

Figure 4: Comparison between the dynamics under Euclidean and cosine asymmetric losses for different initializations in a network with  $M = 2$  output neurons. **(a)** Observed dynamics of the eigenvalues in the two-neuron toy network under three different initializations. Both eigenvalues always converge to 1 regardless of the initialization. **(b)** Same as **(a)**, but for the cosine distance. Under different initializations, the two eigenvalues converge to arbitrary, but equal, values.

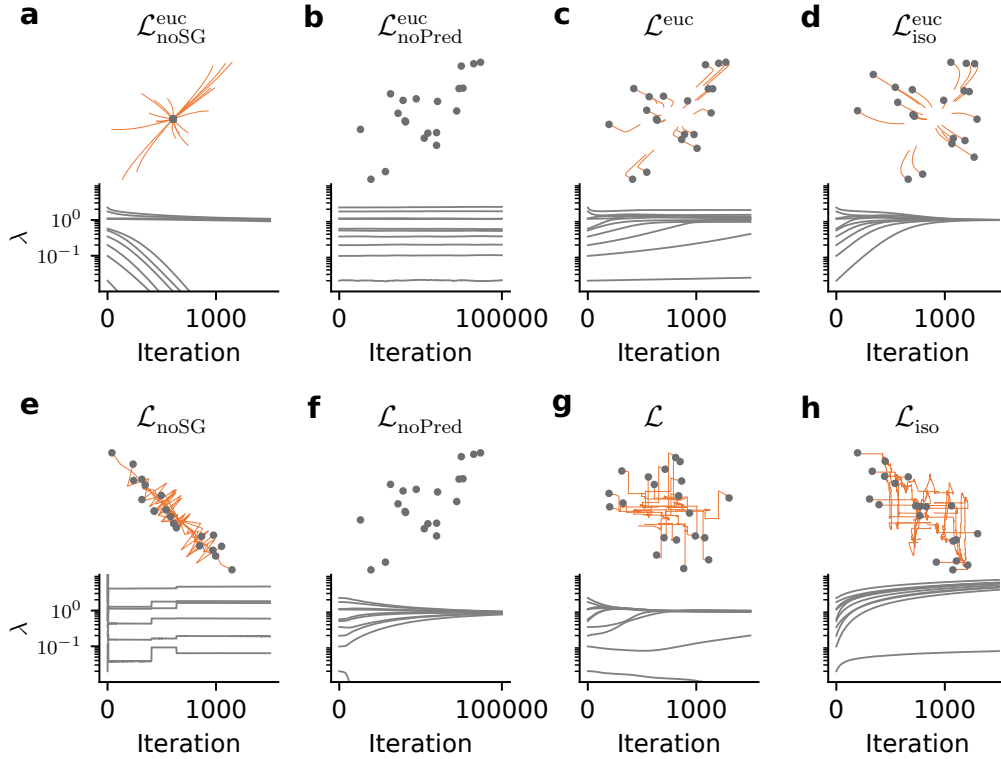


Figure 5: Same as Fig. 2 but with a ReLU nonlinearity on the embeddings. We observe learning dynamics qualitatively similar to the linear network.

## 399 B Proofs

400 **Lemma 1.** (Euclidean and cosine losses in the eigenspace of the predictor) *Let  $W_P$  be a linear*  
 401 *predictor according to DirectPred with eigenvalues  $\lambda_m$ , and  $\hat{\mathbf{z}}$  the representations expressed in the*  
 402 *predictor's eigenbasis, the asymmetric losses  $\mathcal{L}^{\text{euc}}$  and  $\mathcal{L}$  can be expressed as:*

$$\mathcal{L}^{\text{euc}} = \frac{1}{2} \sum_m^M |\lambda_m \hat{z}_m^{(1)} - \text{SG}(\hat{z}_m^{(2)})|^2, \quad (1)$$

$$\mathcal{L} = - \sum_m^M \frac{\lambda_m \hat{z}_m^{(1)} \text{SG}(\hat{z}_m^{(2)})}{\|D\hat{\mathbf{z}}^{(1)}\| \|\text{SG}(\hat{\mathbf{z}}^{(2)})\|}. \quad (2)$$

403 *Proof.* Under DirectPred, the predictor is a symmetric matrix with eigendecomposition  $W_P =$   
 404  $UDU^\top$ . Since  $U$  is an orthogonal matrix, we also have  $UU^\top = I$  so that we can simplify the losses  
 405 as follows:

$$\begin{aligned} \mathcal{L}^{\text{euc}} &= \frac{1}{2} \|W_P \mathbf{z}^{(1)} - \text{SG}(\mathbf{z}^{(2)})\|^2 \\ &= \frac{1}{2} \|UDU^\top \mathbf{z}^{(1)} - \text{SG}(UU^\top \mathbf{z}^{(2)})\|^2 \\ &= \frac{1}{2} \|D\hat{\mathbf{z}}^{(1)} - \text{SG}(\hat{\mathbf{z}}^{(2)})\|^2 \\ &= \frac{1}{2} \sum_m^M |\lambda_m \hat{z}_m^{(1)} - \text{SG}(\hat{z}_m^{(2)})|^2 \end{aligned}$$

$$\begin{aligned} \mathcal{L} &= - \frac{(W_P \mathbf{z}^{(1)})^\top \text{SG}(\mathbf{z}^{(2)})}{\|W_P \mathbf{z}^{(1)}\| \|\text{SG}(\mathbf{z}^{(2)})\|} \\ &= - \frac{(\mathbf{z}^{(1)})^\top UDU^\top \text{SG}(\mathbf{z}^{(2)})}{\|UDU^\top \mathbf{z}^{(1)}\| \|\text{SG}(UU^\top \mathbf{z}^{(2)})\|} \\ &= - \frac{(\hat{\mathbf{z}}^{(1)})^\top D \text{SG}(\hat{\mathbf{z}}^{(2)})}{\|D\hat{\mathbf{z}}^{(1)}\| \|\text{SG}(\hat{\mathbf{z}}^{(2)})\|} \\ &= - \sum_m^M \frac{\lambda_m \hat{z}_m^{(1)} \text{SG}(\hat{z}_m^{(2)})}{\|D\hat{\mathbf{z}}^{(1)}\| \|\text{SG}(\hat{\mathbf{z}}^{(2)})\|}, \end{aligned}$$

406 where we used the fact that  $U$  is orthogonal and therefore does not change the Euclidean norm.  
 407  $\hat{\mathbf{z}} = U^\top \mathbf{z}$  is the representation rotated into the eigenbasis.  $\square$

