

Efficient Separation of Convolutional Image Mixtures*

Sarit Shwartz, Yoav Y. Schechner, and Michael Zibulevsky

Dept. Electrical Engineering, Technion - Israel Inst. Tech., Haifa 32000, Israel
psarit@tx.technion.ac.il, {yoav, mzib}@ee.technion.ac.il

Abstract. Convolutional mixtures of images are common in photography of semi-reflections. They also occur in microscopy and tomography. Their formation process involves focusing on an object layer, over which defocused layers are superimposed. Blind source separation (BSS) of convolutional image mixtures by direct optimization of mutual information is very complex and suffers from local minima. Thus, we devise an efficient approach to solve these problems. Our method is fast, while achieving high quality image separation. The convolutional BSS problem is converted into a set of instantaneous (pointwise) problems, using a short time Fourier transform (STFT). Standard BSS solutions for instantaneous problems suffer, however, from scale and permutation ambiguities. We overcome these ambiguities by exploiting a parametric model of the defocus point spread function. Moreover, we enhance the efficiency of the approach by exploiting the sparsity of the STFT representation as a prior.

1 Introduction

Typical blind source separation (BSS) methods seek separation when the mixing process is unknown. However, loose prior knowledge regarding the mixing process often exists, due to its physical origin. In particular, this process can be represented by a parametric form, rather than a trivial representation of raw numbers. For example, consider convolutional image mixtures caused by defocus blur. This blur can be parameterized, yet the parameters' values are unknown. Such mixtures occur in tomography and microscopy [1, 2]. They also occur in semi-reflections [1], e.g., from a glass window: a scene imaged behind the semi-reflector is superimposed on a reflected scene [3, 4]. Each scene is at a different distance from the camera, thus differently defocus blurred in the mixtures.

We claim that BSS can benefit from such a parametrization, as it makes the estimation more efficient while helping to alleviate ambiguities. In the case of semi-reflections, our goal is to decompose the mixed and blurred images into

* This research has been supported in parts by the "Dvorah" Fund of the Technion and by the HASSIP Research Network Program HPRN-CT-2002-00285, sponsored by the European Commission. The research was carried out in the Ollendorff Minerva Center. Minerva is funded through the BMBF. Yoav Schechner is a Landau Fellow-supported by the Taub Foundation, and an Alon Fellow.

the separate scene layers, by minimizing the mutual information (MI) of the estimated objects. An attempt by Ref. [1] used exhaustive search, hence being computationally prohibitive. Ref. [5] attempted convolutional image separation by minimization of higher order cumulant. That method suffers from a scale ambiguity: the sources are reconstructed up to an unknown filter. Moreover, the method's complexity increases fast with the support of the separation kernel.

The complexity of convolutional source separation has been reduced in the domain of acoustic signals, by using frequency methods [6, 7]. There, BSS is decomposed into several small pointwise problems by applying a short-time-Fourier transform (STFT). Then, standard BSS tools are applied to each of the STFT channels. However, these tools suffer from fundamental ambiguities, which may reduce the overall separation quality. Ref. [8] suggested that these ambiguities can be overcome by nonlinear operations in the image domain. However, this method encountered performance problems when simulated over natural images.

We show that these problems can be efficiently solved by exploiting a parametric model for the unknown blur. Moreover, we use the sparsity of STFT coefficients to yield a practically unique solution, which is derived fast. The algorithm is demonstrated in simulations of semi-reflected natural scenes.

2 Problem Formulation

Let $\{s_1, \dots, s_K\}$ be a set of K independent sources. Each source is of the form $s_k = s_k(\mathbf{x})$, $k = 1, \dots, K$, where $\mathbf{x} = (x, y)$ is a two dimensional (2D) spatial coordinate vector in the case of images. Let $\{m_1, \dots, m_K\}$ be a set of K measured signals, each of which is a linear mixture of a convolved version of the sources

$$m_i(\mathbf{x}) = a_{i1} * s_1(\mathbf{x}) + \dots + a_{iK} * s_K(\mathbf{x}) \quad , i = 1, \dots, K \quad (1)$$

