

1 We would like to thank all reviewers for their careful reading, thoughtful comments, and overall
 2 positive assessment! We address the questions raised by reviewers below.

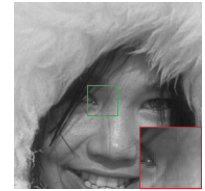
3 **Real-world application [Reviewer #1, #2, #3, #4].** To show the real world applicability of our
 4 theoretical framework, we consider the **local adaptive image denoising** task, where the noise
 5 levels in different parts of the images can be different. More specifically,

6 **(i) Dataset.** We split BSD500 dataset (400 images) [1] into a training set (100 images) and a test
 7 set (300 images). Gaussian noises are added to each pixel with noise levels depending on image
 8 local smoothness, making the noise levels on edges lower than non-edge regions. The task is to
 9 restore the original image from the noisy version $X \in [0; 1]^{180 \times 180}$.

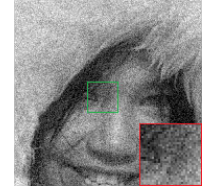
10 **(ii) Architecture.** We designed a hybrid architecture $\text{Alg}^k(E(X; Y))$ where Alg^k is a k -step
 11 unrolled minimization algorithm to the ℓ_2 -regularized reconstruction objective $E(X; Y) :=$
 12 $\frac{1}{2}kY + g(X) + \frac{1}{2} \sum_{i,j} [F(X)]_{i,j} Y_{i,j}^2$, and the residual $g(X)$ and position-wise regu-
 13 larization coefficient $F(X)$ are both DnCNN networks as in [2]. The optimization objective,
 14 $E(X; Y)$, is quadratic in Y , which follows the settings our theory focused on.

15 **(iii) Generalization gap.** We instantiate the hybrid architecture into different models using GD
 16 and NAG algorithms with different unrolled steps k . Each model is trained with 3000 epochs,
 17 and the *generalization gaps* between training and test errors are reported in Fig. 1. The results
 18 also show good consistency with our theory, where stabler algorithm (GD) can generalize better
 19 given *over-/under-*parameterized neural module, and for the *about-right* parameterization case,
 20 the generalization gap behaviors similar to $\text{Stab}(k) \sim \text{Cvg}(k)$. We will conduct more experimental
 21 trials and provide figures with smoother curves and error bars in our revised paper.

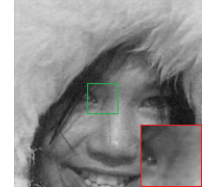
22 **(iv) Visualization.** To show that the learned hybrid model has a good performance in this real
 23 application, we include a visualization of the original, noisy, and denoised images.



(a) original image



(b) noisy image



(c) denoised by $\text{GD}^{12}(E(X; \cdot))$

24 **Generality of problem setting [Reviewer #1, #2, #3, #4].**

25 We acknowledged that our theoretical analysis is performed
 26 under a simplified problem setting, but we'd like to clarify
 27 a few points to avoid confusion.

28 We assume $E(x; y)$ is quadratic in y but it can depend on
 29 *any way* in the input x (i.e., Q can be any neural network).
 30 This can cover many real applications. For example, the
 31 above image denoising task, and many other data reconstruc-
 32 tion problems can be cast into quadratic energy minimization.

33 Even though in the paper we only present the results for
 34 using GD and NAG algorithms as the reasoning module, the main theorems which state the relation between the learning
 35 behavior and algorithm properties can be applied to *any optimization algorithm* as long as corresponding properties of
 36 the algorithm are provided. [Reviewer #3].

37 Our analysis framework can be extended to more general settings where the neural network module outputs a suitable
 38 strongly convex energy function. In fact the key component of our approximation and generalization analysis is the
 39 Lipschitz stability of the maps between the energy function and the exact minimizer, which can be ensured for general
 40 convex optimization problems if suitable regularization terms are introduced in the energy functions. We aim to analyze
 41 this general setting in future research.

42 **Other questions.**

43 **Space shrinks to a single function. [Reviewer #2]** We sincerely thank Reviewer #2 for his/her very detailed comments.
 44 Regarding why Alg^k contains only a single function, this comes from the convergence guarantee of the
 45 algorithms. That is, when the step size is in the stable region, the optimization error will decrease in each iteration
 46 and gradually decrease to 0 when $k \rightarrow \infty$. Therefore, for every step size η , $\text{Alg}^k = \text{Opt}$, the exact minimizer.

47 **Generalized to DP? [Reviewer #3]** Thanks Reviewer #3 for bringing in this interesting question, which is also what
 48 we want to address in our future work. Our analysis for RNN/GNN can potentially be adapted to the case where
 49 RNN/GNN are used to learn problems requiring DP, since RNN/GNN can present those operations in DP. However, as
 50 explained in our theory, one may not able to obtain as a tight bound as GD/NAG due to the difficulty of analyzing RNN.
 51 Furthermore, in the case of DP, the notion of convergence with respect the number of step k is different: the k is a fixed
 52 number of stages needed to run the DP iterations to solve the optimization. This will require more research.

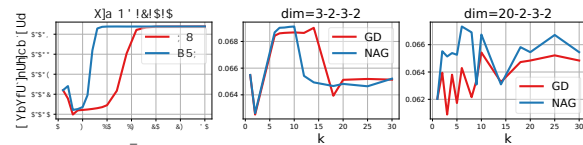


Figure 1: Generalization gap. Each k corresponds to a separately trained model. Left (under-parameterized): f is a DnCNN with 3 channels and 2 hidden layers and $g = 0$. Middle (about-right): both f and g have 3 channels and 2 hidden layers. Right (over-parameterized): f has 20 channels.

53 [1] Pablo Arbelaez, Michael Maire, Charless Fowlkes, and Jitendra Malik. Contour detection and hierarchical image segmentation.
 54 *IEEE transactions on pattern analysis and machine intelligence*, 33(5):898–916, 2010.

55 [2] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of
 56 deep51cnn for image denoising. *IEEE Transactions on Image Processing*, 26(7):3142–3155, 2017.