

421 **A Proofs**

422 **A.1 Proof of Theorem 4.5**

423 *Proof.* We first rewrite  $\mathcal{L}_{\text{tri}}$  as a matrix decomposition objective

$$\begin{aligned} \mathcal{L}_{\text{tri}}(f) &= -2\mathbb{E}_{x,x^+} f(x)^\top S f(x^+) + \mathbb{E}_x \mathbb{E}_{x^-} (f(x)^\top S f(x^-))^2 \\ &= \sum_{x,x'} \left( \frac{A_{xx'}^2}{D_{xx} D_{x'x'}} + D_{xx} D_{x'x'} (f(x)^\top S f(x'))^2 - 2A_{xx'} f(x)^\top S f(x') \right) + \text{const} \quad (16) \\ &= \|\bar{A} - F S F^\top\|^2. \end{aligned}$$

424 According to the Eckart-Young Theorem [Eckart and Young, 1936], the optimal solutions  $F^*, S^*$   
425 satisfy

$$F^* S^* (F^*)^\top = U^k \Sigma (V^k)^\top,$$

426 where  $\Sigma \in \mathbb{R}^{k \times k}$  is a diagonal matrix with the  $k$ -largest eigenvalues of  $\bar{A}$  and  $U \in \mathbb{R}^{N \times k}$  contains  
427 the corresponding eigenvectors of the  $k$ -largest eigenvalues. When the regularizer  $\mathcal{L}_{\text{Dec}}$  is minimized,  
428  $F^*$  satisfy  $(F^*)^\top F^* = I$ . In the next step, we prove the uniqueness of the optimal solution.

429 We denote  $H = F^* \Sigma (F^*)^\top$ . As  $(F^*)^\top F^* = I$ , we obtain  $HH^\top = F^* S^* (S^*)^\top (F^*)^\top$ . If  $\zeta, \sigma$  are  
430 a pair of eigenvector and eigenvalue of  $HH^\top$ , we have

$$\begin{aligned} HH^\top \zeta &= F^* S^* (S^*)^\top (F^*)^\top \zeta = \sigma \zeta, \\ S^* (S^*)^\top (F^*)^\top \zeta &= \sigma (F^*)^\top \zeta, \\ S^* (S^*)^\top ((F^*)^\top \zeta) &= \sigma ((F^*)^\top \zeta). \end{aligned} \quad (17)$$

431 So the eigenvalues of  $HH^\top$  are the eigenvalues of  $S^* (S^*)^\top$ . As the positive eigenvalues of  $HH^\top$   
432 are uniquely determined and  $S^*$  has a descending order,  $S^*$  is also determined and  $S^* = \Sigma$ .

433 We note that  $HH^\top = F^* S^* (S^*)^\top (F^*)^\top$ , i.e.,  $HH^\top F^* = F^* S^* (S^*)^\top$ , which means that the  $k$   
434 columns of  $F^*$  are the eigenvectors of  $HH^\top$  and the corresponding eigenvalues are  $\sigma_1 \cdots \sigma_k$ . As  
435  $HH^\top$  only has  $k$  different non-negative eigenvalues  $\sigma_1, \dots, \sigma_k$ , the eigenspace of each eigenvalue  
436 is one-dimensional. When we consider the real number space, any two eigenvectors  $\zeta_i, \zeta'_i$  of the  
437 same eigenvalue  $\sigma_i$  satisfy  $\zeta_i = c \zeta'_i$ . As  $(F^*)^\top F^* = I$ , we obtain  $c = \pm 1$ . As  $f(x) = \frac{1}{\sqrt{D_{xx}}} F_x$ , we  
438 obtain

$$f_j^*(x) = \pm \frac{1}{\sqrt{D_{xx}}} (U_x^k)_j, S^* = \text{diag}(\sigma_1, \dots, \sigma_k), \quad (18)$$

439 □

440 **A.2 Proof of Theorem 5.1**

441 We first introduce a lemma which theoretically guarantees the generalization performance of spectral  
442 contrastive learning.

443 **Lemma A.1** ([HaoChen et al., 2021]). *For the optimal solutions to spectral contrastive learning*  
444 *(SCL), we have*

$$\mathcal{E}(f_{\text{SCL}}^*) \leq \mathcal{O}\left(\frac{\alpha}{1 - \sigma_{k+1}}\right),$$

445 where we denote  $\alpha$  as the probability that the natural samples and augmented views have different  
446 labels, i.e.,  $\alpha = \mathbb{E}_{\bar{x} \sim \mathcal{P}_u} \mathbb{E}_{x \sim \mathcal{A}(\cdot|\bar{x})} \mathbb{1}[y(\bar{x}) \neq y(x)]$  and  $\sigma_{k+1}$  as the  $(k+1)$ -th largest eigenvalue of  
447 the normalized adjacent matrix  $\bar{A}$ .