Here $*$ denotes convolution and $a_{ik}(\mathbf{x})$, $k = 1, \dots, K$, are linear spatially invariant filters. Denote $\{\hat{s}_1, \dots, \hat{s}_K\}$ as the set of the reconstructed sources. Reconstruction is done by applying a linear operator \mathbf{W} on $\{m_1, \dots, m_K\}$. Each of the reconstructed sources is of the form

$$\hat{s}_k(\mathbf{x}) = w_{k1} * m_1(\mathbf{x}) + \dots + w_{kK} * m_K(\mathbf{x}) \quad , k = 1, \dots, K \quad (2)$$

where $w_{ik}(\mathbf{x})$ are linear spatially invariant filters. Our goal is: given only the measured signals $\{m_1, \dots, m_K\}$, find a linear separation operator \mathbf{W} that inverts the mixing process, thereby separating the sources. The mixing process is inverted by finding \mathbf{W} that minimizes the MI of $\{\hat{s}_1, \dots, \hat{s}_K\}$.

MI is expressed by using the marginal entropies $\mathcal{H}_{\hat{s}_k}$ and the joint entropy of the estimated sources $\mathcal{H}_{\hat{s}_1, \hat{s}_2}$ as $\mathcal{I}_{\hat{s}_1, \hat{s}_2} = \sum_{k=1}^K \mathcal{H}_{\hat{s}_k} - \mathcal{H}_{\hat{s}_1, \dots, \hat{s}_K}$. However, estimation of the joint entropy may be unreliable. It can be avoided if the mixtures are pointwise, rather than convolutional. In pointwise mixtures, the separation operator \mathbf{W} is a simple matrix, termed the separation matrix. In this case, the MI can be expressed as (see for example Ref. [9])

$$\mathcal{I}(\hat{s}_1, \hat{s}_2) = -\log |\det(\mathbf{W})| + \sum_{k=1}^K \mathcal{H}_{\hat{s}_k} \quad (3)$$

It is desirable to do the same for convolutive mixtures. However, if \mathbf{W} is a convolutive operator, Eq. (3) does not hold. We note that expressions similar to (3) have been developed for convolutive mixtures [10] assuming spatially white sources. Nevertheless, algorithms based on these expressions suffer from whitening of the separated sources, corrupting the estimation severely both in acoustic and in imaging applications.

3 Efficient Separation of Convolutive Image Mixtures

We may use Eq. (3) in convolutive mixtures, despite the fact that it is valid only in pointwise mixtures. This is achieved by decomposing the convolutive optimization problem into several smaller ones, which are apparently independent of each other. This approach is inspired by frequency domain algorithms developed for acoustic signals [6, 7]. Nevertheless, this approach has its own fundamental limitations, which are discussed and solved in Secs. 4 and 5.

We apply STFT¹ to the data. Denote $\boldsymbol{\omega} = (\omega_x, \omega_y)$ as the index vector of the frequency variable of the 2D STFT. Assuming that the STFT window size is larger than the effective width² of the blur kernel [6], Eq. (1) becomes

$$m_i(\boldsymbol{\omega}, \mathbf{x}) \approx a_{i1}(\boldsymbol{\omega})s_1(\boldsymbol{\omega}, \mathbf{x}) + \dots + a_{iK}(\boldsymbol{\omega})s_K(\boldsymbol{\omega}, \mathbf{x}), \quad i = 1, \dots, K, \quad (4)$$

since convolution becomes a multiplication in this domain.

Eq. (4) exposes a fundamental problem in cases of energy-preserving convolution operators. In such operators $a_{ik}(\boldsymbol{\omega}) \rightarrow 1$ as $\boldsymbol{\omega} \rightarrow 0$ (the overall light energy over the image area is invariant to the convolution). This occurs in defocus blur, since change of focus does not cause light attenuation, only a different spread of the light energy across the sensor area [1, 2]. As $a_{ik}(\boldsymbol{\omega}) \rightarrow 1$, Eq. (4) becomes

$$m_i(\boldsymbol{\omega}, \mathbf{x}) \approx s_1(\boldsymbol{\omega}, \mathbf{x}) + \dots + s_K(\boldsymbol{\omega}, \mathbf{x}), \quad i = 1, \dots, K. \quad (5)$$

This is a singular set of equations. Therefore, low spatial frequencies are not well reconstructed. Note that this has nothing to do with the ICA problem. Even if the blur kernels a_{ik} are *perfectly known*, the reconstruction is ill-conditioned in the low-frequency bands [1, 2]. Keeping in mind this matter, we continue with the blind estimation process. Note that at each sub-band $\boldsymbol{\omega}$, Eq. (4) expresses a pointwise mixture of sub-band images. At each frequency channel, the mixed sources can be separated by simple ICA optimization. Then, all the separated sources from all the frequency channels may be combined by inverse STFT.

