

# Multi-reference alignment in high dimensions: sample complexity and phase transition

Elad Romanov <sup>\*1</sup>, Tamir Bendory <sup>†2</sup>, and Or Ordentlich<sup>‡1</sup>

<sup>1</sup>School of Computer Science and Engineering, The Hebrew University, Jerusalem, Israel

<sup>2</sup>School of Electrical Engineering, Tel Aviv University, Tel Aviv, Israel

## Abstract

Multi-reference alignment entails estimating a signal in  $\mathbb{R}^L$  from its circularly-shifted and noisy copies. This problem has been studied thoroughly in recent years, focusing on the finite-dimensional setting (fixed  $L$ ). Motivated by single-particle cryo-electron microscopy, we analyze the sample complexity of the problem in the high-dimensional regime  $L \rightarrow \infty$ . Our analysis uncovers a phase transition phenomenon governed by the parameter  $\alpha = L/(\sigma^2 \log L)$ , where  $\sigma^2$  is the variance of the noise. When  $\alpha > 2$ , the impact of the unknown circular shifts on the sample complexity is minor. Namely, the number of measurements required to achieve a desired accuracy  $\varepsilon$  approaches  $\sigma^2/\varepsilon$  for small  $\varepsilon$ ; this is the sample complexity of estimating a signal in additive white Gaussian noise, which does not involve shifts. In sharp contrast, when  $\alpha \leq 2$ , the problem is significantly harder and the sample complexity grows substantially quicker with  $\sigma^2$ .

## 1 Introduction

We study the sample complexity of the multi-reference alignment (MRA) model: the problem of estimating a signal from its circularly-shifted and noisy copies. Specifically, let  $X \sim \mathcal{N}(0, I)$  be an  $L$ -dimensional vector with i.i.d. standard normal entries. We collect  $n$  independent measurements of random cyclic shifts of  $X$ , corrupted by additive white Gaussian noise:

$$Y_i = R_{\ell_i} X + \sigma Z_i, \quad i = 1, \dots, n, \quad (1.1)$$

where  $R_\ell$  denotes a cyclic shift, namely,  $(R_\ell X)_j = X_{(j+\ell) \bmod L}$  for all  $j = 0, \dots, L-1$ ,  $Z_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, I)$ , and  $\ell_i \stackrel{i.i.d.}{\sim} \text{Uniform}(\{0, \dots, L-1\})$  are statistically independent of  $X$ . Given the measurements  $Y^n = (Y_1, \dots, Y_n)$ , one is interested in constructing an estimator  $\hat{X} = \hat{X}(Y^n)$  of the signal. Importantly, the unknown shifts  $\ell_1, \dots, \ell_n$ —while their estimation might be a means to an end—are nuisance variables.

This paper focuses on the high-dimensional regime, where the dimension of the signal grows indefinitely  $L \rightarrow \infty$ . In this setting, we wish to characterize the relations between the number of measurements  $n$ , the length of each observation  $L$ , and the noise level  $\sigma^2$  that allow estimating  $X$  to

---

\*E-mail: elad.romanov@mail.huji.ac.il

†E-mail: bendory@tauex.tau.ac.il

‡E-mail: or.ordentlich@mail.huji.ac.il

a prescribed accuracy. This is in contrast to previous works, surveyed in Section 3, which analyzed the interplay between  $n$  and  $\sigma$ , while considering a fixed  $L$ .

It is important to note that given the measurements, there is no way to distinguish between  $X$  and its cyclic shift since  $P_{Y^n|X=x} = P_{Y^n|X=R_1x} = \dots = P_{Y^n|X=R_{L-1}x}$ . Therefore, we can only estimate the orbit of  $X$  under the group of circular shifts  $\mathbb{Z}_L$ . Accordingly, we use the following distortion measure

$$\rho(X, \hat{X}) = \frac{1}{L} \min_{\ell=0, \dots, L-1} \|X - R_\ell \hat{X}\|^2. \quad (1.2)$$

In the sequel, we loosely say that we aim to estimate  $X$  rather than its orbit, and refer to  $\mathbb{E}\rho(X, \hat{X})$  as the MSE.

**Sample complexity** Our goal in this paper is to characterize the smallest possible number of measurements required to achieve a desired MSE in terms of the dimension  $L$  and the noise level  $\sigma^2$ . To that end, we define the smallest MSE attainable by any estimator as

$$\text{MSE}_{\text{MRA}}^*(L, \sigma^2, n) := \inf_{\hat{X}} \mathbb{E}\rho(X, \hat{X}(Y^n)), \quad (1.3)$$

and the sample complexity of the MRA problem

$$n_{\text{MRA}}^*(L, \sigma^2, \varepsilon) := \min \{n : \text{MSE}_{\text{MRA}}^*(L, \sigma^2, n) \leq \varepsilon\}. \quad (1.4)$$

We define the signal-to-noise ratio (SNR) by

$$\text{SNR} := \frac{\mathbb{E}\|X\|^2}{\sigma^2} = \frac{L}{\sigma^2}. \quad (1.5)$$

This definition is consistent with previous works which considered a fixed  $L$  and  $\sigma \rightarrow \infty$ , implying  $\text{SNR} \rightarrow 0$ ; see Section 3.

The asymptotics in our model turn out to be particularly interesting when the dimension, the noise level, and the SNR are simultaneously large. In particular, it will be convenient to parametrize the noise variance by

$$\sigma^2(\alpha) = \frac{L}{\alpha \log L} \iff \alpha = \frac{L}{\sigma^2 \log L} = \frac{\text{SNR}}{\log L}. \quad (1.6)$$

Accordingly, we define  $\text{MSE}_{\text{MRA}}^*(L, \alpha, n) := \text{MSE}_{\text{MRA}}^*(L, \sigma^2(\alpha), n)$  and  $n_{\text{MRA}}^*(L, \alpha, \varepsilon) := n_{\text{MRA}}^*(L, \sigma^2(\alpha), \varepsilon)$ .

**Motivation.** The MRA model is mainly motivated by single-particle cryo-electron microscopy (cryo-EM)—a leading technology to constitute the 3-D structure of biological molecules. In its most simplified version, the cryo-EM problem involves reconstructing a 3-D structure from its multiple noisy tomographic projections, taken after the structure has been rotated by an unknown 3-D rotation. In analogy, in the MRA model (1.1) the signal  $X$  is measured after an unknown circular shift. In Theorem 2.3, we extend the basic model to include a projection; we refer to this model as the projected MRA model. This projection plays the role, to some extent, of the tomographic projection in cryo-EM. Section 7 discusses further potential extensions.

The correspondence between MRA and cryo-EM, while admittedly not perfect, has motivated an extensive study of the MRA problem in recent years. For example, the resolution limitations of MRA were analyzed in [BJL<sup>+</sup>20] in order to draw an analogy to the achievable resolution of cryo-EM—a crucial aspect from a biological standpoint. More relevant to this work, in [BBSK<sup>+</sup>17, PWB<sup>+</sup>19, BCLS20, APS18], the sample complexity of the MRA and cryo-EM models were analyzed for a fixed dimension  $L$ . Remarkably, it was shown that in the low noise regime (small  $\sigma$ ), the number of measurements should scale like  $\sigma^2$ , while in the high noise regime (large  $\sigma$ )  $n$  must increase with  $\sigma^6$ ; see further discussion in Section 3.

Our high-dimensional analysis is motivated by the size of modern cryo-EM datasets. In a typical cryo-EM experiment, the number of measurements and the dimension of the 3-D structure are of the same order of a few millions. For example, a 3-D structure of size  $200 \times 200 \times 200$  voxels resulting in 8,000,000 parameters to be estimated. In fact, high-dimensional statistical analysis has been already proven to be effective for cryo-EM data processing. For example, a covariance estimation technique based on high-dimensional analysis (the so-called spiked model) has significantly improved image denoising [BZS16].

**Information-theoretic background and asymptotic notation.** The analysis of this work is greatly based on information-theoretic notions and techniques. For completeness, we review the relevant definitions in Appendix A.

We also repeatedly use asymptotic notation. For sequences  $a = a(L)$  and  $b = b(L)$ , we write  $a(L) = O(b(L))$  if there exists a constant  $C > 0$  such that  $a(L) \leq Cb(L)$  for all  $L$ . Similarly,  $a(L) = \Omega(b(L))$  means  $a(L) \geq Cb(L)$ . Occasionally, we use  $a(L) = O_\beta(b(L))$  to signify explicitly that  $C$  depends on some parameter  $\beta$ . The notation  $a(L) = o(b(L))$  means  $a(L)/b(L) \rightarrow 0$  as  $L \rightarrow \infty$ . In particular, if  $a(L) = o(1)$  then  $a(L) \rightarrow 0$  asymptotically. Similarly,  $a(L) = \omega(b(L))$  means  $a(L)/b(L) \rightarrow \infty$ .

The code to reproduce the figures is publicly available at <https://github.com/TamirBendory/high-dimensional-mra-bounds>.<sup>1</sup>

## 2 Main results and discussion

**Phase transition.** This work focuses on the asymptotic setting where  $L$  tends to infinity. Our first main finding is that in this asymptotic limit there is a transition in terms of the behavior of the sample complexity. For  $\alpha > 2$ , the MRA problem is essentially as easy as estimating a signal in additive white Gaussian noise (AWGN), with no random shifts. More precisely, for sufficiently small distortion  $\varepsilon$ , the sample complexity tends to the sample complexity of estimating a signal in AWGN,  $n_{\text{AWGN}}^*(L, \alpha, \varepsilon) = \lceil (\frac{1}{\varepsilon} - 1) \sigma^2(\alpha) \rceil$ , which behaves as  $\frac{\sigma^2(\alpha)}{\varepsilon}$  for small  $\varepsilon$ . In sharp contrast, for  $\alpha \leq 2$  the problem becomes substantially harder.

---

<sup>1</sup>Our expectation-maximization implementation is based on the code of [BBM<sup>+</sup>17].

**Theorem 2.1** *The sample complexity of the MRA model (1.1) obeys:*

1. *For any  $\alpha > 2$  we have*

$$\lim_{\varepsilon \rightarrow 0} \lim_{L \rightarrow \infty} \frac{n_{MRA}^*(L, \alpha, \varepsilon)}{\sigma^2(\alpha)/\varepsilon} = \lim_{\varepsilon \rightarrow 0} \lim_{L \rightarrow \infty} \frac{n_{MRA}^*(L, \alpha, \varepsilon)}{n_{AWGN}^*(L, \alpha, \varepsilon)} = 1.$$

2. *For any  $\alpha \leq 2$  and any  $\varepsilon < 1$  we have*

$$\lim_{L \rightarrow \infty} \frac{n_{MRA}^*(L, \alpha, \varepsilon)}{\sigma^2(\alpha)/\varepsilon} = \lim_{L \rightarrow \infty} \frac{n_{MRA}^*(L, \alpha, \varepsilon)}{n_{AWGN}^*(L, \alpha, \varepsilon)} = \infty.$$

In part 1 of Theorem 2.1, the lower bound  $\frac{n_{MRA}^*(L, \alpha, \varepsilon)}{n_{AWGN}^*(L, \alpha, \varepsilon)} \geq 1$  is trivial: estimating in the MRA model is harder than estimating a signal in AWGN (namely, when the shifts are known). A small subtlety is that the distortion measure  $\mathbb{E}\rho(X, \hat{X})$  is a bit weaker than the standard definition of MSE,  $\mathbb{E}\|X - \hat{X}\|^2$ , as it allows for any cyclic shift. However, we show in Section 5 that, as expected, this has a vanishing effect for large  $L$ . In order to show that  $\lim_{\varepsilon \rightarrow 0} \lim_{L \rightarrow \infty} \frac{n_{MRA}^*(L, \alpha, \varepsilon)}{n_{AWGN}^*(L, \alpha, \varepsilon)} \leq 1$  we introduce an algorithm that for any  $\alpha > 2$  requires about  $\sigma^2(\alpha)/\varepsilon$  samples to achieve  $\mathbb{E}\rho(X, \hat{X}) \leq \varepsilon$ , provided that  $\varepsilon$  is sufficiently small and  $L$  is sufficiently large. The sole purpose of the estimation procedure is establishing an upper bound; its computational complexity is exponential in  $L$  and thus the procedure is far from being efficient. More specifically, it is based on a two-step procedure. First, we construct a  $\delta$ -net that, by definition, contains a member close to  $X$  and look for the most likely candidate within that net given the measurements. Second, we use this candidate in order to determine almost all shifts  $\hat{\ell}_i$ , and then estimate the signal by alignment and averaging  $\hat{X} = \frac{1}{n} \sum_{i=1}^n R_{-\hat{\ell}_i} Y_i$ . The details are given in Section 6.

In order to establish part 2 of Theorem 2.1, we show that for  $\alpha \leq 2$  the mutual information (MI)  $I(X; Y)$  between  $X$  and a single MRA measurement grows with  $L$  significantly slower than  $I(X; X + \sigma Z)$ , as in estimating a signal in AWGN. The details are given in Section 5.

Although our results are asymptotic in  $L$ , the phase transition at  $\alpha = 2$  predicted by Theorem 2.1 is evident already for relatively small  $L$ . Figure 1 presents the root MSE (RMSE) as a function of  $\alpha$  for different values of  $L$ . We take our estimator  $\hat{X}$  to be the output of the expectation-maximization (EM) algorithm [DLR77, BBM<sup>+</sup>17], which is the standard choice for MRA; see details in Section 3. For large values of  $L$  and large  $\alpha$ , the error of EM tends to that of estimating a signal in AWGN, implying that it detects the shifts accurately. For smaller values of  $\alpha$ , the error grows rapidly, especially when  $\alpha < 2$ .

**Connection with template matching.** At this point, the reader may wonder what is the intuitive interpretation of  $\alpha = 2$ . To answer this question we now introduce the *template matching problem*, which is studied in detail in Section 4. In this problem, we are given  $X$  and one MRA measurement  $Y = R_\ell X + Z$ , where  $X$ ,  $R_\ell$  and  $Z$  are distributed as above, and our goal is to recover the shift  $R_\ell$ . We will see that in the asymptotic setting,  $\alpha = 2$  is the critical threshold for this problem. That is, the error probability in recovering  $R_\ell$  from  $(X, Y)$  approaches 0 for all  $\alpha > 2$ , and approaches 1 for all  $\alpha < 2$ .

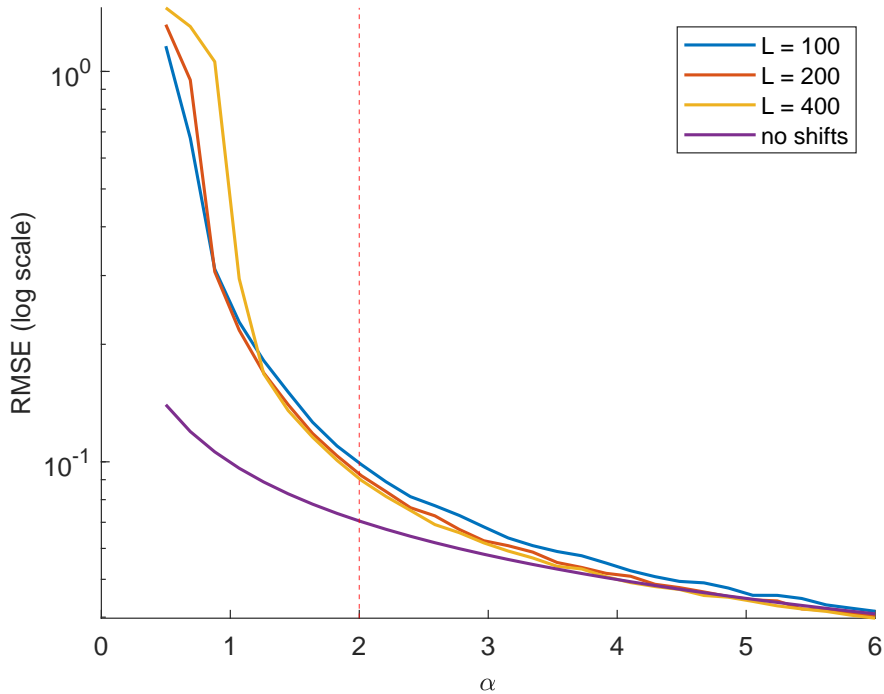


Figure 1: The RMSE of EM (averaged over 100 trials) as a function of  $\alpha$  for different values of  $L$ . The number of measurements was set to be  $n(L) = 100L/\log(L)$ . For large values of  $\alpha$ , the error reduces to the error of estimating a signal in AWGN,  $\sqrt{\frac{\sigma^2}{\sigma^2+n}} = \frac{1}{\sqrt{1+100\alpha}}$ , suggesting that EM performs as if the shifts were known. For small values of  $\alpha$ , and in particular  $\alpha < 2$ , the error rapidly increases.

In the MRA problem, recovering the shifts is harder, as we do not have access to  $X$ . We nevertheless show that for  $\alpha > 2$ , given enough measurements, it is possible to recover a fraction approaching 1 of the shifts correctly. On the other hand, recovering a large fraction of the shifts correctly for  $\alpha < 2$  is impossible since it is impossible even in the template matching model. Intuitively, if we cannot recover almost all shifts, the attained MSE should be much worse than in estimating a signal in AWGN, which means that the sample complexity should be much higher for  $\alpha < 2$ . Our bounds in Section 5 formalize this intuition.

To illustrate the phase transition for template matching, we conducted a “genie-aided” experiment, presented in Figure 2. In this experiment, we use the true  $X$  (the “genie”) in order to estimate the shifts by  $\hat{\ell}_i = \arg \max_{\ell \in \{0, \dots, L-1\}} \langle R_\ell X, Y_i \rangle$ . Then, we estimate the signal by aligning the measurements and averaging  $\hat{X} = \frac{1}{n} \sum_{i=1}^n R_{-\hat{\ell}_i} Y_i$ . For large values of  $\alpha$ , the recovery error converges to the error of estimating a signal in AWGN. For smaller  $\alpha$  values, and in particular  $\alpha < 2$ , the recovery error rapidly increases.

**Tighter lower bound for the low SNR regime.** Theorem 2.1 shows that for all  $\alpha \leq 2$  and fixed  $\varepsilon < 1$  the shifts make a difference: the sample complexity with unknown shifts (i.e., the MRA

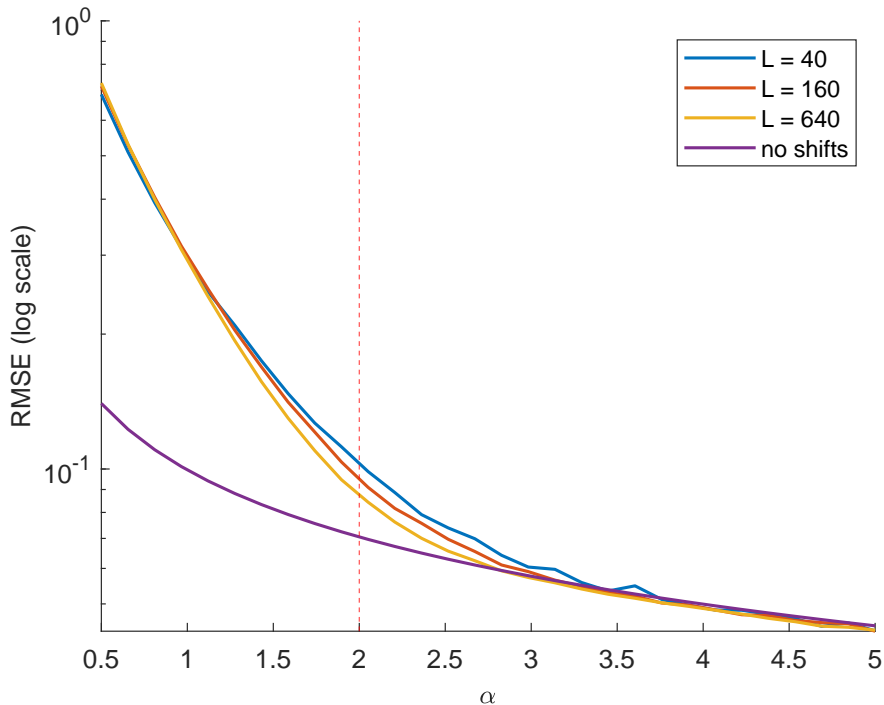


Figure 2: A “genie-aided” experiment: the true  $X$  is used to estimate the shifts  $\hat{\ell}_1, \dots, \hat{\ell}_n$ , as in the template matching problem, and then the signal is estimated by aligning all measurements and averaging  $\hat{X} = \frac{1}{n} \sum_{i=1}^n R_{-\hat{\ell}_i} Y_i$ . The figure presents the RMSE (averaged over 50 trials) as a function of  $\alpha$  for different values of  $L$ . The number of measurements was set to be  $n(L) = 100L/\log(L)$ . For large values of  $\alpha$ , the error reduces to the error of estimating a signal in AWGN (i.e., when the shifts are known)  $\sqrt{\frac{\sigma^2}{\sigma^2+n}} = \frac{1}{\sqrt{1+100\alpha}}$ . For small values of  $\alpha$ , and in particular  $\alpha < 2$ , the template matching error quickly increases.

problem) is  $\omega(\sigma^2(\alpha)/\varepsilon)$ , and is therefore substantially greater than the sample complexity when the shifts are known. For  $\alpha < 1$ , we were able to prove a much stronger lower bound on the sample complexity.

