# HETEROGENEOUS MULTIREFERENCE ALIGNMENT: A SINGLE PASS APPROACH

*Nicolas Boumal,*[1,2]    *Tamir Bendory,*[2]    *Roy R. Lederman,*[2]    *Amit Singer* [1,2]

Mathematics Department[1] and PACM,[2] Princeton University, Princeton, NJ, USA

## ABSTRACT

Multireference alignment (MRA) is the problem of estimating a signal from many noisy and cyclically shifted copies of itself. In this paper, we consider an extension called *heterogeneous* MRA, where $K$ signals must be estimated, and each observation comes from one of those signals, unknown to us. This is a simplified model for the heterogeneity problem notably arising in cryo-electron microscopy. We propose an algorithm which estimates the $K$ signals without estimating either the shifts or the classes of the observations. It requires only one pass over the data and is based on low-order moments that are invariant under cyclic shifts. Given sufficiently many measurements, one can estimate these invariant features averaged over the $K$ signals. We then design a smooth, non-convex optimization problem to compute a set of signals which are consistent with the estimated averaged features. We find that, in many cases, the proposed approach estimates the set of signals accurately despite non-convexity, and conjecture the number of signals $K$ that can be resolved as a function of the signal length $L$ is on the order of $\sqrt{L}$.

*Index Terms*— Multireference alignment, bispectrum, non-convex optimization, expectation-maximization, Gaussian mixture models, cryo-EM, heterogeneity

## 1. INTRODUCTION

Multireference alignment (MRA) seeks to estimate a signal from numerous noisy and cyclically shifted copies of itself. This problem serves as a model for scientific problems in structural biology [1, 2, 3, 4], radar [5, 6] and image processing [7, 8, 9]. Algorithmic and statistical properties of MRA have been analyzed in [10, 11, 12, 13, 14].

In this paper, we consider *heterogeneous* MRA, where more than one signal must be estimated from noisy and shifted observations. This investigation is motivated in part by applications in single particle cryo-electron microscopy (cryo-EM) and X-ray free electron lasers (XFEL). These imaging techniques are used to map the three-dimensional structure of molecules at near-atomic resolution from either two-dimensional noisy tomographic projections (cryo-EM) or diffraction images (XFEL), taken at unknown viewing directions [15, 16, 17]. It is typical in those applications to acquire a large number of very noisy observations. Heterogeneity arises when more than one type of molecule (or conformation) appears in a sample. Heterogeneous MRA is a simplified one-dimensional model of this situation, where the rotations are replaced by cyclic shifts. The tomographic projection from three-dimensional objects to two-dimensional images is not modeled in this formulation of MRA.

One of the opportunities in cryo-EM and XFEL—compared to X-ray crystallography—is that they potentially allow to esti-

mate multiple conformations of molecules observed together in heterogeneous mixtures. Achieving this capability is one of the key technological challenges [18], and the subject of many recent works [19, 20, 21, 22, 23, 24, 25].

This motivates us to consider heterogeneity in MRA, where $K \geq 1$ signals must be recovered from unlabeled, shifted, noisy observations. It has been shown in [12] that such a mixture of signals can be estimated from fifth-order moments using a tensor decomposition algorithm. In this paper, we propose a method which, empirically, estimates all $K$ signals simultaneously using only the third-order moments, same as what is necessary for the homogeneous case [11].

In a nutshell, following [10], for each observation, we compute features (moments) which are invariant under cyclic shift. Namely, we compute the mean, power spectrum and bispectrum. Averaging these features over all observations yields an estimator for the averaged invariant features. We then set up a smooth, non-convex optimization problem designed to recover the unknown signals from the averaged features. Our numerical study demonstrates that this approach performs well for a broad range of parameters, with random initialization, despite the non-convexity of the optimization problem.

Importantly, our approach bypasses the need to estimate the latent variables of the problem, namely, the unknown shifts and classes of the individual observations. Furthermore, it naturally works in single-pass streaming mode and is parallelizable. We show that for large data and low signal-to-noise ratio (SNR), this can be about as accurate and much faster than a standard alternative, namely, expectation maximization (EM).

We note that signal estimation based on invariant features is an old idea in signal processing [26] and cryo-EM [27, 28]. Furthermore, the idea of estimating more than one object from mixed measurements was recently used for Gaussian mixtures [29, 30] and for mixed low-rank matrix sensing [31]. Our main contribution is to demonstrate how a non-convex optimization approach to heterogeneous MRA resolves multiple signals in a single pass over the data at low SNR.