To describe the ICA optimization, denote $\mathbf{W}(\boldsymbol{\omega})$ as the separation matrix at channel $\boldsymbol{\omega}$. In addition, denote $\mathcal{I}^\omega(\hat{s}_1, \hat{s}_2)$ and $\mathcal{H}_{\hat{s}_k}^\omega$ as the MI and marginal entropies of the estimated sources at channel $\boldsymbol{\omega}$, respectively. Then, similarly to Eq. (3), the MI of the estimated sources at each channel is given by

$$\min_{\mathbf{w}(\boldsymbol{\omega})} \left\{ -\log |\det[\mathbf{W}(\boldsymbol{\omega})]| + \sum_{k=1}^K \hat{\mathcal{H}}_{\hat{s}_k}^\omega \right\}, \quad (6)$$

¹ This operation is also termed as a *windowed Fourier transform*, which may be more appropriate for spatial coordinates as we use.

² A discussion regarding the STFT window width is given in Sec. 7.

where $\hat{\mathcal{H}}_{\hat{s}_k}^\omega$ is an estimator of the channel entropy of an estimated source. Hence, using this factorization, MI minimization of a convolutional mixture is expected to be both more accurate and more efficient to obtain.

Sparse Separation in the STFT Domain

Now, we exploit image statistics in order to achieve a computationally efficient solution for the sub-problems in each frequency channel. As shown in [11], sparsity of sources is a strong prior that can be exploited to achieve a very efficient separation. It is known from studies of image statistics (see for example [12]) that sub-band images are *sparse signals*. Motivated by [11, 13], their quasi-maximum likelihood blind separation can be achieved via minimization of

$$\min_{\mathbf{w}(\omega)} \left\{ -\log |\det[\mathbf{W}(\omega)]| + (1/N) \sum_{k=1}^K \sum_{n=1}^N |\hat{s}_k(\omega, n)| \right\}. \quad (7)$$

Here n indexes the STFT shift (out of a total of N). This enables relative Newton optimization [14], which enhances the efficiency of sparse source separation.

4 Inherent Problems

The frequency representation brings efficiency of pointwise separation. With it, however, come fundamental ambiguities that are common in pointwise problems. The *permutation ambiguity* implies that the separated sub-band images appear at each channel in a random permutation. Some sub-band images corresponding to the “first” estimated source may actually belong to the “second” estimated source. When the channels are transformed back to the image domain using the inverse STFT, the reconstructed images can suffer from crosstalk. Even though source separation was achieved in each channel independently, distinct sub-band images from different sources are combined in the reconstruction.

In addition, the scale of different channels is unknown due to *scale ambiguity*, leading to imbalance between frequency channels. When the estimated channels of a source are transformed back to the image domain using the inverse STFT, the reconstructed image can appear unnatural and suffer from artifacts.

Moreover, the performance in each frequency channel is frequency dependent. Typically, there are a few frequency channels with good separation, a few channels with very poor separation and the rest of the channels have mediocre separation quality. There are several reasons for this phenomenon. One reason is related to the different sparsity of different frequency channels [15].

5 Inter-channel Knowledge Transfer

In this section we bypass the permutation and scale ambiguities by exploiting a prior about the unknown convolutional process. Blur caused by optical defocus can be parameterized [16]. As an example, consider a rough parametric model:

a simple 2D Gaussian kernel with different widths in the x and y directions [1]. Denote $\xi_{i,k} = [\xi_{i,k,x}, \xi_{i,k,y}]$ as the vector of the unknown blur parameters of the blur kernel of source k at image i and

$$G_{\xi_{i,k}}(\omega) = \exp[-\omega_x^2/(2\xi_{i,k,x}^2)] \exp[-\omega_y^2/(2\xi_{i,k,y}^2)] \quad (8)$$

as the filter which preserves light energy. In addition to defocus, let us incorporate attenuation $g_{i,k}$ of each source k into any mixture i .

