

1 We thank all the reviewers. Our novelty is appreciated by R1, R3, and R4, the superior performance is recognized by R1,
2 R2, and R4, and the writing clarity is accredited by R1, R3, and R4. Below we addressed all questions and concerns.

3 **R1: shared β .** Following the settings in Table 1, the model with shared β drops PSNR from 34.87 dB to 34.85 dB.

4 **R1: cardinality & performance.** Increasing cardinality can constantly improve performance when model is under-
5 parameterized. Since our ablation study is based on small models, the cardinality is limited by the number of channels.
6 We do observe performance regression point when reducing the size of training dataset.

7 **R2: cardinality.** The most primitive formulation of cardinality groups is that we start from single convolution layer
8 with only c/d ($/$ denotes numeric division) output nodes and build sparse representation with k groups upon it as
9 Section 3.2, then each sparsity group has $c/(dk)$ nodes and the MLP has k outputs. The cardinality groups repeat
10 building this structure for d times and concatenate output nodes of all the d convolutions. In section 3.3, we move the
11 concatenation operation from convolution output nodes to convolution kernels for formulation simplification.

12 **R2: visualization.** Dark blue and light yellow reflect the most certain regions of sparsity group selection. In the second
13 column of Fig. 5, dark blue more likely co-occurs with brown and black texture in original RGB images and light
14 yellow co-occurs with green textures. In the last column, dark blue and light yellow occur only within the region of tree
15 and sky, lion, and squirrel, and especially concentrate in foreground in the second and third images. We think it is more
16 related to semantic, since foreground can be automatically identified and receptive field grows as stacking convolution
17 layers. We will remove the statement about semantic if you still think it is not conclusive.

18 **R3: sparse coding & LISTA.** We focus on exploring sparsity of hidden representation in deep neural networks instead
19 of implanting deep structures in sparse coding approaches, such as LISTA. Although LISTA variants could be unrolled
20 into similar formulation as deep models, the iterative shrinkage algorithms solve sparse coding with several folds of
21 computation cost than a single convolution layer, which is hard to be applied in efficient deep models without parameter
22 sharing across layers. In contrast, our method improves the efficiency of deep models, meanwhile increases the number
23 of distinct parameters. We will revise the statement in line 33 as "cannot be efficiently solved" to reduce ambiguity.

24 **R3: sparsity by ReLU.** As we discussed in Section 2.3, sparsity by ReLU has been addressed in [10]. ReLU thresholds
25 negatives to zeros and has a gap with sparsity definition in Eq. 1, where the number of non-zero values should be much
26 less than representation dimensionality. Also, ReLU can only determine sparsity for given vectors, but our method
27 predicts the sparsity by a side MLP and saves the computation for unnecessary elements. Experiments show our method
28 constantly outperforms networks only with ReLU activations.

29 **R3: group sparse coding.** Our method with multi-level sparsity constraints by ReLU and Softmax shares a similar
30 idea with mixed L2 and L1 norms in group sparse coding. Furthermore, enforcing the constraints into model structure
31 instead of additional objective functions is more favorable for efficiency and performance in deep models. We will add
32 more detailed connection and comparison between group sparse coding and our method in the future version.

33 **R3: guarantee of efficiency.** First, our method constantly outperforms baselines over multiple benchmarks including
34 hundreds of images with either less computation and the same model size or the same computation and larger model
35 size. Second, there are many design choices in our methods, and results in the ablation studies coincide with our design
36 and analysis. Third, visualization shows non-trivial sparsity group selection.

37 **R3&R2: state-of-the-art baselines.** We focus on the importance of sparse representation of deep models and propose
38 a very general approach to achieve sparsity within individual convolution layer, instead of a single state-of-the-art
39 model. Hence, all the baselines are residual networks with only convolutional and activation layers. Other advanced
40 techniques, such as multi-scale (U-net) and self-attention (non-local), are perpendicular to our method and can be easily
41 combined with ours for superior performance.

42 **R4: coupled dictionary.** The coupled dictionary is introduced in [3], where encoding dictionary D_1 and decoding
43 dictionary D_2 are independent two dictionaries. In previous method [1], D_1 and D_2 are the same one.

44 **R4: computation time.** We follow the protocol as Squeeze-
45 and-Excitation networks using CPU inference time to address
46 the latency for embedded device applications. We evaluate
47 the forward running time of a single convolution layer of our
48 implementation in PyTorch. The running time is the average of
49 100 times on quad-core Intel Core i5-2500 at 3.30GHz. We set batch size as 100, input and output channels as 64, patch
50 size as 64. The time for additional computation is around 5-10%. Our implementation can be improved, since SE-Net
51 takes 3ms over 164ms and our approach is more efficient by only modifying convolution kernels instead of features.

Computation time (in milliseconds).						
Group size	N/A	2	4	8	16	
fixed model size	1.96	1.17	0.77	0.56	0.70	
fixed FLOPs	1.96	2.12	2.12	2.08	2.08	

52 **R4: patch-based approach.** We will add the suggested paper and more methods with fair comparisons including
53 patch-based deep models, such as channel attention and self-attention models.