

APPENDIX

A TRAINING STATISTICS

We provide some examples demonstrating the effectiveness of introducing linear and batch normalization layers in the projector in Fig. 8. This contributes to prevent excessing increase of mean or variance. An alternative and promising approach to the linear and batch normalization layers is to penalize the norm of the representations directly in the objective. We leave this to future investigation.

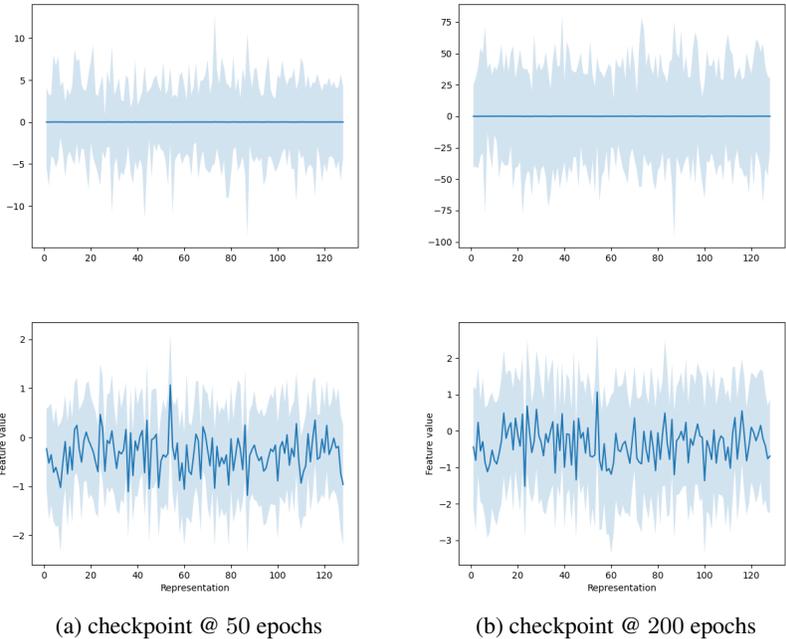


Figure 8: Example of mean and standard deviation statistics for the representation features obtained by the backbone network trained on SVHN data. Statistics are computed for a batch of size of 100 samples. Results corresponds to checkpoints for the projector ($c = 4096$) without (**top**) and with (**bottom**) linear and batch normalization layers. Linear and batch normalization layers contribute to stabilize the training by avoiding mean or variance increase.

B DISCUSSION ON FINITE CAPACITY

It is important to mention that the global minima for the *FALCON* objective might not be reached when using a backbone network of finite and small capacity. In this case, the avoidance of representation and cluster collapses can still be guaranteed when the invariance and the matching prior losses are both minimized. Indeed, we observe that for representation collapse $p_{ij} = p_j$ for all $i \in [n], j \in [c]$ (i.e. the outputs of the overall network are constant with respect to their inputs) and that the corresponding minimum value of the objective is given by the following formula

$$\mathcal{L}_{FALCON}(\mathcal{D}) = \beta H(\mathbf{p}) + CE(\mathbf{q}, \mathbf{p})$$

where the first addend arises from the invariance loss, whereas the second one arises from the matching prior one. Notably, the two terms cannot be minimized at the same time due to their competitive nature. For instance, in the case of uniform \mathbf{q} , the solution of $\mathbf{p} = \mathbf{q}$ is a minimum for the matching prior loss but not for the invariance one (this is actually a saddle point, as corresponding to the maximum for the entropy term in the above equation).

Cluster collapse occurs whenever $\exists j, k \neq j \in [c]$ such that for all $i \in [n], p_{ij} \leq p_{ik}$. The minimization of the invariance loss forces the whole network to make low entropy predictions,

864 whereas the minimization of the matching prior loss forces to distribute these predictions across all
865 codes according to q . Hence, when both losses are minimized cluster collapse is avoided.
866

867 C MINIMA OF THE *FALCON* LOSS

869 *Proof.* We recall here the loss
870

$$871 \mathcal{L}_{FALCON}(\mathcal{D}) = -\frac{\beta}{n} \sum_{i=1}^n \sum_{j=1}^c p_{ij} \log p'_{ij} - \sum_{j=1}^c q_j \log \frac{1}{n} \sum_{i=1}^n p_{ij}$$

874 and prove all optimality conditions. Before doing that, we observe that the loss is convex w.r.t. \mathbf{P}
875 when \mathbf{P}' is fixed, as the first addend is a sum of linear terms, whereas the second addend is a sum
876 of convex terms. Similarly, we observe that convexity holds w.r.t. \mathbf{P}' when \mathbf{P} is fixed by exploiting
877 the same reasoning. However, it is important to mention that the loss is not convex globally. This
878 can be shown firstly by computing the Hessian of the first addend w.r.t. both \mathbf{P} and \mathbf{P}' and secondly
879 by observing that the Hessian is not positive semi-definite (we skip the tedious calculation of the
880 Hessian).