Assume that in each acquired image, one of the layers is focused,³ i.e. $G_{\xi_{k,k}} = 1$. Define $\mathbf{A}(\omega)$ as the mixing operator in frequency channel ω .

$$\mathbf{A}(\omega) = \begin{bmatrix} 1 & g_{1,2}G_{\xi_{1,2}}(\omega) & \dots & \dots \\ g_{2,1}G_{\xi_{2,1}}(\omega) & 1 & \vdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \dots & \dots & g_{K,K-1}G_{\xi_{K,K-1}}(\omega) & 1 \end{bmatrix}. \quad (9)$$

Thus, the separation matrix in each channel is parameterized by $\xi_{i,k}$ and $g_{i,k}$ and is of the form $\mathbf{W}(\omega) = [\mathbf{A}(\omega)]^{-1}$. Note that the parameter $\xi_{i,k}$ and $g_{i,k}$ are the *same for all frequency channels*. Hence, there is a small number of actual unknown blur variables. On the other hand, there is a large number of frequency channels upon which the estimation of these variables can be based.

As we explain in Sec. 5.2, we can automatically select three channels ω^a , ω^b and ω^c , that yield the best separation results according to a ranking criterion. Define $\tilde{\mathbf{A}}(\omega^a) = [\mathbf{W}(\omega^a)]^{-1}$ and similarly $\tilde{\mathbf{A}}(\omega^b)$ and $\tilde{\mathbf{A}}(\omega^c)$. Let $\tilde{a}_{i,k}$ be the coefficients of $\tilde{\mathbf{A}}$. Then, for each blur kernel, we calculate the unknown blur parameters $\xi_{i,k}$ and $g_{i,k}$ by solving the following set of equations:

$$\begin{cases} g_{i,k}G_{\xi_{i,k}}(\omega^a) = \tilde{a}_{i,k}(\omega^a)/\tilde{a}_{i,i}(\omega^a) \\ g_{i,k}G_{\xi_{i,k}}(\omega^b) = \tilde{a}_{i,k}(\omega^b)/\tilde{a}_{i,i}(\omega^b) \\ g_{i,k}G_{\xi_{i,k}}(\omega^c) = \tilde{a}_{i,k}(\omega^c)/\tilde{a}_{i,i}(\omega^c) \end{cases}, \quad (10)$$

We solve this set to find the parameters $\xi_{i,k}$ and $g_{i,k}$, thus deriving the blur and attenuation parameters based on those few selected channels.⁴

Now, we can use these parameters and Eq. (8) to calculate $g_{i,k}G_{\xi_{i,k}}(\omega)$ for all the frequency channels. This directly yields the separation operator \mathbf{W} for *all the frequency channels*. We invert the mixing process by using this \mathbf{W} . It may be possible to achieve higher accuracy by representing each blur kernel using parametric models other than Gaussian, requiring more parameters. This would require selection of additional channels.

³ We stress that we seek layer *separation* rather than *deblurring*. Therefore, if source k is defocused in all the images, we denote the least defocused version of source k as the effective source we aim to reconstruct. Then, we denote $G_{\xi_{i,k}}(\omega)$ as the relative defocus filter between the effective source and the defocused source at image i .

⁴ One might suggest optimizing the MI directly over the parameters $g_{i,k}$ and $\xi_{i,j}$. However, this optimization scheme is not necessarily convex. A detailed discussion on this issue is given in [15].

5.1 Separation of Semi-reflections

Section 5 describes a parametric model for mixtures of blurred images. It consists of an attenuation factor $g_{i,k}$ and an energy preserving filter $G_{\xi_{i,k}}$. However, in common applications such as semi-reflections [1] or widefield optical sectioning [2], no attenuation accompanies the change of focus. Hence, $g_{i,k} = 1$ for all i, k . For each signal, each source is affected only by two parameters in the Gaussian model. Thus, only two channels are needed to solve for the unknown $\xi_{i,k}$. Moreover, in the special case of semi-reflections, we have only two sources. Therefore, the mixing operator and the separation operator are reduced to

$$\mathbf{A}(\omega) = \begin{bmatrix} 1 & G_{\xi_{1,2}} \\ G_{\xi_{2,1}}(\omega) & 1 \end{bmatrix}, \quad \mathbf{W}(\omega) = \begin{bmatrix} 1 & -G_{\xi_{1,2}} \\ -G_{\xi_{2,1}}(\omega) & 1 \end{bmatrix} \{\det(|\mathbf{A}(\omega)|)\}^{-1}. \quad (11)$$