**Theorem 2.2** *For any  $0 < \varepsilon < 1$  there exists a constant  $c = c(\varepsilon)$  such that for all  $0 < \alpha < 1$*

$$n_{MRA}^*(L, \alpha, \varepsilon) > cL^{2-\alpha}. \quad (2.1)$$

Theorems 2.1 and 2.2 are proved in Section 5.

**The sample complexity of the projected MRA model.** Recall that MRA serves as a toy model of the cryo-EM reconstruction problem. An additional complication arising in cryo-EM is a fixed tomographic projection, a line integral, also known as the X-ray transform. To account for this effect, we extend our basic model (1.1) to the *projected multi-reference alignment problem*

(PMRA) model:<sup>2</sup>

$$Y_i = \pi_S R_{\ell_i} X + \sigma Z_i. \quad (2.2)$$

Here,  $\pi_S : \mathbb{R}^L \rightarrow \mathbb{R}^{L'}$  is matrix projecting a vector in  $\mathbb{R}^L$  to  $\mathbb{R}^{L'}$  by keeping only the coordinates that belong to a subset  $S \subset [L]$  of size  $L' \leq L$  and discarding the rest, and  $Z_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, I)$  are  $L'$ -dimensional i.i.d. Gaussian vectors. We assume that  $S$  is fixed and known to the estimator. As in MRA without the projection, the goal is to reconstruct  $X$  up to a circular shift, that is, produce an estimate  $\hat{X}$  such that  $\mathbb{E}\rho(X, \hat{X})$  is as small as possible.

We study the PMRA problem in an asymptotic setting where  $L, L', \sigma^2 \rightarrow \infty$  simultaneously. It makes sense to adopt a slightly different scaling for the noise in PMRA, as

$$\sigma^2 = \sigma_{\text{PMRA}}^2(\alpha) = \frac{L'}{\alpha \log(L)}. \quad (2.3)$$

The reason for this particular scaling will be made clear from the analysis: the numerator is the total signal energy available in a single measurement,  $\mathbb{E}\|\pi_S R_{\ell_i} X\|^2 = L'$ ; the  $\log(L)$  factor is  $\log$  the size of the group of shifts. In Section 7 we provide some remarks as to how to extend our results to other groups. Similarly to our notation for the MRA model, we denote the smallest attainable MSE in the PMRA model as  $\text{MSE}_{\text{PMRA}}^*(L, \alpha, n)$ , and the sample complexity as  $n_{\text{PMRA}}^*(L, \alpha, \varepsilon)$ .

**Theorem 2.3** *Suppose that  $\sigma_{\text{PMRA}}^2(\alpha)$  is scaled as in (2.3), and  $L, L' \rightarrow \infty$ , so that  $L' \leq L$  and  $L' = \omega(\log(L))$  (that is,  $L$  grows strictly less than exponentially fast in  $L'$ ). The sample complexity of the PMRA model (2.2) obeys the following lower bounds:*

1. For any  $\alpha > 2$  and  $0 < \varepsilon < 1$  we have that

$$n_{\text{PMRA}}^*(L, \alpha, \varepsilon) \geq \frac{L}{L'} \left( \frac{1}{\varepsilon} - 1 \right) \sigma_{\text{PMRA}}^2(\alpha) (1 + o(1)). \quad (2.4)$$

2. For any  $\alpha \leq 2$  and  $0 < \varepsilon < 1$  we have that

$$n_{\text{PMRA}}^*(L, \alpha, \varepsilon) = \omega \left( \frac{L}{L'} \frac{\sigma_{\text{PMRA}}^2(\alpha)}{\varepsilon} \right). \quad (2.5)$$

The proof of the theorem relies heavily on the proof of Theorem 2.1 and is sketched in Section 5.4. We conjecture that at high SNR ( $\alpha > 2$ ), the lower bound given in Theorem 2.3 is in fact tight at very low MSE (formally  $\varepsilon \rightarrow 0$ , as in Theorem 2.1).

### 3 Prior art

The multi-reference alignment problem was introduced by [BCSZ14], and fully formulated in [BCLS20]. The general MRA model reads

$$Y_i = T_i(g_i \circ X) + \sigma Z_i, \quad i = 1, \dots, n, \quad (3.1)$$

---

<sup>2</sup>We mention that other projected MRA models were studied in [BBSK<sup>+</sup>17, BJL<sup>+</sup>20].

where  $g_i$  is a random element of a compact group  $G$  (drawn from a possibly unknown distribution over  $G$ ) acting on a vector space  $X \in \mathbb{X}$ , and  $T_i, i = 1, \dots, n$ , are known linear operators. If  $T_i = I$  for all  $i$ ,  $g_i$  are drawn uniformly from the group of cyclic shifts  $\mathbb{Z}_L$ , and  $X \sim \mathcal{N}(0, I)$ , then (3.1) reduces to the MRA model (1.1). This model can be thought of as a special case of a Gaussian mixture model, where all centers are connected through a group action (i.e., a cyclic shift). If  $T_i = \pi_S$  for all  $i$ , we get the projected MRA model (2.2). In cryo-EM—the main motivation of this work— $G$  is the group of 3-D rotations  $SO(3)$ ,  $\mathbb{X}$  is the space of 3-D “band-limited” functions (that is, functions that can be expanded by finitely many basis functions), and  $T_i$  encodes the (fixed) tomographic projection, as well as other linear effects, such as the microscope’s point spread function (which varies across images) and sampling [Sin18, BBS20].

The sample complexity of the MRA model (1.1), in the minimax sense, was first studied in [BRW17, PWB<sup>+</sup>19]. The focus of these works, as well as the rest of the works mentioned in this section, is on the regime where the noise level  $\sigma$  and number of measurements  $n$  diverge, while the dimension of each measurement  $L$  is fixed, implying  $\text{SNR} \rightarrow 0$ . These results were extended to the general MRA model (3.1) by [BBSK<sup>+</sup>17] and [APS18] (the latter generalizes the framework proposed in [ABL<sup>+</sup>18]). These papers constitute an intimate connection between the MRA model and the method of moments—a classical estimation technique. Let  $\bar{d}$  be the lowest order moment that distinguishes two different signals (signals that are not in the same orbit) given a specific MRA model (namely, fixed  $T_i, \mathbb{X}$ , and a distribution over  $G$ ). Then, unless  $n \cdot \text{SNR}^{\bar{d}} \rightarrow \infty$ , the MSE is bounded from below. More informally, the moments determine the optimal (minimax) estimation rate of the problem. For example, for the MRA model (1.1) it is known that the third moment determines a generic signal uniquely (in this work we only consider normal i.i.d. signals that fall into this category), i.e.,  $\bar{d} = 3$ , and thus  $n \cdot \text{SNR}^3 \gg 1$  is a necessary condition. Remarkably, this phenomenon was observed empirically in context of cryo-EM early on by Signworth [Sig98].

In this work we discover that the statistical properties of MRA in high-dimensions, at least for  $X \sim \mathcal{N}(0, I)$ , are not characterized by moments, but rather by the parameter  $\alpha$  that balances the noise level and the dimension (1.6). In particular, in our setting  $\text{SNR} = \alpha \log L$  diverges, rather than  $\text{SNR} \rightarrow 0$  as in previous works. In this sense, our results imply that the “low SNR” regime is not only  $\text{SNR} \rightarrow 0$ , and actually extends into unbounded values of SNR provided that it grows slowly enough with  $L$ .

From the algorithmic perspective, two main computational frameworks were applied to MRA problems. The first approach is based on expectation-maximization (EM)—a popular heuristic to maximize the posterior distribution [DLR77]. EM is the most popular and successful methodology to elucidate high-resolution 3-D structures using cryo-EM [Sch12, BBS20], and it was successfully applied to a variety of MRA setups [BBM<sup>+</sup>17, BBLS18, ABL<sup>+</sup>18, MBB<sup>+</sup>19, B JL<sup>+</sup>20]. A recent work [FSWW20] studies the likelihood landscape for the general MRA model (3.1), where  $G$  is a discrete group and  $T_i = I$ . The latter paper shows that when the dimension is fixed and the SNR is sufficiently high, the log likelihood has certain favorable features from an optimization perspective; their results give a compelling argument for why EM seems to give good performance for MRA in high SNR. In [Bru19], it is shown that usually maximum likelihood achieves the parametric rate  $\rho(X, \hat{X}_{\text{MLE}}) \sim 1/n$ , although in some cases the rate can be  $\sim 1/\sqrt{n}$ .

The second algorithmic framework is based on the method of moments. This approach has an appealing property: it requires only one pass over the measurements, and thus its computa-



tional load is relatively low, unless  $L$  is large [BBM<sup>+</sup>17, BBLS18, ABL<sup>+</sup>18, MBB<sup>+</sup>19, PWB<sup>+</sup>19, PSB19]. In addition, as mentioned, it achieves the optimal estimation rate when  $L$  is fixed and  $\text{SNR} \rightarrow 0$ . Consequently, a variety of moment-based algorithms were proposed. For example, the authors of [PWB<sup>+</sup>19] suggest estimating the third-order tensor moment of the signal  $T^{(3)} = L^{-1} \sum_{\ell=0}^{L-1} (R_\ell X)^{\otimes 3}$ , from which  $X$  can be recovered by Jenrich’s method [Har70, LRA93]. Using the robustness analysis of [GVX14], they were able to show that  $n = O(\varepsilon^{-1} \sigma^6 \text{poly}(L))$  samples suffice to achieve  $\rho(X, \widehat{X}) \leq \varepsilon$  with constant probability. This bound depends polynomially on both the dimensional and on the inverse smallest DFT coefficient of  $X$ ; when  $X \sim \mathcal{N}(0, I)$ , one can verify that typically all the DFT coefficients of  $X$  are greater than  $\Omega(L^{-1/2})$ . The  $\text{poly}(L)$  dependence is not computed explicitly, but to the best of our understanding, the analysis of [GVX14] provides a significantly worse dimensional scaling than the  $\Omega(L^2)$  in our lower bound (as  $\alpha \rightarrow 0$ ). Another work [BBM<sup>+</sup>17] studies recovery by bispectrum inversion, which is equivalent to the third-order moment if the distribution of shifts is uniform. They argue that when  $L$  is fixed, the sample complexity should scale like  $O(\sigma^6)$ , hiding an implicit dependence on  $L$ . The method of moments was also applied to cryo-EM and related technologies, see for example [Kam80, DZS15, LBB<sup>+</sup>18, SKK<sup>+</sup>20], as well as to additional MRA setups [APS17, ALS19, HL19].

A recent work [KB20] establishes an enticing connection between likelihood-based techniques and the method of moments for the general MRA model (3.1) for fixed  $L$ ,  $\text{SNR} \rightarrow 0$ , and  $T_i = I$ . Specifically, it was shown that likelihood optimization in the low SNR regime reduces to a sequence of moment matching problems. In addition, the method of moments is also closely-related to invariant theory and thus tools from the latter field can be applied to analyze MRA models; see in particular [BBSK<sup>+</sup>17].

## 4 Phase transition of template matching

Suppose that the shifts  $R_{\ell_i}$  are all known. In this scenario, estimating the signal is easy: one needs to align each observation  $R_{\ell_i}^{-1} y_i$  and average out the noise. Therefore, if possible, it makes sense to try and estimate the shifts. In this section, we study the problem of estimating a shift when the signal is assumed to be known (which is not the case in MRA); we refer to this problem as *template matching*. Specifically, suppose that one has access to a signal, a “template”  $X \in \mathbb{R}^L$ , and observes a single sample  $Y = R_\ell X + \sigma Z$ , where  $X \sim \mathcal{N}(0, I)$ ,  $R_\ell \sim \text{Uniform}(\{0, \dots, L-1\})$  is a random uniform shift,  $Z \sim \mathcal{N}(0, I)$ , and  $R_\ell$ ,  $Z$  and  $X$  are mutually independent. The goal, then, is to recover  $R_\ell$  from  $X$  and  $Y$ .<sup>3</sup>

While the template matching problem seems to be significantly easier than the MRA problem, we show a surprising phenomenon: in high dimensions, template matching and MRA share the exact same phase transition point. In particular, it turns out that in high dimensions, under our parameterization  $\sigma^2(\alpha)$ , which amounts to  $L/\sigma^2 = \alpha \log(L)$ , the template matching problem displays a *sharp recoverability threshold*. That is: (i) whenever  $\alpha > 2$ , the random shift can be recovered with error probability  $p_e \rightarrow 0$  as  $L \rightarrow \infty$ ; (ii) whenever  $\alpha < 2$ , the shift cannot be consistently recovered, and in fact for any estimator,  $p_e \rightarrow 1$ .

Observe that the optimal estimator (in the sense of maximum a posteriori probability) for  $R_\ell$  is

---

<sup>3</sup>A more general setting, where  $X$  is not necessarily Gaussian, and  $R_\ell X$  goes through some general channel, not necessarily Gaussian, was studied by Wang, Hu, and Shayevitz [WHS17], but under different asymptotics.

given by:

$$\widehat{R}_{\text{MAP}} = \underset{\ell'}{\operatorname{argmin}} \|X - R_{\ell'}^{-1}Y\|^2 = \underset{\ell'}{\operatorname{argmax}} \frac{\langle X, R_{\ell'}^{-1}Y \rangle}{\|X\|^2}. \quad (4.1)$$

Denote its error probability by

$$p_e = \Pr\left(R_{\ell} \neq \widehat{R}_{\text{MAP}}\right). \quad (4.2)$$

We start by establishing that with overwhelming probability, the template  $X$  is “incoherent”, in the sense that the correlations  $\langle X, R_{\ell'}X \rangle / \|X\|^2$  are very small, unless  $\ell' = 0$ . The lemma is proved in Appendix B.

**Lemma 4.1** *For  $\kappa > 0$ , let  $\mathcal{A}(\kappa)$  be the event that*

$$|L^{-1}\|X\|^2 - 1| < \kappa \quad \text{and} \quad \max_{\ell' \neq 0} L^{-1} |\langle X, R_{\ell'}X \rangle| \leq \kappa,$$

and let  $\overline{\mathcal{A}(\kappa)}$  be its complement. Then,

$$\Pr(\overline{\mathcal{A}(\kappa)}) \leq 2L \exp(-cL \min(\kappa, \kappa^2)),$$

for a universal constant  $c > 0$ . In particular, one can choose a sequence  $\kappa = \kappa_L$  such that  $\kappa \rightarrow 0$  sufficiently slowly, for example,  $\kappa = CL^{-1/2} \log(L)$  for  $C > 0$  large enough, so that  $\Pr(\mathcal{A}_L(\kappa_L)) = 1 - o(1)$ .

Let

$$\Theta_{\ell'} = \frac{\langle X, R_{\ell'}^{-1}Y \rangle}{\|X\|^2} = \frac{\langle X, R_{\ell-\ell'}X \rangle}{\|X\|^2} + \frac{\sigma \langle X, R_{\ell'}^{-1}Z \rangle}{\|X\|^2}, \quad (4.3)$$

and

$$W_{\ell'} = \|X\|^{-1} \langle X, R_{\ell'}^{-1}Z \rangle. \quad (4.4)$$

Recalling that  $\widehat{R}_{\text{MAP}} = \underset{\ell'}{\operatorname{argmax}} \Theta_{\ell'}$ , and plugging  $\sigma^2 = (\alpha \log(L))^{-1}L$ , Lemma 4.1 implies that with high probability,

$$\Theta_{\ell'} = \begin{cases} 1 + (1 + o(1)) \frac{1}{\sqrt{\alpha \log(L)}} \cdot W_{\ell} & \text{if } \ell' = \ell, \\ o(1) + (1 + o(1)) \frac{1}{\sqrt{\alpha \log(L)}} \cdot W_{\ell'} & \text{if } \ell' \neq \ell. \end{cases} \quad (4.5)$$

Since for every  $\ell'$ ,  $W_{\ell'} \sim \mathcal{N}(0, 1)$ , it is obvious that  $\Theta_{\ell} \xrightarrow{p} 1$  as  $L \rightarrow \infty$ . Thus, to analyze the error of the MAP estimator, it simply remains to understand the behavior of  $\max_{\ell'} W_{\ell'}$ . To this end, we recall the following three results. We start with a well-known fact about the maximum of i.i.d. standard Gaussians:

**Lemma 4.2** *Let  $Z_1, \dots, Z_L$  be i.i.d  $\mathcal{N}(0, 1)$  random variables. Then, as  $L \rightarrow \infty$ ,*

$$\mathbb{E} \left[ \max_{\ell} Z_{\ell} \right] / \sqrt{2 \log(L)} \rightarrow 1.$$

The upper bound  $\mathbb{E}[\max_{\ell} Z_{\ell}] \leq \sqrt{2 \log(L)}$  is elementary, and holds even when  $Z_1, \dots, Z_L$  are not independent. The proof follows from  $\mathbb{E} \max_{\ell} Z_{\ell} \leq \beta^{-1} \log \mathbb{E} \max_{\ell} e^{\beta Z_{\ell}} \leq \beta^{-1} \log \mathbb{E} \sum_{\ell=1}^L e^{\beta Z_{\ell}} = \beta/2 + \beta^{-1} \log(L)$ , which holds for all  $\beta > 0$ ; now take  $\beta = \sqrt{2 \log(L)}$ . The proof of the matching lower bound, on the other hand, is more involved and follows from results in extreme value theory, see, for instance, Example 1.1.7 in [DHF07]. We also use the following “quantitative” version of the Sudakov-Fernique inequality:

**Lemma 4.3 (Theorem 2.2.5 in [AT09])** *Let  $(X_1, \dots, X_L)$  and  $(Y_1, \dots, Y_L)$  be Gaussian vectors so that  $\mathbb{E}[X_i] = \mathbb{E}[Y_i]$  for all  $i$ . Set*

$$\gamma_{i,j}^X = \mathbb{E}(X_i - X_j)^2, \quad \gamma_{i,j}^Y = \mathbb{E}(Y_i - Y_j)^2,$$

and  $\gamma = \max_{i,j} |\gamma_{i,j}^X - \gamma_{i,j}^Y|$ . Then

$$\left| \mathbb{E} \left[ \max_i X_i \right] - \mathbb{E} \left[ \max_i Y_i \right] \right| \leq \sqrt{2\gamma \log(L)}.$$

To get concentration around the mean, we use (a simple case of) the Borell-TIS inequality:

**Lemma 4.4** *Let  $(X_1, \dots, X_L)$  be a Gaussian vector, and set  $\sigma^2 = \max_i \mathbb{E}[X_i^2]$ . Then*

$$\Pr \left( \left| \max_i X_i - \mathbb{E} \left[ \max_i X_i \right] \right| \geq t \right) \leq 2e^{-t^2/2\sigma^2}.$$

See, e.g., [AT09, Theorem 2.1.1] (there only a one sided bound is stated; the other side follows the same way). The following is now an immediate corollary of Lemmas 4.1, 4.2, 4.3 and 4.4:

**Theorem 4.5 (Sharp threshold for template matching)** *If  $\alpha > 2$ , then  $p_e \rightarrow 0$  as  $L \rightarrow \infty$ . Conversely, if  $\alpha < 2$ , then  $p_e \rightarrow 1$ .*

**Proof.** We start by estimating  $\mathbb{E} \max_{\ell'} W_{\ell'}$ . Choose  $\kappa = o(1)$  such that the event  $\mathcal{A}(\kappa)$  of Lemma 4.1 holds with probability  $1 - o(1)$ . Conditioned on  $X$ ,  $\{W_{\ell'}\}_{\ell'=0, \dots, L-1}$  is a centered Gaussian vector, with covariance

$$C_{i,j}(X) = \mathbb{E}[W_i W_j | X] = \|X\|^{-2} \langle R_i X, R_j X \rangle,$$

whereby under  $\mathcal{A}$ ,  $|C_{i,j}(X) - \delta_{i,j}| = o(1)$ .