## 2. MRA WITHOUT HETEROGENEITY

We begin by introducing the homogeneous MRA model. Let $x \in \mathbb{R}^L$ be the unknown signal and let $R_r$ be the cyclic shift operator: $(R_r x)[n] = x[n - r]$, with all indices considered modulo $L$. We are given $N$ measurements:

$$y_j = R_{r_j} x + \varepsilon_j, \quad j = 1, \dots, N, \quad (2.1)$$

where $\varepsilon_j \sim \mathcal{N}(0, \sigma^2 I)$ is i.i.d. white Gaussian noise. Our goal is to estimate $x$ up to shift in a high-noise regime where the unknown shifts $r_j$ could not be recovered reliably even if $x$ were known (see Cramér–Rao lower bounds in [32]).

Following recent work [10], we turn to *invariant features*: moments of the signal which are invariant under shifts. We denote by $\hat{x}$ the discrete Fourier transform (DFT) of $x$, with

$\hat{x}[k] = \sum_{n=0}^{L-1} x[n]e^{-2\pi ink/L}$. A shift by $r$ adds phase to the DFT: $\widehat{(R_r x)}[k] = \hat{x}[k]e^{-\frac{2\pi ir}{L}k}$. Using this fact, it is easy to verify that the mean, power spectrum, and bispectrum of $x$, defined by the formulae

$$\mu_x := \hat{x}[0]/L, \tag{2.2}$$

$$P_x[k] := \hat{x}[k]\overline{\hat{x}[k]} = |\hat{x}[k]|^2, \tag{2.3}$$

$$B_x[k,\ell] := \hat{x}[k]\overline{\hat{x}[\ell]}\hat{x}[\ell-k], \tag{2.4}$$

respectively, are invariant to shifts.

Simple expectation computations show the average invariant features of the measurements converge to the invariant features of the signal (up to bias terms) as $N \to \infty$, allowing to estimate them:

$$M_1 := \frac{1}{N}\sum_{j=1}^{N}\mu_{y_j} \to \mu_x, \tag{2.5}$$

$$M_2 := \frac{1}{N}\sum_{j=1}^{N}P_{y_j} \to P_x + \sigma^2 L\mathbf{1}, \tag{2.6}$$

$$M_3 := \frac{1}{N}\sum_{j=1}^{N}B_{y_j} \to B_x + \mu_x \cdot \sigma^2 L^2 A, \tag{2.7}$$

where $\mathbf{1}$ is a vector of all-ones and $A \in \mathbb{R}^{L \times L}$ is a zero matrix except for $A[0,0] = 3$ and 1's on the remaining entries of the diagonal and the first row and column (if working with complex signals, subtract one from the first column of $A$).

The variance on $M_1$ scales like $O(\sigma^2/N)$. Because they square and cube the noise, variance on $M_2$ and $M_3$, respectively, scales like $O(\sigma^4/N)$ and $O(\sigma^6/N)$, with cross-terms contributing additional $O(\sigma^2/N)$ variance, relevant only at high SNR. Thus, provided $N/(\sigma^2+\sigma^6)$ is large enough (which is necessary for MRA [11]), the invariant features $\mu_x, P_x$ and $B_x$ can be estimated reliably. Various algorithms were proposed in [10] to recover $x$ from these moments. In the following section, we show how similar principles can be harnessed for the heterogeneous MRA model.

## 3. HETEROGENEITY VIA MIXED INVARIANT FEATURES

In this work, we extend the invariant features approach to heterogeneous MRA. In this setup, $K$ unknown signals $x_1,\ldots,x_K \in \mathbb{R}^L$ must be estimated (we assume they are distinct even up to shift), from the $N$ observations

$$y_j = R_{r_j}x_{v_j} + \varepsilon_j, \quad j = 1,\ldots,N, \tag{3.1}$$

where classes $v_j$ as well as shifts $r_j$ are unknown, and $\varepsilon_j$ is i.i.d. white Gaussian noise of variance $\sigma^2$ as before. We assume $v_j$'s are drawn i.i.d. from a (possibly unknown) mixing probability $w \in \Delta_K = \{w \in \mathbb{R}^K : w \geq 0 \text{ and } \sum_k w[k] = 1\}$ (the simplex): $w[k]$ indicates the proportion of measurements which come from class $k$. Without the shifts, this model reduces to the well-studied Gaussian mixture model (GMM) with known, diagonal noise covariance, for which low-order moment methods have been studied [29, 30].