448 Then we construct the generalization guarantee of tri-contrastive learning.

449 *Proof.* Following the proof of Theorem 4.5, we know that the optimal solutions learned by triCL are

$$\begin{aligned} F^* &= U^k, \\ S^* &= \text{diag}(\sigma_1, \dots, \sigma_k). \end{aligned}$$

450 So we know that the optimal encoder of triCL satisfies,  $\forall x \in \mathcal{D}$

$$f^*(x) = \frac{1}{\sqrt{D_{xx}}} (U_x^k)^\top.$$

451 Compared with the optimal solutions of spectral contrastive learning (Eq. 2), we know

$$(\text{diag}(\sigma_1, \dots, \sigma_k)R)^\top f_{triCL}^*(x) = f_{SCL}^*(x), \quad (19)$$

452 where  $f_{triCL}^*, f_{SCL}^*$  denote the optimal solutions of tri-contrastive learning and spectral contrastive  
 453 learning. As  $\text{diag}(\sigma_1, \dots, \sigma_k)R$  is an invertible matrix, we then prove that the invertible matrix can  
 454 be absorbed in the linear probing. We denote  $\text{diag}(\sigma_1, \dots, \sigma_k)R$  as  $Q$  and we denote the linear  
 455 classifier as  $B$ , i.e.,  $g(f(x)) = f(x)^\top B$ . For a linear classifier  $B$ , let  $\tilde{B} = BQ^{-1}$ . We then obtain  
 456  $f_{triCL}^*(x)^\top \tilde{B} = f_{SCL}^*(x)^\top B$ .

457 So

$$\mathcal{E}(f_{triCL}^*) = \mathcal{E}(f_{SCL}^*).$$

458 With lemma A.1, we have

$$\mathcal{E}(f_{triCL}^*) \leq \mathcal{O}\left(\frac{\alpha}{1 - \sigma_{k+1}}\right).$$

459

□

### 460 A.3 Proof of Theorem 5.2

461 *Proof.* Based on the proof of Theorem 4.5, we know that the  $t$ -th dimension of the optimal solutions  
 462 satisfies

$$\begin{aligned} F^* &= U_t^k, \\ S_t^* &= \text{diag}(\sigma_1, \dots, \sigma_k)_t. \end{aligned}$$

463 With the analysis in Eckart-Young theorem [Eckart and Young, 1936], we have

$$\begin{aligned} \|\bar{A} - F_t^* S_t^* (F_t^*)^\top\|_F^2 &= \|\bar{A} - U_t^k \text{diag}(\sigma_1, \dots, \sigma_k)_t (U_t^k)^\top\|_F^2 \\ &= \sum_{i=1}^{t-1} \sigma_i^2 + \sum_{i=t+1}^k \sigma_i^2. \end{aligned}$$

464 As  $\sigma_i$  is the  $i$ -th largest eigenvalues of  $\bar{A}$ , so

$$\|\bar{A} - F_1^* S_1^* (F_1^*)^\top\|_F^2 \leq \dots \leq \|\bar{A} - F_k^* S_k^* (F_k^*)^\top\|_F^2.$$

465 Following Eq 16, we obtain

$$\mathcal{L}_{triCL}(f_t, S_t) = \|\bar{A} - F_t^* S_t^* (F_t^*)^\top\|_F^2 + \text{const},$$

466 we obtain

$$\mathcal{L}_{triCL}(f_1^*, S_1^*) \leq \dots \leq \mathcal{L}_{triCL}(f_k^*, S_k^*).$$

467

□

### 468 A.4 Feature Identifiability of Asymmetric Tri-contrastive Learning

469 We first extend the augmentation graph to an asymmetric form. The asymmetric augmentation graph  
 470 is defined over the set of all samples with its adjacent matrix denoted by  $P_O$ . In the augmentation  
 471 graph, each node corresponds to a sample, and the weight of the edge connecting two nodes  $x_A$  and  
 472  $x_B$  is equal to the probability that they are selected as a positive pair, i.e.,  $(P_O)_{x_a, x_b} = \mathcal{P}_O(x_a, x_b)$ .  
 473 And we denote  $\bar{P}_O$  as the normalized adjacent matrix of the augmentation graph, i.e.,  $(\bar{P}_O)_{x_a, x_b} =$   
 474  $\frac{\mathcal{P}_O(x_a, x_b)^2}{\mathcal{P}_A(x_a)\mathcal{P}_B(x_b)}$ .