Let  $(\tilde{W}_1, \dots, \tilde{W}_{L-1})$  be i.i.d  $\mathcal{N}(0, 1)$  random variables. By Lemmas 4.2 and 4.3, conditioned on  $X$  and under  $\mathcal{A}$ ,

$$\mathbb{E}[\max_{\ell'} W_{\ell'} | X, \mathcal{A}] = \mathbb{E}[\max_{\ell'} \tilde{W}_{\ell'}] + o(\sqrt{\log(L)}) = \sqrt{(2+o(1)) \log(L)}.$$

Lemma 4.4 gives us a uniform (in  $X$ ) concentration inequality, conditioned on  $X$  and under  $\mathcal{A}$ ,

$$\Pr \left( \left| \max_{\ell'} W_{\ell'} - \sqrt{2 \log(L)} \right| \geq \sqrt{\varepsilon \log(L)} \mid X, \mathcal{A} \right) \leq 2L^{-(\varepsilon+o(1))/2},$$

so that

$$\Pr \left( \left| \max_{\ell'} W_{\ell'} - \sqrt{2 \log(L)} \right| \geq \sqrt{\varepsilon \log(L)} \right) \leq 2L^{-(\varepsilon+o(1))/2} + \Pr(\bar{\mathcal{A}}) = o_{\varepsilon}(1).$$

Thus, we have shown that  $\max_{\ell'} W_{\ell'} / \sqrt{2 \log(L)} \xrightarrow{p} 1$ . Using equation (4.5), we deduce that  $\Theta_{\ell} \xrightarrow{p} 1$  whereas  $\max_{\ell' \neq \ell} \Theta_{\ell'} \xrightarrow{p} \sqrt{2/\alpha}$ . Since  $\widehat{R}_{\text{MAP}} = \operatorname{argmax}_{\ell'} \Theta_{\ell'}$ , we conclude that  $p_e \rightarrow 0$  when  $\alpha > 2$  and  $p_e \rightarrow 1$  when  $\alpha < 2$ . ■

**A remark on the relation between template matching and synchronization.** In the MRA model, one does not have access to the true template and thus needs to estimate the relative shifts based solely on the data; this problem is referred to as *synchronization*.

For simplicity, let us assume we are given two measurements  $Y_1 = X + \sigma Z_1$  and  $Y_2 = R_{\ell} X + \sigma Z_2$ , and would like to estimate  $R_{\ell}$  (recall that  $X$  is unknown). The optimal (MAP) estimator is  $\widehat{R}_{\text{syn}} = \operatorname{argmax}_{\ell'} \Pr(R_{\ell'} | Y_1, Y_2)$ . It is straightforward to show that

$$\begin{aligned} \widehat{R}_{\text{syn}} &= \operatorname{argmax}_{\ell'} \langle Y_1, R_{\ell'}^{-1} Y_2 \rangle = \operatorname{argmax}_{\ell'} \langle (X + \sigma Z_1), R_{\ell'}^{-1} (R_{\ell} X + \sigma Z_2) \rangle \\ &= \operatorname{argmax}_{\ell'} \left\{ \langle X, R_{\ell-\ell'} X \rangle + \sigma \langle X, R_{\ell'}^{-1} Z_2 \rangle + \sigma \langle X, R_{\ell-\ell'}^{-1} Z_1 \rangle + \sigma^2 \langle Z_1, R_{\ell'}^{-1} Z_2 \rangle \right\}. \end{aligned}$$

In order for this to consistently return the true relative shift  $R_{\ell}$ , one needs to ensure that the “noise” term,

$$\sigma \langle X, R_{\ell'}^{-1} Z_2 \rangle + \sigma \langle X, R_{\ell-\ell'}^{-1} Z_1 \rangle + \sigma^2 \langle Z_1, R_{\ell'}^{-1} Z_2 \rangle$$

is small compared to  $\|X\|^2 \sim L$ . The “typical” size of the first two terms is  $\sigma \langle X, R_{\ell'}^{-1} Z_2 \rangle + \sigma \langle X, R_{\ell-\ell'}^{-1} Z_1 \rangle \sim \sigma \sqrt{L}$ , whereas the third is  $\sigma^2 \langle Z_1, R_{\ell'}^{-1} Z_2 \rangle \sim \sigma^2 \sqrt{L}$ , and is therefore the dominant one for large  $\sigma$ . Thus, to succeed with non-vanishing probability, we need that  $\sigma^2 \sqrt{L} \lesssim L$ , that is,  $\sigma^2 \lesssim \sqrt{L}$ . In the regime we are interested in, the noise level is  $\sigma^2 \sim L / \log(L)$ , and this turns out to be far too large.

We mention in passing that if many measurements are available, one can leverage the redundancy in the data to recover the true relative shifts in challenging environments; see for example [Sin11, SS11, Bou16, PWBM18, RG19].

## 5 Sample complexity lower bounds

### 5.1 The information-theoretic method for estimation lower bounds

We employ a standard information-theoretic method of obtaining estimation error lower bounds, via rate-distortion theory (see e.g. [PW19]). We refer the reader to Appendix A for a basic review of the information-theoretic definitions and facts we use in this section. Let  $\widehat{X}$  be an estimator of  $X$  from the measurements  $Y^n = (Y_1, \dots, Y_n)$ , which achieves expected error (“distortion”)

$$\mathbb{E} \rho(X, \widehat{X}) = L^{-1} \mathbb{E} \min_{\ell=0, \dots, L-1} \|X - R_{\ell}^{-1} \widehat{X}\|^2 \leq \varepsilon. \quad (5.1)$$

Since the estimator depends only on the measurements, and not on  $X$ , the triplet  $X - Y^n - \widehat{X}$  constitutes a Markov chain. Hence, by the data processing inequality (Proposition A.3.3) we have that  $I(X; \widehat{X}) \leq I(X; Y^n)$ . We lower-bound the left-hand side by the *rate distortion function* (RDF)  $R(\cdot)$  associated with the source  $X \sim \mathcal{N}(0, I)$ , and distortion measure  $\rho(\cdot, \cdot)$ :  $I(X; \widehat{X}) \geq R(\varepsilon)$  where

$$R(\varepsilon) = \min_{P_{W|X}: \mathbb{E} \rho(X, W) \leq \varepsilon} I(X; W).$$

Note that the minimization here is over conditional distributions  $P_{W|X}$ , or equivalently, over joint distributions  $P_{X,W}$  whose  $X$ -marginal is  $P_X$ —in our case  $\mathcal{N}(0, I)$ —obeying the average distortion constraint  $\mathbb{E}\rho(X, W) \leq \varepsilon$ . Combining the upper and lower bounds on  $I(X; \hat{X})$ , we have

$$R(\varepsilon) \leq I(X; Y^n), \quad (5.2)$$

and we shall next derive a lower bound for  $R(\varepsilon)$  in terms of  $\varepsilon$ .

## 5.2 A lower bound on the rate-distortion function

We start by obtaining a lower bound on the RDF. While the RDF problem for a Gaussian source under MSE distortion measure is classical, the MSE up to the best alignment (the distortion measure we consider) is somewhat non-standard. Obtaining a precise expression for the true RDF seems difficult, but a simple lower bound can be obtained as follows.

**Proposition 5.1** *For an  $L$  dimensional i.i.d. Gaussian vector  $X \sim \mathcal{N}(0, I)$ , and distortion measure  $\rho(\cdot, \cdot)$  as defined in (1.2), the rate distortion function satisfies*

$$R(\varepsilon) \geq \frac{L}{2} \log \left( \frac{1}{\varepsilon} \right) - \log(L).$$

**Proof.** By definition of the rate distortion function, to establish the claim we need to show that for any conditional distribution (“test-channel”)  $P_{W|X}$  that satisfies the constraint  $\mathbb{E}\rho(X, W) \leq \varepsilon$ , where  $\rho(X, W) = L^{-1} \min_{\ell=0, \dots, L-1} \|X - R_\ell^{-1}W\|^2$ , it holds that  $I(X; W) \geq \frac{L}{2} \log \left( \frac{1}{\varepsilon} \right) - \log(L)$ . To that end, let  $R = R(X, W) = \operatorname{argmin}_{\ell \in [0, \dots, L-1]} \|X - R_\ell W\|$  be the difference minimizing shift. By the chain law of MI (Proposition A.3.2),

$$I(X; W) = I(X; W, R) - I(X; R|W) \geq I(X; W, R) - \log(L), \quad (5.3)$$

where we used  $I(X; R|W) \leq H(R|W) \leq \log(L)$ ; the former follows from the definition of MI and non-negativity of entropy (Proposition A.1.1), and the latter follows from Proposition A.1.2 as the random variable  $R$  can take at most  $L$  values. Recall that  $L^{-1}\mathbb{E}\|X - RW\|^2 \leq \varepsilon$  by definition of  $R$ . We therefore have that

$$I(X; RW) \geq \min_{P_{W'|X}: L^{-1}\mathbb{E}\|X - W'\|^2 \leq \varepsilon} I(X; W') = \frac{L}{2} \log \left( \frac{1}{\varepsilon} \right),$$

where in the second equality we have used the well-known expression for the quadratic Gaussian rate distortion function (Proposition A.4). Thus, using the data processing inequality (Proposition A.3.3), we have

$$I(X; W, R) \geq I(X; RW) \geq \frac{L}{2} \log \left( \frac{1}{\varepsilon} \right).$$

Substituting this into (5.3) establishes the claim. ■

Combining Proposition 5.1 with equation (5.2), we obtain:

**Corollary 5.2** Suppose that  $X \sim \mathcal{N}(0, I)$  is an  $L$  dimensional i.i.d. Gaussian vector,  $\hat{X}$  is any estimator of  $X$  from  $Y_1, \dots, Y_n$ , and  $\rho(\cdot, \cdot)$  is as defined in (1.2). Then

$$\mathbb{E}\rho(X, \hat{X}) \geq \exp\left(-\frac{2I(X, Y^n) + 2\log(L)}{L}\right) = \exp(-2L^{-1} \cdot I(X, Y^n) + o(1)).$$

Equivalently,

$$\text{MSE}_{MRA}^*(L, \alpha, n) \geq \exp\left(-\frac{2I(X, Y^n) + 2\log(L)}{L}\right) = \exp(-2L^{-1} \cdot I(X, Y^n) + o(1)).$$

### 5.3 Upper bounds on the mutual information

In light of Corollary 5.2, an upper bound on the MI  $I(X; Y^n)$  provides a lower bound on the expected error of any estimator of  $X$  from  $Y^n = (Y_1, \dots, Y_n)$ .

We start with the rather trivial observation that the MI between the signal  $X$  and the measurements  $Y^n$  is smaller than the MI in a problem where there are no random shifts, which is equal to  $\frac{L}{2} \log(1 + n\sigma^{-2})$ . The next lemma formalizes this intuition and quantifies the MI difference between the two problems.

**Lemma 5.3** The mutual information between the signal  $X$  and measurements  $Y_1, \dots, Y_n$  is

$$I(X; Y^n) = \frac{L}{2} \log(1 + n\sigma^{-2}) - I(R^n; X|Y^n), \quad (5.4)$$

where  $R^n = (R_{\ell_1}, \dots, R_{\ell_n})$ . In particular,  $I(X; Y^n) \leq \frac{L}{2} \log(1 + n\sigma^{-2})$ .

**Proof.** Let  $\tilde{Y}_i = R_{\ell_i}^{-1} Y_i = X + \sigma R_{\ell_i}^{-1} Z_i$ . We may write

$$\begin{aligned} I(X; Y^n) &= I(X; Y^n, R^n) - I(X; R^n|Y^n) \\ &= I(X; \tilde{Y}^n, R^n) - I(X; R^n|Y^n) \\ &= I(X; \tilde{Y}^n) + I(X; R^n|\tilde{Y}^n) - I(X; R^n|Y^n), \end{aligned}$$

where the first and third equalities follow by the chain rule for MI (Proposition A.3.2), and the second follows from Proposition A.3.4, and the fact that the mapping  $(Y^n, R^n) \mapsto (\tilde{Y}^n, R^n)$  is invertible. By the fact that the Gaussian distribution is rotation invariant, and in particular  $R_{\ell_i}^{-1} Z \sim \mathcal{N}(0, I)$ , we have that  $R^n$  is statistically independent of  $(X, \tilde{Y}^n)$ , and consequently

$$I(X; R|\tilde{Y}^n) = H(R|\tilde{Y}^n) - H(R|\tilde{Y}^n, X) = H(R) - H(R) = 0,$$

where the first equality follows by definition of conditional mutual information and the second by Proposition A.3.5. It remains to compute  $I(X; \tilde{Y}^n)$ . To this end, note that  $P_{\tilde{Y}^n|X=x} = \mathcal{N}^{\otimes n}(x, \sigma^2 I)$ , that is,  $X$ , and  $\tilde{Y}^n$  have the same joint distributed as  $X$  and  $(X + \sigma Z_1, \dots, X + \sigma Z_n)$ , i.e., as  $n$  measurements of a signal in AWGN. It is well known that the sample average  $\frac{1}{n} \sum_{i=1}^n X + \sigma Z_i$  is a sufficient statistic of  $(X + \sigma Z_1, \dots, X + \sigma Z_n)$  for  $Y$ . We therefore have that

$$\begin{aligned} I(X; \tilde{Y}^n) &= I(X; X + \sigma Z_1, \dots, X + \sigma Z_n) = I\left(X; X + \frac{\sigma}{n} \sum_{i=1}^n Z_i\right) \\ &= I(X; X + \mathcal{N}(0, (\sigma^2/n)I)) = \frac{L}{2} \log(1 + n\sigma^{-2}), \end{aligned} \quad (5.5)$$

where the last equality follows from Proposition A.3, 6. ■

Combining Corollary 5.2 and Lemma 5.3, we obtain the following lower bound, that essentially says the MSE in the MRA model is no better than in estimating a signal in AWGN.

**Corollary 5.4** *The smallest attainable MSE in the MRA model satisfies*

$$\text{MSE}_{MRA}^*(L, \sigma^2, n) \geq \frac{L^{-\frac{2}{L}}}{1 + n\sigma^{-2}} = \frac{1}{1 + n\sigma^{-2}}(1 + o(1)),$$

and the sample complexity satisfies

$$n_{MRA}^*(L, \sigma^2, \varepsilon) \geq \left\lceil \left( \frac{L^{-\frac{2}{L}}}{\varepsilon} - 1 \right) \sigma^2 \right\rceil = n_{AWGN}^*(L, \sigma^2, \varepsilon)(1 + o(1)).$$

Lemma 5.3 tells us that the gap between  $I(X; Y^n)$  and the MI in estimating a signal in AWGN, without the shifts,  $\frac{L}{2} \log(1 + n\sigma^{-2})$ , is  $I(X; R^n | Y^n)$ . This quantity is intimately related to a multi-sample version of the template matching problem, as was considered in Section 4. This connection will be exploited later on, when we derive an upper bound on the single sample MI  $I(X; Y_i)$ .

**Information combining.** Observe that the measurements  $Y_1, \dots, Y_n$  are mutually independent conditioned on  $X$ ; that is, the samples are obtained by passing the same signal  $X$  independently through a memoryless channel  $P_{Y^n|X} = P_{Y|X}^{\otimes n}$ . By Proposition A.3, 5, this implies that

$$I(X; Y^n) \leq \sum_{i=1}^n I(X; Y_i) = nI(X; Y), \quad (5.6)$$

where  $Y = R_\ell X + \sigma Z$  is a single measurement in the MRA model. Substituting (5.6) into Corollary 5.2, yields the following.

**Proposition 5.5** *The smallest attainable MSE in the MRA model satisfies*

$$\text{MSE}_{MRA}^*(L, \sigma^2, n) \geq L^{-\frac{2}{L}} \exp\left(-n \frac{2}{L} I(X; Y)\right) = \exp\left(-n \frac{2}{L} I(X; Y)\right) (1 + o(1)),$$

and the sample complexity satisfies

$$n_{MRA}^*(L, \sigma^2, \varepsilon) \geq \frac{L}{2} \cdot \frac{\log\left(\frac{1}{\varepsilon}\right) - \frac{2 \log(L)}{L}}{I(X; Y)} = \log\left(\frac{1}{\varepsilon}\right) \cdot \frac{L}{2I(X; Y)} (1 + o(1)),$$

where  $Y = R_\ell X + \sigma Z$  is a single measurement in the MRA model.

It is important to emphasize at this point that the bound in (5.6) becomes very loose for  $n$  sufficiently large. Indeed, Lemma 5.3 implies that  $I(X; Y^n)$  should scale at best logarithmically, rather than linearly, with  $n$ . Consequently, the lower bound on  $\text{MSE}_{MRA}^*(L, \sigma^2, n)$  in Proposition 5.5 decreases exponentially fast with  $n$ , whereas we know from Corollary 5.4 that it cannot decrease faster than the parametric rate of  $1/n$  as in estimating a signal in AWGN. Despite its grossly wrong dependence on  $n$ , the upper bound  $I(X; Y^n) \leq nI(X; Y)$  does suffice to say something non-trivial

about the sample complexity of the problem. As seen from Proposition 5.5: in order for the estimation error to be strictly bounded away from one, one needs at least  $\Omega(L \cdot I(X; Y)^{-1})$  samples. We will see that this rather “naïve” analysis is already enough to accurately separate between a “high SNR” and a “low SNR” regime, where the behavior of the MRA problem is qualitatively different. Intuitively, as the measurements  $Y_1, \dots, Y_n$  are only dependent through the random variable  $X$ , if  $n$  is so small that it is impossible to learn much about  $X$  from  $Y^n$ , the dependence between  $Y_1, \dots, Y_n$  must be weak. Thus, in that regime, ignoring this dependence and bounding  $I(X; Y^n) \leq nI(X; Y)$  is a rather accurate estimate.

The problem of obtaining a stronger bound on multi-sample MI  $I(X; Y^n)$  in terms of the single-sample MI  $I(X; Y)$  is an instance of a so-called *information combining* problem. Several problems of this type have been studied in the information theory literature, mostly dealing with binary channels [SSZ05, LHHH05]. In our case, we believe this problem to be quite hard, at least in the low SNR regime, and thus we could not obtain a tighter bound. Deriving such bounds can yield stronger lower bounds on  $\text{MSE}_{\text{MRA}}^*(L, \alpha, n)$  in the low-SNR regime ( $\alpha < 2$ ) than the ones we obtain here using the simple bound  $I(X; Y^n) \leq nI(X; Y)$ .

**Roadmap.** We will devote the rest of this section to deriving upper bounds on  $I(X; Y)$ . These bounds, together with Proposition 5.5, will immediately imply lower bounds on the MSE and the sample complexity. In particular, we will derive two bounds, using different methods, that will be effective at two SNR regimes.

- We estimate the mutual information using Jensen’s inequality to facilitate the computation of several expectations. One could expect this method to give somewhat tight results when  $I(X; Y)$  is very small, and indeed, we shall see that when  $0 < \alpha < 1$ , we obtain a bound  $I(X; Y) = O(L^{\alpha-1})$ , which tends to 0 as  $L \rightarrow \infty$ . For  $\alpha \geq 1$ , the obtained bound will turn out to be too loose.
- In Lemma 5.3 we have found that  $I(X; X + \sigma Z) - I(X; Y) = I(X, R_\ell | Y)$ . We lower bound this gap using a Fano-like inequality, which in the case  $\alpha < 2$  amounts to “quantifying” how well  $R_\ell$  can be estimated from  $X$  and  $Y$ , in a somewhat more precise sense than Theorem 4.5 (which tells us that in this case, the error is  $p_e = 1 - o(1)$ ). This will allow us to show that when  $\alpha < 2$ ,  $I(X; Y) = o(\log(L))$ . We will not, however, be able to recover the estimate in the case of  $0 < \alpha < 1$  using this method.

### 5.3.1 MI bound at very low SNR ( $\alpha < 1$ )

We first express  $I(X; Y)$  in the following way:

**Lemma 5.6** *Suppose that  $X \sim \mathcal{N}(0, I)$ ,  $Z \sim \mathcal{N}(0, I)$ , and  $R \sim \text{Uniform}(\{R_0, \dots, R_{L-1}\})$  are mutually independent. Then,*

$$I(X; Y) = \frac{L}{2} \log(1 + \sigma^{-2}) - L\sigma^{-2} + \mathbb{E}_{X, Z} \left[ \log \mathbb{E}_R \exp \left( \frac{1}{\sigma^2} \langle X + \sigma Z, RX \rangle \right) \right].$$

**Proof.** Write  $I(X; Y) = h(Y) - h(Y|X)$ . Note that for any shift  $R_\ell$ ,  $R_\ell X \sim \mathcal{N}(0, I)$  and therefore  $Y \sim \mathcal{N}(0, (1 + \sigma^2)I)$ ; this means that  $Y = R_\ell X + \sigma Z$  is independent of  $R_\ell$ . The differential entropy of  $Y$  is  $h(Y) = h(\mathcal{N}(0, (1 + \sigma^2)I)) = \frac{L}{2} \log(2\pi e) + \frac{L}{2} \log(1 + \sigma^2)$ , by Proposition A.1.3.