408 **Lemma 2.** (General learning dynamics of representations) *Assuming that a given loss  $\mathcal{L}$  is optimized*  
 409 *by gradient descent on the parameters of a neural network with the empirical NTK  $\hat{\Theta}$  and learning*  
 410 *rate  $\eta$ , then the representations  $\hat{\mathbf{z}}$  evolve according to the dynamics:*

$$\frac{d\hat{\mathbf{z}}}{dt} = -\eta \hat{\Theta}_t(\mathbf{x}, \mathcal{X}) \nabla_{\hat{\mathbf{z}}} \mathcal{L} \quad , \quad (6)$$

411 *Proof.* Let  $\boldsymbol{\theta}$  be the parameters of the neural network. Then we obtain the representational dynamics  
 412 using the chain rule in the continuous-time gradient-flow setting [1]:

$$\begin{aligned} \frac{d\hat{\mathbf{z}}}{dt} &= \nabla_{\boldsymbol{\theta}} \hat{\mathbf{z}} \frac{d\boldsymbol{\theta}}{dt} \\ &= \nabla_{\boldsymbol{\theta}} \hat{\mathbf{z}} (-\eta \nabla_{\boldsymbol{\theta}} \mathcal{L}) \\ &= \nabla_{\boldsymbol{\theta}} \hat{\mathbf{z}} \left( -\eta \nabla_{\boldsymbol{\theta}} \hat{\mathbf{Z}}^\top \nabla_{\hat{\mathbf{z}}} \mathcal{L} \right) \\ &= -\eta \hat{\Theta}_t(\mathbf{x}, \mathcal{X}) \nabla_{\hat{\mathbf{z}}} \mathcal{L} \quad . \end{aligned}$$

413 Note, that structurally these dynamics are the same as the embedding space dynamics [1, 2] but  
 414 merely expressed in the predictor eigen basis.  $\square$

415 We proceed by proving the following Lemma which we will use in our proofs of Theorems 1  
 416 and 2.

417 **Lemma 3.** *The NTK for a linear network is invariant under orthogonal transformations of the*  
 418 *network output.*

419 *Proof.* We first note that for a linear network, the parameters  $\boldsymbol{\theta}$  are just the feedforward weights  $W$ .  
 420 Therefore, for any orthogonal transformation  $U$  of the network output:

$$\begin{aligned} \hat{\mathbf{z}} &= U^\top f(\mathbf{x}) = U^\top W \mathbf{x} \\ \Rightarrow \nabla_{\boldsymbol{\theta}} \hat{\mathbf{z}} &= \nabla_W \hat{\mathbf{z}} = \nabla_W (U^\top W \mathbf{x}) = \mathbf{x}^\top \otimes U^\top, \end{aligned} \quad (18)$$

421 where  $\otimes$  is the Kronecker product resulting from the fact that every input vector component appears  
 422 in the update once for each output component.

423 We now study  $\hat{\Theta}_t(\mathcal{X}, \mathcal{X})$ , the transformed empirical NTK (cf. Lemma 2). The  $(M \times M)$  diagonal  
 424 blocks in the full  $(M|\mathcal{D}| \times M|\mathcal{D}|)$  empirical NTK  $\hat{\Theta}_t(\mathcal{X}, \mathcal{X})$  correspond to single samples and the  
 425 off-diagonal blocks are cross-terms between samples, where  $|\mathcal{D}|$  denotes the size of the training  
 426 dataset and  $M$  the dimension of the outputs. We can develop a generic expression for each  $(M \times M)$   
 427 block  $\hat{\Theta}_t(\mathbf{x}_i, \mathbf{x}_j)$  corresponding to the interactions between samples  $i$  and  $j$  as:

$$\begin{aligned} \hat{\Theta}_t(\mathbf{x}_i, \mathbf{x}_j) &= \nabla_W \hat{\mathbf{z}}_i \nabla_W \hat{\mathbf{z}}_j^\top \\ &= (\mathbf{x}_i^\top \otimes U^\top) (\mathbf{x}_j^\top \otimes U^\top)^\top \\ &= (\mathbf{x}_i^\top \otimes U^\top) (\mathbf{x}_j \otimes U) \\ &= (\mathbf{x}_i^\top \mathbf{x}_j) \otimes (U^\top U) \\ &= (\mathbf{x}_i^\top \mathbf{x}_j) \otimes I_M \\ &= (\mathbf{x}_i^\top \mathbf{x}_j) I_M. \end{aligned} \quad (19)$$

428 where we have used the fact that  $(A \otimes B)^\top = A^\top \otimes B^\top$  and  $(A \otimes B)(C \otimes D) = AC \otimes BD$ . Here,  
 429  $I_M$  is the identity matrix of size  $M$ . Noting that Eq. (19) is unchanged when  $U$  is just the identity  
 430 matrix completes the proof.  $\square$

431 **Theorem 1.** (Representational dynamics under  $\mathcal{L}^{\text{euc}}$ ) For a linear network with i.i.d Gaussian  
 432 inputs learning with  $\mathcal{L}^{\text{euc}}$ , the representational dynamics of each mode  $m$  independently follows the  
 433 gradient of the loss  $-\nabla_{\hat{\mathbf{z}}} \mathcal{L}^{\text{euc}}$ . More specifically, the dynamics uncouple and follow a system of  $M$   
 434 independent differential equations:

$$\frac{d\hat{\mathbf{z}}_{m,t}^{(1)}}{dt} = -\eta \frac{\partial \mathcal{L}^{\text{euc}}}{\partial \hat{\mathbf{z}}_{m,t}^{(1)}}(t) = \eta \lambda_{m,t} \left( \hat{\mathbf{z}}_{m,t}^{(2)} - \lambda_{m,t} \hat{\mathbf{z}}_{m,t}^{(1)} \right), \quad (8)$$

435 which, after taking the expectation over augmentations, becomes:

$$\frac{d\hat{\mathbf{z}}_{m,t}}{dt} = \eta \lambda_{m,t} (1 - \lambda_{m,t}) \hat{\mathbf{z}}_{m,t} \quad . \quad (9)$$

436 *Proof.* For a linear network with weights  $W \in \mathbb{R}^{M \times N}$ , we have from Lemma 3 that the empirical  
 437 NTK  $\hat{\Theta}(\mathcal{X}, \mathcal{X})$  in the orthogonal eigenbasis is equal to the empirical NTK  $\Theta(\mathcal{X}, \mathcal{X})$  in the original  
 438 basis. Furthermore from the proof for the lemma (see Eq. (19) above), each  $(M \times M)$  block of the  
 439 full  $(M|\mathcal{D}| \times M|\mathcal{D}|)$  empirical NTK is given by:

$$\hat{\Theta}_t(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j) I_M. \quad (20)$$