881 **Invariance.** We observe that \mathbf{P}' appears only in the first addend of \mathcal{L}_{FALCON} and that this addend can
882 be equivalently rewritten in the following way:

$$883 -\frac{\beta}{n} \sum_{i=1}^n \sum_{j=1}^c p_{ij} \log p'_{ij} = -\frac{\beta}{n} \sum_{i=1}^n \sum_{j=1}^c p_{ij} \log p_{ij} - \frac{\beta}{n} \sum_{i=1}^n \sum_{j=1}^c p_{ij} \log \frac{p'_{ij}}{p_{ij}} \\ 884 = \frac{\beta}{n} \sum_{i=1}^n H(\mathbf{p}_i) + \frac{\beta}{n} \sum_{i=1}^n KL(\mathbf{p}_i \parallel \mathbf{p}'_i) \quad (6)$$

889 where $H(\cdot)$, $KL(\cdot)$ are the entropy and Kullback-Leibler divergence, respectively. Therefore min-
890 imizing \mathcal{L}_{FALCON} w.r.t. \mathbf{P}' is equivalent to minimizing Eq. 6. The solution is given by $\mathbf{p}_i = \mathbf{p}'_i$,
891 $\forall i \in [n]$, thus proving the invariance condition.

892 **Extrema.** We first leverage the invariance condition, $\mathbf{p}_i = \mathbf{p}'_i$, $\forall i \in [n]$, and rewrite \mathcal{L}_{FALCON}
893 accordingly:
894

$$895 \mathcal{L}_{FALCON}(\mathcal{D}) = -\frac{\beta}{n} \sum_{i=1}^n \sum_{j=1}^c p_{ij} \log p_{ij} - \sum_{j=1}^c q_j \log \frac{1}{n} \sum_{i=1}^n p_{ij} \quad (7)$$

899 We observe that the loss in Eq. 7 is convex w.r.t. \mathbf{P} . Therefore, we can obtain its optimality
900 conditions, by deriving the closed-form solutions for the minima of the second addend in Eq. 7
901 and then constraining the optimization of the first addend with these solutions and deriving the
902 corresponding minima.

903 Let's start by considering the following constrained convex minimization problem, obtained from the
904 first addend in Eq. 7 with n, β being dropped as being constant for the optimization:
905

$$906 \min_{\mathbf{P}} - \sum_{i=1}^n \sum_{j=1}^c p_{ij} \log p_{ij} \\ 907 \text{s.t.} \quad \sum_{j=1}^c p_{ij} = 1, \quad \forall i \in [n] \\ 908 \epsilon \leq p_{ij} \leq 1 - \epsilon(c-1), \quad \forall i \in [n], j \in [c], \quad (8)$$

914 and the corresponding Lagrangian with multipliers $\mathbf{\Lambda}, \mathbf{\Delta} \in \mathbb{R}_+^{n \times c}$, $\boldsymbol{\nu} \in \mathbb{R}^n$ is:

$$915 \mathcal{L}_1(\mathbf{P}; \mathbf{\Lambda}, \mathbf{\Delta}, \boldsymbol{\nu}) \equiv - \sum_{i=1}^n \sum_{j=1}^c p_{ij} \log p_{ij} + \sum_{i=1}^n \nu_i \left(\sum_{j=1}^c p_{ij} - 1 \right) +$$

$$+ \sum_{i=1}^n \sum_{j=1}^c [\lambda_{ij}(\epsilon - p_{ij}) + \delta_{ij}(p_{ij} - 1 + \epsilon(c-1))] \quad (9)$$

We observe that the Lagrangian is constructed so as to satisfy the following relation

$$-\sum_{i=1}^n \sum_{j=1}^c p_{ij} \log p_{ij} \geq \mathcal{L}_1(\mathbf{P}; \mathbf{\Lambda}, \mathbf{\Delta}, \mathbf{\Omega}, \nu) \quad (10)$$

Let's maximize \mathcal{L}_1 w.r.t. \mathbf{P} by setting $\nabla_{p_{ij}} \mathcal{L}_1 = 0$. This leads to the following closed-form expression:

$$p_{ij}^* = e^{-1-\lambda_{ij}+\nu_i+\delta_{ij}} \quad \forall i \in [n], j \in [c] \quad (11)$$

By evaluating \mathcal{L}_1 at the solutions in Eq. [11](#), we obtain the Lagrange dual function

$$\mathcal{L}_1(\mathbf{P}^*; \mathbf{\Lambda}, \mathbf{\Delta}, \mathbf{\Omega}, \nu) = n + \sum_{i=1}^n \left\{ -\nu_i + \sum_{j=1}^c [\lambda_{ij}\epsilon - \delta_{ij}(1 - \epsilon(c-1))] \right\} \quad (12)$$

The Lagrange multipliers in Eq. [12](#) depend on the values of \mathbf{P}^* through the Karush-Kuhn-Tucker (KKT) conditions. We distinguish two main cases for \mathbf{P}^* , each leading to different evaluation of the Lagrange dual function:

- *Case 1.* When all probability values touch their extrema, such as

$$\forall i \in [n], \exists! j \in [c], \forall k \in [c] \text{ with } k \neq j \text{ s.t. } p_{ij}^* = 1 - \epsilon(c-1) \text{ and } p_{ik}^* = \epsilon$$