Let us now write the conditional differential entropy explicitly. The conditional density of  $Y$  given  $X$  is  $p_{Y|X}(y|x) = \mathbb{E}_R \left[ (2\pi\sigma^2)^{-L/2} \exp\left(-\frac{1}{2\sigma^2}\|y - Rx\|^2\right) \right]$  for uniform  $R$ . The conditional entropy is then simply

$$\begin{aligned} h(Y|X) &= \mathbb{E}_{X,Y} \left[ -\log p_{Y|X}(Y|X) \right] \\ &= \frac{L}{2} \log(2\pi\sigma^2) - \mathbb{E}_{X,Y} \left[ \log \mathbb{E}_R \exp\left(-\frac{1}{2\sigma^2}\|Y - RX\|^2\right) \right] \\ &= \frac{L}{2} \log(2\pi\sigma^2) - \mathbb{E}_{X,Y} \left[ \log \mathbb{E}_R \exp\left(-\frac{1}{2\sigma^2} (\|Y\|^2 + \|X\|^2 - 2\langle Y, RX \rangle)\right) \right] \\ &= \frac{L}{2} \log(2\pi\sigma^2) + \frac{L + (1 + \sigma^2)L}{2\sigma^2} - \mathbb{E}_{X,Y} \left[ \log \mathbb{E}_R \exp\left(\frac{1}{\sigma^2}\langle Y, RX \rangle\right) \right]. \end{aligned}$$

We can write  $Y = R'(X + \sigma Z)$ , where  $R'$  is another uniform shift (independent of  $X, Y, R$ ); here we used the orthogonal invariance of  $Z \sim \mathcal{N}(0, I)$ . Since  $R$  is uniformly distributed,

$$\begin{aligned} \mathbb{E}_{X,Z,R} \left[ \log \mathbb{E}_R \exp\left(\frac{1}{\sigma^2}\langle R'(X + \sigma Z), RX \rangle\right) \right] &= \mathbb{E}_{X,Z,R} \left[ \log \mathbb{E}_R \exp\left(\frac{1}{\sigma^2}\langle (X + \sigma Z), (R')^{-1}RX \rangle\right) \right] \\ &= \mathbb{E}_{X,Z} \left[ \log \mathbb{E}_R \exp\left(\frac{1}{\sigma^2}\langle (X + \sigma Z), RX \rangle\right) \right], \end{aligned}$$

that is, we can “drop”  $R'$ . The claimed formula now readily follows. ■

The following proposition is the main estimate of this section. The proof uses some properties of the spectrum of  $R_\ell$ , stated and proved in Appendix C.

**Proposition 5.7** *We have the following upper bound on the single sample MI:*

$$I(X; Y) \leq \log\left(1 + L^{-1}e^{\sigma^{-2}L}\right) + O(\sigma^{-4}L).$$

*In particular, if  $\sigma^{-2}L = \alpha \log(L)$  for  $0 < \alpha < 1$ , then the MI asymptotically vanishes as  $L \rightarrow \infty$  with  $I(X; Y) \leq L^{-1+\alpha}(1 + o(1))$ .*

**Proof.** By the concavity of the log function, we always have  $\mathbb{E}_W \log(W) \leq \log(\mathbb{E}W)$ . Thus,

$$\begin{aligned} \mathbb{E}_{X,Z} \left[ \log \mathbb{E}_R \exp\left(\frac{1}{\sigma^2}\langle X + \sigma Z, RX \rangle\right) \right] &\leq \mathbb{E}_X \left[ \log \mathbb{E}_{Z,R} \exp\left(\frac{1}{\sigma^2}\langle X + \sigma Z, RX \rangle\right) \right] \\ &= \mathbb{E}_X \left[ \log \mathbb{E}_R \exp\left(\frac{1}{\sigma^2}\langle X, RX \rangle + \frac{1}{2\sigma^2}\|RX\|^2\right) \right] \\ &= \mathbb{E}_X \left[ \log \mathbb{E}_R \exp\left(\frac{1}{\sigma^2}\langle X, RX \rangle + \frac{1}{2\sigma^2}\|X\|^2\right) \right] \\ &= \frac{1}{2}\sigma^{-2}L + \mathbb{E}_X \left[ \log \mathbb{E}_R \exp\left(\frac{1}{\sigma^2}\langle X, RX \rangle\right) \right] \\ &\leq \frac{1}{2}\sigma^{-2}L + \log \mathbb{E}_{R,X} \exp\left(\frac{1}{\sigma^2}\langle X, RX \rangle\right). \end{aligned}$$

Plugging into the expression in Lemma 5.6, we get

$$I(X; Y) \leq \frac{L}{2} \log(1 + \sigma^{-2}) - \frac{1}{2}L\sigma^{-2} + \log \mathbb{E}_{R,X} \exp\left(\frac{1}{\sigma^2}\langle X, RX \rangle\right).$$

Note that as  $L, \sigma^2 \rightarrow \infty$ , already  $\frac{L}{2} \log(1 + \sigma^{-2}) - \frac{1}{2}L\sigma^{-2} = O(\sigma^{-4}L)$ . Observe that  $\langle X, RX \rangle = \langle X, R^\top X \rangle = \frac{1}{2}\langle X, (R + R^\top)X \rangle$ . By Lemma C.1, all the matrices  $R_\ell + R_\ell^\top$  are diagonalized by some orthonormal basis with eigenvalues  $\{2 \cos(\frac{2\pi}{L}k\ell)\}_{k=0}^{L-1}$ . By the orthogonal invariance of  $X \sim \mathcal{N}(0, I)$ , there are i.i.d.  $W_{k,\ell} \sim \mathcal{N}(0, 1)$  such that for all  $\ell$ ,

$$\sigma^{-2}\langle X, R_\ell X \rangle = \sigma^{-2} \sum_{k=0}^{L-1} \cos\left(\frac{2\pi}{L}k\ell\right) W_{k,\ell}^2.$$

Recall that the moment generating function of a  $\chi^2$  random variable is

$$\mathbb{E}_{W \sim \mathcal{N}(0,1)}[e^{tW^2}] = (1 - 2t)^{-1/2} \quad \text{for } t > 1/2.$$

Therefore, assuming  $\sigma^2$  is sufficiently large (e.g.,  $\sigma^2 > 2$ ),

$$\begin{aligned} \log \mathbb{E}_{R,X} \exp\left(\frac{1}{\sigma^2}\langle X, RX \rangle\right) &= \log \left[ L^{-1} \sum_{\ell=0}^{L-1} \prod_{k=0}^{L-1} \left(1 - 2\sigma^{-2} \cos\left(\frac{2\pi}{L}k\ell\right)\right)^{-1/2} \right] \\ &= \log \sum_{\ell=0}^{L-1} e^{\psi_\ell} - \log(L), \end{aligned}$$

where

$$\psi_\ell = -\frac{1}{2} \sum_{k=0}^{L-1} \log\left(1 - 2\sigma^{-2} \cos\left(\frac{2\pi}{L}k\ell\right)\right).$$

Expanding the log function to first order around 1 and using Lemma C.1, for large values of  $L$  and  $\sigma^2$ , we get

$$\psi_\ell = \sum_{k=0}^{L-1} \sigma^{-2} \cos\left(\frac{2\pi}{L}k\ell\right) + O(\sigma^{-4}L) = \begin{cases} \sigma^{-2}L + O(\sigma^{-4}L) & \text{if } \ell = 0, \\ O(\sigma^{-4}L) & \text{otherwise.} \end{cases}$$

Thus, we have the estimate

$$\begin{aligned} \log \sum_{\ell=0}^{L-1} e^{\psi_\ell} - \log(L) &= \log\left(\frac{1}{L}e^{\sigma^{-2}L + O(\sigma^{-4}L)} + \frac{L-1}{L}e^{O(\sigma^{-4}L)}\right) \\ &= \log\left(1 + L^{-1}e^{\sigma^{-2}L}\right) + O(\sigma^{-4}L), \end{aligned}$$

from which the claimed result immediately follows.  $\blacksquare$

Observe that for  $\alpha > 1$ , Proposition 5.7 gives an upper bound of the order  $I(X; Y) = O(\log(L))$ . It will turn out that when  $\alpha > 2$ , this is indeed the right order of magnitude. However, for  $1 < \alpha \leq 2$  the bound is too loose, and in fact  $I(X; Y) = o(\log(L))$ .

### 5.3.2 MI bound using template matching

We start from Lemma 5.3 which gives, for  $n = 1$  and  $Y = RX + \sigma Z$ ,  $I(X; Y) = \frac{L}{2} \log(1 + \sigma^{-2}) - I(R; X|Y)$ . We make the important observation that  $R$  and  $Y$  are independent; indeed, regardless

of  $R$ , it holds that  $Y|R \sim \mathcal{N}(0, (1 + \sigma^2)I)$ . We remark, however, that when  $n > 1$ ,  $Y^n$  is not independent of  $R^n$ . We can therefore use Proposition A.1.5, and Proposition A.1.2 to write

$$I(R; X|Y) = H(R|Y) - H(R|X, Y) = H(R) - H(R|X, Y) = \log(L) - H(R|X, Y),$$

so that

$$I(X; Y) = \frac{L}{2} \log(1 + \sigma^{-2}) - \log(L) + H(R|X, Y). \quad (5.7)$$

The following is now an immediate consequence of Fano's inequality (Proposition A.2) and Theorem 4.5.

**Proposition 5.8** *Suppose that  $\sigma^{-2}L = \alpha \log(L)$  with  $\alpha > 2$ . Then,*

$$\begin{aligned} I(X; Y) &= \frac{L}{2} \log(1 + \sigma^{-2}) - (1 + o(1)) \log(L) \\ &= \left( \frac{\alpha}{2} - 1 + o(1) \right) \log(L) + O(\sigma^{-4}L). \end{aligned}$$

**Proof.** We estimate  $H(R|X, Y)$ . Clearly,  $H(R|X, Y) \geq 0$  by non-negativity of entropy (Proposition A.1.1). As for an upper bound, by Fano's inequality (Proposition A.2), for any estimator  $\widehat{R}$  of  $R$  from  $X, Y$ , the error probability  $p_e = \Pr(R \neq \widehat{R})$  satisfies

$$H(R|X, Y) \leq \log 2 + p_e \log(L).$$

By Theorem 4.5,  $\widehat{R}_{\text{MAP}}$  has error  $p_e \rightarrow 0$ , which means that  $H(R|X, Y) = o(1) \cdot \log(L) = o(\log(L))$ . Plugging this into equation (5.7) and expanding  $\frac{L}{2} \log(1 + \sigma^{-2}) = \frac{\alpha}{2} \log(L) + O(\sigma^{-4}L)$ , we obtain the desired estimate for  $I(X; Y)$ . ■

When  $\alpha < 2$  we have  $p_e \rightarrow 1$ , so that it is no longer true that  $H(R|X, Y) = o(\log(L))$ . Indeed, since  $I(X; Y) = (\alpha/2 - 1) \log(L) + O(\sigma^{-4}L) + H(R|X, Y)$ , we must have that  $H(R|X, Y) \geq (1 - \alpha/2 - o(1)) \log(L)$ , since the MI is non-negative. While, indeed, in this regime  $R$  cannot be recovered from  $X, Y$ , we can still obtain a non-trivial upper bound (of the form  $c(\alpha) \log(L)$  for some  $c(\alpha) < 1$ ) on the conditional entropy  $H(R|X, Y)$ ; the idea is that given  $X, Y$ , we can form a relatively small list that contains  $R$  with high probability.

Our goal, then, is to non-trivially upper bound  $H(R|X, Y)$  in the regime  $\alpha \leq 2$  where  $p_e \not\rightarrow 0$ . Let  $\tau > 0$ , and denote by  $\mathcal{S}_\tau$  the set of  $\tau$ -likely shifts:

$$\mathcal{S}_\tau = \left\{ R' : \frac{\langle X, (R')^{-1}Y \rangle}{\|X\|^2} \geq 1 - \tau \right\}. \quad (5.8)$$

The analysis of Section 4 tells us that for any  $\tau > 0$ , the true shift  $R$  belongs with high probability to the set  $\mathcal{S}_\tau$ . Moreover, when  $\alpha > 2$  (and  $\tau > 0$  is a sufficiently small constant), in fact with high probability  $\mathcal{S}_\tau = \{R\}$ . When  $\alpha \leq 2$  this will no longer be the case; nonetheless, we show that  $|\mathcal{S}_\tau|$  is with high probability significantly smaller than  $L$ . This means that given  $X$  and  $Y$ , we can produce a list of likely candidates for  $R$  which is much smaller than the entire group of shifts. The following lemma is proved in Section D.

**Lemma 5.9** *Let  $\kappa, \tau, \delta > 0$ . Set  $M = L^{1-\frac{1}{2}\alpha(1-\kappa)(1-\tau-\frac{\kappa}{1-\kappa})^2+\delta}$ , and assume that  $\alpha \leq 2$ . Then*

$$\Pr(R \notin \mathcal{S}_\tau \text{ or } |\mathcal{S}_\tau| > M) \leq 2Le^{-cL \min(\kappa, \kappa^2)} + L^{-\frac{1}{2}\alpha(1-\kappa)(1-\tau-\frac{\kappa}{1-\kappa})^2} + 2L^{-\delta}, \quad (5.9)$$

where  $c > 0$  is the universal constant of Lemma 4.1.

Lemma 5.9 implies that there are slowly decaying sequences  $\tau = \tau_L = o(1)$ ,  $\delta = \delta_L = o(1)$  such that the event

$$\mathcal{B} = \left\{ R \in \mathcal{S}_{\tau_L} \text{ and } |\mathcal{S}_{\tau_L}| \leq L^{1-\frac{1}{2}\alpha+\delta_L} \right\}$$

holds with high probability of  $\Pr(\mathcal{B}) = 1 - o(1)$ . We use this to bound the conditional entropy  $H(R|X, Y)$ , and obtain a bound on the MI:

**Proposition 5.10** *Suppose that  $\alpha \leq 2$ . Then,*

$$I(X; Y) = o(\log(L)).$$

**Proof.** We upper bound the conditional entropy  $H(R|X, Y)$  using a ‘‘Fano-like’’ argument. Let  $E$  be the indicator for the event  $\mathcal{B}$  above. Since  $E$  is completely deterministic given  $(R, X, Y)$ , we have that  $H(E|R, X, Y) = 0$  by Proposition A.1.1 and by the chain rule of entropy (Proposition A.1.4) we have

$$\begin{aligned} H(R|X, Y) &= H(R|X, Y) + H(E|R, X, Y) \\ &= H(R, E|X, Y) \\ &= H(E|X, Y) + H(R|X, Y, E) \\ &\leq H(E) + H(R|X, Y, E = 1) \Pr(E = 1) + H(R|X, Y, E = 0) \Pr(E = 0), \end{aligned}$$

where we have bounded  $H(E|X, Y) \leq H(E)$  using Proposition A.1.5, and expanded  $H(R|X, Y, E)$  according to the definition of conditional entropy, averaging only with respect to  $E$ .

Now, given that  $E = 1$ , we know that  $R$  belongs to  $\mathcal{S}_{\tau_L}$ , which has size  $|\mathcal{S}_{\tau_L}| \leq M = L^{1-\frac{1}{2}\alpha+\delta_L}$ . Hence,  $H(R|X, Y, E = 1) \leq \log(M) = (1 - \frac{1}{2}\alpha + \delta_L) \log(L)$  by Proposition A.1.2, and by the same reason  $H(R|X, Y, E = 0) \leq \log(L)$ . By definition,  $\Pr(E = 1) = \Pr(\mathcal{B}) = 1 - o(1)$ , and  $H(E) \leq \log(2)$  by Proposition A.1.2. Thus,  $H(R|X, Y) \leq (1 - \frac{1}{2}\alpha + o(1)) \log(L)$ . Plugging this into Eq. (5.7),

$$\begin{aligned} I(X; Y) &= \frac{L}{2} \log(1 + \sigma^{-2}) - \log(L) + H(R|X, Y) \\ &= \left( \frac{\alpha}{2} - 1 + o(1) \right) \log(L) + O(\sigma^{-4}L) + \left( 1 - \frac{\alpha}{2} + o(1) \right) \log(L) \\ &= o(\log(L)) + O(\sigma^{-4}L), \end{aligned}$$

as claimed. ■

**Remark 5.11** *One might wonder if the argument above (if carried out delicately enough) can match the estimate  $I(X; Y) = O(L^{-1+\alpha})$  we have already seen for  $\alpha < 1$ . Unfortunately, the bound  $\Pr(|\mathcal{S}_\tau| \geq M) \leq 2L^{-\delta}$  (using Markov’s inequality; see the proof of Lemma 5.9 in Section D) is already too crude for that purpose: since we need to choose  $\delta = o(1)$ , the  $o(1)$  correction above must decay slower than  $L^{-c}$  (for any  $c > 0$ ).*

### 5.3.3 Proof of main results

We are ready to prove Theorem 2.2 and the sample complexity lower bounds of Theorem 2.1.

#### Proof of Theorems 2.1 (lower bounds) and 2.2.

- Theorem 2.1,  $\alpha > 2$  (lower bound): Corollary 5.4 immediately implies that  $\lim_{\varepsilon \rightarrow 0} \lim_{L \rightarrow \infty} \frac{n_{\text{MRA}}^*(L, \alpha, \varepsilon)}{\sigma^2/\varepsilon} \geq 1$ .
- Theorem 2.1,  $\alpha \leq 2$ : Combining Proposition 5.5 and Proposition 5.10, give  $n_{\text{MRA}}^*(L, \alpha, \varepsilon) = \omega\left(\frac{L}{\log(L)} \log(1/\varepsilon)\right)$ , which is  $\omega(\sigma^2/\varepsilon)$  for fixed  $\varepsilon$ .
- Theorem 2.2,  $\alpha < 1$ : Combining Proposition 5.5 and Proposition 5.7 yield  $n_{\text{MRA}}^*(L, \alpha, \varepsilon) = \Omega(L^{2-\alpha} \log(1/\varepsilon))$ .

The proof of the upper bound  $\lim_{\varepsilon \rightarrow 0} \lim_{L \rightarrow \infty} \frac{n_{\text{MRA}}^*(L, \alpha, \varepsilon)}{\sigma^2/\varepsilon} \leq 1$  for  $\alpha > 2$  (item (1) of Theorem 2.1) appears in Section 6.

### 5.4 Projected MRA

In this section, we sketch a proof of Theorem 2.3. Recall that in the PMRA model, the measurements  $Y_1, \dots, Y_n \in \mathbb{R}^{L'}$  have the form

$$Y_i = \pi_S R_{\ell_i} X + \sigma Z_i,$$

where  $X \sim \mathcal{N}(0, I)$  is  $L$ -dimensional,  $Z_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, I)$  are  $L'$ -dimensional, and  $\pi_S : \mathbb{R}^L \rightarrow \mathbb{R}^{L'}$  is the projection onto the coordinates in  $S \subset [L]$ , with  $|S| = L'$ . Here, the set  $S$  is fixed across all samples, and is a priori known.

As before, we are interested in asymptotics as  $L, L', \sigma^2 \rightarrow \infty$  simultaneously. In the PMRA, we parameterize the noise as  $\sigma^2 = \frac{L'}{\alpha \log(L)}$ ; this is smaller than how we scaled  $\sigma^2$  in MRA by a factor of  $L'/L$ . The numerator  $L'$  comes from the total “signal energy” that each measurement sees:  $\mathbb{E}\|\pi_S R_{\ell_i} X\|^2 = L'$ , whereas the  $\log(L)$  factor is  $\log$  the size of the group of shifts (and therefore is the same as in MRA).

In the interest of space, we only provide a brief sketch for the proof of Theorem 2.3. We essentially follow the steps of the proof of Theorem 2.1, outlining what modifications need to be made for the argument to work for the PMRA model.

**Template matching.** The MAP estimator is given by

$$\widehat{R}_{\text{MAP}} = \operatorname{argmax}_{\ell'} \frac{\langle \pi_S R_{\ell'} X, Y \rangle}{L'} = \operatorname{argmax}_{\ell'} \left\{ \frac{\langle \pi_S R_{\ell'} X, \pi_S R_{\ell} X \rangle}{L'} + \frac{\langle \pi_S R_{\ell'} X, \sigma Z \rangle}{L'} \right\}.$$

One can prove, as in Lemma 4.1, that with high probability

$$\max_{\ell, \ell'} \left| \frac{\langle \pi_S R_{\ell'} X, \pi_S R_{\ell} X \rangle}{L'} - \mathbb{1}_{\{\ell = \ell'\}} \right| = o(1)$$

holds. Note that the assumption that  $L$  is not too large with respect to  $L'$  (strictly less than exponential in  $L'$ ) is *essential* here: following the proof of Lemma 4.1, we can obtain a concentration bound of the form

$$\Pr\left(\left|(L')^{-1}\langle\pi_S R_{\ell'} X, \pi_S R_{\ell} X\rangle - \mathbb{1}_{\{\ell=\ell'\}}\right| > \kappa\right) \leq \exp(-cL' \min(\kappa, \kappa^2)),$$

which needs to beat a union bound over all indices  $\ell, \ell'$ . Having shown that, we can compare the maximum of the noise term to the maximum of a sequence of standard Gaussians (using Lemmas 4.2, 4.3 and 4.4), to deduce

$$\max_{\ell'=0,\dots,L-1} \frac{\langle\pi_S R_{\ell'} X, \sigma Z\rangle}{L'} \approx \frac{1}{\sqrt{\alpha \log(L)}} \left\langle \frac{\pi_S R_{\ell'} X}{\|\pi_S R_{\ell'} X\|}, Z \right\rangle \approx \sqrt{\frac{2}{\alpha}}.$$

Since  $\frac{\langle\pi_S R_{\ell} X, \pi_S R_{\ell} X\rangle}{L'} \approx 1$ , we conclude that the MAP estimator is successful consistently when  $\alpha > 2$  and fails consistently when  $\alpha < 2$ .