440 where  $I_M \in \mathbb{R}^{M \times M}$  is the identity. Eq. (20) gives the total effective interaction between the samples  
 441  $i$  and  $j$  from the dataset. For high-dimensional inputs  $\mathbf{x}$  drawn from an i.i.d standard Gaussian  
 442 distribution, we have  $\mathbf{x}_i^\top \mathbf{x}_j \approx \delta_{ij}$  by the central limit theorem. Therefore, in the special case of  
 443 a linear network with Gaussian i.i.d inputs, the representational dynamics (Lemma 2) simplify as  
 444 follows:

$$\begin{aligned} \frac{d\hat{\mathbf{z}}_{i,t}^{(1)}}{dt} &= -\eta \hat{\Theta}_t(\mathbf{x}_i, \mathcal{X}) \nabla_{\hat{\mathbf{z}}} \mathcal{L} \\ &= -\eta \hat{\Theta}_t(\mathbf{x}_i, \mathbf{x}_i) \nabla_{\hat{\mathbf{z}}_i} \mathcal{L} - \eta \sum_{j \neq i} \hat{\Theta}_t(\mathbf{x}_i, \mathbf{x}_j) \nabla_{\hat{\mathbf{z}}_j} \mathcal{L} \\ &= -\eta (\mathbf{x}_i^\top \mathbf{x}_i) \nabla_{\hat{\mathbf{z}}_i} \mathcal{L} - \eta \sum_{j \neq i} (\mathbf{x}_i^\top \mathbf{x}_j) \nabla_{\hat{\mathbf{z}}_j} \mathcal{L} \\ &= -\eta \nabla_{\hat{\mathbf{z}}_i} \mathcal{L} \quad . \end{aligned} \quad (21)$$

445 While the assumption of Gaussian i.i.d inputs is quite restrictive, we offer a generalizing interpretation  
 446 here. Specifically, the above argument also holds when the inputs  $\mathbf{x}$  are not all mutually orthogonal,  
 447 but fall into  $P$  orthogonal clusters in the input dataset. Then, we would have  $\mathbf{x}_i^\top \mathbf{x}_j = \delta_{p_i=p_j}$  where  
 448  $p_i$  is the “label” of the cluster corresponding to sample  $i$ . If  $\mathcal{P}_i$  is the number of all the samples with  
 449 the same label  $p_i$ , then Eq. (21) would simply be scaled to give  $\frac{d\hat{\mathbf{z}}_{i,t}^{(1)}}{dt} = -\eta \mathcal{P}_i \nabla_{\hat{\mathbf{z}}_i} \mathcal{L}$ .

450 For brevity, we proceed with the simplest case Eq. (21) in which every input is orthogonal. For  $\mathcal{L}^{\text{euc}}$ ,  
 451 the representational gradient  $\nabla_{\hat{\mathbf{z}}_i} \mathcal{L}$  is then given by:

$$\nabla_{\hat{\mathbf{z}}_i} \mathcal{L}^{\text{euc}} = \left( D_t \hat{\mathbf{z}}_{i,t}^{(1)} - \hat{\mathbf{z}}_{i,t}^{(2)} \right) D_t$$

452 Noting that  $D_t$  is just a diagonal matrix containing the eigenvalues  $\lambda_{m,t}$  and dropping the sample  
 453 subscript  $i$  for notational ease, we obtain for the  $m$ -th component of  $\nabla_{\hat{\mathbf{z}}} \mathcal{L}^{\text{euc}}$ :

$$\frac{\partial \mathcal{L}^{\text{euc}}}{\partial \hat{\mathbf{z}}_{m,t}} = \lambda_{m,t} (\lambda_{m,t} \hat{\mathbf{z}}_{m,t}^{(1)} - \hat{\mathbf{z}}_{m,t}^{(2)}) \quad .$$

454 Substituting this result in Eq. (21) gives us Eq. (8), the expression we were looking for. Finally,  
 455 introducing  $\hat{\mathbf{z}}_{m,t} \equiv \mathbb{E}[\hat{\mathbf{z}}_{m,t}^{(1)}] = \mathbb{E}[\hat{\mathbf{z}}_{m,t}^{(2)}]$  as the expectation over augmentations, we find that each  
 456 eigenmode evolves independently in expectation value as:

$$\begin{aligned} \mathbb{E} \left[ \frac{d\hat{\mathbf{z}}_{m,t}^{(1)}}{dt} \right] &= \frac{d\hat{\mathbf{z}}_{m,t}}{dt} = \eta \lambda_m \left( \mathbb{E}[\hat{\mathbf{z}}_{m,t}^{(2)}] - \lambda_m \mathbb{E}[\hat{\mathbf{z}}_{m,t}^{(1)}] \right) \\ &= \eta \lambda_{m,t} (1 - \lambda_{m,t}) \hat{\mathbf{z}}_{m,t} \quad . \end{aligned}$$

457 □

**Theorem 2.** (Representational dynamics under  $\mathcal{L}$ ) For a linear network with i.i.d Gaussian inputs trained with  $\mathcal{L}$ , the dynamics follow a system of  $M$  coupled differential equations:

$$\frac{d\hat{z}_m^{(1)}}{dt} = \eta \frac{\lambda_m}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \sum_{k \neq m} \lambda_k \left( \lambda_k (\hat{z}_k^{(1)})^2 \hat{z}_m^{(2)} - \lambda_m \hat{z}_m^{(1)} \hat{z}_k^{(1)} \hat{z}_k^{(2)} \right), \quad (10)$$

and, in the regime where eigenvalues are of comparable magnitude, the expected update over augmentations is well approximated by:

$$\frac{d\hat{z}_m}{dt} \approx \eta \lambda_m \cdot \mathbb{E} \left[ \frac{\hat{z}_m^2}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} \left[ \frac{\hat{z}_m}{\|\hat{z}\|} \right] \cdot \sum_{k \neq m} \lambda_k (\lambda_k - \lambda_m), \quad (11)$$

*Proof.* We can retrace the steps from the proof for Theorem 1 until Eq. (21):

$$\frac{d\hat{z}_{i,t}^{(1)}}{dt} = -\eta \nabla_{\hat{z}_i} \mathcal{L}.$$

$\nabla_{\hat{z}_i} \mathcal{L}$  is a vector of dimension  $M$ . Ignoring the subscripts  $i$  and  $t$  for simplicity, and focusing on the  $m$ -th component of  $\nabla_{\hat{z}_i} \mathcal{L}$ , we get:

$$\begin{aligned} \mathcal{L} &= - \sum_m^M \frac{\lambda_m \hat{z}_m^{(1)} \text{SG}(\hat{z}_m^{(2)})}{\|D\hat{z}^{(1)}\| \|\text{SG}(\hat{z}^{(2)})\|} \\ \Rightarrow \frac{\partial \mathcal{L}}{\partial \hat{z}_m^{(1)}} &= - \frac{\lambda_m \hat{z}_m^{(2)}}{\|D\hat{z}^{(1)}\| \|\hat{z}^{(2)}\|} + \frac{\sum_k \lambda_k \hat{z}_k^{(1)} \hat{z}_k^{(2)}}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \cdot \lambda_m^2 \hat{z}_m^{(1)} \\ &= - \frac{\lambda_m}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \left[ \|D\hat{z}^{(1)}\|^2 \hat{z}_m^{(2)} - \lambda_m \hat{z}_m^{(1)} \left( \sum_k \lambda_k \hat{z}_k^{(1)} \hat{z}_k^{(2)} \right) \right] \\ &= - \frac{\lambda_m}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \left[ \left( \sum_k \lambda_k^2 (\hat{z}_k^{(1)})^2 \right) \hat{z}_m^{(2)} - \lambda_m \hat{z}_m^{(1)} \left( \sum_k \lambda_k \hat{z}_k^{(1)} \hat{z}_k^{(2)} \right) \right] \\ &= - \frac{\lambda_m}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \sum_{k \neq m} \lambda_k \left( \lambda_k (\hat{z}_k^{(1)})^2 \hat{z}_m^{(2)} - \lambda_m \hat{z}_m^{(1)} \hat{z}_k^{(1)} \hat{z}_k^{(2)} \right) \\ \Rightarrow \frac{d\hat{z}_m^{(1)}}{dt} &= -\eta \frac{\partial \mathcal{L}}{\partial \hat{z}_m^{(1)}} \\ &= \frac{\eta \lambda_m}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \sum_{k \neq m} \lambda_k \left( \lambda_k (\hat{z}_k^{(1)})^2 \hat{z}_m^{(2)} - \lambda_m \hat{z}_m^{(1)} \hat{z}_k^{(1)} \hat{z}_k^{(2)} \right), \end{aligned}$$

proving Eq. (10). Assuming sufficiently small augmentations,  $\hat{z}_k^{(1)}$  and  $\hat{z}_k^{(2)}$  carry the same sign, and the net sign of both terms inside the parenthesis is fully determined by  $\gamma_m \equiv \text{sign}(\hat{z}_m^{(1)})$ . Hence, we may write:

$$\frac{d\hat{z}_m^{(1)}}{dt} = \frac{\eta \lambda_m \gamma_m}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \sum_{k \neq m} \left( \lambda_k^2 (\hat{z}_k^{(1)})^2 |\hat{z}_m^{(2)}| - \lambda_m \lambda_k |\hat{z}_m^{(1)}| |\hat{z}_k^{(1)}| |\hat{z}_k^{(2)}| \right).$$

It is useful to separate out  $\gamma_m$  in this manner because every other term in the expression is now non-negative. Then  $\text{sign}(\gamma_m \cdot \frac{d\hat{z}_m}{dt}) = \text{sign}(\hat{z}_m \cdot \frac{d\hat{z}_m}{dt})$  tells us whether  $\hat{z}_m$  tends to increase or decrease in magnitude, as we have argued in the main text.

**Asymptotic analysis.** To get a handle on how the different eigenvalues influence each other, we consider two important limiting cases. First, we consider the asymptotic regime dominated by one eigenvalue, and show that it tends towards a more symmetric solution in which the gap between different eigenvalues decreases. Second, we derive asymptotic expressions for the near-uniform regime in which all eigenvalues are comparable in size and show that this solution tends toward the uniform solution (cf. Eq. (11)).

477 To facilitate our analysis, we define each mode's relative contribution  $\chi_m \equiv \frac{|\hat{z}_m|}{\|\hat{z}\|}$  and evaluate  
 478 Eq. (10) taking the expectation value over augmentations:

$$\begin{aligned}
 \mathbb{E} \left[ \frac{d\hat{z}_m^{(1)}}{dt} \right] &= \eta \lambda_m \sum_{k \neq m} \left( \lambda_k^2 \cdot \mathbb{E} \left[ \frac{(\hat{z}_k^{(1)})^2 \hat{z}_m^{(2)}}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \right] - \lambda_m \lambda_k \cdot \mathbb{E} \left[ \frac{\hat{z}_m^{(1)} \hat{z}_k^{(1)} \hat{z}_k^{(2)}}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \right] \right) \\
 \frac{d\hat{z}_m}{dt} &= \eta \lambda_m \sum_{k \neq m} \left( \lambda_k^2 \cdot \mathbb{E} \left[ \frac{\hat{z}_k^2}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} \left[ \frac{\hat{z}_m}{\|\hat{z}\|} \right] - \lambda_m \lambda_k \cdot \mathbb{E} \left[ \frac{\hat{z}_m \hat{z}_k}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} \left[ \frac{\hat{z}_k}{\|\hat{z}\|} \right] \right) \\
 &= \eta \lambda_m \gamma_m \sum_{k \neq m} \left( \lambda_k^2 \cdot \mathbb{E} \left[ \chi_k^2 \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [\chi_m] - \lambda_m \lambda_k \cdot \mathbb{E} \left[ \chi_m \chi_k \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [\chi_k] \right).
 \end{aligned} \tag{22}$$

479 In the second equality we used the fact that the expectation value taken over augmentations is  
 480 conditioned on the input sample, which makes them conditionally independent.

481 **One dominant eigenvalue.** First, we consider the low-rank regime in which one eigenvalue  
 482 dominates. Without loss of generality, we assume  $\lambda_1 \gg \lambda_k \forall k \neq 1$ . We then have:

$$\begin{aligned}
 \chi_1 &\sim 1 \\
 \chi_k &\sim \epsilon \quad (0 < \epsilon \ll 1) \quad \forall \quad k \neq 1
 \end{aligned}$$

483 Plugging these values into Eq. (22) gives the following dynamics for the dominant eigenmode:

$$\begin{aligned}
 \frac{d\hat{z}_1}{dt} &\approx \eta \lambda_1 \gamma_1 \sum_{k \neq 1} \left( \lambda_k^2 \cdot \mathbb{E} \left[ \epsilon^2 \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [1] - \lambda_1 \lambda_k \cdot \mathbb{E} \left[ \epsilon \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \mathbb{E} [\epsilon] \right) \\
 &= \eta \lambda_1 \gamma_1 \sum_{k \neq 1} \left( \lambda_k^2 \epsilon^2 \cdot \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [1] - \lambda_1 \lambda_k \epsilon^2 \cdot \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \right) \\
 &= \eta \lambda_1 \gamma_1 \epsilon^2 \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \sum_{k \neq 1} \lambda_k (\lambda_k - \lambda_1) \quad .
 \end{aligned}$$

484 These updates are always opposite in sign to the representation component, which corresponds to  
 485 decaying dynamics for the leading eigenmode because  $\gamma_1 \frac{d\hat{z}_1}{dt} < 0$ .