By the KKT conditions (i.e. complementary slackness), we have that $\lambda_{ij} = 0$ and $\delta_{ik} = 0$, whereas $\lambda_{ik} \geq 0, \delta_{ij} \geq 0$. By substituting these conditions in Eq. [12](#), we obtain that

$$\mathcal{L}_1(\mathbf{P}^*; \mathbf{\Lambda}, \mathbf{\Delta}, \mathbf{\Omega}, \nu)|_{\{\lambda_{ij}=\delta_{ik}=0\}} = n + \sum_{i=1}^n \left\{ -\nu_i - \delta_{ij}(1 - \epsilon(c-1)) + \sum_{k \neq j} \lambda_{ik}\epsilon \right\} \quad (13)$$

By taking into account also Eq. [11](#), we have that $\forall i \in [n], \exists! j \in [c], \forall k \in [c]$

$$\delta_{ij} = 1 - \nu_i + \log(1 - \epsilon(c-1)) \text{ and } \lambda_{ik} = -1 + \nu_i - \log \epsilon \quad (14)$$

And by substituting Eq. [14](#) into Eq. [13](#), we obtain that

$$\mathcal{L}_1(\mathbf{P}^*; \mathbf{\Lambda}, \mathbf{\Delta}, \mathbf{\Omega}, \nu)|_{\{\lambda_{ij}=\delta_{ik}=0\} \text{ and Eq. [14](#)\}} = -n(1 - \epsilon(c-1)) \log(1 - \epsilon(c-1)) - n\epsilon(c-1) \log \epsilon \quad (15)$$

- *Case 2.* When all probability values never touch the highest extrema, such as

$$\forall i \in [n], j \in [c], \text{ s.t. } p_{ij}^* < 1 - \epsilon(c-1)$$

By KKT conditions, we have that $\delta_{ij} = 0$. By substituting these conditions in Eq. [12](#), we obtain that

$$\mathcal{L}_1(\mathbf{P}^*; \mathbf{\Lambda}, \mathbf{\Delta}, \mathbf{\Omega}, \nu)|_{\{\delta_{ij}=0\}} = n + \sum_{i=1}^n \left\{ -\nu_i + \sum_{j=1}^c \lambda_{ij}\epsilon \right\} \quad (16)$$

which always satisfies the inequality

$$\mathcal{L}_1(\mathbf{P}^*; \mathbf{\Lambda}, \mathbf{\Delta}, \mathbf{\Omega}, \nu)|_{\{\delta_{ij}=0\}} \geq \mathcal{L}_1(\mathbf{P}^*; \mathbf{\Lambda}, \mathbf{\Delta}, \mathbf{\Omega}, \nu)|_{\{\lambda_{ij}=\delta_{ik}=0\}} \quad (17)$$

and therefore also

$$\mathcal{L}_1(\mathbf{P}^*; \mathbf{\Lambda}, \mathbf{\Delta}, \mathbf{\Omega}, \nu)|_{\{\delta_{ij}=0\}} \geq \mathcal{L}_1(\mathbf{P}^*; \mathbf{\Lambda}, \mathbf{\Delta}, \mathbf{\Omega}, \nu)|_{\{\lambda_{ij}=\delta_{ik}=0\} \text{ and Eq. [14](#)\}} \quad (18)$$

Finally, we observe that the objective of the optimization problem of Eq. 8 evaluated at the solutions of *Case 1* is

$$-\sum_{i=1}^n \sum_{j=1}^c p_{ij} \log p_{ij} = \mathcal{L}_1(\mathbf{P}^*; \mathbf{\Lambda}, \mathbf{\Delta}, \mathbf{\Omega}, \boldsymbol{\nu})|_{\{\lambda_{i,j}=\delta_{i,k}=0\}} \text{ and Eq. 14} \quad (19)$$

And by leveraging also the result in Eq. 18 we can state that the solutions of *Case 1* are the global minima of the objective in Eq. 8. Thus concluding the proof for the extrema condition.

Matched prior. We consider the minimization of the second addend in Eq. 7 subject to the extrema condition

$$\begin{aligned} \min_{\mathbf{P}} & -\sum_{j=1}^c q_j \log \frac{1}{n} \sum_{i=1}^n p_{ij} \\ \text{s.t.} & \sum_{j=1}^c p_{ij} = 1, \quad \forall i \in [n] \\ & p_{ij} \in \{\epsilon, 1 - \epsilon(c-1)\}, \quad \forall i \in [n], j \in [c], \end{aligned} \quad (20)$$

Let's define $\tilde{p}_j \equiv \frac{1}{n} \sum_{i=1}^n p_{ij}$ for all $j \in [c]$ and observe that $\sum_{j=1}^c \tilde{p}_j = 1$ and $\epsilon \leq \tilde{p}_j \leq 1 - \epsilon(c-1)$. Therefore, we can rewrite the problem in Eq. 20 equivalently

$$\begin{aligned} \min_{\mathbf{P}} & -\sum_{j=1}^c q_j \log \tilde{p}_j \\ \text{s.t.} & \sum_{j=1}^c \tilde{p}_j = 1, \\ & \epsilon \leq \tilde{p}_j \leq 1 - \epsilon(c-1), \quad \forall j \in [c], \end{aligned} \quad (21)$$

Now, we observe that the optimization objective satisfies the following equality

$$-\sum_{j=1}^c q_j \log \tilde{p}_j = H(\mathbf{q}) + KL(\mathbf{q} \parallel \tilde{\mathbf{p}}) \quad (22)$$

The minimum for Eq. 22 is obtained at $\mathbf{q} = \tilde{\mathbf{p}}$ and this solution satisfies the constraints in Eq. 21 because $\epsilon \leq q_j \leq 1 - \epsilon(c-1)$ for all $j \in [c]$ (indeed we can always choose ϵ to satisfy the inequality), thus being the global optimum. In other words, we have that $\frac{1}{n} \sum_{i=1}^n p_{ij} = q_j$ for all $j \in [c]$.