**Lower bound at high SNR ( $\alpha > 2$ ).** The lower bound on the sample complexity follows from applying Corollary 5.2 with the following easy bound on the multi-sample MI  $I(X; Y^n)$ :

$$I(X; Y^n) \leq \frac{L}{2} \log\left(1 + \frac{L'}{L} n \sigma^{-2}\right). \quad (5.10)$$

The idea for proving (5.10) is as follows. Suppose that the shifts  $R_{\ell_1}, \dots, R_{\ell_n}$  were all known. Each measurement  $Y_i$  contains noisy measurements of  $L'$  out of  $L$  coordinates of  $X$ , and note that if we knew the shifts, we would also know to which coordinate of  $X$  each coordinate of  $Y_i$  corresponds. For each coordinate  $i \in [L]$ , let  $n_i$  be the total number of (noisy) measurements of  $X_i$  available across all samples  $Y_1, \dots, Y_n$ . Thus, assuming the shifts are given and known, we can think of the problem as follows: we have  $L$  independent standard (one dimensional) Gaussians,  $X_1, \dots, X_L$ ; for each  $i$ , we measure  $n_i$  measurements of  $X_i$  through an AWGN. Thus,

$$\begin{aligned} I(X; Y^n | R^n = r^n) &\leq \sum_{i=1}^L \frac{1}{2} \log(1 + n_i(r^n) \sigma^{-2}) \leq \frac{L}{2} \log\left(1 + \frac{n_1(r^n) + \dots + n_L(r^n)}{L} \sigma^{-2}\right) \\ &= \frac{L}{2} \log\left(1 + \frac{L'}{L} n \sigma^{-2}\right), \end{aligned}$$

where the second inequality follows from convexity. Averaging over all possible shifts  $r^n$ ,  $I(X; Y^n) \leq I(X; Y^n | R^n) \leq \frac{L}{2} \log\left(1 + \frac{L'}{L} n \sigma^{-2}\right)$ , as claimed.

**Lower bound at low SNR ( $\alpha \leq 2$ ).** We can reiterate the Fano-type argument of Proposition 5.10 without substantial modifications. The single-sample MI from equation (5.7) now becomes

$$I(X; Y) = \frac{L'}{2} \log(1 + \sigma^{-2}) - \log(L) + H(R|X, Y).$$

Lemma 5.9 goes through almost verbatim with

$$\mathcal{S}_\tau = \left\{ R' : \frac{\langle\pi_S R' X, Y\rangle}{\|\pi_S R' X\|^2} \geq 1 - \tau \right\}.$$

instead of the definition given in (5.8), and with the first term in the left-hand-side of (5.9) decaying exponentially fast in  $L'$ , rather than  $L$ . Thus, by the same argument as in the proof of Proposition 5.10, we bound

$$H(R|X, Y) \leq \left(1 - \frac{\alpha}{2} + o(1)\right) \log(L).$$

Expanding  $\frac{L'}{2} \log(1 + \sigma^{-2}) = \frac{L'\sigma^{-2}}{2} + O(L'\sigma^{-4})$  and plugging  $\sigma^2 = L'/(\alpha \log(L))$ , we conclude that  $I(X; Y) = o(\log(L))$ . Combining with Proposition 5.5,

$$n_{\text{PMRA}}^*(L, \alpha, \varepsilon) \geq \log\left(\frac{1}{\varepsilon}\right) \cdot \frac{L}{2I(X; Y)}(1 + o(1)) = \omega\left(\frac{L}{\log(L)}\right) = \omega\left(\frac{L}{L'} \cdot \sigma^2\right).$$

### 5.5 Some remarks on the capacity of the MRA channel

One can think of the model  $Y = RX + \sigma Z$  as a communication channel whose input is  $X$  and output is  $Y$ . A natural question in information theory, then, is to find the *capacity* of this channel, defined as

$$C_{\text{MRA}}(L, \sigma^2) = \max_{P_X: \mathbb{E}\|X\|^2 \leq L} I(X; Y),$$

where the optimization is over all input distributions  $X$  obeying a mean power constraint  $\mathbb{E}\|X\|^2 \leq L$ . The channel capacity is a central quantity in information theory, and characterizes exactly the fundamental limits of data transmission over this channel: in each channel use, one can at best transmit reliably  $C_{\text{MRA}}$  nats of information.

Determining the capacity of the additive white Gaussian channel  $Y = X + \sigma Z$  is a classical problem. It is well-known that

$$C_{\text{AWGN}}(L, \sigma^2) = \frac{L}{2} \log(1 + \sigma^{-2}),$$

and the capacity-achieving distribution is i.i.d Gaussian  $X \sim \mathcal{N}(0, I)$ . It is easy to see that  $C_{\text{MRA}} \leq C_{\text{AWGN}}$ . Indeed, note that  $Y = RX + \sigma Z \stackrel{d}{=} R(X + \sigma Z)$  (by rotation invariance), hence by the data processing inequality (Proposition A.3.3), applied to the Markov chain  $X - (X + \sigma Z) - R(X + \sigma Z)$ , we get

$$I(X; X + \sigma Z) \geq I(X; R(X + \sigma Z)) = I(X; Y),$$

from which  $C_{\text{AWGN}} \geq C_{\text{MRA}}$  follows. At this point, one naturally wonders: (i) Can something non-trivial be said about the ratio  $C_{\text{MRA}}/C_{\text{AWGN}}$ ; in particular, when is it approximately one (say as  $L, \sigma^2 \rightarrow \infty$ )? (ii) What is the capacity achieving input distribution for the MRA channel? In particular, is  $X \sim \mathcal{N}(0, I)$  the capacity achieving input distribution at some (every?) SNR regime?

At very high SNR, namely  $\sigma^{-2}L = \omega(\log(L))$ , equation (5.7) tells us that an i.i.d Gaussian input is “essentially” capacity achieving: if  $X \sim \mathcal{N}(0, I)$ , then

$$I(X; Y) \geq C_{\text{AWGN}} - \log(L) = \frac{L}{2} \log(1 + \sigma^{-2}) - \log(L),$$

and the loss of information,  $\log(L)$  nats, is negligible compared to  $\frac{L}{2} \log(1 + \sigma^{-2})$ .

At very low SNR, however, it turns out that an i.i.d input distribution is very much suboptimal. Consider the input distribution  $X \sim \mathcal{N}(0, \mathbf{1}\mathbf{1}^\top)$ , that is, we allocate the entire power budget on

the direction  $\mathbf{1}/\sqrt{L} = (1/\sqrt{L}, \dots, 1/\sqrt{L})$ . Since all the coordinates of  $X$  are the same, the signal is completely invariant to the shifts, meaning that  $X = RX$  exactly. In that case,

$$I(X; Y) = I(X; X + \sigma Z) = \frac{1}{2} \log(1 + \sigma^{-2}L),$$

so that under *extremely low SNR*, where  $\sigma^{-2}L < 1$  is a constant but small number, we have  $I(X; Y) = \frac{1}{2}\sigma^{-2}L - O(\sigma^{-4}L^2)$ . We can also expand  $C_{\text{AWGN}} = \frac{1}{2}L\sigma^{-2} + O(\sigma^{-4}L)$ , so that  $I(X; Y)$  matches  $C_{\text{AWGN}}$  to leading order in the SNR. On the other hand, recall that for an i.i.d input distribution  $X \sim \mathcal{N}(0, I)$ , we have seen that if the SNR is  $\sigma^{-2}L < \log(L)$  then already  $I(X; Y) = o(1)$ . Thus, i.i.d inputs are highly suboptimal at low SNR.

Determining the channel capacity and the capacity-achieving input distribution, inbetween the extreme SNR regimes  $\sigma^{-2}L = \omega(\log(L))$  and  $\sigma^{-2}L = o(1)$ , looks like an interesting but quite challenging task. An i.i.d input  $\mathcal{N}(0, I)$  has the advantage that it utilizes optimally the available degrees of freedom ( $L$ , the dimension); its disadvantage is that it does not play well with the random shift, in that the signals  $R_\ell X$  are very different to one another. On the other hand, the input  $\mathcal{N}(0, \mathbf{1}\mathbf{1}^\top)$  mitigates best the negative effect of the random shift (it is not affected by it at all), but this is done at the expense of the available degrees of freedom (one instead of  $L$ ). It is interesting to find out how the capacity achieving distribution balances delicately between these two effects.

## 6 Sample complexity upper bound for $\alpha > 2$ via brute-force template matching

In this section we propose a recovery algorithm for the high SNR regime  $\alpha > 2$ , which essentially matches our  $\Omega(L/\log L)$  lower bound on the sample complexity. Our goal here is not to propose a new MRA algorithm, but rather to establish a matching upper bound on the *statistical difficulty* of the problem; that is, we are studying the fundamental information-theoretic (rather than computational) limits of MRA. This is an important distinction because previous papers conjectured that a natural extension of the MRA model, called heterogeneous MRA, suffers from a fundamental computational-statistical gap [BBL18, BBSK<sup>+</sup>17]. In particular, the proposed algorithm is computationally intractable, and involves a brute-force search on an exponentially sized set of candidates. Moreover, our approach is tailored to the case  $\alpha > 2$ , which is exactly the SNR regime where template matching is statistically possible.

**Outline of our algorithm.** Before diving into the technical details of our proposed scheme, we give a brief outline of the approach. The estimation algorithm works in two stages. Suppose we are given  $n$  independent samples. We divide them into two subsamples of sizes  $n_1$  and  $n_2$ ,  $n_1 + n_2 = n$ . We do this so to ensure that the estimator  $\hat{Q}$  produced in step 1 is statistically independent of the additive noise in the samples used for step 2. This simplifies our analysis considerably. The two stages performed by the algorithm are the following.

1. *Brute-force search for a template:* In the first stage, we use the first  $n_1$  samples to find some direction  $\hat{Q} \in \mathbb{S}^{L-1}$  (here  $\mathbb{S}^{L-1}$  is the unit sphere in  $\mathbb{R}^L$ ) such that  $\hat{Q}$  is sufficiently well-aligned with some shift of the true signal, that is,  $\max_\ell L^{-1/2} \langle X, R_\ell^{-1} \hat{Q} \rangle \geq 1 - \eta$ , where  $\eta = \eta(\alpha)$  is small. To do this, we consider a fine-enough cover of the sphere,  $\mathcal{N} \subset \mathbb{S}^{L-1}$ , and take  $\hat{Q} \in \mathcal{N}$



as the minimizer of a certain score:  $\hat{Q} = \operatorname{argmin}_{Q \in \mathcal{N}} \sum_{i=1}^{n_1} s_i(Q)$ , where  $s_i(Q)$  is computed from the  $i$ -th sample  $Y_i$ . Minimizing  $\sum_{i=1}^{n_1} s_i(Q)$  over  $\mathbb{S}^{L-1}$  boils down to a brute-force search over the cover, whose size is exponential in  $L$ . Hence, this algorithm is not efficient. In principle, one could take at this point  $\sqrt{L}\hat{Q} \approx \|X\|\hat{Q}$  as an estimator for  $X$ . Unfortunately, the MSE of this estimator decays at a suboptimal rate with respect to the number of samples  $n$ ; this is remedied by the second step.

2. *Alignment and averaging:* Using  $\hat{Q}$  from the previous step, we perform template matching on the remaining  $n_2$  samples  $Y_1, \dots, Y_{n_2}$  in order to estimate their shifts relative to  $\hat{Q}$ :

$$\hat{R}_{\ell_i} = \operatorname{argmax}_{\ell} \langle Y_i, R_{\ell} \hat{Q} \rangle.$$

The final estimator for  $X$  is then the average of the aligned measurements:

$$\hat{X} = \frac{1}{n_2} \sum_{i=1}^{n_2} \hat{R}_{\ell_i}^{-1} Y_i.$$

All the missing technical details are provided in the next two sections. To ease the reading, the proofs of all lemmas are given in Appendix E.

**Main result of this section.** The main result of this section is the following:

**Proposition 6.1** *Suppose that  $\alpha > 2$ , fix  $\varepsilon > 0$ , and let  $n, L \rightarrow \infty$ . Then, there exists some  $c(\alpha) > 0$  depending on  $\alpha$  such that if*

$$n_1 = c(\alpha)\sigma^2, \quad n_2 = (1 + o(1))\frac{\sigma^2}{\varepsilon},$$

*then the estimator  $\hat{X}$  returned by our algorithm satisfies  $\rho(X, \hat{X}) \leq \varepsilon$  with probability  $1 - o(1)$ .*

Note that when  $\varepsilon > 0$  is small, the sample complexity is dominated by  $n_2$ :

$$n = c(\alpha)\sigma^2 + (1 + o(1))\frac{\sigma^2}{\varepsilon} \approx (1 + o(1))\frac{\sigma^2}{\varepsilon},$$

and thus almost independent of the constant  $c(\alpha)$ . Proposition 6.1 should be compared with the optimal achievable MSE for estimating a signal in AWGN, without the shifts  $L^{-1}\mathbb{E}\|X - \hat{X}_{\text{MMSE}}\|^2 = \frac{\sigma^2}{\sigma^2 + n}$ .

**Proof of Theorem 2.1 (upper bound)** The upper bound for  $\alpha > 2$  follows readily from Proposition 6.1. To show this, we construct a new estimator  $[\hat{X}]$  as follows:  $[\hat{X}] = \hat{X}$  if  $\|\hat{X}\| \leq 10\sqrt{L}$  and  $[\hat{X}] = 0$  otherwise. Note that under the high-probability event  $\|X\| \leq 2\sqrt{L}$ , necessarily  $\rho(X, [\hat{X}]) \leq \rho(X, \hat{X})$ . Write

$$\mathbb{E}\rho(X, [\hat{X}]) = \mathbb{E}\left[\rho(X, [\hat{X}])\mathbf{1}_{\|X\| \leq 2\sqrt{L}}\right] + \mathbb{E}\left[\rho(X, [\hat{X}])\mathbf{1}_{\|X\| > 2\sqrt{L}}\right].$$

Under  $\|X\| \leq 2\sqrt{L}$ , the random variable  $\rho(X, [\hat{X}])$  is bounded, hence by Proposition 6.1,

$$\mathbb{E}\left[\rho(X, [\hat{X}])\mathbf{1}_{\|X\| \leq 2\sqrt{L}}\right] \leq \varepsilon + o(1).$$

As for the other term,

$$\mathbb{E} \left[ \rho(X, [\widehat{X}]) \mathbf{1}_{\|X\| > 2\sqrt{L}} \right] \leq \mathbb{E} \left[ L^{-1/2} (\|X\| + 10L^{1/2}) \mathbf{1}_{\|X\| > 2\sqrt{L}} \right] \leq 6\mathbb{E} \left[ L^{-1/2} \|X\| \mathbf{1}_{L^{-1/2}\|X\| > 2} \right] = o(1).$$

Thus,  $[\widehat{X}]$  uses  $n = [(1 + o(1))/\varepsilon + c(\alpha)] \sigma^2$  samples and achieves  $\mathbb{E}\rho(X, [\widehat{X}]) \leq \varepsilon + o(1)$ , so that

$$\limsup_{L \rightarrow \infty} \frac{n_{\text{MRA}}^*(L, \alpha, \varepsilon)}{\sigma^2/\varepsilon} \leq 1 + O_\alpha(\varepsilon).$$

**Class of “nice signals.”** Before getting to the details of the algorithm, in the analysis that follows, it is convenient to treat the signal  $X$  as fixed and belonging some class of “nice” signals. Specifically, we require that: (i) the signal is sufficiently uncorrelated with its shifts, in that  $L^{-1}\langle X, R_\ell X \rangle \approx 0$  for all  $\ell \neq 0$ , and its norm is concentrated around  $L^{-1}\|X\|^2 \approx 1$ ; (ii) The Fourier (DFT) coefficients of  $X$  are uniformly bounded.

Let  $f_0, \dots, f_{L-1} \in \mathbb{C}^L$  be the DFT basis vectors, that is,  $(f_\ell)_j = L^{-1/2} e^{\frac{2\pi i}{L} \ell j}$ , and  $\mathcal{F} \in U(L)$  be the matrix whose columns are  $f_0, \dots, f_{L-1}$ , so that  $\mathcal{F}^* X \in \mathbb{C}^L$  are the Fourier coefficients of  $X$  (here  $\mathcal{F}^*$  denotes the Hermitian conjugate of  $\mathcal{F}$ .) For  $\kappa > 0$ , we formally consider the set

$$\mathbb{X}_\kappa = \left\{ X \in \mathbb{R}^L \quad : \quad \max_\ell |L^{-1}\langle X, R_\ell X \rangle - \mathbf{1}_{\{\ell=0\}}| \leq \kappa, \quad \text{and} \quad \|\mathcal{F}^* X\|_\infty \leq \sqrt{10 \log(L)} \right\}, \quad (6.1)$$

where  $\mathbf{1}_{\{\ell=0\}} = 1$  when  $\ell = 0$  and is zero otherwise. We take  $\kappa = o(1)$  sufficiently large so to ensure that when  $X \sim \mathcal{N}(0, I)$ , the constraint  $\max_\ell |L^{-1}\langle X, R_\ell X \rangle - \mathbf{1}_{\{\ell=0\}}| \leq \kappa$  holds with probability  $1 - o(1)$  as  $L \rightarrow \infty$ ; by Lemma 4.1, we may choose  $\kappa = c \log(L)/\sqrt{L}$  for  $c > 0$  a large enough constant. Let  $\mathbb{X}$  be the set corresponding to such choice. To lighten the notation, we will not keep track of  $\kappa$  explicitly, instead referring to all vanishing terms as  $o(1)$ . For the other constraint, the exact bound  $\|\mathcal{F}^* X\|_\infty \leq \sqrt{10 \log(L)}$  is somewhat arbitrary, in that 10 can be replaced with any constant greater than 4. The following is quite immediate at this point:

**Lemma 6.2** *Suppose that  $X \sim \mathcal{N}(0, I)$ . Then,  $\Pr(X \notin \mathbb{X}) = o(1)$ .*

We note that it is likely that without assuming that the estimation is over a class of “nice” signals (for example, the class  $\mathbb{X}_\kappa$ ), the situation changes. On that note, we mention the work [Bru19], where it is shown that there are signals  $X$  for which the MLE only attains the rate  $\rho(X, \widehat{X}_{\text{MLE}}) \sim n^{-1/2}$ .

## 6.1 Step 1: Brute force template matching

Recall that our intermediate goal here is to find a direction  $\widehat{Q} \in \mathbb{S}^{L-1}$  such that  $\max_\ell L^{-1/2} \langle X, R_\ell^{-1} \widehat{Q} \rangle \geq 1 - \eta$ , where  $\eta > 0$  is some desired accuracy level. Since, assuming  $X \in \mathbb{X}$ , for any  $Q \in \mathbb{S}^{L-1}$ ,

$$\left\| \frac{X}{\|X\|} - R_\ell^{-1} Q \right\|^2 = 2 - 2 \left\langle \frac{X}{\|X\|}, R_\ell^{-1} Q \right\rangle = 2 - 2L^{-1/2} \langle X, R_\ell^{-1} Q \rangle + o(1),$$

then taking  $\mathcal{N}$  to be a  $\sqrt{\eta}$ -cover of  $\mathbb{S}^{L-1}$ , it must contain some  $Q \in \mathcal{N}$  with  $L^{-1/2} \langle Q, R_\ell^{-1} X \rangle \geq 1 - \frac{1}{2}\eta + o(1)$ . It is well known that one can find a cover of the sphere which is not too large:

**Lemma 6.3** [Lemma 5.13 in [vH14]] *There exists an  $\sqrt{\eta}$ -cover  $\mathcal{N}$  of  $\mathbb{S}^{L-1}$  of size  $|\mathcal{N}| \leq (3/\sqrt{\eta})^L$ . That is, there exists a set  $\mathcal{N} \subset \mathbb{S}^{L-1}$  of size  $|\mathcal{N}| \leq (3/\sqrt{\eta})^L$ , such that  $\forall X \in \mathbb{S}^{L-1}, \exists Q \in \mathcal{N}$  with  $\|X - Q\| \leq \sqrt{\eta}$ .*

For each  $Q \in \mathcal{N}$ , we define its per-sample score:

$$s_i(Q) = s_i^\eta(Q) = \mathbb{1} \left[ \max_{\ell} L^{-1/2} \langle Y_i, R_\ell^{-1} Q \rangle \geq 1 - \frac{3}{4} \eta \right],$$

and the total score  $s(Q) = \sum_{i=1}^{n_1} s_i(Q)$ ,  $n_1$  being the number of samples allocated for this step. That is,  $s(Q)$  is the number of samples  $Y_i$  such that  $L^{-1/2} \langle Q, R_\ell^{-1} Y_i \rangle \geq 1 - \frac{3}{4} \eta$  for some  $\ell$ . The returned estimator is then simply

$$\hat{Q} = \operatorname{argmax}_{Q \in \mathcal{N}} s(Q).$$

Note that  $s_i(\cdot)$  could be thought of as a discontinuous proxy for the log-likelihood (restricted to  $X \in \mathbb{S}^{L-1}$ ):  $\log P(Y_i|X) = \log \sum_{\ell=0}^{L-1} \exp\left(\frac{1}{\sigma^2} \langle X, R_\ell^{-1} Y_i \rangle\right) + \text{constant}$ . When  $\sigma$  is small, the log-likelihood is essentially dominated by  $\max_{\ell} \sigma^{-2} \langle X, R_\ell^{-1} Y_i \rangle$ . Maximizing the likelihood is computationally more straightforward (in the sense that this is a continuous optimization problem, no need to quantize the domain as we do); however, analyzing the MLE directly appears to be difficult [FSWW20, KB20].

