

1 We thank all the reviewers for their constructive comments and useful suggestions.

2 **Q (R1): "Comparison with other methods like encoder" & "why do we need this technique"**

3 **A:** This is a very important point that we need to clarify in our paper. GD inversion is not a straw man here: almost all
4 the prior work on using generative models for solving inverse problems (see reference [3,4,10-14] in our paper) uses
5 gradient descent as the main inversion technique so improving upon that is significant. There are also methods that
6 train encoders or even end-to-end reconstruction methods from measurements, but are harder to train and are tuned to a
7 specific inverse problem as opposed to a general method, see [3]. We will expand on this in the paper.

8 As compared to GD-based methods, our algorithm is much more efficient. GD costs on average 1.2 minutes to compress
9 one single image (momentum makes the process faster but still slower than our method). While our algorithm takes
10 only 0.5 second. See appendix for time comparisons.

11 We have compared the performance of GD with our method under different network architectures, i.e. different levels
12 of expansiveness.

13 Our paper focuses on fundamental theoretical results, since there are very few in this area. In practice our technique can
14 also be used as a good initialization for gradient-based methods. We verified this for DCGAN. See reply for R3.

15 **Q (R1): On reductions *from* known NP-hard problem A:** You are correct, of course. We fixed the confusing
16 expression.

17 **Q (R1): whether the binary case is representative for general case**

18 **A:** Thank you for your comment. We were able to extend our proof from binary to general real-valued inputs. We
19 achieve this by constraining an intermediate layer to be binary, using an additional output that enforces an intermediate
20 layer to have binary values. Combining this with our previous argument establishes that it is NP-hard to invert a 4-layer
21 network with real-valued inputs.

22 Specifically, we design a network $f : \mathbb{R}^k \rightarrow \mathbb{R}^2$ as follows: After input layer $\mathbf{z} \in \mathbb{R}^k$, we add 2 ReLU layers to make
23 sure the output of second hidden layer $\mathbf{u} \in [-1, 1]$, i.e. $\mathbf{u} = \min\{\max\{\mathbf{z}, -1\}, 1\}$.¹ Afterwards, we copy the entire
24 network we used for the binary proof to layer 3 (the original m hidden nodes as layer 3's first m nodes) and to the
25 output layer \mathbf{o} (the original scalar output as the first observation o_1) but add 2 more nodes to layer 3 and one node for
26 output. The previous binary argument makes sure that if \mathbf{u} has to be binary, we could solve 3SAT. Meanwhile, we let
27 the two additional nodes on layer 3 to be $a = \sum_i \max\{u_i, 0\}$ and $b = \sum_i \min\{u_i, 0\}$ and the second observation o_2
28 to be $a + b$, which actually satisfies $a + b \equiv \sum_i |u_i|$. At inference time, we let this node o_2 to be equal to k . In this way
29 to test for exact recovery, $\sum_{i=1}^k |u_i| = k$ and one has to let each u_i to be +1 or -1.

30 Therefore we show that for a 4-layer real network, it is NP-hard to determine if it could be exactly recovered for a given
31 observation.

32 **Q (R2): Comparison with 'invertibility of convolutional neural networks' or other RIP properties**

33 **A:** We will add the discussion with (Gilbert et al.) and (Bourrier et al.) as suggested. Our work is substantially different
34 from (Gilbert et al.) since they still work on linear mappings (convolutional layers without activation functions), while
35 our work is targeting ReLU or LeakyReLU activations.

36 Thank you for the pointer of l_∞ -RIP property. Typically RIP is for fat matrices with structured input while we deal with
37 tall matrices. Also we only need the lower bound side of the RIP, so we redefined the condition in the paper. But indeed
38 the transpose of the weight matrices should satisfy lower bound side of l_∞ -RIP. We will add the discussions properly in
39 the revised version.

40 **Q (R2): the hypothesis that m_i coordinates need to be bounded away from zero is not natural**

41 **A:** For the l_∞ case, we have shown that for random matrices, we do not explicitly need this requirement, as shown
42 in Corollary 1. When the network is expansive and with random weights, there will be enough mass on the positive
43 observations with high probability.

44 **Q (R3): Additional Experiments on DCGAN:**

45 Thank you for your suggestion. We trained a DCGAN architecture using MNIST data. We used 3 convolutional layers
46 with ReLU activations that are expanding and the last layer being convolutional with sigmoid activation. We used
47 projected gradient descent (PGD, gradient descent on the last layer, and projection over the first 3 layers) for a denoising
48 task. We use our algorithm as an initialization for the projection step. Over multiple runs, we compare inversion using
49 1) GD, 2) GD with momentum, 3) PGD with random initialization (for the projection step), and 4) PGD with our
50 scheme for initialization. The average relative error for each method was: 1) 0.26, 2) 0.25, 3) 0.17, and 4) 0.088. This
51 additional experiment shows the benefits of our method for convolutional architectures. We will include more details in
52 the final paper.

¹Here $\max\{\mathbf{z}, -1\}$ could be achieved through $\text{ReLU}(\mathbf{z} + 1) - 1$, and similarly for the min operation.