486 For all other modes we have:

$$\begin{aligned}
 \frac{d\hat{z}_{m \neq 1}}{dt} &\approx \eta \lambda_m \gamma_m \sum_{k \notin \{m, 1\}} \left( \lambda_k^2 \epsilon^2 \cdot \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [\epsilon] - \lambda_m \lambda_k \epsilon^2 \cdot \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [\epsilon] \right) \\
 &\quad + \eta \lambda_m \gamma_m \left( \lambda_1^2 \cdot \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [\epsilon] - \lambda_m \lambda_1 \epsilon \cdot \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [1] \right) \\
 &= \eta \lambda_m \gamma_m \epsilon \cdot \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \left( \lambda_1 (\lambda_1 - \lambda_m) + \epsilon^2 \sum_{k \notin \{m, 1\}} \lambda_k (\lambda_k - \lambda_m) \right) \\
 &\approx \eta \lambda_m \gamma_m \lambda_1 \epsilon \cdot \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] (\lambda_1 - \lambda_m) \quad ,
 \end{aligned}$$

487 so that  $\gamma_m \frac{d\hat{z}_m}{dt} > 0$ , i.e., the updates have the *same* sign as the representation component, which  
 488 corresponds to growth dynamics. In other words: The dominant eigenvalue ‘‘pulls all the other  
 489 eigenvalues up,’’ a form of implicit cooperation between the eigenmodes. We also note that the  
 490 non-dominant eigenmodes increase at a rate proportional to  $\epsilon$ , whereas the dominant eigenmode  
 491 decreases at a slower rate proportional to  $\epsilon^2$ . Thus, for sensible initializations with at least one  
 492 large and many small eigenvalues, the modes will tend toward an equilibrium at some non-zero  
 493 intermediate value, without a dominant mode. Next we study this other limiting case in which all  
 494 eigenvalues are of similar size.

495 **Near-uniform regime.** To study the dynamics in a near-uniform regime, we note that all  $\chi_m$  are  
 496 of order  $\mathcal{O}(1)$  in  $\hat{z}_m$ , whereas the eigenvalues  $\lambda_m$  are of order  $\mathcal{O}(\hat{z}_m^2)$ . In this setting, the effect of  
 497 the eigenvalue terms  $\lambda_m$  on the dynamics is stronger than the  $\chi_m$  terms which are bounded between  
 498 0 and 1. With a sufficiently high-dimensional representation, all  $\chi_m$  terms will be centered around  
 499  $1/\sqrt{M}$ . Based on these observations, we may make the simplifying assumption that the contributions  
 500 are all approximately equal, i.e,  $\chi_i = \chi$  for all  $i$ . Substituting this value in Eq (22) gives:

$$\frac{d\hat{z}_m}{dt} = \eta \lambda_m \gamma_m \cdot \mathbb{E} \left[ \chi^2 \frac{\|\hat{\mathbf{z}}\|^2}{\|D\hat{\mathbf{z}}\|^3} \right] \cdot \mathbb{E} [\chi] \cdot \sum_{k \neq m} \lambda_k (\lambda_k - \lambda_m) \quad . \quad (23)$$

501 Finally, substituting for  $\chi$ , which by assumption are all approximately equal:

$$\chi \approx \chi_m = \frac{|\hat{z}_m|}{\|\hat{\mathbf{z}}\|} \quad ,$$

502 and absorbing back the sign from  $\gamma_m$ , we obtain the approximate dynamics in Eq (11):

$$\frac{d\hat{z}_m}{dt} \approx \eta \lambda_m \cdot \mathbb{E} \left[ \frac{\hat{z}_m^2}{\|D\hat{\mathbf{z}}\|^3} \right] \cdot \mathbb{E} \left[ \frac{\hat{z}_m}{\|\hat{\mathbf{z}}\|} \right] \cdot \sum_{k \neq m} \lambda_k (\lambda_k - \lambda_m)$$

503

□

## 504 C Derivation of idealized learning dynamics for different loss variations

### 505 C.1 Removing the stop-grad from the Euclidean loss $\mathcal{L}^{\text{euc}}$

506 Omitting the stop-grad operator from  $\mathcal{L}^{\text{euc}}$  gives:

$$\begin{aligned}\mathcal{L}_{\text{noSG}}^{\text{euc}} &= \frac{1}{2} \|W_P \mathbf{z}^{(1)} - \mathbf{z}^{(2)}\|^2 \\ &= \frac{1}{2} \sum_m^M |\lambda_m \hat{z}_m^{(1)} - \hat{z}_m^{(2)}|^2 \quad .\end{aligned}$$

507 Tracing the steps to prove Theorem 1 and assuming Gaussian i.i.d inputs for a linear network, we  
508 write:

$$\begin{aligned}\frac{\partial \mathcal{L}_{\text{noSG}}^{\text{euc}}}{\partial \hat{z}_m} &= \frac{\partial \mathcal{L}_{\text{noSG}}^{\text{euc}}}{\partial \hat{z}_m^{(1)}} + \frac{\partial \mathcal{L}_{\text{noSG}}^{\text{euc}}}{\partial \hat{z}_m^{(2)}} \\ &= \left( \lambda_m \hat{z}_m^{(1)} - \hat{z}_m^{(2)} \right) \lambda_m - \left( \lambda_m \hat{z}_m^{(1)} - \hat{z}_m^{(2)} \right) \\ &= \left( \lambda_m \hat{z}_m^{(1)} - \hat{z}_m^{(2)} \right) (\lambda_m - 1) \\ \Rightarrow \frac{d\hat{z}_m}{dt} &= -\eta \mathbb{E} \left[ \frac{\partial \mathcal{L}_{\text{noSG}}^{\text{euc}}}{\partial \hat{z}_m} \right] \\ &= -\eta \left( \lambda_m \mathbb{E}[\hat{z}_m^{(1)}] - \mathbb{E}[\hat{z}_m^{(2)}] \right) (\lambda_m - 1) \\ &= -\eta (1 - \lambda_m)^2 \hat{z}_m \quad ,\end{aligned}$$

509 which results in decaying representations and thus collapse.

### 510 C.2 Removing the stop-grad from the Cosine loss $\mathcal{L}$

511 Following the same arguments as above, omitting the stop-grad operator from  $\mathcal{L}$  gives:

$$\begin{aligned}\mathcal{L}_{\text{noSG}} &= -\frac{(W_P \mathbf{z}^{(1)})^\top \mathbf{z}^{(2)}}{\|W_P \mathbf{z}^{(1)}\| \|\mathbf{z}^{(2)}\|} \\ \Rightarrow \frac{\partial \mathcal{L}_{\text{noSG}}}{\partial \hat{z}_m} &= \frac{-\lambda_m}{\|D\hat{\mathbf{z}}^{(1)}\|^3 \|\hat{\mathbf{z}}^{(2)}\|} \sum_{k \neq m} \left( \lambda_k^2 (\hat{z}_k^{(1)})^2 \hat{z}_m^{(2)} - \lambda_m \lambda_k \hat{z}_m^{(1)} \hat{z}_k^{(1)} \hat{z}_k^{(2)} + \lambda_k^2 \lambda_m (\hat{z}_k^{(1)})^3 - \lambda_k \hat{z}_m^{(2)} \hat{z}_k^{(2)} \hat{z}_k^{(1)} \right) \\ &\quad + \frac{-\lambda_m}{\|D\hat{\mathbf{z}}^{(1)}\|^3 \|\hat{\mathbf{z}}^{(2)}\|} \left( \lambda_m^3 (\hat{z}_m^{(1)})^3 - \lambda_m (\hat{z}_m^{(2)})^2 \hat{z}_m^{(1)} \right) \quad ,\end{aligned}$$