We start by showing that there are only a few shifts  $\ell$  such that  $L^{-1/2} \langle X, R_\ell^{-1} Q \rangle$  are all large.

**Lemma 6.4** *Suppose that  $X \in \mathbb{X}$ . For  $Q \in \mathbb{S}^{L-1}$ , let*

$$N_Q(h) = \left| \left\{ \ell : L^{-1/2} |\langle X, R_\ell^{-1} Q \rangle| \geq h \right\} \right|.$$

*Then,  $N_Q(h) \leq h^{-2} \|\mathcal{F}^* X\|_\infty^2 \leq h^{-2} \cdot 10 \log(L)$ .*

We next show that if  $\max_{\ell} L^{-1/2} \langle X, R_\ell^{-1} Q \rangle$  is small, then with high probability the score  $s(Q)$  is not large.

**Lemma 6.5** *Assume that  $X \in \mathbb{X}$ ,  $\alpha > 2$ ,  $\eta < 1 - \sqrt{2/\alpha}$ , and  $L$  is large enough so that  $\log(L) \leq L^{3\eta^2\alpha/128}$ . Suppose that  $Q \in \mathbb{S}^{L-1}$  is such that  $\max_{\ell} L^{-1/2} \langle X, R_\ell^{-1} Q \rangle \leq 1 - \eta$ , then*

$$\Pr(s(Q) \geq n_1/2) \leq \left[ 16 \left( 2 + \frac{640}{\left(1 - \sqrt{\frac{2}{\alpha}}\right)^2} \right) L^{-\eta^2\alpha/128} \right]^{n_1/2}.$$

Next, we prove that if  $\max_{\ell} \langle X, R_\ell^{-1} Q \rangle$  is sufficiently large, then  $s(Q)$  is large with high probability.

**Lemma 6.6** *Assume that  $X \in \mathbb{X}$ ,  $\alpha > 2$ , and  $L$  is large enough so that  $L^{\eta^2\alpha/64} \geq 4$ . Suppose that  $Q \in \mathbb{S}^{L-1}$  is such that  $\max_{\ell} \langle X, R_\ell^{-1} Q \rangle \geq 1 - 5\eta/8$ . Then,*

$$\Pr(s(Q) < n_1/2) \leq e^{-n_1/32}.$$

We are now ready to conclude the analysis of Step 1 of our algorithm.

**Proposition 6.7** Assume that  $X \in \mathbb{X}$ ,  $\alpha > 2$ , and  $\eta < 1 - \sqrt{2/\alpha}$ . Then, there is constant  $c > 0$ , such that whenever

$$n_1 \geq c \frac{L \log(1/\eta)}{\alpha \eta^2 \log(L)} = c \frac{\sigma^2 \log(1/\eta)}{\eta^2},$$

the vector  $\widehat{Q} = \operatorname{argmax}_{Q \in \mathcal{N}} s(Q)$  satisfies  $\max_{\ell} \langle X, R_{\ell}^{-1} Q \rangle \geq 1 - \eta$  with probability  $1 - o(1)$  as  $n_1, L \rightarrow \infty$ . In fact, the error probability decays exponentially fast with  $n_1$ .

**Proof.** As argued in the beginning of this section, the  $\sqrt{\eta}$ -cover  $\mathcal{N}$  contains some  $Q \in \mathbb{S}^{L-1}$  such that  $L^{-1/2} \langle X, R_{\ell}^{-1} Q \rangle \geq 1 - \eta/2 - o(1) \geq 1 - 5\eta/8$  for some  $\ell$ . By Lemma 6.6, with probability greater than  $1 - e^{-n_1/32}$ , this vector has score  $s(Q) \geq n_1/2$ . It therefore suffices to show that with high probability, all the vectors  $Q \in \mathcal{N}$  that are bad, meaning that  $\max_{\ell} L^{-1/2} \langle X, R_{\ell}^{-1} Q \rangle < 1 - \eta$ , have score  $s(Q) < n_1/2$ . By Lemmas 6.3 and 6.5,

$$\begin{aligned} \Pr(\exists \text{bad } Q \in \mathcal{N} : s(Q) \geq n_1/2) &\leq |\mathcal{N}| \cdot \Pr(s(Q) \geq n_1/2 \mid Q \text{ is bad}) \\ &\leq (9/\eta)^{L/2} \cdot \left[ 16 \left( 2 + \frac{640}{\left(1 - \sqrt{\frac{2}{\alpha}}\right)^2} \right) L^{-\eta^2 \alpha / 128} \right]^{n_1/2} \\ &\leq \left( C(\alpha) e^{-c_1 \eta^2 \alpha \log(L) + c_2 \frac{L}{n} \log(1/\eta)} \right)^{n_1}, \end{aligned}$$

where  $c_1, c_2 > 0$  are absolute constants, and  $C(\alpha)$  depends on  $\alpha$ . Then, this probability tends to 0 as  $n_1, L \rightarrow \infty$  (exponentially fast in  $n_1$ ) whenever  $n_1 \geq c \frac{L \log(1/\eta)}{\alpha \eta^2 \log(L)}$  for some other  $c > 0$ . ■

Note that at this point we could take  $\widehat{X} = L^{1/2} \cdot \widehat{Q}$  as an estimator for  $X$ , so that

$$\rho(X, \widehat{X}) = \min_{\ell} \|L^{-1/2} X - R_{\ell}^{-1} Q\|^2 \leq 2\eta + o(1),$$

holds with high probability. For *fixed*  $\eta$ , this estimator indeed captures the correct dimensional scaling of the sample complexity, namely, that  $n = O(L/(\alpha \log L))$  samples are sufficient to get non-trivial alignment error. However, its dependence on  $\eta$  is seemingly quite bad: for estimating a signal in AWGN, without the shifts, the optimal dependence on  $\eta$  should look like  $O(L/(\alpha \log L) \cdot \eta^{-1})$ , rather than the much worse  $O(L/(\alpha \log L) \cdot \eta^{-2} \log(1/\eta))$  we were able to show. In the next section, we see how to achieve this “correct” rate by essentially recovering the shifts on all but a vanishing fraction of the samples, and averaging the properly aligned measurements.

## 6.2 Step 2: Achieving optimal MSE decay rate by alignment and averaging

Suppose that one has access to a known template  $Q \in \mathbb{S}^{L-1}$ , such that  $\langle X, Q \rangle \geq 1 - \eta$ . Since  $L^{-1} \|X\|^2 = 1 + o(1)$ , this is the same as having  $\|L^{-1/2} X - Q\|^2 \leq 2\eta + o(1)$ , and since  $\max_{\ell \neq 0} L^{-1} |\langle X, R_{\ell} X \rangle| = o(1)$ , we see that for any  $\ell \neq 0$ ,

$$\|L^{-1/2} R_{\ell} X - Q\| \geq \|L^{-1/2} [R_{\ell} X - X]\| - \|L^{-1/2} X - Q\| \geq \sqrt{2} - \sqrt{2\eta} - o(1).$$

In particular, we see that when  $\sqrt{2\eta} < \sqrt{2} - \sqrt{2\eta}$ , that is,  $\eta < 1/4$  (and  $L$  is sufficiently large), there is a *unique*  $\ell$  (specifically,  $\ell = 0$ ) such that  $\|L^{-1/2} X - R_{\ell} Q\|^2 \leq 2\eta + o(1)$ . In that case, the idea of matching a sample  $Y_i = R_{\ell_i} X + \sigma Z$  against the template  $Q$  becomes well-posed, in the sense that its desired outcome is clear: we would like to recover the shift  $R_{\ell_i}$ .

**Lemma 6.8** Assume that  $X \in \mathbb{X}$  and  $\alpha > 2$ . Let  $Y = R_\ell X + \sigma Z$ , and suppose that  $Q \in \mathbb{S}^{L-1}$  is independent of  $Y$  and satisfies  $\max_{\ell'} L^{-1/2} \langle X, R_{\ell'}^{-1} Q \rangle \geq 1 - \eta$ , where

$$\sqrt{\eta} < \frac{1}{2}(1 - \sqrt{2/\alpha}).$$

Denote the maximizing shift by  $\ell^*$ . Let  $\hat{\ell} = \operatorname{argmax}_{\ell'} \langle Y, R_{\ell'} Q \rangle$ . Then

$$\Pr \left( \hat{\ell} \neq \ell - \ell^* \right) \leq 2L^{-\frac{1}{2}\alpha(1/2-1/\sqrt{2\alpha}-\sqrt{\eta})^2+o(1)}.$$

Given Lemma 6.8, we propose the following estimation strategy. Suppose we would like to estimate  $X$  up to error  $\rho(X, \hat{X}) \leq \varepsilon < 1$ . Fix some  $\eta > 0$  with  $\sqrt{\eta} < (1 - \sqrt{2/\alpha})/2$  (for concreteness, say  $\eta = (1 - \sqrt{2/\alpha})^2/16$ ). We first apply the algorithm of Step 1 (Section 6.1) to obtain  $\hat{Q} \in \mathbb{S}^{L-1}$  such that  $\max_{\ell} \langle X, R_{\ell}^{-1} \hat{Q} \rangle \geq 1 - \eta$ . Assuming that  $n_1 \geq \frac{c \log(1/\eta)}{\eta^2} \sigma^2 = c_\eta \sigma^2$ , we are successful with probability  $1 - o(1)$ . Let  $\ell^*$  be such that  $\langle X, R_{\ell^*}^{-1} \hat{Q} \rangle \geq 1 - \eta$ . Next, for  $n_2$  new independent samples, we compute for each measurement  $\hat{\ell}_i = \operatorname{argmax}_{\ell} \langle Y_i, R_{\ell} \hat{Q} \rangle$  and return the aligned sample average:

$$\hat{X} = \frac{1}{n_2} \sum_{i=1}^{n_2} R_{\hat{\ell}_i}^{-1} Y_i. \quad (6.2)$$

Lemma 6.8 tells us that we should expect most of the aligned measurements  $R_{\hat{\ell}_i}^{-1} Y_i$  to be well-aligned with  $R_{\ell^*}$ , that is,  $R_{\hat{\ell}_i}^{-1} Y_i = R_{\ell^*} X + \mathcal{N}(0, \sigma^2 I)$ . This means that,  $\hat{X} \approx R_{\ell^*} X + \mathcal{N}(0, (\sigma^2/n_2)I)$ , hence  $\|R_{\ell^*} X - \hat{X}\|^2 \approx \sigma^2/n_2$ , which is smaller than  $\varepsilon$  if  $n_2 \geq \sigma^2/\varepsilon$ . We make this argument precise below:

**Proposition 6.9** Assume that  $X \in \mathbb{X}$  and  $\alpha > 2$ . Fix  $\varepsilon > 0$  and some  $\eta < \frac{1}{2}(1 - \sqrt{2/\alpha})^2$ . Let  $\hat{Q} \in \mathbb{S}^{L-1}$  be the output of Step 1 (run with a tuning parameter  $\eta$  and  $n_1$  samples). Let  $\hat{X}$  be as in equation (6.2), computed from  $n_2$  new samples. Suppose that  $n_1, n_2, L \rightarrow \infty$  with

$$n_1/\sigma^2 \rightarrow \gamma_1, \quad n_2/\sigma^2 \rightarrow \frac{\gamma_2}{\varepsilon},$$

where  $\gamma_1$  and  $\gamma_2$  are constants satisfying

$$\gamma_1 = \gamma_1(\eta) \geq \frac{c \log(1/\eta)}{\eta^2}, \quad \gamma_2 > 1,$$

( $c$  being the universal constant from Proposition 6.7). Then,

$$\Pr \left( \rho(X, \hat{X}) \leq \varepsilon \right) \rightarrow 1.$$

Proposition 6.1 now immediately follows from Lemma 6.2 and Proposition 6.9.

## 7 Conclusions and extensions

In this work we have studied the sample complexity of the MRA problem in the limit of large  $L$ . In this regime, we have shown that the parameter  $\alpha = \frac{\sigma^2 \log L}{L}$  plays a crucial role in characterizing the best attainable performance of any estimator.

As mentioned above, the MRA model is primarily motivated by the cryo-EM technology to constitute the 3-D structure of biological molecules. In the cryo-EM literature, it was shown that it is effective to assume that the molecule was drawn from a Gaussian prior with decaying power spectrum [Sch12]. In addition, the 3-D rotations are usually not distributed uniformly over the group  $SO(3)$ . We now discuss briefly how these different aspects can be potentially incorporated into our framework.

**Prior on the signal.** Our model assumes a Gaussian i.i.d. prior on the signal  $X$  to be reconstructed. While this assumption lends itself to a relatively clean analysis, and allows to compare our bounds on  $n_{\text{MRA}}^*(L, \alpha, \varepsilon)$  to the simple benchmark  $n_{\text{AWGN}}^*(L, \alpha, \varepsilon)$ , many of our results can be generalized to treat other priors on  $X$ . In particular, all of our sample complexity lower bounds are based on lower bounding the mutual information between  $X$  and  $\hat{X}$  under the constraint  $\mathbb{E}[\rho(X, \hat{X})] \leq \varepsilon$  on the one hand, and upper bounding  $I(X; Y^n)$  under the MRA model, on the other hand. In Proposition 5.1 we have relied on the Gaussian rate distortion function to lower bound  $I(X; \hat{X})$  for any estimator that achieves MSE at most  $\varepsilon$ . For  $X$  whose distribution is not  $\mathcal{N}(0, I)$ , we can either compute the corresponding rate distortion function explicitly, or simply apply Shannon’s lower bound  $R(D) \geq h(X) - \frac{1}{2} \log(2\pi e D)$ , see [Ber71]. Our upper bounds on  $I(X; Y^n)$  in the regime  $\alpha > 1$  are based on Lemma 5.3, followed by lower bounding  $I(R^n; X|Y^n)$  using Fano-like arguments. It is easy to see that (5.4) continues to hold, with  $\leq$  instead of  $=$ , for any random variable  $X$  with  $\mathbb{E}\|X\|^2 \leq L$ . Furthermore, the lower bounds on  $I(R^n; X|Y^n)$  we derive in Section 5.3.2 remain valid whenever  $\frac{\|X\|}{L}$  is sufficiently concentrated around 1 and  $\frac{\langle X, R_\ell X \rangle}{L}$  is sufficiently concentrated around 0 for all  $\ell = 1, \dots, L-1$ . In particular, this is the case for (sufficiently light-tailed) i.i.d. zero-mean and unit variance distributions. In light of the discussion above, we see that the parameter  $\alpha = \frac{\sigma^2 \log L}{L}$  is of great importance whenever the random signal  $X$  satisfies the above concentration requirements and has differential entropy proportional to  $L$ .

**Shift distribution.** Assuming uniform prior on the i.i.d. shifts  $R_{\ell_1}, \dots, R_{\ell_n}$  is a worst-case analysis. Indeed, for any given distribution, shifting all measurements again  $R_{u_i} Y_i$ , for  $u_i \stackrel{i.i.d.}{\sim} \text{Uniform}(\{0, \dots, L-1\})$  before feeding them to the estimator leads to (1.1). However, previous works (for fixed  $L$ ) showed that harnessing non-uniformity can make a big difference in the sample complexity [ABL<sup>+</sup>18, SKK<sup>+</sup>20]. With some effort, our upper bounds on  $I(X; Y^n)$  in the regime  $\alpha > 1$  should also extend to treat this case. Here, the main challenge is to generalize Lemma 5.9 to the case of non-uniform distribution, i.e., to find a sharp estimate on the smallest possible size of a list of candidates for the true shift, which contains the true shift with high probability.

**Extension to other groups.** We believe that many aspects of our information-theoretical analysis can be generalized to other (families of) discrete groups, denoted here by  $\mathcal{G}_L$ , which satisfy the following properties (roughly speaking): (i) If  $X$  is suitably generic and  $g \neq h$ , then  $\langle gX, hX \rangle$  is very small - concretely, if  $X \sim \mathcal{N}(0, I)$ , then  $\mathbb{E}[\langle gX, hX \rangle] = 0$ ; (ii) The size of the group  $|\mathcal{G}_L|$  does not grow too fast (strictly less than exponentially fast in  $L$ ). These conditions imply that whenever  $X$  is isotropic and sufficiently light-tailed (e.g., sub-Gaussian),  $\{gX\}_{g \in \mathcal{G}}$  are “almost orthogonal.” The proper noise scaling to consider would then be  $\sigma^2 = \frac{L}{\alpha \log |\mathcal{G}_L|}$ , with  $\alpha = 2$  being the critical noise level—this comes from the fact that  $\max_{g \in \mathcal{G}_L} \langle gX, Z \rangle \approx \sqrt{2 \log |\mathcal{G}_L|}$ . For continuous compact groups, we suspect that one might be able to apply some of our arguments by cleverly discretizing

the suitable group action. Carrying out a program of this type seems as a promising direction for future research.

## Acknowledgment

E.R. and O.O. are supported in part by the ISF under Grant 1791/17. T.B. is supported in part by NSF-BSF grant no. 2019752, and the Zimin Institute for Engineering Solutions Advancing Better Lives.

## References

- [ABL<sup>+</sup>18] Emmanuel Abbe, Tamir Bendory, William Leeb, João M Pereira, Nir Sharon, and Amit Singer. Multireference alignment is easier with an aperiodic translation distribution. *IEEE Transactions on Information Theory*, 65(6):3565–3584, 2018.
- [ALS19] Yariv Aizenbud, Boris Landa, and Yoel Shkolnisky. Rank-one multi-reference factor analysis. *arXiv preprint arXiv:1905.12442*, 2019.
- [APS17] Emmanuel Abbe, João M Pereira, and Amit Singer. Sample complexity of the boolean multireference alignment problem. In *2017 IEEE International Symposium on Information Theory (ISIT)*, pages 1316–1320. IEEE, 2017.
- [APS18] Emmanuel Abbe, João M Pereira, and Amit Singer. Estimation in the group action channel. In *2018 IEEE International Symposium on Information Theory (ISIT)*, pages 561–565. IEEE, 2018.
- [AT09] Robert J Adler and Jonathan E Taylor. *Random fields and geometry*. Springer Science & Business Media, 2009.
- [BBL18] Nicolas Boumal, Tamir Bendory, Roy R Lederman, and Amit Singer. Heterogeneous multireference alignment: A single pass approach. In *2018 52nd Annual Conference on Information Sciences and Systems (CISS)*, pages 1–6. IEEE, 2018.
- [BBM<sup>+</sup>17] Tamir Bendory, Nicolas Boumal, Chao Ma, Zhizhen Zhao, and Amit Singer. Bispectrum inversion with application to multireference alignment. *IEEE Transactions on signal processing*, 66(4):1037–1050, 2017.
- [BBS20] Tamir Bendory, Alberto Bartesaghi, and Amit Singer. Single-particle cryo-electron microscopy: Mathematical theory, computational challenges, and opportunities. *IEEE Signal Processing Magazine*, 37(2):58–76, 2020.
- [BBSK<sup>+</sup>17] Afonso S Bandeira, Ben Blum-Smith, Joe Kileel, Amelia Perry, Jonathan Weed, and Alexander S Wein. Estimation under group actions: recovering orbits from invariants. *arXiv preprint arXiv:1712.10163*, 2017.
- [BCLS20] Afonso S Bandeira, Yutong Chen, Roy R Lederman, and Amit Singer. Non-unique games over compact groups and orientation estimation in cryo-EM. *Inverse Problems*, 36(6):064002, 2020.

- [BCSZ14] Afonso S Bandeira, Moses Charikar, Amit Singer, and Andy Zhu. Multireference alignment using semidefinite programming. In *Proceedings of the 5th conference on Innovations in theoretical computer science*, pages 459–470. ACM, 2014.
- [Ber71] Toby Berger. *Rate distortion theory: A mathematical basis for data compression*. Prentice-Hall, 1971.
- [BJL<sup>+</sup>20] Tamir Bendory, Ariel Jaffe, William Leeb, Nir Sharon, and Amit Singer. Super-resolution multi-reference alignment. *arXiv preprint arXiv:2006.15354*, 2020.
- [Bou16] Nicolas Boumal. Nonconvex phase synchronization. *SIAM Journal on Optimization*, 26(4):2355–2377, 2016.
- [Bru19] Victor-Emmanuel Brunel. Learning rates for gaussian mixtures under group action. In *Conference on Learning Theory*, pages 471–491, 2019.
- [BRW17] Afonso S Bandeira, Philippe Rigollet, and Jonathan Weed. Optimal rates of estimation for multi-reference alignment. *arXiv preprint arXiv:1702.08546*, 2017.
- [BZS16] Tejal Bhamre, Teng Zhang, and Amit Singer. Denoising and covariance estimation of single particle cryo-EM images. *Journal of structural biology*, 195(1):72–81, 2016.
- [CT12] Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley & Sons, 2012.
- [DHF07] Laurens De Haan and Ana Ferreira. *Extreme value theory: an introduction*. Springer Science & Business Media, 2007.
- [DLR77] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.
- [DZS15] Jeffrey J Donatelli, Peter H Zwart, and James A Sethian. Iterative phasing for fluctuation X-ray scattering. *Proceedings of the National Academy of Sciences*, 112(33):10286–10291, 2015.
- [FSWW20] Zhou Fan, Yi Sun, Tianhao Wang, and Yihong Wu. Likelihood landscape and maximum likelihood estimation for the discrete orbit recovery model. *arXiv preprint arXiv:2004.00041*, 2020.
- [GVX14] Navin Goyal, Santosh Vempala, and Ying Xiao. Fourier PCA and robust tensor decomposition. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pages 584–593, 2014.
- [Har70] Richard A Harshman. Foundations of the PARAFAC procedure: Models and conditions for an explanatory multimodal factor analysis. 1970.
- [HL19] Matthew Hirn and Anna Little. Wavelet invariants for statistically robust multi-reference alignment. *arXiv preprint arXiv:1909.11062*, 2019.
- [Kam80] Zvi Kam. The reconstruction of structure from electron micrographs of randomly oriented particles. *Journal of Theoretical Biology*, 82(1):15–39, 1980.