The equation system we need to solve in order to estimate $\xi_{1,2}$ and $\xi_{2,1}$ is

$$\begin{cases} -G_{\xi_{1,2}}(\omega^a) = w_{1,2}(\omega^a)/w_{1,1}(\omega^a) \\ -G_{\xi_{1,2}}(\omega^b) = w_{1,2}(\omega^b)/w_{1,1}(\omega^b) \\ -G_{\xi_{2,1}}(\omega^a) = w_{2,1}(\omega^a)/w_{2,2}(\omega^a) \\ -G_{\xi_{2,1}}(\omega^b) = w_{2,1}(\omega^b)/w_{2,2}(\omega^b) \end{cases}. \quad (12)$$

Here, $w_{i,k}$ are the coefficients of matrix $\mathbf{W}(\omega)$ and ω^a, ω^b are the best and second best channels according to the ranking we describe next.⁵

We stress that thanks to this approach of parameter-based inter-channel knowledge transfer, the permutation, scale and sign ambiguities are solved: the sources are not derived in a random order or with inter-channel imbalance, but in a way that must be consistent with the blur model, hence with the image formation process. In addition, the problem of channel and data dependent performance is alleviated, since the separation operator is estimated based on selected channels performing well.

5.2 Selecting a Good Frequency Channels

The parameter estimation method requires *ranking* of the channels. The ranking relies on a quality criterion for the separation (i.e., independence) of \hat{s}_1 and \hat{s}_2 at each frequency channel ω , given the sparsity assumption.

The scatter plot of sparse independent signals has a cross shape aligned with the axes, in the (\hat{s}_1, \hat{s}_2) plane, i.e., most of the samples should have small angles relative to the \hat{s}_1 and \hat{s}_2 axes. Define

$$\chi_{\mathcal{L}_1}^{\omega} = \sum_{k=1}^2 \left(\left\{ \sum_{n=1}^N |\hat{s}_k(\omega, n)| \right\} / \left\{ \sum_{n=1}^N [\hat{s}_k(\omega, n)]^2 \right\} \right). \quad (13)$$

This criterion increases as the samples in the scatter plot deviate from the \hat{s}_1 and \hat{s}_2 axes, and is reduced when each sample n has non-zero values exclusively in

⁵ It might be possible to achieve better estimation by using more than two channels, for example, by solving a non-linear least squares problem.

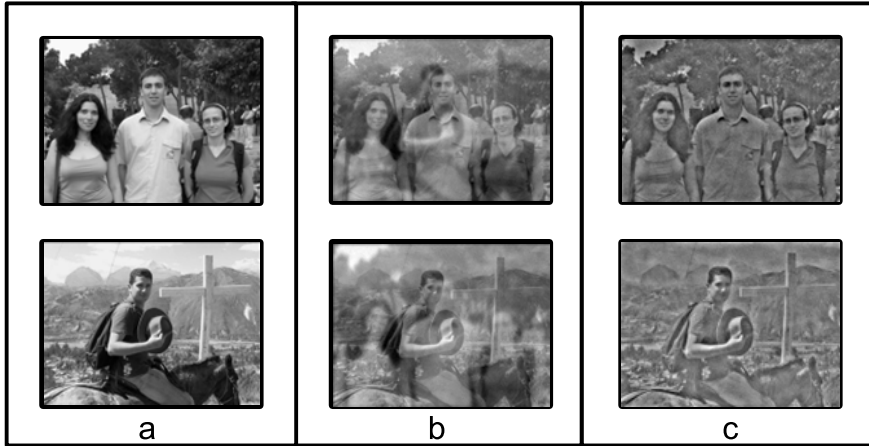


Fig. 1. Simulation results: (a) Two original natural images. (b) The two convolved and mixed images. (c) Reconstructed layers.

\hat{s}_1 or \hat{s}_2 . This closed form expression automatically determines which frequency channels yield the most separated sources, and are thus preferable.

Thus, in our algorithm, we first perform ICA in all the frequency channels. We then calculate $\chi_{\mathcal{Z}_1}^{\omega}$, thus ranking the channels. Then, we select the best channels as those that correspond to the smallest values of the penalty function $\chi_{\mathcal{Z}_1}^{\omega}$. These channels are used in Sec. 5.

