Finally, we assume that \mathbf{w}_i 's are unit-normalized. As this is a special
 382 case of Thm. 1C C4, we know that there exists a constant vector $\boldsymbol{\psi}$ such that:

$$\mathbf{w}_a = \frac{\alpha}{\beta}\mathcal{A}\mathbf{v}_a + \boldsymbol{\psi}, \quad (35)$$

383 for any $a \in \mathcal{C}$. We are going to prove that $\mathcal{O} = \frac{\alpha}{\beta}\mathcal{A}$ is orthogonal and $\boldsymbol{\psi} = \mathbf{0}$.

384 Let $\mathcal{O} = \mathcal{U}^\top \Sigma \mathcal{V}$ be the singular value decomposition (SVD) of \mathcal{O} . Consequently, after premultiplying
 385 with \mathcal{U} , we receive:

$$\mathcal{U}\mathbf{w}_a = \Sigma \mathcal{V}\mathbf{v}_a + \mathcal{U}\boldsymbol{\psi}. \quad (36)$$

386 As orthogonal transformations \mathcal{U} and \mathcal{V} keep their arguments unit-normalized and $\{\mathcal{V}\mathbf{v}_a - \mathcal{V}\mathbf{v}_b\}$ is
 387 still an affine generator system (Lem. 1), we may assume without the loss of generality that

$$\mathbf{w}_a = \Sigma \mathbf{v}_a + \boldsymbol{\psi}, \quad (37)$$

388 for any $a \in \mathcal{C}$, where all \mathbf{v}_a 's and \mathbf{w}_a 's are unit-normalized.

389 Let us assume that $\boldsymbol{\psi} \neq \mathbf{0}$. In that case both sides of (37) can be scaled such that the offset $\boldsymbol{\psi}$ has
 390 unit norm. In this case \mathbf{w}_a 's are no longer on the unit hypersphere, but they instead have a mutual
 391 norm r . Assuming that the diagonal elements of Σ are $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_d)$, this is equivalent to:

$$r^2 = \|\Sigma \mathbf{v}_a + \boldsymbol{\psi}\|^2 = \|\Sigma \mathbf{v}_a\|^2 + 2\langle \Sigma \mathbf{v}_a, \boldsymbol{\psi} \rangle + \|\boldsymbol{\psi}\|^2 \quad (38)$$

$$= \langle \mathbf{v}_a \odot \mathbf{v}_a, \boldsymbol{\sigma} \odot \boldsymbol{\sigma} \rangle + \langle \mathbf{v}_a, 2\boldsymbol{\sigma} \odot \boldsymbol{\psi} \rangle + 1, \quad (39)$$

392 where $[\mathbf{x} \odot \mathbf{y}]_i = x_i y_i$ is the elementwise product. Eq. (39) is equivalent to the following:

$$(\mathbf{v}_a \odot \mathbf{v}_a)^\top (\boldsymbol{\sigma} \odot \boldsymbol{\sigma}) + \mathbf{v}_a^\top (2\boldsymbol{\sigma} \odot \boldsymbol{\psi}) - r^2 = -1. \quad (40)$$

393 Collecting the equations for all $a \in \mathcal{C}$ yields:

$$\mathcal{D} \begin{pmatrix} \boldsymbol{\sigma} \odot \boldsymbol{\sigma} \\ 2\boldsymbol{\sigma} \odot \boldsymbol{\psi} \\ r^2 \end{pmatrix} = -\mathbf{1}_{|\mathcal{C}|}, \quad (41)$$

394 where \mathcal{D} is the following $|\mathcal{C}| \times (2d + 1)$ matrix:

$$\mathcal{D} = \begin{pmatrix} \dots\dots\dots & \dots\dots\dots & \dots \\ (\mathbf{v}_a \odot \mathbf{v}_a)^\top & \mathbf{v}_a^\top & -1 \\ \dots\dots\dots & \dots\dots\dots & \dots \end{pmatrix}. \quad (42)$$

395 By Assum. 2, the left $|\mathcal{C}| \times 2d$ submatrix of \mathcal{D} has full rank of $2d$. Consequently, the solution space
 396 to the more general, linear equation $\mathcal{D}\mathbf{t} = -\mathbf{1}_{|\mathcal{C}|}$, where $\mathbf{t} \in \mathbb{R}^d$, has a dimensionality of at most 1.

397 Using the unit-normality of \mathbf{v}_a 's, we see that $(\mathbf{v}_a \odot \mathbf{v}_a)^\top \mathbf{1}_d = 1$. From this, it follows that the
 398 solutions are exactly the following:

$$\mathbf{t} = \begin{pmatrix} \gamma \cdot \mathbf{1}_d \\ \mathbf{0}_d \\ \gamma + 1 \end{pmatrix}, \quad \text{where } \gamma \in \mathbb{R}. \quad (43)$$

399 Therefore, for any solution of (41) there exists γ such that:

$$\boldsymbol{\sigma} \odot \boldsymbol{\sigma} = \gamma \cdot \mathbf{1}_d \quad (44)$$

$$\boldsymbol{\sigma} \odot \boldsymbol{\psi} = \mathbf{0}_d. \quad (45)$$

400 However, as the original transformation \mathcal{A} was invertible, all singular values σ_i are strictly positive
 401 and, thus, it follows that $\boldsymbol{\psi} = \mathbf{0}$. Technically speaking, this is a contradiction to our initial assumption
 402 that $\boldsymbol{\psi} \neq \mathbf{0}$. All in all, it follows that $\boldsymbol{\psi} = \mathbf{0}$ is the only possibility.