- [KB20] Anya Katsevich and Afonso Bandeira. Likelihood maximization and moment matching in low SNR Gaussian mixture models. *arXiv preprint arXiv:2006.15202*, 2020.
- [LBB<sup>+</sup>18] Eitan Levin, Tamir Bendory, Nicolas Boumal, Joe Kileel, and Amit Singer. 3D ab initio modeling in cryo-EM by autocorrelation analysis. In *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, pages 1569–1573. IEEE, 2018.
- [LHHH05] I. Land, S. Huettinger, P. A. Hoehner, and J. B. Huber. Bounds on information combining. *IEEE Transactions on Information Theory*, 51(2):612–619, 2005.
- [LM00] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338, 2000.
- [LRA93] Sue E Leurgans, Robert T Ross, and Rebecca B Abel. A decomposition for three-way arrays. *SIAM Journal on Matrix Analysis and Applications*, 14(4):1064–1083, 1993.
- [MBB<sup>+</sup>19] Chao Ma, Tamir Bendory, Nicolas Boumal, Fred Sigworth, and Amit Singer. Heterogeneous multireference alignment for images with application to 2D classification in single particle reconstruction. *IEEE Transactions on Image Processing*, 29:1699–1710, 2019.
- [PSB19] Thomas Pumis, Amit Singer, and Nicolas Boumal. The generalized orthogonal Procrustes problem in the high noise regime. *arXiv preprint arXiv:1907.01145*, 2019.
- [PW19] Yury Polyanskiy and Yihong Wu. Lecture notes on information theory. 2019. [http://people.lids.mit.edu/yp/homepage/data/itlectures\\_v5.pdf](http://people.lids.mit.edu/yp/homepage/data/itlectures_v5.pdf).
- [PWB<sup>+</sup>19] Amelia Perry, Jonathan Weed, Afonso S Bandeira, Philippe Rigollet, and Amit Singer. The sample complexity of multireference alignment. *SIAM Journal on Mathematics of Data Science*, 1(3):497–517, 2019.
- [PWBM18] Amelia Perry, Alexander S Wein, Afonso S Bandeira, and Ankur Moitra. Message-passing algorithms for synchronization problems over compact groups. *Communications on Pure and Applied Mathematics*, 71(11):2275–2322, 2018.
- [RG19] Elad Romanov and Matan Gavish. The noise-sensitivity phase transition in spectral group synchronization over compact groups. *Applied and Computational Harmonic Analysis*, 2019.
- [RV13] Mark Rudelson and Roman Vershynin. Hanson-Wright inequality and sub-gaussian concentration. *Electronic Communications in Probability*, 18, 2013.
- [Sch12] Sjors HW Scheres. RELION: implementation of a Bayesian approach to cryo-EM structure determination. *Journal of structural biology*, 180(3):519–530, 2012.
- [Sig98] Fred J Sigworth. A maximum-likelihood approach to single-particle image refinement. *Journal of structural biology*, 122(3):328–339, 1998.
- [Sin11] Amit Singer. Angular synchronization by eigenvectors and semidefinite programming. *Applied and computational harmonic analysis*, 30(1):20–36, 2011.
- [Sin18] Amit Singer. Mathematics for cryo-electron microscopy. *arXiv preprint arXiv:1803.06714*, 2018.

- [SKK<sup>+</sup>20] Nir Sharon, Joe Kileel, Yuehaw Khoo, Boris Landa, and Amit Singer. Method of moments for 3D single particle ab initio modeling with non-uniform distribution of viewing angles. *Inverse Problems*, 36(4):044003, 2020.
- [SS11] Amit Singer and Yoel Shkolnisky. Three-dimensional structure determination from common lines in cryo-EM by eigenvectors and semidefinite programming. *SIAM journal on imaging sciences*, 4(2):543–572, 2011.
- [SSZ05] I. Sutskever, S. Shamai, and J. Ziv. Extremes of information combining. *IEEE Transactions on Information Theory*, 51(4):1313–1325, 2005.
- [vH14] Ramon van Handel. Probability in high dimension. Technical report, PRINCETON UNIV NJ, 2014.
- [WHS17] L. Wang, S. Hu, and O. Shayevitz. Quickest sequence phase detection. *IEEE Transactions on Information Theory*, 63(9):5834–5849, 2017.

## A Information Theoretic Background

In this section we review some basic information theoretic definitions and results that are needed throughout this paper. The proofs of the results below can be found in any textbook on information theory, e.g. [CT12], and are therefore omitted.

For a discrete random variable  $X \sim P_X$  supported on the alphabet  $\mathcal{X}$ , the entropy is defined as

$$H(X) = H(P_X) := \sum_{x \in \mathcal{X}} P_X(x) \log \frac{1}{P_X(x)} = \mathbb{E}_{X \sim P_X} \left[ \log \frac{1}{P_X(X)} \right].$$

For a pair of random variables  $(X, Y) \sim P_{XY}$ , where  $X$  is discrete, the conditional entropy of  $X$  given  $Y$  is defined as

$$H(X|Y) := \mathbb{E}_{y \sim P_Y} [H(X|Y = y)] = \mathbb{E}_{y \sim P_Y} [H(P_{X|Y=y})].$$

Similarly, if  $X$  is a continuous random variable on  $\mathbb{R}^d$  with density  $p_X$ , its differential entropy is defined as

$$h(X) = h(P_X) := \int_{x \in \mathbb{R}^d} p_X(x) \log \frac{1}{p_X(x)} dx = \mathbb{E}_{X \sim P_X} \left[ \log \frac{1}{p_X(X)} \right].$$

For a pair of random variables  $(X, Y) \sim P_{XY}$ , where  $X$  is continuous and has conditional density  $p_{X|Y=y}$  for all  $y \in \mathcal{Y}$ , where  $\mathcal{Y}$  is the alphabet of  $Y$ , the conditional entropy is defined as

$$h(X|Y) = \mathbb{E}_{y \sim P_Y} [h(X|Y = y)] = \mathbb{E}_{y \sim P_Y} [h(P_{X|Y=y})].$$

### Proposition A.1 (Properties of entropy and differential entropy)

1. **Non-negativity of entropy:** For a discrete random variable  $X$  the entropy satisfies  $H(X) \geq 0$ , with equality if and only if  $X$  is deterministic.

2. **Uniform distribution maximizes entropy:** For a discrete random variable  $X$  supported on  $\mathcal{X}$

$$H(X) \leq \log |\mathcal{X}|,$$

and this is attained with equality if and only if  $X \sim \text{Uniform}(\mathcal{X})$ .

3. **Gaussian distribution maximizes differential entropy under second moment constraints:** Suppose that the continuous random variable  $X$  is supported on  $\mathbb{R}^d$ , and has covariance matrix  $\Sigma = \mathbb{E}[(X - \mathbb{E}(X))(X - \mathbb{E}(X))^\top]$ . Then,

$$h(X) \leq \frac{1}{2} \log \left( (2\pi e)^d \det(\Sigma) \right), \quad (\text{A.1})$$

and this is attained with equality if and only if  $X \sim \mathcal{N}(\mu, \Sigma)$  for some  $\mu \in \mathbb{R}^d$ .

4. **Chain rule:** For discrete random variables  $(X, Y) \sim P_{XY}$  we have

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y).$$

For continuous random variables  $(X, Y) \sim P_{XY}$ , we have

$$h(X, Y) = h(X) + h(Y|X) = h(Y) + h(X|Y).$$

5. **Concavity:** The functions  $P_X \mapsto H(P_X)$  and  $P_X \mapsto h(P_X)$  are concave. Consequently, conditioning reduces entropy, that is

$$H(X|Y) \leq H(X)$$

if  $X$  is discrete, and

$$h(X|Y) \leq h(X)$$

if  $X$  is continuous. In both cases, the bounds are attained with equality iff  $X$  and  $Y$  are statistically independent.

We will also make use of Fano's inequality, as stated below.

**Proposition A.2 (Fano's inequality)** Let  $(X, Y) \sim P_{XY}$ , where  $X$  is a discrete random variable supported on  $\mathcal{X}$ . Then, for any estimator  $\hat{X} = \hat{X}(Y)$  of  $X$  from  $Y$ , we have

$$H(X|Y) \leq \log 2 + \Pr(X \neq \hat{X}) \log |\mathcal{X}|.$$

If both  $(X, Y) \sim P_{XY}$  are discrete, the mutual information between  $X$  and  $Y$  is defined as

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X),$$

and if they are both continuous

$$I(X; Y) = h(X) - h(X|Y) = h(Y) - h(Y|X).$$

If one is discrete, say  $X$ , and the other continuous, say  $Y$ , then

$$I(X; Y) = H(X) - H(X|Y) = h(Y) - h(Y|X).$$

For a triplet of random variables  $(X, Y, Z) \sim P_{XYZ}$ , the conditional mutual information is defined as

$$I(X; Y|Z) = \mathbb{E}_{z \sim P_Z} [I(X; Y|Z = z)],$$

where  $I(X; Y|Z = z)$  is the mutual information between  $X$  and  $Y$  under the distribution  $(X, Y) \sim P_{XY|Z=z}$ .

**Proposition A.3 (Properties of Mutual Information)**

1. **Non-negativity of mutual information:**  $I(X; Y) \geq 0$  with equality iff  $X$  and  $Y$  are statistically independent.
2. **Chain rule:** For  $(X, Y, Z) \sim P_{XYZ}$  we have

$$I(X; Y, Z) = I(X; Y) + I(X; Z|Y) = I(X; Z) + I(X; Y|Z).$$

3. **Data processing inequality:** Assume  $X - Y - Z$  is a Markov chain in this order, that is their joint distribution decomposes as  $P_{XYZ} = P_X P_{Y|X} P_{Z|Y}$ , then

$$I(X; Z) \leq I(X; Y).$$

4. **Invertible functions:** For any function  $f : \mathcal{Y} \rightarrow \mathcal{A}$ , where  $\mathcal{A}$  is an arbitrary alphabet, we have  $I(X; f(Y)) \leq I(X; Y)$  with equality if  $f$  is invertible.
5. **Mutual information for memoryless channels:** Let  $(X^n, Y^n) \sim P_{X^n Y^n} = P_{X^n} P_{Y^n|X^n}$  and assume the channel from  $X^n$  to  $Y^n$  is a product channel, that is  $P_{Y^n|X^n} = \prod_{i=1}^n P_{Y_i|X_i}$ . Then

$$I(X^n; Y^n) \leq \sum_{i=1}^n I(X_i; Y_i).$$

This bound is attained with equality if  $P_{X^n} = \prod_{i=1}^n P_{X_i}$ , i.e., if  $X^n$  is memoryless as well.

6. **Gaussian mutual information:** Let  $X, Z \sim \mathcal{N}(0, I)$  be statistically independent  $L$ -dimensional random vectors with i.i.d. standard normal entries. Then

$$I(X; X + \sigma Z) = \frac{L}{2} \log \left( 1 + \frac{1}{\sigma^2} \right).$$

For a random variable  $X \sim P_X$  supported on alphabet  $\mathcal{X}$ , a reconstruction alphabet  $\hat{\mathcal{X}}$  and a distortion measure  $d : \mathcal{X} \times \hat{\mathcal{X}} \rightarrow \mathbb{R}$ , the rate distortion function (RDF) is defined as

$$R(D) = \min_{P_{\hat{X}|X} : \mathbb{E}[d(X, \hat{X})] \leq D} I(X; \hat{X}),$$

where both  $I(X; \hat{X})$  and  $\mathbb{E}[d(X, \hat{X})]$  are evaluated with respect to the joint distribution  $P_X P_{\hat{X}|X}$ . The solution of the optimization problem above for the quadratic Gaussian case is well known, and is summarized in the proposition below.

**Proposition A.4 (Quadratic Gaussian RDF)** *Let  $X \sim \mathcal{N}(0, \sigma^2 I)$  be a random vector in  $\mathbb{R}^d$ ,  $\hat{\mathcal{X}} = \mathbb{R}^d$ , and  $d(x, \hat{x}) = \frac{1}{d} \|x - \hat{x}\|^2$ . Then,*

$$R(D) = \frac{d}{2} \log \left( \frac{\sigma^2}{D} \right). \quad (\text{A.2})$$

*In particular, if  $X \sim \mathcal{N}(0, \sigma^2 I)$  and  $\hat{X}$  is such that  $\frac{1}{d} \mathbb{E} \|X - \hat{X}\|^2 \leq D$ , then*

$$I(X; \hat{X}) \geq \frac{d}{2} \log \left( \frac{\sigma^2}{D} \right).$$

## B Proof of Lemma 4.1

Before getting to the proof, we recall the Hanson-Wright inequality:

**Lemma B.1 (Hanson-Wright inequality for sub-Gaussian random vectors, Theorem 1.1 in [RV13])** *Let  $X$  be a random vector with independent entries such that for all  $i$ ,*

$$\mathbb{E} X_i = 0, \quad \|X_i\|_{\psi_2} \leq K,$$

*where  $\|X_i\|_{\psi_2} = \inf \left\{ s > 0 : \mathbb{E} e^{(X_i/s)^2} \leq 2 \right\}$ . Let  $A$  be any matrix. Then, there is a universal constant  $c > 0$  such that*

$$\Pr \left( \left| X^\top A X - \mathbb{E}(X^\top A X) \right| > t \right) \leq 2 \exp \left[ -c \min \left( \frac{t^2}{K^4 \|A\|_F^2}, \frac{t}{K^2 \|A\|} \right) \right].$$

It is immediate to verify that if  $X \sim \mathcal{N}(0, \sigma^2)$ , then  $\|X\|_{\psi_2} = \sigma/\sqrt{2 \log 2} = c\sigma$ . Also, for any  $\ell$ ,  $\|R_\ell\| = 1$  (since  $R_\ell \in O(L)$ ) and therefore  $\|R_\ell\|_F^2 \leq L$ . Also,

$$\mathbb{E}(\langle X, R_\ell X \rangle) = \text{tr}(R_\ell) = \begin{cases} L & \text{if } \ell = 0, \\ 0 & \text{otherwise.} \end{cases}$$

By the Hanson-Wright inequality, Lemma B.1,

$$\Pr (|\langle X, R_\ell X \rangle - \mathbb{E}(\langle X, R_\ell X \rangle)| \geq L\kappa) \leq 2 \exp(-c \min((L\kappa)^2/L, L\kappa)) = 2 \exp(-cL \min(\kappa, \kappa^2)).$$

The claimed result follows by a union bound.

## C The spectrum of the operators $R_\ell$

We recall some elementary facts about the spectrum of the operators  $R_\ell$ :

**Lemma C.1** *The eigenvalues of  $R_\ell + R_\ell^\top$  are exactly (with multiplicities)  $\lambda_{\ell,k} = 2 \cos\left(\frac{2\pi}{L}\ell k\right)$ ,  $k = 0, \dots, L-1$ . Moreover,*

$$\sum_{k=0}^{L-1} \lambda_{\ell,k} = \begin{cases} L & \text{if } \ell = 0, \\ 0 & \text{otherwise.} \end{cases}$$

**Proof.** Let  $f_k \in \mathbb{C}^L$ ,  $k = 0, \dots, L-1$ , be the DFT basis vectors, namely  $f_{k,j} = L^{-1/2} e^{\frac{2\pi i}{L}kj}$ . It is immediate to verify that  $f_k$  is an eigenvector of  $R_\ell$  with eigenvalue  $\lambda_{\ell,k} = e^{\frac{2\pi i}{L}\ell k}$ :

$$(R_\ell f_k)_j = (f_k)_{j+\ell} = e^{\frac{2\pi i}{L}\ell k} (f_k)_j.$$

Hence,  $(R_\ell + R_\ell^\top)f_k = (e^{\frac{2\pi i}{L}\ell k} + e^{-\frac{2\pi i}{L}\ell k})f_k = 2 \cos\left(\frac{2\pi i}{L}\ell k\right) f_k$ . This means that  $\lambda_{\ell,k}$  are the eigenvalues of  $R_\ell + R_\ell^\top$  as an operator  $\mathbb{C}^L \rightarrow \mathbb{C}^L$ . But since  $R_\ell + R_\ell^\top$  is also diagonalizable over  $\mathbb{R}^L$  by an orthogonal matrix, there also exists a *real* orthonormal eigenbasis  $u_1, \dots, u_L \in \mathbb{R}^L$  with  $(R_\ell + R_\ell^\top)u_k = \lambda_{\ell,k}u_k$ . As for the last claim, it follows from  $\sum_{k=0}^{L-1} \lambda_{\ell,k} = 2\Re\left\{\sum_{k=0}^{L-1} e^{\frac{2\pi i}{L}\ell k}\right\}$ , the right-hand side being  $L$  when  $\ell = 0$  and zero otherwise. ■

## D Proof of Lemma 5.9

Suppose that the event  $\mathcal{A} = \mathcal{A}(\kappa)$  from Lemma 4.1 holds, meaning that  $|L^{-1}\|X\|^2 - 1| \leq \kappa$  and  $\max_{\ell' \neq 0} L^{-1} |\langle X, R_{\ell'} X \rangle| \leq \kappa$ . Observe that

$$R \notin \mathcal{S}_\tau \Leftrightarrow \frac{\sigma\langle X, R^{-1}Z \rangle}{\|X\|^2} < -\tau.$$

Conditioned on  $X$ ,

$$\frac{\sigma\langle X, R^{-1}Z \rangle}{\|X\|^2} \sim \mathcal{N}\left(0, \sigma^2/\|X\|^2\right),$$

and under  $\mathcal{A}$ , this variance is  $\sigma^2/\|X\|^2 = \frac{L}{\|X\|^2 \alpha \log(L)} \leq \frac{1}{\alpha(1-\kappa)\log(L)}$ . Thus,

$$\Pr(R \notin \mathcal{S}_\tau | \mathcal{A}) \leq e^{-\frac{1}{2}\tau^2 \alpha(1-\kappa)\log(L)} = L^{-\frac{1}{2}\tau^2 \alpha(1-\kappa)}.$$

Now, suppose that  $R' \neq R$ . Then

$$\begin{aligned} \Pr(R' \in \mathcal{S}_\tau | \mathcal{A}, R) &= \Pr\left(\frac{\langle X, (R')^{-1}RX \rangle}{\|X\|^2} + \frac{\sigma\langle X, (R')^{-1}Z \rangle}{\|X\|^2} \geq 1 - \tau \mid \mathcal{A}, R\right) \\ &\leq \Pr\left(\frac{\sigma\langle X, (R')^{-1}Z \rangle}{\|X\|^2} \geq 1 - \tau - \frac{\kappa}{1-\kappa} \mid \mathcal{A}, R\right) \\ &\leq L^{-\frac{1}{2}\alpha(1-\kappa)\left(1-\tau-\frac{\kappa}{1-\kappa}\right)^2}, \end{aligned}$$

where we used the fact that under  $\mathcal{A}$ ,  $\frac{\langle X, (R')^{-1}RX \rangle}{\|X\|^2} \leq \kappa/(1-\kappa)$ , and uniformly bounded the variance of  $\frac{\sigma\langle X, (R')^{-1}Z \rangle}{\|X\|^2}$  conditioned on  $X$  and under  $\mathcal{A}$  as before. Since the bound above is uniform in  $R$ , of course,

$$\Pr(R' \in \mathcal{S}_\tau | \mathcal{A}) \leq L^{-\frac{1}{2}\alpha(1-\kappa)\left(1-\tau-\frac{\kappa}{1-\kappa}\right)^2}.$$

Now,

$$\mathbb{E} \left[ |\mathcal{S}_\tau| \mid \mathcal{A} \right] \leq 1 + (L-1) \cdot L^{-\frac{1}{2}\alpha(1-\kappa)(1-\tau-\frac{\kappa}{1-\kappa})^2} \leq 1 + L^{1-\frac{1}{2}\alpha(1-\kappa)(1-\tau-\frac{\kappa}{1-\kappa})^2}.$$

Setting  $M = L^{1-\frac{1}{2}\alpha(1-\kappa)(1-\tau-\frac{\kappa}{1-\kappa})^2+\delta}$ , by Markov's inequality, and assuming  $\alpha \leq 2$ ,

$$\Pr(|\mathcal{S}_\tau| \geq M \mid \mathcal{A}) \leq \frac{1 + L^{1-\frac{1}{2}\alpha(1-\kappa)(1-\tau-\frac{\kappa}{1-\kappa})^2}}{L^{1-\frac{1}{2}\alpha(1-\kappa)(1-\tau-\frac{\kappa}{1-\kappa})^2+\delta}} \leq 2L^{-\delta}.$$

Combining both estimates and taking a union bound,

$$\begin{aligned} \Pr(R \notin \mathcal{S}_\tau \text{ or } |\mathcal{S}_\tau| > M) &\leq \Pr(R \notin \mathcal{S}_\tau \mid \mathcal{A}) + \Pr(|\mathcal{S}_\tau| \geq M \mid \mathcal{A}) + \Pr(\overline{\mathcal{A}}) \\ &\leq L^{-\frac{1}{2}\alpha(1-\kappa)(1-\tau-\frac{\kappa}{1-\kappa})^2} + 2L^{-\delta} + \Pr(\overline{\mathcal{A}}) \\ &\leq L^{-\frac{1}{2}\alpha(1-\kappa)(1-\tau-\frac{\kappa}{1-\kappa})^2} + 2L^{-\delta} + 2Le^{-cL \min(\kappa, \kappa^2)}, \end{aligned}$$

where the last inequality follows from Lemma 4.1.