512 so that, when taking the expectation value over augmentations, the dynamics follow:

$$\begin{aligned}\frac{d\hat{z}_m}{dt} &= -\eta \mathbb{E} \left[ \frac{\partial \mathcal{L}_{\text{noSG}}}{\partial \hat{z}_m} \right] \\ &= \eta \lambda_m \gamma_m \sum_{k \neq m} \lambda_k \left( \lambda_k \cdot \mathbb{E} \left[ \chi_k^2 \frac{\|\hat{\mathbf{z}}\|^2}{\|D\hat{\mathbf{z}}\|^3} \right] \cdot \mathbb{E}[\chi_m] - \lambda_m \cdot \mathbb{E} \left[ \chi_m \chi_k \frac{\|\hat{\mathbf{z}}\|^2}{\|D\hat{\mathbf{z}}\|^3} \right] \cdot \mathbb{E}[\chi_k] \right) \\ &\quad + \eta \lambda_m \gamma_m \sum_{k \neq m} \lambda_k \left( \lambda_m \lambda_k \mathbb{E} \left[ \chi_k^3 \frac{\|\hat{\mathbf{z}}\|^3}{\|D\hat{\mathbf{z}}\|^3} \right] \cdot \mathbb{E} \left[ \frac{1}{\|\hat{\mathbf{z}}\|} \right] - \mathbb{E} \left[ \chi_k \frac{\|\hat{\mathbf{z}}\|}{\|D\hat{\mathbf{z}}\|^3} \right] \cdot \mathbb{E}[\chi_m \chi_k \|\hat{\mathbf{z}}\|] \right) \\ &\quad + \eta \lambda_m^2 \gamma_m \left( \lambda_m^2 \mathbb{E} \left[ \chi_m^3 \frac{\|\hat{\mathbf{z}}\|^3}{\|D\hat{\mathbf{z}}\|^3} \right] \cdot \mathbb{E} \left[ \frac{1}{\|\hat{\mathbf{z}}\|} \right] - \mathbb{E} \left[ \chi_m \frac{\|\hat{\mathbf{z}}\|}{\|D\hat{\mathbf{z}}\|^3} \right] \cdot \mathbb{E}[\chi_m^2 \|\hat{\mathbf{z}}\|] \right) \quad .\end{aligned}$$



513 In the asymptotic regime with dominant eigenvalue  $\lambda_1$ , we get the dynamics:

$$\begin{aligned}
\frac{d\hat{z}_1}{dt} &= \eta\lambda_1\gamma_1 \sum_{k \neq m} \lambda_k \left( \lambda_k \epsilon^2 \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] - \lambda_m \epsilon^2 \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \right) \\
&\quad + \eta\lambda_1\gamma_1 \sum_{k \neq m} \lambda_k \left( \lambda_1 \lambda_m \epsilon^3 \mathbb{E} \left[ \frac{\|\hat{z}\|^3}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} \left[ \frac{1}{\|\hat{z}\|} \right] - \epsilon^2 \mathbb{E} \left[ \frac{\|\hat{z}\|}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [\|\hat{z}\|] \right) \\
&\quad + \eta\lambda_1^2\gamma_1 \left( \lambda_1^2 \cdot \mathbb{E} \left[ \frac{\|\hat{z}\|^3}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} \left[ \frac{1}{\|\hat{z}\|} \right] - \mathbb{E} \left[ \frac{\|\hat{z}\|}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [\|\hat{z}\|] \right) \\
&\approx \eta\lambda_1^4\gamma_1 \cdot \mathbb{E} \left[ \frac{\|\hat{z}\|^3}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} \left[ \frac{1}{\|\hat{z}\|} \right] \\
\frac{d\hat{z}_{m \neq 1}}{dt} &= \eta\lambda_m\gamma_m \sum_{k \notin \{m,1\}} \lambda_k \left( \lambda_k \epsilon^3 \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] - \lambda_m \epsilon^3 \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \right) \\
&\quad + \eta\lambda_m\gamma_m\lambda_1 \left( \lambda_1 \epsilon \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] - \lambda_m \epsilon \mathbb{E} \left[ \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \right) \\
&\quad + \eta\lambda_m\gamma_m \sum_{k \notin \{m,1\}} \lambda_k \left( \lambda_m \lambda_k \epsilon^3 \mathbb{E} \left[ \frac{\|\hat{z}\|^3}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} \left[ \frac{1}{\|\hat{z}\|} \right] - \epsilon^3 \mathbb{E} \left[ \frac{\|\hat{z}\|}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [\|\hat{z}\|] \right) \\
&\quad + \eta\lambda_m\gamma_m\lambda_1 \left( \lambda_m \lambda_1 \mathbb{E} \left[ \frac{\|\hat{z}\|^3}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} \left[ \frac{1}{\|\hat{z}\|} \right] - \epsilon \mathbb{E} \left[ \frac{\|\hat{z}\|}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [\|\hat{z}\|] \right) \\
&\quad + \eta\lambda_m^2\gamma_m \left( \lambda_m^2 \epsilon^3 \mathbb{E} \left[ \frac{\|\hat{z}\|^3}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} \left[ \frac{1}{\|\hat{z}\|} \right] - \epsilon^3 \mathbb{E} \left[ \frac{\|\hat{z}\|}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [\|\hat{z}\|] \right) \\
&\approx \eta\lambda_m^2\gamma_m\lambda_1^2 \cdot \mathbb{E} \left[ \frac{\|\hat{z}\|^3}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} \left[ \frac{1}{\|\hat{z}\|} \right],
\end{aligned}$$

514 Thus, all eigenmodes diverge because  $\gamma_m \frac{d\hat{z}_m}{dt} > 0$ .

515 Similarly, we find divergent dynamics when starting in the near-uniform regime:

$$\begin{aligned}
\frac{d\hat{z}_m}{dt} &= \eta\lambda_m\gamma_m \mathbb{E} \left[ \chi^2 \frac{\|\hat{z}\|^2}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [\chi] \sum_{k \neq m} \lambda_k (\lambda_k - \lambda_m) \\
&\quad + \eta\lambda_m^2\gamma_m \mathbb{E} \left[ \chi^3 \frac{\|\hat{z}\|^3}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} \left[ \frac{1}{\|\hat{z}\|} \right] \sum_k \lambda_k^2 \\
&\quad - \eta\lambda_m\gamma_m \mathbb{E} \left[ \chi \frac{\|\hat{z}\|}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} [\chi^2 \|\hat{z}\|] \sum_k \lambda_k \\
&\approx \eta\lambda_m^2\gamma_m \mathbb{E} \left[ \chi^3 \frac{\|\hat{z}\|^3}{\|D\hat{z}\|^3} \right] \cdot \mathbb{E} \left[ \frac{1}{\|\hat{z}\|} \right] \sum_k \lambda_k^2,
\end{aligned}$$

516 selecting the terms with the highest power in the eigenvalues.