403 Therefore, (37) becomes:

$$\mathbf{w}_a = \Sigma \mathbf{v}_a, \quad (46)$$

404 where all \mathbf{v}_a 's and \mathbf{w}_a 's are unit-normalized. Following the same derivation yields:

$$1 = \|\Sigma \mathbf{v}_a\|^2 = (\mathbf{v}_a \odot \mathbf{v}_a)^\top (\boldsymbol{\sigma} \odot \boldsymbol{\sigma}), \quad (47)$$

405 or, after collecting the equations for all $a \in \mathcal{C}$:

$$\mathcal{B}(\boldsymbol{\sigma} \odot \boldsymbol{\sigma}) = \mathbf{1}_{|\mathcal{C}|}, \quad (48)$$

406 where \mathcal{B} is the $|\mathcal{C}| \times d$ matrix

$$\mathcal{B} = \begin{pmatrix} \dots\dots\dots \\ (\mathbf{v}_a \odot \mathbf{v}_a)^\top \\ \dots\dots\dots \end{pmatrix}. \quad (49)$$

407 By Assum. 2, \mathcal{B} has full rank, thus, there is at most one solution to the equation $\mathcal{B}\mathbf{t} = \mathbf{1}_{|\mathcal{C}|}$. Due to
 408 the unit-normality of \mathbf{v}_a 's, this solution is exactly $\mathbf{t} = \mathbf{1}_d$. However, as the singular values σ_i are all
 409 positive, the only solution to $\boldsymbol{\sigma} \odot \boldsymbol{\sigma} = \mathbf{1}_d$ is $\boldsymbol{\sigma} = \mathbf{1}_d$. This is equivalent to saying that $\mathcal{O} = \frac{\alpha}{\beta} \mathcal{A}$ is
 410 orthogonal.

411 Furthermore, $\mathbf{h} = (\mathcal{A}^\top)^{-1} = \left(\frac{\beta}{\alpha} \mathcal{O}^\top\right)^{-1} = \frac{\alpha}{\beta} \mathcal{O}$.

412

□

413 B Additional experimental results

414 In Tab. 2, we present additional ablation studies exploring the effect of varying the levels of con-
 415 centration for \mathbf{v}_c across the unit hyper-sphere. We do not observe any significant impact on the R^2
 scores from more concentrated cluster centroids \mathbf{v}_c .

Table 2: Identifiability in the synthetic setup. Mean \pm standard deviation across 5 random seeds. Settings that match our theoretical assumptions are \checkmark . We report the R^2 score for linear mappings, $\tilde{\mathbf{z}} \rightarrow \mathbf{z}$ and $\mathbf{w}_i \rightarrow \mathbf{v}_c$ for cases with normalized (o) and unnormalized (a) \mathbf{w}_i . For unnormalized \mathbf{w}_i , we verify that mappings $\tilde{\mathbf{z}} \rightarrow \mathbf{z}$ are orthogonal by reporting the mean absolute error between their singular values and those of an orthogonal transformation.

N	d	$ \mathcal{C} $	$p(\mathbf{v}_c)$	$p(\mathbf{z} \mathbf{v}_c)$	M.	normalized \mathbf{w}_i cases				unnormalized \mathbf{w}_i	
						$R_o^2(\uparrow)$	$\tilde{\mathbf{z}} \rightarrow \mathbf{z}$	$\mathbf{w}_i \rightarrow \mathbf{v}_c$	MAE _o (\downarrow)	$\tilde{\mathbf{z}} \rightarrow \mathbf{z}$	$\mathbf{w}_i \rightarrow \mathbf{v}_c$
10^3	5	100	Uniform	vMF($\kappa=10$)	\checkmark	98.6 ± 0.01	99.9 ± 0.01	0.01 ± 0.00	0.00 ± 0.00	99.0 ± 0.00	99.9 ± 0.00
10^3	5	100	Laplace	vMF($\kappa=10$)	\checkmark	98.7 ± 0.00	99.5 ± 0.00	0.01 ± 0.00	0.00 ± 0.00	99.1 ± 0.00	99.8 ± 0.00
10^3	5	100	Normal	vMF($\kappa=10$)	\checkmark	98.2 ± 0.01	99.2 ± 0.01	0.01 ± 0.00	0.00 ± 0.00	99.2 ± 0.00	99.8 ± 0.00

416

417 **C Acronyms**

418 **CL** Contrastive Learning

419 **DGP** data generating process

420 **ICA** Independent Component Analysis

421 **LVM** latent variable model

422 **SSL** Self-Supervised Learning

423 **vMF** von Mises-Fisher