

1 **R1: For a nonsymmetric matrix ...** We will clarify from the start that our method is designed for symmetric matrices.
 2 **R1: PI is unstable if two eigenvalues are close. So what is the advantage of PI vs the analytic solution?** The
 3 instability of PI is due to inaccurate initialization of the eigenvectors. In our case, the forward pass provides accurate
 4 values using SVD and PI then becomes stable during the backward pass. Another advantage over the analytic solution
 5 what we demonstrate is that **PI yields an upper bound on the gradients' magnitude, which guarantees they will not**
 6 **explode.** As shown in Sec. 2.2, PI is the geometric expansion of the analytic solution: With $\lambda_i \geq \lambda_j$, the term $1/(\lambda_i - \lambda_j)$,
 7 which appears in the analytic solution, is approximated by $1/\lambda_i + 1/\lambda_i (\lambda_j/\lambda_i) + 1/\lambda_i (\lambda_j/\lambda_i)^2 + \dots + 1/\lambda_i (\lambda_j/\lambda_i)^{K-1} \leq$
 8 K/λ_i in PI, where K represents the power iteration number. When $\lambda_i = \lambda_j$, the term $1/(\lambda_i - \lambda_j)$ in the analytic solution
 9 explodes, while the PI gradients are naturally upper bounded by K/λ_i .

10 **R1: A hidden regularizer, ϵ , is introduced. What if you set ϵ to 0 or 10^{-12} .** ϵ is commonly used to stabilize the
 11 computation (e.g., [2] [3] & [r1]). As shown in Eq. 13, ϵ controls the gradients' upper bound and appears in the
 12 denominator. Thus, too small an epsilon, e.g., 10^{-12} , will yield a large upper bound and increase the risk of gradient
 13 explosion. Following standard practice, e.g., [3], ϵ is set to 10^{-4} for all the methods, including SVD and PI, but our
 14 method always achieves 100% success rate while the others do not.

15 **R1: Your method uses some tricks. The comparison seems unfair, since you can also truncate SVD.** For a fair
 16 comparison, we recomputed results using the same tricks to truncate the eigenvalues for our SVD baseline. For PI,
 17 truncation was already used in our submission to mitigate the round-off errors caused by the deflation process. The
 18 results on CIFAR10 are given in the table below. Note that, for matrix dimensions $d > 16$, SVD and PI still fail in all
 19 cases. By contrast, we achieve 100% success rate even when the dimension is as large as 128.

Methods	Evaluation Metrics	$d = 4$	$d = 8$	$d = 16$	$d = 32$	$d = 64$	$d = 128$
SVD	Min Error & Suc. Rate	4.50%, 60%	4.75%, 33.3%	4.65%, 40%	-, 0%	-, 0%	-, 0%
PI	Min Error & Suc. Rate	4.44%, 100%	6.28%, 6.7%	-, 0%	-, 0%	-, 0%	-, 0%
Ours	Min Error & Suc. Rate	4.59%, 100%	4.43%, 100%	4.40%, 100%	4.46%, 100%	4.44%, 100%	4.75%, 100%

20 **R1: Results on CIFAR100 are far away from SOTA performance. It would be nice to compare with [r1].** Our
 21 paper focuses on solving the stability issues of ED, not designing better normalization layers (i.e., PCA & ZCA).
 22 Stability is measured as the success rate of the methods, which, for our purpose, is more important than accuracy. ZCA
 23 and PCA constitute two applications of our method to demonstrate stability. We therefore just used simple backbones
 24 (i.e., ResNet18/50), which translates to accuracies inferior to the SOTA. By contrast, [r1] focuses on designing a better
 25 normalization layer using an iterative normalization method. We nonetheless acknowledge the relevance of this paper,
 26 which we will cite in the final version. Note that [r1] was not published at the time of NeurIPS submission.

27 **R1: I doubt ED is widely used as there are numerical issues. Justify why ED is important for deep learning.**
 28 Indeed, ED has many numerical stability issues, and this is exactly what our paper addresses. Nevertheless, as discussed
 29 in the introduction from Line 14 to 18, ED has been used for image/point matching [6,7,8], second-order pooling [4],
 30 and pose estimation [9]. It has not been well integrated into deep networks because of the numerical instability, and, as
 31 stated by R3, one can expect that our paper will bring insight to this problem and be the basis for many new ideas.

32 **R2: Will the failure cases become more or less common when matrices are large enough?** As shown in the table
 33 above, the baselines' failures become more common as d increases, whereas our method succeeds 100% of the time for
 34 all dimensions. This remains true when increasing d to 128 (twice as many as in the submission), by putting the ZCA
 35 layer on top of a 128-channel conv. layer. The underlying reasons are that, thanks to our use of SVD in the forward pass,
 36 we have more accurate eigen value/vector estimates than the PI baseline, and that, as shown in Eq. 13, the gradients of
 37 our method are always bounded regardless of matrix size while those of the SVD baseline may easily explode.

38 **R2: The convergence behavior looks similar whether using the existing methods and the proposed one.** The
 39 convergence curves shown in Fig. 2 are based solely on the successful cases for the baselines and ignore the failure cases
 40 (see success rates in Table 2). Including these numerous failures would render the baseline curves entirely meaningless.

41 **R2: In Tables 3 and 4, the prediction error is not monotone with the matrix size.** For ZCA in Table 3, with
 42 ResNet18, all values are virtually the same, and with ResNet50, the trend shows that larger d values, which only our
 43 method can handle, give better results. For PCA in Table 4 (a,b), when too much information is preserved, some noise
 44 is kept and the accuracy drops. Conversely, when too much information is removed, some useful signal is discarded and
 45 the accuracy also drops. The right number of dimensions to preserve is dataset dependent and can be determined by
 46 cross validation.

47 **R3: Do we have to do blocks of d ?** Dividing the features into blocks of dimension d is only useful to compare our
 48 approach with the baselines, which only succeed for small matrix dimensions. Given our stabilized ED method, the
 49 blocks become unnecessary, and the largest dimension d in each experiment corresponds to not using blocks.

50 **R3: PCA denoising doesn't look better than batch normalization (BN).** While PCA denoising indeed has
 51 marginal improvement over BN, it is not our main focus. Similarly to ZCA, PCA denoising only is another ap-
 52 plication to demonstrate the stability of our method. The baselines to truly look at in Table 5 are PCA(PI) and
 53 PCA(SVD), which often break down in the training phase. We will emphasize this in the final version.

54 **R3: Is there a way of minimizing the dependency on ZCA?** We will minimize the emphasis on ZCA whitening
 55 and PCA denoising in the abstract and introduction.