

1 We thank the reviewers for their time and their reviews. We address the questions below.

2 **Ambiguities beyond flips and rotations (R3)** As pointed out by R3, the albedo of the hidden scene is fundamentally  
3 ambiguous, as any intensity can be compensated by a reciprocal intensity in the transport matrix. To anchor the solution  
4 colors, we use the common “gray world assumption”, and impose it by a simple chromaticity prior that discourages  
5 large differences between color channels. The color of the observed scene can therefore tint the colors of the hidden  
6 scene solution.

7 The space of ambiguities and potential distortions can be characterized as follows (see Koenderink et al., The Generic  
8 Bilinear Calibration-Estimation Problem). Let  $T_0$  and  $L_0$  be the true underlying factors, the observed video thus being  
9  $Z = T_0 L_0$ . All “valid” factorizations are of the form  $T = T_0 A^\dagger$  and  $L = A L_0$ , where  $A$  is chosen (almost) arbitrarily  
10 and  $A^\dagger$  is its (pseudo)inverse. This can be seen by substituting  $TL = (T_0 A^\dagger)(A L_0) = T_0(A^\dagger A)L_0 = T_0 L_0 = Z$ .

11 The result of any factorization implicitly corresponds to some choice of  $A$  and  $A^\dagger$ . In simple cases, the matrix  $A$  can  
12 represent e.g. a permutation that flips the image, whence  $A^\dagger$  is a flip that restores the original orientation: this case is  
13 illustrated in Figure 4’s conversely flipped matrices. They can also represent complementary color transformations  
14 as discussed above. However, for classical factorization methods, they tend to consist of unstructured “noise” that  
15 scrambles the image-like structure in  $T_0$  and  $L_0$  beyond recognition.

16 Our finding in the paper is that via DIP-based factorization, these transformations instead tend to express continuous and  
17 bijective image warps (and color modulations) that preserve the general image structure. As observed by R3, this does  
18 in practice include more complex distortions than just flips and rotations — see for example the nonlinear stretching of  
19 the cameraman image in Figure 4. In full two dimensions, there is room for more complex distortions, but we still find  
20 that e.g. the relative motions of independent objects often remain readable.

21 **Geometric complexity (R3)** We assume that the scene contains a sufficient amount of geometric complexity to  
22 generate high-frequency features like shadows. This improves the conditioning of the problem, as discussed in the  
23 literature on frequency analysis of light transport effects (see e.g. A Theory of Locally Linear Light Transport by  
24 Mahajan et al.). We will emphasize this in the revised paper.

25 **Comparisons to previous methods (R4)** To our knowledge, no existing work attempts to solve the problem under a  
26 similarly general setup, with no assumptions about the shapes viewed in the scene. Attempts to use standard factorization  
27 methods consistently produce unstructured and scrambled results, analogous to the baselines in Figure 4. An example is  
28 seen in the supplemental video, where an SVD factorization is visualized at time 2:00 - 2:09.

29 To provide a comparison, we generalized a recent algorithm that addresses the closest analogue we could think of, i.e.  
30 blind deconvolution with a classical sparse gradient prior that models natural image statistics. As discussed in Section  
31 5.2, we were unable to obtain competitive results despite fair efforts put into the experiment.

32 Regarding comparisons to other non-line-of-sight methods: active non-line-of-sight methods are outside the scope of  
33 the paper, as these techniques assume fundamentally different imaging modalities (usually static hidden scenes, actively  
34 probed over an extended period of time). Similarly, a recent passive non-line-of-sight technique by Bouman et al.  
35 (Turning Corners into Cameras: Principles and Methods, ICCV 2017) assumes a specific scene geometry with clearly  
36 defined “corners” and focuses on near-invisible signals, while our method assumes the geometry and reflectances of the  
37 relay objects are unknown.

38 **Validity of reconstructions for machine vision tasks (R4)** When factorizing light transport using traditional factor-  
39 ization techniques, there is no guarantee that the two factors correspond to the true visible and hidden scene, or even to  
40 any plausible image signal. The question raised by R4 is whether the factors reconstructed using DIP actually correspond  
41 to the true visible and hidden scenes, or are just arbitrary natural-looking images. Figure 6 of the paper provides a  
42 partial answer to that question. We show that for controlled scenes such as the Disks sequence, the reconstructed signal  
43 is clearly not arbitrary, but closely matches the ground truth sequence. Even for more complicated sequences (Hands),  
44 this correspondence still appears to hold.