475 Similar to the symmetric form, we then rewrite  $\mathcal{L}_{\text{tri}}$  as a matrix decomposition objective

$$\begin{aligned}
\mathcal{L}_{\text{tri}}(f_A, f_B, S) &= -2\mathbb{E}_{x_a, x_b} f_A(x_a)^\top S f_B(x_b) + \mathbb{E}_{x_a^-, x_b^-} (f_A(x_a^-)^\top S f_B(x_b^-))^2 \\
&= \sum_{x_a, x_b} \left( \frac{\mathcal{P}_O(x_a, x_b)^2}{\mathcal{P}_A(x_a)\mathcal{P}_B(x_b)} + \mathcal{P}_A(x_a)\mathcal{P}_B(x_b) (f_A(x_a)^\top S f_B(x_L))^2 \right. \\
&\quad \left. - 2\mathcal{P}_O(x_a, x_b) f_A(x_a)^\top S f_B(x_L) \right) + \text{const} \\
&= \|\bar{P}_O - F_A S F_B^\top\|^2.
\end{aligned}$$

476 According to the Eckart-Young Theorem [Eckart and Young, 1936], the optimal solutions  $F_A^*, S^*, F_B^*$   
477 satisfy

$$F_A^* S^* (F_B^*)^\top = U^k \Sigma (V^k)^\top,$$

478 where  $\Sigma \in \mathbb{R}^{k \times k}$  is a diagonal matrix with the  $k$ -largest eigenvalues of  $\bar{P}_O$  and  $U \in \mathbb{R}^{N_A \times k}$  contains  
479 the corresponding eigenvectors of the  $k$ -largest eigenvalues. When the regularizer  $\mathcal{L}_{\text{Dec}}$  is minimized,  
480  $F_A^*$  and  $F_B^*$  satisfy  $(F_A^*)^\top F_A^* = I$ ,  $(F_B^*)^\top F_B^* = I$ . In the next step, we prove the uniqueness of the  
481 optimal solution.

482 We denote  $H = F_A^* \Sigma F_B^*$ , and we obtain  $HH^\top = F_A^* S^* (S^*)^\top (F_A^*)^\top$ . If  $\zeta, \sigma$  are a pair of  
483 eigenvector and eigenvalue of  $HH^\top$ , we have

$$\begin{aligned}
HH^\top \zeta &= F_A^* S^* (S^*)^\top (F_A^*)^\top \zeta = \sigma \zeta, \\
S^* (S^*)^\top (F_A^*)^\top \zeta &= \sigma (F_A^*)^\top \zeta, \\
S^* (S^*)^\top ((F_A^*)^\top \zeta) &= \sigma ((F_A^*)^\top \zeta).
\end{aligned} \tag{20}$$

484 So the eigenvalues of  $HH^\top$  are the eigenvalues of  $S^* (S^*)^\top$ . As the positive eigenvalues of  $HH^\top$   
485 are uniquely determined and  $S^*$  has an increasing order,  $S^*$  is also determined and  $S^* = \Sigma$ .

486 We note that  $HH^\top = F_A^* S^* (S^*)^\top (F_A^*)^\top$ , i.e.,  $HH^\top F_A^* = F_A^* S^* (S^*)^\top$ , which means that the  $k$   
487 columns of  $F_A^*$  are the eigenvectors of  $HH^\top$  and the corresponding eigenvalues are  $\sigma_1 \cdots \sigma_k$ . As  
488  $HH^\top$  only has  $k$  different non-negative eigenvalues  $\sigma_1, \cdots, \sigma_k$ , the eigenspace of each eigenvalue  
489 is one-dimensional. When we consider the real number space, any two eigenvectors  $\zeta_i, \zeta'_i$  of the  
490 same eigenvalue  $\sigma_i$  satisfy  $\zeta_i = c \zeta'_i$ . As  $(F_A^*)^\top F_A^* = I$ , we obtain  $c = +1$ . Then we eliminate the  
491 ambiguity of sign following Eq 11 and  $F_A^*$  is unique. Similarly,  $F_B^*$  is also unique. So the optimal  
492 solution of  $\mathcal{L}_{\text{triCLIP}}$  is unique.