Finally, recall that $I_{max}(j) \equiv \{i \in [n] : p_{ij} = 1 - \epsilon(c-1)\}, \forall j \in [c]$, which identifies all elements having the highest possible value of probability in \mathbf{P} . We observe that

$$\begin{aligned} \sum_{i=1}^n p_{ij} &= \sum_{i \in I_{max}(j)} p_{ij} + \sum_{i \notin I_{max}(j)} p_{ij} \\ &= \sum_{i \in I_{max}(j)} (1 - \epsilon(c-1)) + \sum_{i \notin I_{max}(j)} \epsilon \quad (\text{by Extrema condition}) \\ &= |I_{max}(j)|(1 - c\epsilon) + n\epsilon \end{aligned}$$

By the condition $\frac{1}{n} \sum_{i=1}^n p_{ij} = q_j$ and the above relation we have that

$$|I_{max}(j)|(1 - c\epsilon) + n\epsilon = nq_j, \quad \forall j \in [c]$$

or equivalently that

$$|I_{max}(j)| = \left(\frac{q_j - \epsilon}{1 - c\epsilon} \right) n \quad (23)$$

Now, for the case of uniform prior, Eq. 23 becomes

$$q_j = \frac{1}{c} \implies |I_{max}(j)| = \frac{n}{c}, \quad \forall j \in [c] \quad (24)$$

This concludes the proof for the matching prior condition.

Finally the global minimum value of the *FALCON* objective can be obtained by dividing Eq. 15 by n and adding the entropy term (as for the result obtained by the matched prior condition). This concludes the proof of the Lemma. \square

D EMBEDDING THEOREM

Proof. Recall the extrema condition from Lemma 1, that is

$$\forall i \in [n], \exists! j \in [c], \forall k \in [c] \text{ with } k \neq j \text{ s.t. } p_{ij} = 1 - \epsilon(c-1) \text{ and } p_{ik} = \epsilon$$

Moreover, due to orthogonality of \mathbf{W} we can express the *Span* condition, i.e. $\mathbf{h}_i = \sum_{j'=1}^c \alpha_{ij'} \mathbf{w}_{j'}$ for all $i \in [n]$ with $\alpha_{ij} \in \mathbb{R}$. This fact leads us to the following equation

$$p_{ij} = \frac{e^{\mathbf{w}_j^T \mathbf{h}_i / \tau}}{\sum_{j''=1}^c e^{\mathbf{w}_{j''}^T \mathbf{h}_i / \tau}} \stackrel{\text{Span}}{=} \frac{e^{\alpha_{ij} f / \tau}}{\sum_{j'=1}^c e^{\alpha_{ij'} f / \tau}} \quad \forall i \in [n], j \in [c] \quad (25)$$

Combining the extrema condition with Eq. 25 gives us a system of equations for each $i \in [n]$

$$\begin{cases} \frac{e^{\alpha_{ij} f / \tau}}{\sum_{j'=1}^c e^{\alpha_{ij'} f / \tau}} = 1 - \epsilon(c-1) \\ \frac{e^{\alpha_{ik} f / \tau}}{\sum_{j'=1}^c e^{\alpha_{ij'} f / \tau}} = \epsilon \end{cases} \quad \forall k \neq j$$

By taking the logarithm on both sides of the two equations and resolving the above system, the solution is equal to

$$\begin{aligned} \alpha_{ik} &= \alpha_{ij} - \frac{\tau}{f} \log \left(\frac{1 - \epsilon(c-1)}{\epsilon} \right) \\ &= \alpha_{ij} - \frac{1}{\sqrt{n}} \quad \forall k \neq j \end{aligned} \quad (26)$$

where the last equality holds due to the choice $\tau = f / (\sqrt{n} \log((1 - \epsilon(c-1))/\epsilon))$. Using Eq. 26 in the *Span* condition gives us the following result

$$\forall i \in [n], \exists! j \in [c] \text{ s.t. } \mathbf{h}_i = \alpha_{ij} \mathbf{w}_j + \left(\alpha_{ij} - \frac{1}{\sqrt{n}} \right) \sum_{k \neq j} \mathbf{w}_k \quad (27)$$