Our goal in heterogeneous MRA is to esimate the signals $x_1,\ldots,x_K$ (up to shifts and ordering) and possibly to estimate $w$. As before, we are not interested in $r_j$ or $v_j$ of individual measurements. We propose to do this based on $M_1, M_2$ and $M_3$ (2.5)–(2.7), which are now *mixed* invariant features. Expectation computations

yield:

$$M_1 \to \sum_{k=1}^{K} w[k]\mu_{x_k}, \tag{3.2}$$

$$M_2 \to \sum_{k=1}^{K} w[k]P_{x_k} + \sigma^2 L\mathbf{1}, \tag{3.3}$$

$$M_3 \to \sum_{k=1}^{K} w[k]\left(B_{x_k} + \mu_{x_k} \cdot \sigma^2 L^2 A\right). \tag{3.4}$$

Having computed $M_1, M_2$ and $M_3$ in $O(NL^2)$ flops (parallelized over $N$), we search for $K$ signals and possibly for a mixing density $w$ which best agree with the data in a least-squares sense. The weights below proceed from a crude approximation of the variances of the individual terms, where it is assumed the power spectra of the unknown signals are roughly constant and close to the value $P$ (Appendix A) (a common factor $\frac{N}{\sigma^2 L}$ was suppressed):

$$\min_{\substack{\tilde{x}_1,\ldots,\tilde{x}_K \in \mathbb{R}^L \\ \tilde{w} \in \Delta_K}} \left|\sum_{k=1}^{K} \tilde{w}[k]L\mu_{\tilde{x}_k} - LM_1\right|^2$$

$$+ \frac{1}{\sigma^2 L + 2P}\left\|\sum_{k=1}^{K} \tilde{w}[k]P_{\tilde{x}_k} + \sigma^2 L\mathbf{1} - M_2\right\|_2^2$$

$$+ \frac{1}{\sigma^4 L^2 + 3P^2}\left\|\sum_{k=1}^{K} \tilde{w}[k]B_{\tilde{x}_k} + M_1 \cdot \sigma^2 L^2 A - M_3\right\|_F^2. \tag{3.5}$$

We simplified the third-order term somewhat by substituting $M_1$ for $\sum_k \tilde{w}[k]\mu_{\tilde{x}_k}$. The cost function is smooth in all variables, but it is non-convex. We use Manopt [33] to optimize it from random initializations. This toolbox allows to turn the simplex $\Delta_K$ into a Riemannian manifold [34], then to run a trust-region algorithm. The cost and its gradient can be computed in $O(KL^2)$ flops—independent of $N$. Because of non-convexity, the algorithm could converge to suboptimal points. In Section 4, we observe that this is rarely the case in practice for a wide range of parameters.

Importantly, our approach relies only on invariant features up to third order—a concise summary of the data as soon as $N \gg L$—and these can be estimated accurately as long as $N/(\sigma^2 + \sigma^6)$ is large enough. This in turn implies that the noise levels can be arbitrarily high, provided sufficiently many measurements are available.

## 4. NUMERICAL EXPERIMENTS

We conduct a few numerical experiments solving (3.5). In all our experiments, signals $x_1,\ldots,x_K$ are generated with i.i.d. standard Gaussian entries. For two signals $x$ and $\tilde{x}$, we define a cyclic-shift invariant distance as

$$\text{dist}(x,\tilde{x}) = \min_{s \in \{0,\ldots,L-1\}} \|R_s x - \tilde{x}\|_2. \tag{4.1}$$

This is computed in $O(L \log L)$ flops with FFTs. An estimator $\tilde{\mathbf{x}} = (\tilde{x}_1,\ldots,\tilde{x}_K)$ for $\mathbf{x} = (x_1,\ldots,x_K)$ is defined up to ordering. Thus, we define the permutation invariant distance:

$$\text{dist}(\mathbf{x},\tilde{\mathbf{x}})^2 = \min_{\pi \in S_K} \sum_{k=1}^{K} \text{dist}(x_k, \tilde{x}_{\pi(k)})^2. \tag{4.2}$$

Optimization over $S_K$ (permutations over $K$ elements) is solved via the Hungarian algorithm in $O(K^3)$ operations. The relative estimation errors we report below are given by:

$$\text{relative\_error}(\mathbf{x}, \tilde{\mathbf{x}}) = \frac{\text{dist}(\mathbf{x}, \tilde{\mathbf{x}})}{\sqrt{\sum_{k=1}^{K} \|x_k\|_2^2}}. \qquad (4.3)$$

If the mixing probabilities $w$ are estimated by $\tilde{w}$, given an optimal permutation $\pi$ in (4.2), we report the estimation error as a total variation distance over the simplex:

$$\text{TV\_dist}(w, \tilde{w}) = \frac{1}{2} \sum_{k=1}^{K} |w_k - \tilde{w}_{\pi(k)}|. \qquad (4.4)$$

This value is between 0 and 1.