## 493 B Experimental Details

### 494 B.1 Experiment Details of Section 6.1

495 We first generate a random matrix  $A$  with size  $5000 \times 3000$ , and make sure that it does not contain  
496 multiple eigenvectors (which is easy to satisfy). For the matrix factorization problem  $\|A - FG^\top\|_F^2$ ,  
497 we apply off-the-shelf algorithms and repeat this process ten times. We then calculate the mean and  
498 variance of the  $l_2$  pairwise distance between the obtained solutions of  $F$ . For the trifactorization  
499 objective  $\|A - FSG^\top\|_F^2$ , we use SVD to obtain an initial solution, and apply the sign identification  
500 procedure to determine the sign of each eigenvector. Similarly, we also repeat this process ten times  
501 and calculate the mean and variance of the  $l_2$  pairwise distance between different solutions.

### 502 B.2 Experiment Details of Section 6.2

503 **Pretraining Setups.** For different evaluation tasks (k-NN, linear evaluation, image retrieval), we use  
504 the same pretrained models. We adopt ResNet-18 as the backbone. For CIFAR-10 and CIFAR-100,  
505 the projector is a two-layer MLP with hidden dimension 2048 and output dimension 256. And for  
506 ImageNet-100, the projector is a two-layer MLP with hidden dimension 4096 and output dimension  
507 256. We pretrain the models with batch size 256 and weight decay 0.0001. For CIFAR-10 and  
508 CIFAR-100, we pretrain the models for 200 epochs. While for ImageNet-100, we pretrain the models  
509 for 400 epochs. We use the cosine anneal learning rate scheduler and set the initial learning rate to  
510 0.4 on CIFAR-10, CIFAR-100, and 0.3 on ImageNet-100.

511 As the importance matrix is learned on the projection layer, we conduct the downstream tasks on the  
 512 features encoded by the complete networks (containing both the backbones and the projectors).

513 **The Distribution of the Importance Matrix.** When observing the distribution of feature importance  
 514 discovered by the importance matrix  $S$ , we first apply the softmax activation functions on the  
 515 diagonal values of  $S$  and sort different rows of  $S$  by the descending order of corresponding diagonal  
 516 values in  $S$ . We denote the non-negative ordered diagonal values of  $S$  as  $(s_1, \dots, s_k)$ . When  
 517 we present the distribution of them in Figure 2(a), we normalize the diagonal values and obtain

$$518 \left( s_1 / \sum_{i=1}^k s_i, \dots, s_k / \sum_{i=1}^k s_i \right).$$

519 **The K-NN Accuracy on Selected Dimensions.** For k-NN evaluation on 10 selected dimensions, we  
 520 do not finetune the models. We sort the dimensions of  $f(x)$  by the descending order of corresponding  
 521 diagonal values in the importance matrix. The k-NN is conducted on the standard split of CIFAR-10,  
 522 CIFAR-100 and ImageNet-100 and the predicted label of samples is decided by the 10 nearest  
 523 neighbors.

524 **Linear Evaluation on Selected Dimensions.** We train the linear classifier on 20 dimensions of the  
 525 frozen networks for 30 epochs during the linear evaluation. We set batch size to 256 and weight  
 526 decay to 0.0001. For triCL, we sort the dimensions by descending order of the importance matrix.  
 527 And for SCL, we randomly choose 20 dimensions.

## 528 C More Extensions of Tri-contrastive Learning

529 In this section, we apply tri-contrastive learning to another representative contrastive learning objec-  
 530 tive: the non-contrastive loss [Grill et al., 2020, Chen and He, 2021].

531 Besides contrastive learning, non-contrastive learning is another popular self-supervised framework  
 532 that throws the negative samples in contrastive learning and learns the meaningful representations  
 533 only by aligning the positive pairs. Taking the state-of-the-art algorithm BYOL [Grill et al., 2020] as  
 534 an example, they use an MSE loss:

$$\mathcal{L}_{\text{MSE}}(f, g) = 2 - 2 \cdot \mathbb{E}_{x, x^+} \frac{g(x)^\top f(x^+)}{\|g(x)\|_2 \cdot \|f(x^+)\|_2}, \quad (21)$$

535 where  $g(x)$  and  $f(x)$  are two different networks to avoid the feature collapse. Then we consider  
 536 adapting the tri-term loss to the non-contrastive learning, i.e.,

$$\mathcal{L}_{\text{triMSE}}(f, g) = 2 - 2 \cdot \mathbb{E}_{x, x^+} \frac{g(x)^\top S f(x^+)}{\|g(x)\|_2 \cdot \|f(x^+)\|_2} + \|\mathbb{E}_x g(x) g(x)^\top - I\|^2. \quad (22)$$

537 It is noticed that BYOL utilizes the stop-gradient technique on the target network  $f$  and it is updated  
 538 by exponential moving average. So we only calculate the feature decorrelation loss on the online  
 539 network  $g$ .