Note that the α_{ij} could potentially take any value in $\alpha_{ij} \in \mathbb{R}$. This is not allowed as embeddings are normalized by design choice (cf. Eq. 1). Indeed, the norm of the embeddings can be rewritten to exploit Eq. 27

$$\begin{aligned} \|\mathbf{h}_i\|_2^2 &= \mathbf{h}_i^T \mathbf{h}_i \\ &\stackrel{\text{Eq. 27}}{=} c f \alpha_{ij}^2 - \frac{2(c-1)f}{\sqrt{n}} \alpha_{ij} + \frac{f}{n} (c-1) \end{aligned} \quad (28)$$

and by equating Eq. 28 to the fact that embeddings are normalized $\|\mathbf{h}_i\|_2^2 = \frac{f}{n}$ for all $i \in [n]$ we obtain the following quadratic equation

$$\alpha_{ij}^2 - \frac{2(c-1)f}{\sqrt{n}} \alpha_{ij} + \frac{f}{n} (c-1) - \frac{f}{n} = 0$$

whose solutions are given by

$$\alpha_{ij} = \left\{ \begin{array}{l} \frac{1}{\sqrt{n}} \\ \left(1 - \frac{2}{c}\right) \frac{1}{\sqrt{n}} \end{array} \right.$$

This concludes the proof. \square

E DIAGONAL COVARIANCE

Proof. Recall from Theorem [1](#) that

$$\forall i \in [n], \exists! j \in [c] \text{ s.t. } \mathbf{h}_i = \alpha_{ij} \mathbf{w}_j + \left(\alpha_{ij} - \frac{1}{\sqrt{n}} \right) \sum_{k \neq j} \mathbf{w}_k$$

By assumption $\alpha_{ij} = \frac{1}{\sqrt{n}}$ and therefore

$$\forall i \in [n], \exists! j \in [c] \text{ s.t. } \mathbf{h}_i = \frac{1}{\sqrt{n}} \mathbf{w}_j \quad (29)$$

meaning that the rows of \mathbf{H} are equal up to a constant to the codes in the dictionary and that they span the same space of the columns of \mathbf{W} , namely the whole embedding space. We can therefore express \mathbf{H} as linear combination of \mathbf{W} .

Without loss of generality, we can always define \mathbf{H} so as to ensure that nearby rows are associated to the same codes in the dictionary. Therefore, by combining this with Eq. [29](#) we have that

$$\mathbf{H} = \mathbf{A}^T \mathbf{W}^T$$

with

$$\mathbf{A} = \begin{pmatrix} \frac{1}{\sqrt{n}} \mathbf{1}_{n/c}^T & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \frac{1}{\sqrt{n}} \mathbf{1}_{n/c}^T & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \frac{1}{\sqrt{n}} \mathbf{1}_{n/c}^T \end{pmatrix} \in \mathbb{R}^{c \times n}$$

where $\mathbf{1}_{n/c}$ is a vector containing n/c ones (whose size follows due to the assumption on uniformity of \mathbf{q}). Importantly, matrix \mathbf{A} satisfies the following property

$$\mathbf{A} \mathbf{A}^T = \frac{1}{c} \mathbf{I} \quad (30)$$

Therefore, we have that

$$\begin{aligned} \mathbf{H}^T \mathbf{H} &= \mathbf{W} \mathbf{A} \mathbf{A}^T \mathbf{W}^T \\ &\stackrel{\text{Eq. } \a href="#">30}{=} \frac{1}{c} \mathbf{W} \mathbf{W}^T \\ &= \mathbf{I} \end{aligned}$$

where the last equality simply follows by the orthogonality condition $\mathbf{W}^T \mathbf{W} = f \mathbf{I}$ and the fact that \mathbf{W} is a square matrix ($c = f$). Indeed, we have that

$$\begin{aligned} \mathbf{W}^T \mathbf{W} &= f \mathbf{I} \\ \mathbf{W} \mathbf{W}^T \mathbf{W} &= f \mathbf{I} \mathbf{W} \\ (\mathbf{W} \mathbf{W}^T) \mathbf{W} &= (f \mathbf{I}) \mathbf{W} \\ \mathbf{W} \mathbf{W}^T &= f \mathbf{I} \end{aligned}$$

thus concluding the proof. \square

F GENERALIZATION TO SUPERVISED LINEAR DOWNSTREAM TASK

We first observe that by the results of Theorem [1](#) and the uniformity of \mathbf{q} , \mathbf{H} has full rank. Moreover, considering that \mathbf{H} is a function of \mathbf{Z} through the first layer of the projector in Eq. [1](#), \mathbf{Z} must be also full rank. As a consequence,

$$\mathbf{Z}^T \mathbf{Z} \text{ has full rank. (Full Rank Property)} \quad (31)$$

Now, we recall an existing result for generalization to supervised downstream tasks from [Shwartz-Ziv et al. \(2023\)](#) (Section 6.1) and demonstrate that the *Full Rank Property* reduces the generalization error.