## E Proofs of Section 6

### E.1 Proof of Lemma 6.2

Recall that  $\kappa$  was chosen so that the first constraint holds with probability  $1 - o(1)$ . All that remains, then, is to show that  $\|\mathcal{F}^*X\|_\infty \leq \sqrt{10 \log(L)}$  holds with high probability. Let  $f_\ell \in \mathbb{C}^L$  be the  $\ell$ -th DFT basis vector, so that  $(\mathcal{F}^*X)_\ell = \langle X, f_\ell \rangle$ . Observe that the real and imaginary parts of  $(\mathcal{F}^*X)_\ell$  are both Gaussians, with variances bounded by 1. Hence,

$$\begin{aligned} \Pr(|(\mathcal{F}^*X)_\ell|^2 > 10 \log(L)) &\leq \Pr(|\Re(\mathcal{F}^*X)_\ell|^2 > 5 \log(L)) + \Pr(|\Im(\mathcal{F}^*X)_\ell|^2 > 5 \log(L)) \\ &\leq 4e^{-\frac{5}{2} \log(L)} = 4L^{-5/2}, \end{aligned}$$

so  $\Pr(\|\mathcal{F}^*X\|_\infty > \sqrt{10 \log(L)}) \leq L \cdot 4L^{-5/2} = 4L^{-3/2} = o(1)$ .

### E.2 Proof of Lemma 6.4

Bounding  $\mathbb{1}[|X| \geq a] \leq \frac{|X|}{a}$ , as in the proof of Markov's inequality, we have

$$N_Q(h) = \sum_{\ell=0}^{L-1} \mathbb{1} \left[ L^{-1} |\langle X, R_\ell^{-1}Q \rangle|^2 \geq h^2 \right] \leq h^{-2} L^{-1} \sum_{\ell=0}^{L-1} |\langle X, R_\ell^{-1}Q \rangle|^2.$$

We may write

$$L^{-1} \sum_{\ell=0}^{L-1} \langle X, R_\ell^{-1}Q \rangle^2 = Q^\top \left( L^{-1} \sum_{\ell=0}^{L-1} (R_\ell X)(R_\ell X)^\top \right) Q \leq \|\mathcal{M}(X)\|,$$

where  $\mathcal{M}(X)$  is the operator

$$\mathcal{M}(X) = L^{-1} \sum_{\ell=0}^{L-1} (R_\ell X)(R_\ell X)^\top.$$

It is convenient to write  $\mathcal{M}(X)$  in terms of the DFT basis  $f_0, \dots, f_{L-1}$

$$\begin{aligned}\mathcal{M}(X) &= L^{-1} \sum_{\ell=0}^{L-1} \sum_{k,j=0}^{L-1} e^{\frac{2\pi}{L}\ell(k-j)} \langle X, f_k \rangle \langle X, f_j \rangle^* f_k f_j^* \\ &= \sum_{k=0}^{L-1} |\langle X, f_k \rangle|^2 f_k f_k^*,\end{aligned}$$

which means that the eigenvalues of  $\mathcal{M}(X)$  are exactly the magnitudes of the fourier coefficients of  $X$ , squared. In particular,  $\|\mathcal{M}(X)\| = \|\mathcal{F}^* X\|_\infty^2 \leq 10 \log(L)$ .

### E.3 Proof of Lemma 6.5

Note that  $s_1(Q), \dots, s_n(Q)$  are i.i.d Bernoulli-distributed. Write

$$\begin{aligned}\Pr(s_i(Q) = 1) &= \Pr\left(\exists \ell : L^{-1/2} \langle Y_i, R_\ell^{-1} Q \rangle \geq 1 - \frac{3}{4}\eta\right) \\ &= \Pr\left(\exists \ell : L^{-1/2} [\langle X, R_\ell^{-1} Q \rangle + \sigma \langle Z, R_\ell^{-1} Q \rangle] \geq 1 - \frac{3}{4}\eta\right) \\ &= \Pr\left(\exists \ell : L^{-1/2} \langle X, R_\ell^{-1} Q \rangle + \frac{\langle Z, R_\ell^{-1} Q \rangle}{\sqrt{\alpha \log(L)}} \geq 1 - \frac{3}{4}\eta\right).\end{aligned}$$

Let

$$\mathcal{L}(Q) = \left\{ \ell : L^{-1/2} \langle X, R_\ell^{-1} Q \rangle \geq 1 - \sqrt{\frac{2}{\alpha}} - \frac{7\eta}{8} \right\}$$

be the set of shifts for which  $L^{-1/2} \langle X, R_\ell^{-1} Q \rangle$  is somewhat large. For  $S \subset [L]$ , set

$$p(S) = \Pr\left(\exists \ell \in S : L^{-1/2} \langle X, R_\ell^{-1} Q \rangle + \frac{\langle Z, R_\ell^{-1} Q \rangle}{\sqrt{\alpha \log(L)}} \geq 1 - \frac{3}{4}\eta\right),$$

so that  $\Pr(s_i(Q) = 1) \leq p(\mathcal{L}(Q)) + p(\overline{\mathcal{L}(Q)})$ . Since

$$\mathbb{E} \max_{\ell \in \mathcal{L}(Q)} \langle Z, R_\ell^{-1} Q \rangle \leq \sqrt{2 \log |\mathcal{L}(Q)|} \leq \sqrt{2 \log(L)},$$

(since each  $\langle Z, R_\ell^{-1} Q \rangle \sim \mathcal{N}(0, 1)$ ; see comment after Lemma 4.2), we apply Lemma 4.4 to get

$$\begin{aligned}p(\overline{\mathcal{L}(Q)}) &\leq \Pr\left(\exists \ell \in \overline{\mathcal{L}(Q)} : \frac{\langle Z, R_\ell^{-1} Q \rangle}{\sqrt{\alpha \log(L)}} \geq \sqrt{2/\alpha} + \eta/8\right) \\ &\leq \Pr\left(\max_{\ell \in \overline{\mathcal{L}(Q)}} \langle Z, R_\ell^{-1} Q \rangle \geq \mathbb{E} \left[ \max_{\ell \in \mathcal{L}(Q)} \langle Z, R_\ell^{-1} Q \rangle \right] + \frac{1}{8}\eta\sqrt{\alpha \log(L)}\right) \\ &\leq 2e^{-\frac{1}{2}(\eta/8)^2 \alpha \log(L)} = 2L^{-\eta^2 \alpha / 128}.\end{aligned}$$



For the other term,

$$p(\mathcal{L}(Q)) \leq |\mathcal{L}(Q)| \Pr \left( \frac{\langle Z, R_\ell^{-1} Q \rangle}{\sqrt{\alpha \log(L)}} \geq \eta/4 \right) \leq |\mathcal{L}(Q)| \cdot e^{-\frac{1}{2}(\eta/4)^2 \alpha \log(L)} = |\mathcal{L}(Q)| \cdot L^{-\eta^2 \alpha / 32}.$$

By Lemma 5.9,

$$|\mathcal{L}(Q)| \leq \frac{10 \log(L)}{\left(1 - \sqrt{\frac{2}{\alpha}} - \frac{7\eta}{8}\right)^2} \leq \frac{640 \log(L)}{\left(1 - \sqrt{\frac{2}{\alpha}}\right)^2},$$

where we also used  $\eta < 1 - \sqrt{2/\alpha}$ . Combining,

$$\Pr(s_i(Q) = 1) \leq 2L^{-\eta^2 \alpha / 128} + \frac{640 \log(L)}{\left(1 - \sqrt{\frac{2}{\alpha}}\right)^2} L^{-\eta^2 \alpha / 32} \leq \left(2 + \frac{640}{\left(1 - \sqrt{\frac{2}{\alpha}}\right)^2}\right) L^{-\eta^2 \alpha / 128},$$

where we used the assumption that  $L$  is large enough so that  $\log(L) \leq L^{3\eta^2 \alpha / 128}$ . We use

$$\Pr(\text{Binom}(n_1, p) \geq k) = \sum_{t=k}^{n_1} \binom{n_1}{t} p^t (1-p)^{n_1-t} \leq p^k \sum_{t=k}^{n_1} \binom{n_1}{t} \leq 2^{n_1} p^k.$$

Since  $s(Q) \sim \text{Binom}(n_1, \Pr(s_i(Q) = 1))$ ,

$$\Pr(s(Q) \geq n_1/2) \leq \left[ 16 \left(2 + \frac{640}{\left(1 - \sqrt{\frac{2}{\alpha}}\right)^2}\right) L^{-\eta^2 \alpha / 128} \right]^{n_1/2}$$

as claimed.

#### E.4 Proof of Lemma 6.6

Let  $\ell$  be such that  $\langle X, R_\ell^{-1} Q \rangle \geq 1 - 5\eta/8$ . Then

$$\Pr(s_i(Q) = 0) \leq \Pr \left( \frac{\langle Z, R_\ell^{-1} Q \rangle}{\sqrt{\alpha \log(L)}} < (5/8 - 3/4)\eta \right) \leq L^{-\eta^2 \alpha (5/8 - 3/4)^2} = L^{-\eta^2 \alpha / 64} \leq 1/4.$$

Thus, using Hoeffding's inequality,

$$\Pr(s(Q) < n_1/2) \leq \Pr(\text{Ber}(3/4, n_1) < n_1/2) \leq e^{-n_1/32}.$$

#### E.5 Proof of Lemma 6.8

To simplify the notation, assume without loss of generality that  $\ell = \ell^* = 0$ . By the discussion above, for any  $\ell' \neq 0$ ,  $L^{-1/2} \langle X, R_{\ell'} Q \rangle \leq 1 - (\sqrt{2} - \sqrt{2\eta})^2 / 2 + o(1) = \sqrt{4\eta} - \eta + o(1)$ . For any  $\tau$ ,

$$\begin{aligned} \Pr(\widehat{\ell} \neq 0) &\leq \Pr \left( L^{-1/2} \langle Y, Q \rangle < \tau \quad \text{or} \quad \exists \ell' \neq 0 : L^{-1/2} \langle Y, R_{\ell'} Q \rangle \geq \tau \right) \\ &\leq \Pr \left( L^{-1/2} \langle Y, Q \rangle < \tau \right) + \Pr \left( \max_{\ell' \neq 0} L^{-1/2} \langle Y, R_{\ell'} Q \rangle \geq \tau \right). \end{aligned}$$

Suppose that  $\tau \leq 1 - \eta$ . We may bound

$$\Pr\left(L^{-1/2}\langle Y, Q \rangle < \tau\right) \leq \Pr\left(L^{-1/2}\sigma\langle Z, Q \rangle < \tau - (1 - \eta)\right) \leq L^{-\frac{1}{2}\alpha(\tau - 1 + \eta)^2}.$$

Using Lemma 4.3 and assuming  $\tau \geq \sqrt{4\eta} - \eta + \sqrt{2/\alpha} + o(1)$ , we may also bound

$$\begin{aligned} \Pr\left(\max_{\ell \neq 0} L^{-1/2}\langle Y, R_\ell Q \rangle \geq \tau\right) &\leq \Pr\left(\max_{\ell \neq 0} L^{-1/2}\sigma\langle Z, R_\ell Q \rangle \geq \tau - \left(\sqrt{4\eta} - \eta + o(1)\right)\right) \\ &\leq L^{-\frac{1}{2}\alpha\left(\tau - (\sqrt{4\eta} - \eta + o(1)) - \sqrt{2/\alpha}\right)^2}. \end{aligned}$$

We would now like to choose  $\sqrt{4\eta} - \eta + \sqrt{2/\alpha} < \tau < 1 - \eta$  so to maximize  $\min\left(|\tau - (1 - \eta)|, |\tau - (\sqrt{4\eta} - \eta + \sqrt{2/\alpha})|\right)$  observe that this interval is non-empty exactly iff  $\sqrt{4\eta} < 1 - \sqrt{2/\alpha}$ . The best  $\tau$  is then simply the midpoint,  $\tau_* = 1/2 - \eta + 1/\sqrt{2\alpha} + \sqrt{\eta}$ , which gives

$$\Pr(\widehat{\ell} \neq 0) \leq 2L^{-\frac{1}{2}\alpha(1/2 - 1/\sqrt{2\alpha} - \sqrt{\eta})^2 + o(1)}.$$

## E.6 Proof of Proposition 6.9

Let  $\widehat{Q} \in \mathbb{S}^{L-1}$  be the output of Step 1. Let  $\mathcal{V}_1$  be the event that  $\max_\ell \langle X, R_\ell^{-1}\widehat{Q} \rangle \geq 1 - \eta$ , and call the maximizing shift  $\ell^*$ . By Proposition 6.7,  $\Pr(\mathcal{V}_1) = 1 - o(1)$ .

Let  $Y_1, \dots, Y_{n_2}$  be  $n_2$  new samples (independent of those used for Step 1), and let  $\mathcal{I} \subset [n_2]$  be the set of misaligned samples, namely,  $\mathcal{I} = \left\{i \in [n_2] : \widehat{\ell}_i \neq \ell_i - \ell^*\right\}$ . We start by providing a high-probability bound on  $|\mathcal{I}|$ . Lemma 6.8 tells us that conditioned on  $\mathcal{V}_1$ , the random variables  $\mathbf{1}_{\{i \in \mathcal{I}\}}$  are i.i.d Bernoullis with  $\Pr(\mathbf{1}_{\{i \in \mathcal{I}\}} = 1) = p \leq 2L^{-\frac{1}{2}\alpha(1/2 - 1/\sqrt{2\alpha} - \sqrt{\eta})^2 + o(1)}$  (the exponent being strictly negative by our requirements on  $\alpha, \eta$ ), thus  $|\mathcal{I}| \sim \text{Binom}(n_2, p)$ . By Bernstein's inequality,

$$\Pr\left(|\mathcal{I}| \geq pn_2 + t \mid \mathcal{V}_1\right) \leq \exp\left(-\frac{\frac{1}{2}t^2}{np(1-p) + \frac{1}{3}t}\right).$$

Note that the right hand side is  $o(1)$  whenever  $t = t(L)$  is such that  $t \rightarrow \infty$  and  $t = \omega(\sqrt{n_2 p})$  as  $n_2, L \rightarrow \infty$ . Thus, there is some  $c = c(\alpha, \eta) > 0$  such that for  $K = L^{-c}n_2$ , the event  $|\mathcal{I}| \leq K$  holds with high probability.

Let  $\widehat{\mu} = \frac{1}{n_2} \sum_{i=1}^{n_2} R_{\widehat{\ell}_i}^{-1} R_{\ell_i} X$  and  $\widehat{W} = \frac{1}{n_2} \sum_{i=1}^{n_2} R_{\widehat{\ell}_i}^{-1} Z_i$ , so that  $\widehat{X} = \widehat{\mu} + \sigma\widehat{W}$ . We decompose the error,

$$\begin{aligned} L^{-1/2}\|R_{\ell^*} X - \widehat{X}\| &\leq L^{-1/2}\|R_{\ell^*} X - \widehat{\mu}\| + L^{-1/2}\sigma\|\widehat{W}\| \leq L^{-1/2}\frac{2|\mathcal{I}|}{n_2}\|X\| + L^{-1/2}\sigma\|\widehat{W}\| \\ &= \frac{2|\mathcal{I}|}{n_2}(1 + o(1)) + L^{-1/2}\sigma\|\widehat{W}\|. \end{aligned}$$

We have already argued that with high probability  $\frac{2|\mathcal{I}|}{n_2} = o(1)$ ; it therefore remains to show that for the appropriate choice of  $n_2$ , the bound  $L^{-1/2}\sigma\|\widehat{W}\| \leq \sqrt{\varepsilon}$  holds with probability  $1 - o(1)$ .

Observe that conditioned on  $|\mathcal{I}| \leq K$ ,  $R_{\ell^*}^{-1}\widehat{W}$  can be written as

$$R_{\ell^*}^{-1}\widehat{W} = \frac{1}{n_2} \sum_{i=1}^{n_2} R_i Z_i,$$

where  $R_i \neq R_{\ell_i}^{-1}$  for at most  $K$  indices. Note that the estimated shifts  $R_{\widehat{\ell}_i}$  generally depend on the noise  $Z_i$ , and therefore we cannot simply conclude that  $R_{\widehat{\ell}_i} Z_i \sim \mathcal{N}(0, I)$ , which would have meant that  $R_{\ell^*}^{-1}\widehat{W} \sim \mathcal{N}(0, n_2^{-1}I)$ . We need to use a slightly more elaborate argument to overcome this difficulty.

For a subset  $S \subset [n]$ ,  $|S| = K$ ,  $S = \{i_1, \dots, i_K\}$ , and shifts  $\mathbf{R} = (R_1, \dots, R_K)$ , define

$$W(S, \mathbf{R}) = \frac{1}{n_2} \sum_{i \notin S} R_{\ell_i}^{-1} Z_i + \frac{1}{n_2} \sum_{j=1}^K R_j Z_{i_j}.$$

Conditioned on the high-probability event  $|\mathcal{I}| \leq K$ , we have

$$\|\widehat{W}\| = \|R_{\ell^*} \widehat{W}\| \leq \max_{|S|=K, \mathbf{R} \in [L]^K} \|W(S, \mathbf{R})\|,$$

where the maximization is over all possible subsets  $S$  of size  $K$  and shifts  $R_1, \dots, R_K$ . It is therefore enough to show that  $L^{-1/2}\sigma \cdot \max_{|S|=k, \mathbf{R} \in [L]^K} \|W(S, \mathbf{R})\| \leq \sqrt{\varepsilon}$  holds with probability  $1 - o(1)$ . Since the shifts  $R_{\ell_i}$  are independent of the noise  $Z_i$ , for every fixed  $S$  and  $\mathbf{R}$  we have  $W(S, \mathbf{R}) \sim \mathcal{N}(0, n_2^{-1}I)$ . Therefore, by a union bound,

$$\Pr \left( L^{-1/2}\sigma \cdot \max_{|S|=k, \mathbf{R} \in [L]^K} \|W(S, \mathbf{R})\| > \sqrt{\varepsilon} \right) \leq n_2^K L^K \Pr \left( \|G\|^2 > \varepsilon \cdot \frac{L}{\sigma^2} \cdot n_2 \right),$$

where  $\mathbb{R}^L \ni G \sim \mathcal{N}(0, I)$ , hence  $\|G\|^2$  is a standard  $\chi^2$ -distributed random variable with  $L$  degrees of freedom, and we bounded  $\binom{n_2}{K} \leq n_2^K$  for the number of possible choices of  $S$ . Using the tail bound of [LM00, Lemma 1]:

$$\Pr_{G \sim \mathcal{N}(0, I)} \left( \|G\|^2 \geq L + 2\sqrt{Lx} + 2x \right) \leq e^{-x}.$$

Plugging in any  $x_0 = x_0(L)$  such that  $(n_2 L)^K e^{-x_0} = o(1)$ , that is,  $x_0 = \omega(K \log(n_2 L))$ , we obtain

$$\Pr \left( L^{-1/2}\sigma \cdot \max_{|S|=k, \mathbf{R} \in [L]^K} \|W(S, \mathbf{R})\| > \left( \frac{\sigma^2}{n_2} \left[ 1 + 2\sqrt{x_0/L} + 2x_0/L \right] \right)^{1/2} \right) = o(1),$$

hence the condition

$$n_2 \geq \frac{\sigma^2}{\varepsilon} \left[ 1 + 2\sqrt{x_0/L} + 2x_0/L \right]$$

suffices. Since  $K = L^{-c}n_2$ , if moreover  $n_2 = o(L)$  then  $x_0 = o(L)$ , hence  $n_2 = \gamma_2 \sigma^2 / \varepsilon$  for any  $\gamma_2 > 1$  would suffice for large enough  $L$ .