6 Demonstration

The method was simulated using two natural images of size 122×162 pixels (Fig. 1a) as the two scene layers. The blur kernels we used are Gaussians with parameter vectors $\xi_{1,2} = [1, 2]$ and $\xi_{2,1} = [2, 1]$ pixels. We did not use attenuation coefficients because in photography of real semi-reflections, the image layers are only blurred but not attenuated by change of focus. We added i.i.d Gaussian noise with standard deviation of ~ 2.5 gray levels to the convolved and mixed images. The resulting mixed and noisy images are shown in Fig. 1b. Separation was performed using STFT having 13×13 frequency channels. The separation results are presented in Fig. 1c. The resulting images are indeed well separated. There is no visible crosstalk between the images. The contrast of the reconstructed images is reduced compared to the original images. This stems from inherent ill-conditioning of the mixing matrix at low frequencies (see Sec. 4), i.e., this is not associated with the blindness of the separation problem.

7 Discussion

The convolutive image separation algorithm has currently a single parameter to tweak: the width of the STFT window. It can affect the separation results, and

the optimal size somewhat depends on the acquired images. As mentioned in Sec. 3, it must be larger than the effective width of the blur kernel. On the other hand, a very wide window can degrade the sparsity of the sub-band images. A detailed discussion is given in [15]. We determined the window width by trial and error, but we believe this can be automated. For example, multi-window STFT may be used, followed by selection of the the best window width using the criterion described in Sec. 5.2. This requires further research.

References

1. Schechner, Y.Y., Kiryati, N., Basri, R.: Separation of transparent layers using focus. *Int. J. Computer Vision* **89** (2000) 25–39
2. Macias-Garza, F., Bovik, A.C., Diller, K.R., Aggarwal, S.J., Aggarwal, J.K.: The missing cone problem and low-pass distortion in optical serial sectioning microscopy. In: *Proc. ICASSP*. Volume 2. (1988) 890–893
3. Schechner, Y.Y., Shamir, J., Kiryati, N.: Polarization and statistical analysis of scenes containing a semi-reflector. *J. Opt. Soc. America A* **17** (2000) 276–284
4. Bronstein, A.M., Bronstein, M.M., Zibulevsky, M., Zeevi, Y.Y.: Sparse ICA for blind separation of transmitted and reflected images. *Intl. J. Imaging Science and Technology* **15**(1) (2005) 84–91
5. Castella, M., Pesquet, J.C.: An iterative blind source separation method for convolutional mixtures of images. In: *Proc. ICA2004*. (2004) 922–929
6. Parra, L., Spence, C.: Convolutional blind separation of non-stationary sources. *IEEE Trans. on Speech and Audio Processing* **8** (2000) 320–327
7. Smaragdakis, P.: Blind separation of convolved mixtures in the frequency domain. *Neurocomputing* **22** (1998) 21–34
8. Kasprzak, W., Okazaki, A.: Blind deconvolution of timely-correlated sources by homomorphic filtering in Fourier space. In: *Proc. ICA2003*. (2003) 1029–34
9. Hyvärinen, A., Karhunen, J., Oja, E.: *Independent component analysis*. John Wiley and Sons, NY (2001)
10. Pham, D.T.: Contrast functions for blind source separation and deconvolution of sources. In: *Proc. ICA2001*. (2001) 37–42
11. Zibulevsky, M., Pearlmutter, B.A.: Blind source separation by sparse decomposition in a signal dictionary. *Neural Computations* **13**(4) (2001) 863–882
12. Simoncelli, E.P.: Statistical models for images: Compression, restoration and synthesis. In: *Proc. IEEE Asilomar Conf. Sig. Sys. and Computers*. (1997) 673–678
13. Pham, D.T., Garrat, P.: Blind separation of a mixture of independent sources through a quasi-maximum likelihood approach. *IEEE Trans. Sig. Proc.* **45**(7) (1997) 1712–1725
14. Zibulevsky, M.: Blind source separation with relative newton method. In: *Proc. ICA2003*. (2003) 897–902
15. Shwartz, S., Schechner, Y.Y., Zibulevsky, M.: Efficient blind separation of convolutional image mixtures. Technical report, CCIT No. 553, Dep. Elec. Eng., Technion Israel Inst.Tech. (2005)
16. Born, M., Wolf, E.: *Principles of optics*. Pergamon, Oxford (1975)