Indeed, consider a classification problem with r classes. Given an unlabeled dataset \mathcal{D} , used for training *FALCON*, with the corresponding unknown ground truth labels $\mathbf{Y}_{\mathcal{D}} \in \mathbb{R}^{n \times r}$ and a supervised dataset $\mathcal{S} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m$, with \mathbf{y}_i being the rows of the label matrix $\mathbf{Y}_{\mathcal{S}} \in \mathbb{R}^{m \times r}$, define $\mathbf{Z} \in \mathbb{R}^{n \times f}$ and $\bar{\mathbf{Z}} \in \mathbb{R}^{m \times f}$ the representations obtained by feeding datasets \mathcal{D} and \mathcal{S} , respectively, through the backbone network g . Moreover, define

$$\begin{aligned} \mathbf{P}_{\mathcal{D}} &\equiv \mathbf{I} - \mathbf{Z}(\mathbf{Z}^T \mathbf{Z})^\dagger \mathbf{Z}^T \\ \mathbf{P}_{\mathcal{S}} &\equiv \mathbf{I} - \bar{\mathbf{Z}}(\bar{\mathbf{Z}}^T \bar{\mathbf{Z}})^\dagger \bar{\mathbf{Z}}^T \end{aligned}$$

where symbol \cdot^\dagger denotes the pseudo-inverse. Now, suppose we train a linear classifier with parameters $\mathbf{U} \in \mathbb{R}^{f \times r}$ on the latent representations obtained from dataset \mathcal{S} through the following supervised loss

$$\ell_{\mathbf{x}, \mathbf{y}}(\mathbf{U}) \equiv \|g(\mathbf{x})\mathbf{U} - \mathbf{y}\|_2^2 + \|\mathbf{U}\|_F$$

Then, we can state the following theorem

Th. 1 (restated from [Shwartz-Ziv et al. \(2023\)](#)). $\forall \delta > 0$ with probability at least $1 - \delta$, we have that

$$\begin{aligned} \mathbb{E}_{\mathbf{x}, \mathbf{y}}\{\ell_{\mathbf{x}, \mathbf{y}}(\mathbf{U})\} &\leq \frac{1}{n} \sum_{i=1}^n \|g(\mathbf{x}_i) - g(\mathbf{x}'_i)\|_2 + \frac{2}{m} \mathbb{E}_{\mathcal{D}, \xi} \left\{ \sup_g \sum_{i=1}^n \xi_i \|g(\mathbf{x}_i) - g(\mathbf{x}'_i)\|_2 \right\} + \\ &+ \frac{2}{\sqrt{n}} \|\mathbf{P}_{\mathcal{D}} \mathbf{Y}_{\mathcal{D}}\|_F + \frac{1}{\sqrt{m}} \|\mathbf{P}_{\mathcal{S}} \mathbf{Y}_{\mathcal{S}}\|_F + \text{const}(n, m) \end{aligned} \quad (32)$$

where ξ is a vector of i.i.d. Rademacher random variables.

Therefore, the expected supervised loss in Eq. 32 can be reduced by minimizing its upper bound. Note that the first addend in Eq. 32 is minimized by the *FALCON* loss, whereas the second addend is also statistically minimized when n is large. The third addend refers to the contribution term for the classification on the unlabeled data. While ground truth $\mathbf{Y}_{\mathcal{D}}$ is unknown, this addend can be minimized by exploiting the following relation

$$\|\mathbf{P}_{\mathcal{D}} \mathbf{Y}_{\mathcal{D}}\|_F \leq \|\mathbf{P}_{\mathcal{D}}\|_F \|\mathbf{Y}_{\mathcal{D}}\|_F$$

Indeed, note that in order to minimize the left-hand side of the inequality, it suffices to minimize the term $\|\mathbf{P}_{\mathcal{D}}\|_F$, which occurs when $\mathbf{Z}^T \mathbf{Z}$ has maximum rank. This is our case due to the *Full Rank Property*. Finally, by the same argument used for the third term in Eq. 32, we can minimize the fourth one by having $\bar{\mathbf{Z}}^T \bar{\mathbf{Z}}$ with maximum rank. This condition holds because $\mathbf{Z}^T \mathbf{Z}$ and $\bar{\mathbf{Z}}^T \bar{\mathbf{Z}}$ concentrate to each other by concentration inequalities (cf. [Shwartz-Ziv et al. \(2023\)](#) for more details).

To summarize, minimizing the *FALCON* loss ensures that we reduce the invariance of representations to data augmentations and increase the rank of the representation covariance. This leads to a decrease of the generalization error as from the result of Theorem 1.

G BLOCK-DIAGONAL ADJACENCY

Proof. The proof follows step by step the one for the diagonal covariance except for the fact that

$$\mathbf{H}\mathbf{H}^T = \mathbf{A}^T \mathbf{W}^T \mathbf{W} \mathbf{A} \quad \underbrace{=}_{\mathbf{W}^T \mathbf{W} = f\mathbf{I}} \quad f \mathbf{A}^T \mathbf{A} = f \mathbf{B}_{\mathbf{A}}$$

where

$$\mathbf{B}_{\mathbf{A}} \equiv \mathbf{A}^T \mathbf{A} = \begin{pmatrix} \mathbf{1}_{\frac{n}{c} \times \frac{n}{c}} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \frac{1}{n} \mathbf{1}_{\frac{n}{c} \times \frac{n}{c}} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \frac{1}{n} \mathbf{1}_{\frac{n}{c} \times \frac{n}{c}} \end{pmatrix} \in \mathbb{R}^{n \times n}$$

and $\mathbf{1}_{\frac{n}{c} \times \frac{n}{c}}$ is a matrix of ones. This concludes the proof. \square

Table 3: Resnet architecture. Conv2D(A,B,C) applies a 2d convolution to input with B channels and produces an output with C channels using stride (1, 1), padding (1, 1) and kernel size (A, A).