517 Thus, omission of stop-grad precludes successful representation learning for both the Euclidean and  
518 the cosine loss, but due to different mechanisms. Euclidean loss yields collapse, whereas the Cosine  
519 loss succumbs to run-away activity.

### 520 C.3 Removing the predictor from the Euclidean loss $\mathcal{L}^{\text{euc}}$

521 To analyze the representational dynamics in the absence of the predictor network, we consider  
522  $\mathcal{L}_{\text{noPred}}^{\text{euc}}$ :

$$\begin{aligned}
\mathcal{L}_{\text{noPred}}^{\text{euc}} &= \frac{1}{2} \|\mathbf{z}^{(1)} - \text{SG}(\mathbf{z}^{(2)})\|^2 \\
&= \frac{1}{2} \sum_m^M |\hat{z}_m^{(1)} - \text{SG}(\hat{z}_m^{(2)})|^2.
\end{aligned}$$

523 The dynamics resulting from this loss function are a special case of the dynamics derived in Theorem 1  
 524 with all the eigenvalues equal to one ( $\lambda_k = 1$ ). In particular Eq. (8) becomes:

$$\frac{d\hat{z}_{m,t}^{(1)}}{dt} = -\eta \frac{\partial \mathcal{L}_{\text{noPred}}^{\text{euc}}}{\partial \hat{z}_{m,t}^{(1)}}(t) = \eta \left( \hat{z}_{m,t}^{(2)} - \hat{z}_{m,t}^{(1)} \right),$$

525 which evaluates to 0 under expectation over augmentations. Hence there is no learning without the  
 526 predictor.

#### 527 C.4 Isotropic losses for equalized convergence rates

528 In Expressions (9) and (11) we see that the overall learning dynamics have a quadratic dependence  
 529 on the eigenvalues with a root near collapsed solutions, which causes these modes to learn slower.  
 530 We reasoned that this anisotropy could be detrimental for learning. To address this issue, we sought  
 531 to derive alternative loss functions that encourage isotropic learning dynamics for all modes.

##### 532 C.4.1 Euclidean IsoLoss.

533 We start by deriving an IsoLoss function for the Euclidean case  $\mathcal{L}^{\text{euc}}$ . To avoid the unwanted quadratic  
 534 dependence, we first note that we would like to arrive at the following expression for the dynamics:

$$\frac{d\hat{z}_m}{dt} = \eta (1 - \lambda_m) \hat{z}_m \quad .$$

535 By recalling the Euclidean loss and corresponding dynamics:

$$\mathcal{L}^{\text{euc}} = \frac{1}{2} \sum_m^M |\lambda_m \hat{z}_m^{(1)} - \text{SG}(\hat{z}_m^{(2)})|^2 \Rightarrow \frac{d\hat{z}_m}{dt} = \eta \lambda_m (1 - \lambda_m) \hat{z}_m \quad ,$$

536 we note that the leading  $\lambda_m$  term has no influence on the overall sign of the dynamics, and is  
 537 introduced by the second step in the chain rule:

$$\frac{\partial \mathcal{L}^{\text{euc}}}{\partial \hat{z}_m^{(1)}} = (\lambda_m \hat{z}_m^{(1)} - \hat{z}_m^{(2)}) \cdot \frac{\partial}{\partial \hat{z}_m^{(1)}} (\lambda_m \hat{z}_m^{(1)} - \hat{z}_m^{(2)}) \quad .$$

538 Based on this realization we see that this second step needs to be modified. To that end, we start with  
 539 the desired derivative:

$$\frac{\partial \mathcal{L}_{\text{iso}}^{\text{euc}}}{\partial \hat{z}_m^{(1)}} = (\lambda_m \hat{z}_m^{(1)} - \hat{z}_m^{(2)}) \cdot \frac{\partial}{\partial \hat{z}_m^{(1)}} (\hat{z}_m^{(1)} - \hat{z}_m^{(2)}) \quad ,$$

540 and see that several loss functions are possible. The one we have reported in Eq. (15) we derived by  
 541 applying an appropriate stop-gradient while integrating:

$$\frac{\partial \mathcal{L}_{\text{iso}}^{\text{euc}}}{\partial \hat{z}_m^{(1)}} = (\hat{z}_m^{(1)} + \lambda_m \hat{z}_m^{(1)} - \hat{z}_m^{(2)} - \hat{z}_m^{(1)}) \cdot \frac{\partial}{\partial \hat{z}_m^{(1)}} (\hat{z}_m^{(1)} - \hat{z}_m^{(2)}) \quad .$$

542 to give:

$$\mathcal{L}_{\text{iso}}^{\text{euc}} = \frac{1}{2} \sum_m^M |\hat{z}_m^{(1)} - \text{SG}(\hat{z}_m^{(2)} + \hat{z}_m^{(1)} - \lambda_m \hat{z}_m^{(1)})|^2$$

543 Another alternative loss with the same desired isotropic learning dynamics, but using a different  
 544 placement of the stop-gradient operators, is given by:

$$\mathcal{L}_{\text{iso}}^{\text{euc}} = \sum_m^M \text{SG} \left( \lambda_m \hat{z}_m^{(1)} - \hat{z}_m^{(2)} \right) \cdot \left( \hat{z}_m^{(1)} - \text{SG}(\hat{z}_m^{(2)}) \right)$$

##### 545 C.4.2 Cosine Similarity IsoLoss.

546 Since most practical SSL approaches rely on cosine similarity, which suffers from a similar anisotropy  
 547 of the learning dynamics, we sought to find IsoLosses in this setting. With the same goal as above,  
 548 we would like to arrive at the dynamics:

$$\frac{d\hat{z}_m}{dt} = \eta \frac{\hat{z}_m^{(2)}}{\|D\hat{z}^{(1)}\| \|\hat{z}^{(2)}\|} - \eta \frac{\sum_k \lambda_k \hat{z}_k^{(1)} \hat{z}_k^{(2)}}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \cdot \lambda_m \hat{z}_m^{(1)}$$

549 starting from the cosine loss and corresponding dynamics:

$$\mathcal{L} = - \sum_m^M \frac{\lambda_m \hat{z}_m^{(1)} \text{SG}(\hat{z}_m^{(2)})}{\|D\hat{z}^{(1)}\| \|\hat{z}^{(2)}\|} \quad (24)$$

$$\Rightarrow \frac{d\hat{z}_m}{dt} = \eta \frac{\lambda_m \hat{z}_m^{(2)}}{\|D\hat{z}^{(1)}\| \|\hat{z}^{(2)}\|} - \eta \frac{\sum_k \lambda_k \hat{z}_k^{(1)} \hat{z}_k^{(2)}}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \cdot \lambda_m^2 \hat{z}_m^{(1)} \quad (25)$$

550 The IsoLoss in this case can be derived by noting how  $\lambda_m$  arises in each of the two terms in Eq. (25),  
551 and engineering an alternative loss function corresponding to each term separately.

