



Figure 1: (a) Some of the features extracted by LFOICA. (b) std. of the results of each method and the p-value of Welch’s  $t$ -test between each baseline method and our LFOICA.

1 We would like to thank all reviewers for the constructive comments. We will fully address all the review concerns in the  
 2 revision. The citation indices below are consistent with those in the main paper.

3 **To R#1: MoG and sparsity have been studied.** It is true that MoG and sparsity have been extensively studied in  
 4 existing ICA methods. Aside from the two references suggested by the reviewer, we also cited several representative  
 5 works in this direction ([14,15,16]). In fact, one major motivation of the proposed LFOICA is that the classical MoG-  
 6 based ICA methods use maximum likelihood to estimate the parameters, which involve computationally intractable  
 7 posterior inference. In our practical considerations, we discussed in the paper that LFOICA can be further specified with  
 8 additional constraints, such as MoG and sparsity, in a likelihood-free way that enjoys both statistical and computational  
 9 efficiency. We will further discuss the suggested references in the "practical consideration" subsection.

10 **Identifiability and convergence of our model.** The identifiability of the mixing matrix  $\mathbf{A}$  in our model ( $\mathbf{x} =$   
 11  $G_{\mathbf{A},\theta}(\mathbf{z}) = \mathbf{A}[f_{\theta_1}(z_1), \dots, f_{\theta_d}(z_d)]^\top$ ) follows the identifiability results for OICA [29], which is summarized in  
 12 the following Theorem. **Theorem 1** Given two OICA models  $\mathbf{x} = G_{\mathbf{A},\theta}(\mathbf{z})$  and  $\mathbf{x}' = G_{\mathbf{A}',\theta'}(\mathbf{z}')$  that specify  
 13 distributions  $\mathbb{P}(\mathbf{x})$  and  $\mathbb{P}(\mathbf{x}')$ , respectively. Under the condition that none of the distributions of  $f_i(\mathbf{z}_i)$  is Gaussian, if  
 14  $MMD(\mathbb{P}(\mathbf{x}), \mathbb{P}(\mathbf{x}')) = 0$ , then  $\mathbf{A}' = \mathbf{A}\mathbf{P}_p\mathbf{S}_p$ , where  $\mathbf{P}_p$  is a  $p \times p$  column permutation matrix and  $\mathbf{S}_p$  is a  $p \times p$   
 15 scaling matrix. The proof is almost the same as that of Theorem 3 in [29], except that in order to guarantee  $\mathbb{P}(\mathbf{x}) = \mathbb{P}(\mathbf{x}')$ ,  
 16 we use  $MMD = 0$  while [29] uses maximum likelihood (KL divergence). Given the identifiability results, the  
 17 estimated mixing matrix converges to the scaled and permuted version of the true mixing matrix and so does the  
 18 source distributions. With additional constraints on  $\mathbf{A}$ , the permutation and scaling ambiguity can be further reduced.  
 19 The proofs to this are provided in the original papers of the two causal discovery problems evaluated in our paper.  
 20 **The parameters in our MLPs (i.e.,  $\theta$ ) are not identifiable ( $\theta \neq \theta'$ ), but we do not need the identifiability of  $\theta$  to**  
 21 **perform certain downstream tasks, such as the two causal discovery tasks mentioned in our paper.**

22 **GAN-based methods to solve nonlinear ICA [21].** Although both our work and [21] use a GAN style approach to  
 23 solve ICA, they are largely different to each other. **First**, the main purpose of [21] is to recover the ICs instead of how the  
 24 ICs are mixed (i.e., the mixing matrix). It models the mixing and unmixing procedure implicitly with an encoder-decoder  
 25 architecture. As a consequence of non-linearity, there is no guarantee for identifiability. In contrast, we concentrate  
 26 on the mixing matrix estimation for causal discovery purpose. **Second**, the encoder-decoder architecture in [21] cannot  
 27 be easily extended for OICA because the posterior of ICs cannot be modeled by a deterministic encoder. **Third**,  
 28 the adversarial training target of LFOICA and [21] are different. While [21] aims at matching the joint distribution  
 29 and product of marginal distribution of the recovered ICs (this is also how [21] makes the components independent),  
 30 LFOICA is trained to match the distributions of the generated mixtures and true mixtures. And the estimated ICs by  
 31 LFOICA are naturally independent because they are generated from independent latent noises with separate networks.

32 **To R#2: Describe the details of standard overcomplete ICA.** Thanks for the suggestion. We will add an algorithmic  
 33 description that summarizes maximum likelihood (the EM-type) algorithms for OICA in the Supplementary Material.

34 **To R#3: The generative model does not have an inference algorithm.** Since the main focus of our work is OICA  
 35 and its applications to causal discovery, we put significant effort on mixing matrix estimation instead of latent variable  
 36 inference. Nevertheless, this is a valuable suggestion and we’ll consider extending our LFOICA.

37 **Comparison to score-matching ICA.** We conduct additional experiments with the suggested score matching ICA  
 38 (SMICA) under the same conditions. For each experiment setting in Table 1 in the paper, MSE of the recovered mixing  
 39 matrix obtained by SMICA are 3.52e-2, 6.70e-2, 8.11e-3 and 1.94e-2 respectively, which is not as good as LFOICA.

40 **Scalability to higher-dimensional data such as images.** LFOICA is originally designed for causal discovery, which  
 41 is different from image feature extraction. To demonstrate the generality of the method to higher dimensional problems,  
 42 we apply it to extract features from image patches. The features form a dictionary which can be used with the recovered  
 43 ICs for image restoration tasks, such as image denoising. Example features are shown in the Figure 1(a). Patch size is  
 44 16x16 and the number of atoms is set to 400. We use DCT dictionary for initialization.

45 **"empirical estimator".** It refers to Maximum Mean Discrepancy (MMD), which is mentioned in line 83.

46 **"s.e.m.s? How significant?".** Here we interpret "s.e.m.s" as standard deviation of the MSE values in Table 1 of our  
 47 paper and report it as well as the p-values by conduct Welch’s  $t$ -test between baseline works and our LFOICA in Figure  
 48 1(b) (this figure corresponds to Table 1 in our paper). As one can see, the differences are significant enough.