Name	Layer	Res. Layer
Block 1	Conv2D(3,3,F)	Conv2D(1,3,F) no padding
	LeakyRELU(0.2)	
	Conv2D(3,F,F)	
	AvgPool2D(2)	
Sum		
Block 2	LeakyRELU(0.2)	
	Conv2D(3,F,F)	
	LeakyRELU(0.2)	
	Conv2D(3,F,F)	
Block 3	AvgPool2D(2)	
	LeakyRELU(0.2)	
	Conv2D(3,F,F)	
	LeakyRELU(0.2)	
Block 4	Conv2D(3,F,F)	
	LeakyRELU(0.2)	
	Conv2D(3,F,F)	
	AvgPool2D(all)	

Table 4: Hyperparameters (in terms of optimizer and data augmentation) used in SVHN, CIFAR-10 and CIFAR-100 experiments.

Class	Name param.	SVHN	CIFAR-10	CIFAR-100
Data augment.	Color jitter prob.	0.1	0.1	0.1
	Gray scale prob.	0.1	0.1	0.1
	Random crop	Yes	Yes	Yes
	Additive Gauss. noise (std)	0.03	0.03	0.03
	Random horizontal flip	No	Yes	Yes
Optimizer	Batch size	64	64	64
	Epochs	20	200	200
	Adam β_1	0.9	0.9	0.9
	Adam β_2	0.999	0.999	0.999
	Learning rate	$1e - 4$	$1e - 4$	$1e - 4$

H EXPERIMENTAL DETAILS ON SVHN, CIFAR10 AND CIFAR100

Training. We used a ResNet-8 (details are provided in Table 3). We consider the hyperparameters in Table 4 for training. Beta is chosen to ensure both losses are minimized, cf. Table 5.

Evaluation. For linear probe evaluation, we followed standard practice by removing the projector head and train a linear classifier on the backbone representation. We train the classifier with Adam optimizer for 100 epochs and learning rate equal to $1e - 2$.

I ADDITIONAL RESULTS ON DICTIONARY SIZE

We provide additional visualization results for the covariance and adjacency matrices on SVHN and CIFAR-10, cf. Figs. 9, 10. Moreover, we add the analysis of generalization on downstream tasks on SVHN and CIFAR-100 varying the size of the dictionary in Figs 11, 12.

Table 5: Values of β hyperparameter. This is chosen from the range $\{0.01, 0.05, 0.1, 0.25, 0.5, 1, 2.5, 5, 10\}$ to ensure that both losses are minimized.

Dictionary Size	10	128	256	512	1024	2048	4096	8192	16384
SVHN	0.5	0.5	0.5	0.25	0.1	0.1	0.1	0.1	0.1
CIFAR-10	0.5	0.5	0.5	0.25	0.1	0.1	0.1	0.1	0.1
CIFAR-100	0.5	0.5	0.5	0.25	0.1	0.1	0.1	0.1	0.1

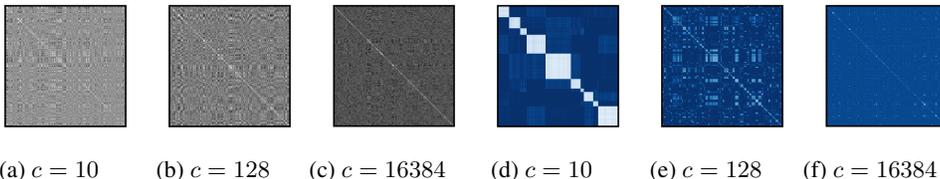


Figure 9: Realization of embedding covariance (**left**) and adjacency matrices (**right**) for the whole SVHN test dataset. Increasing c reduces the value of the off-diagonal elements of the covariance, thus contributing to increase the decorrelation of features (cf. Corollary 2). Moreover, increasing c has the effect to reduce the block sizes of the adjacency matrix (cf. Corollary 3).

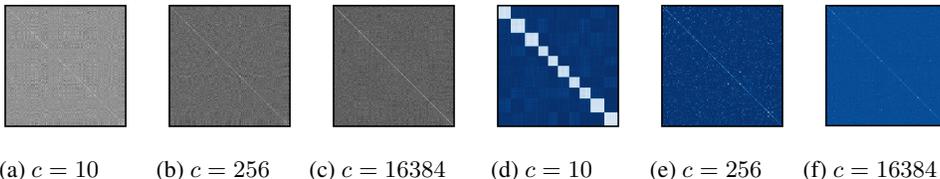


Figure 10: Realization of embedding covariance (**left**) and adjacency matrices (**right**) for the whole CIFAR-100 test dataset. Increasing c reduces the value of the off-diagonal elements of the covariance, thus contributing to increase the decorrelation of features (cf. Corollary 2). Moreover, increasing c has the effect to reduce the block sizes of the adjacency matrix (cf. Corollary 3).

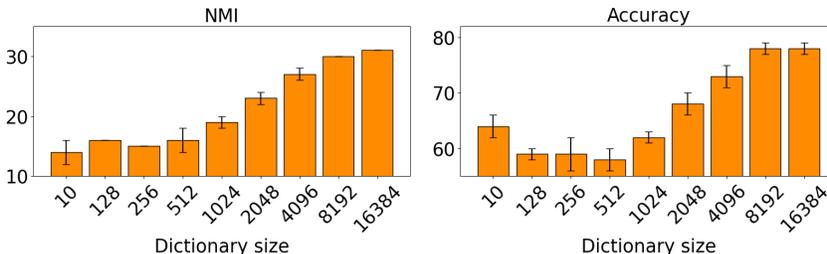


Figure 11: Analysis of downstream generalization for different values of dictionary size on SVHN dataset.