552 In the first term,  $\lambda_m$  arises from the partial derivative of the numerator  $\lambda_m \hat{z}_m^{(1)} \text{SG}(\hat{z}_m^{(2)})$  in the original  
553 loss (Eq. (24)). This can be remediated by using  $\hat{z}_m^{(1)} \text{SG}(\hat{z}_m^{(2)})$  as the numerator instead.

554 In the second term in Eq. (25),  $\lambda_m^2$  arises from the partial derivative of  $\|D\hat{z}^{(1)}\| = \sqrt{\sum_k (\lambda_k \hat{z}_k^{(1)})^2}$  in  
555 the denominator. We can reduce  $\lambda_m^2$  to  $\lambda_m$  by instead taking the partial derivative of  $\|D^{1/2}\hat{z}^{(1)}\| =$   
556  $\sqrt{\sum_k (\lambda_k^{1/2} \hat{z}_k^{(1)})^2}$ .

557 Putting these insights together, we arrive at the desired partial derivative:

$$\begin{aligned} \frac{\partial \mathcal{L}_{\text{iso}}}{\partial \hat{z}_m^{(1)}} &= \frac{-1}{\|D\hat{z}^{(1)}\| \|\hat{z}^{(2)}\|} \cdot \frac{\partial \hat{z}_m^{(1)} \hat{z}_m^{(2)}}{\partial \hat{z}_m^{(1)}} + \frac{\sum_k \lambda_k \hat{z}_k^{(1)} \hat{z}_k^{(2)}}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \cdot \frac{1}{2} \frac{\partial \lambda_m (\hat{z}_m^{(1)})^2}{\partial \hat{z}_m^{(1)}} \\ &= \frac{-1}{\|D\hat{z}^{(1)}\| \|\hat{z}^{(2)}\|} \cdot \frac{\partial (\hat{z}^{(1)})^\top \hat{z}^{(2)}}{\partial \hat{z}_m^{(1)}} + \frac{\sum_k \lambda_k \hat{z}_k^{(1)} \hat{z}_k^{(2)}}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \cdot \frac{1}{2} \frac{\partial \|D^{1/2}\hat{z}^{(1)}\|^2}{\partial \hat{z}_m^{(1)}} \quad , \end{aligned}$$

558 and the integrated IsoLoss in eigenspace:

$$\mathcal{L}_{\text{iso}} = -(\hat{z}^{(1)})^\top \text{SG} \left( \frac{\hat{z}^{(2)}}{\|D\hat{z}^{(1)}\| \|\hat{z}^{(2)}\|} \right) + \frac{1}{2} \text{SG} \left( \frac{(D\hat{z}^{(1)})^\top \hat{z}^{(2)}}{\|D\hat{z}^{(1)}\|^3 \|\hat{z}^{(2)}\|} \right) \|D^{1/2}\hat{z}^{(1)}\|^2 \quad .$$

559 Rotating all terms back to the original space gives the desired IsoLoss for Cosine similarity as reported  
560 (Eq. (17)):

$$\mathcal{L}_{\text{iso}} = -(\mathbf{z}^{(1)})^\top \text{SG} \left( \frac{\mathbf{z}^{(2)}}{\|W_P \mathbf{z}^{(1)}\| \|\mathbf{z}^{(2)}\|} \right) + \frac{1}{2} \text{SG} \left( \frac{(W_P \mathbf{z}^{(1)})^\top \mathbf{z}^{(2)}}{\|W_P \mathbf{z}^{(1)}\|^3 \|\mathbf{z}^{(2)}\|} \right) \|W_P^{1/2} \mathbf{z}^{(1)}\|^2 \quad .$$

## D Experimental details

**Self-supervised pretraining.** We used the CIFAR-10, CIFAR-100 [3], STL-10 [4], and TinyImageNet [5] datasets for self-supervised learning with a ResNet-18 [6] encoder and the SimCLR set of transformations [7]. We also adopted several modifications of ResNet-18 and the augmentation set which have been proposed to deal with the low resolution of the images in these datasets [7]. The ResNet modifications comprise using  $3 \times 3$  convolutional kernels instead of  $7 \times 7$  kernels and skipping the first max-pooling operation. The modifications to the standard SimCLR augmentations are excluding the blur transformation and using a weaker color jitter strength of 0.5. The configurations we used for each dataset are summarized in Table 3. We used BatchNorm in the backbone and the projector MLP in the hidden layer for all methods. For BYOL, we included BatchNorm also in the hidden layer of the predictor MLP.

As stated in the main text, we used SGD with learning rate 0.1, momentum 0.9 and weight decay  $4 \times 10^{-4}$ . Furthermore, we used a warmup period of 10 epochs for the learning rate followed by a cosine decay schedule and a batch size of 512. For the EMA, we started with  $\tau_{\text{base}} = 0.996$  and increased  $\tau_{\text{EMA}}$  to 1 with a cosine schedule exactly following the configuration reported in [8]. For DirectPred, we used  $\alpha = 0.5$ ,  $\tau = 0.3$  for the moving average estimate of the correlation matrix updated at every step, and clipped the eigenvalues of the correlation matrix at  $10^6$ .

Table 3:

	CIFAR-10	CIFAR-100	STL-10	TinyImageNet
Resolution	$32 \times 32$	$32 \times 32$	$96 \times 96$	$64 \times 64$
Kernel size	$3 \times 3$	$3 \times 3$	$7 \times 7$	$3 \times 3$
First max-pool	No	No	Yes	Yes
Blur	No	No	Yes	No
Color jitter	0.5	0.5	1.0	0.5

**Linear evaluation protocol.** We reported the held-out classification accuracy on the test sets for CIFAR-10/100 and STL-10, and the validation set for TinyImageNet, after training the linear classifier on frozen features for all labeled examples available in each training set.

## References

- [1] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide Neural Networks of Any Depth Evolve as Linear Models Under Gradient Descent. *Advances in Neural Information Processing Systems*, 32:8572–8583, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/0d1a9651497a38d8b1c3871c84528bd4-Abstract.html>.
- [2] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- [3] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images, 2009.
- [4] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223. JMLR Workshop and Conference Proceedings, 2011.
- [5] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.

601 [8] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena  
602 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,  
603 et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural*  
604 *information processing systems*, 33:21271–21284, 2020.