J ADDITIONAL ANALYSIS ON COLLAPSES

We provide additional results for the collapses on SVHN and CIFAR100. Specifically, in Fig. 13 we show the analysis of dimensional collapses, whereas in Fig. 14 we show the one for intracluster collapse.

K EXPERIMENTAL DETAILS ON IMAGENET-100

Training. We used a ViT-small backbone network and train it for 100 epochs with learning rate equal to $5e - 4$ and batch-size per GPU equal to 64 on a node with 8 NVIDIA A100 GPUs. Beta is selected

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

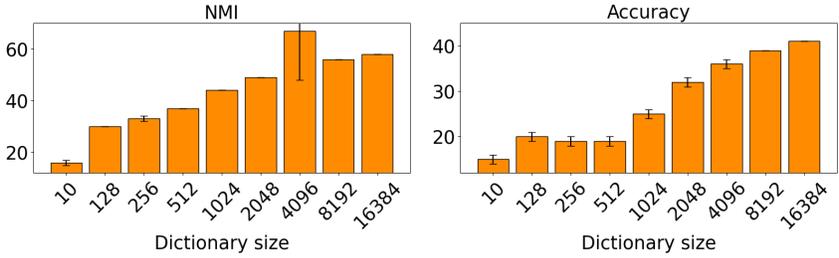


Figure 12: Analysis of downstream generalization for different values of dictionary size on CIFAR-100 dataset.

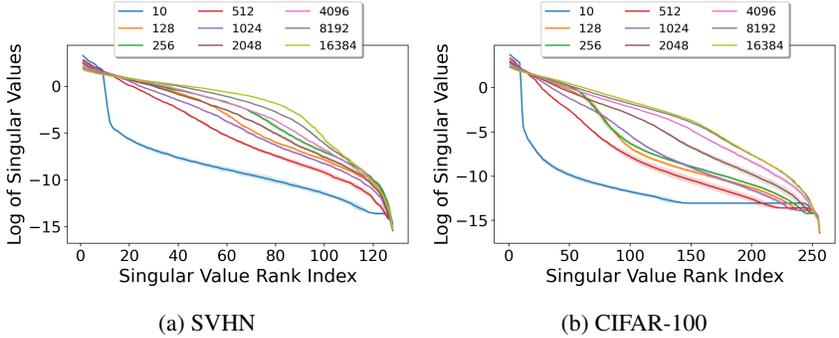


Figure 13: Dimensional collapse analysis on test data for different size of dictionary. Results are averaged over 5 training runs obtained from random initialization seeds.

from a smaller subset of values $\{0.1, 0.25, 0.5\}$ (given the more expensive nature of the experiments) to ensure both losses are minimized and chosen being equal to 0.25.

Evaluation. For linear probe evaluation, we use the DINO codebase and train the classifier with Adam optimizer (Caron et al., 2021).

L PRACTICAL IMPLEMENTATION OF THE LOSS

We observed training instability when using the larger backbone on ImageNet-100. The issue is due to some dictionary codes being unused during the initial training phase (cluster collapse), making the KL matching prior term infinity. Indeed, we have that

$$\begin{aligned} \mathcal{L}_{FALCON}(\mathcal{D}) &= \beta CE(\mathbf{p}, \mathbf{p}') + CE(\mathbf{q}, \mathbf{p}) \\ &\propto \beta CE(\mathbf{p}, \mathbf{p}') + KL(\mathbf{q}, \mathbf{p}) \end{aligned}$$

In practice, the reverse KL term is sufficient to avoid the issue:

$$\mathcal{L}_{FALCON}(\mathcal{D}) = \beta CE(\mathbf{p}, \mathbf{p}') + KL(\mathbf{p}, \mathbf{q})$$

1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403

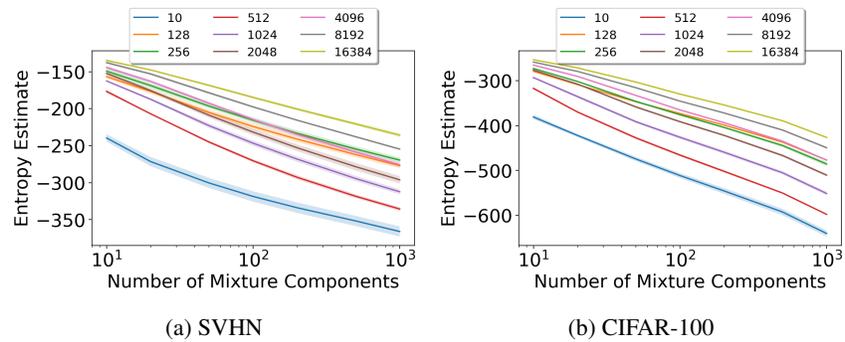


Figure 14: Intracuster collapse analysis on test data for different size of dictionary. Results are averaged over 5 training runs obtained from random initialization seeds.