**Red dots** on Figures 4.1 and 4.2 mark upper bounds on how large $K$ may be as a function of $L$ for demixing to be possible. We reason as follows: for generic real signals, $M_1$ provides 1 real number; $M_2$ provides $\sim \frac{1}{2}L$ distinct real numbers, and $M_3$ provides $\sim \frac{1}{6}L^2$ distinct real numbers (separating real and imaginary parts, and accounting for symmetries; a precise accounting is used for the figures.) A total of $KL$ real numbers must be estimated to recover the signals, and possibly an additional $K - 1$ numbers are required to estimate mixing probabilities. Thus, even if all numbers provided by $M_1, M_2, M_3$ bear independent information, we still need $\sim \frac{1}{6}L^2$ to exceed $KL$ or $KL + (K-1)$. Solving for $K$ yields the displayed bound. In both cases, for large $L$, the bounds behave like $\lceil L/6 \rceil$.

**Experiment 1** explores an infinite data regime, where the mixed invariant features $M_1, M_2, M_3$ are available exactly. The question is then: from these mixed features, can we recover the individual signals $x_1, \ldots, x_K$ by solving (3.5)? If so, how large can we allow $K$ to grow as a function of $L$? Two separate issues are involved: (a) can we solve (3.5) to global optimality despite non-convexity? And (b), do global optimizers coincide with the ground truth? For values of $K$ and $L$ on a grid, we generate ground truth signals once with i.i.d. standard Gaussian entries (the power spectrum $P$ is the constant $L$ in expectation), and their exact mixed features are computed, with uniform mixing probability $w$ (known to the algorithm). Then, for each pair $(L, K)$, we generate 30 random initial guesses for the algorithm and optimize. We declare a global optimum is found if the cost value drops below $10^{-16}$. (In all experiments, the cost function is scaled by $\frac{\sigma^4 L^2 + 3P^2}{2}$.) For each run where optimality is declared, we compute the relative estimation error according to (4.3) and report the worst one for that $(L, K)$—the worst one, because in practice we would not be able to distinguish between global optima. Figure 4.1 leads to the following empirical observation: (i) for $K$ up to approximately $\sqrt{L}$, the local optimization algorithm reliably finds a global optimizer, and it corresponds to the ground truth within numerical errors. There is also a regime where $K$ is so large that recovery is impossible, yet the optimization problem is easily solved. In that regime, the solution to the problem is not unique: the problem is ill-posed.

**Experiment 2** is the same as Experiment 1, except the mixing probabilities $w$ are now random and unknown to the algorithm. For values of $K$ and $L$ on a grid, we generate ground truth signals once, and a mixing probability $w$ as a vector whose entries are i.i.d. uniform in $[0, 1]$, normalized to be a probability density. Their exact mixed invariant moments are computed. Figure 4.2 suggests $w$ can also be recovered, although performance deteriorates.

**Experiment 3** investigates resilience of the algorithm to high levels of noise—see Figure 4.3. We demix $K = 2$ signals of length $L = 50$ from $N = 10^6$ observations. Mixing probabilities $w$ are uniform. Our algorithm initializes $\tilde{w}$ uniform, but optimizes for it
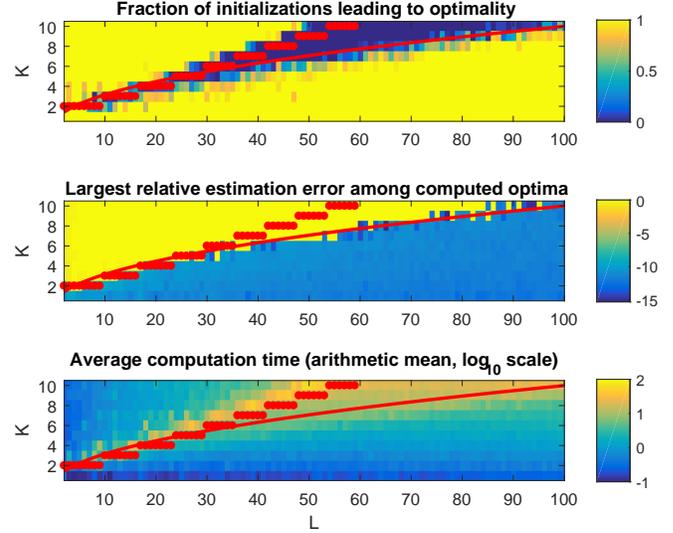


**Fig. 4.1**. Experiment 1 ($N \to \infty$) suggests $K$ up to $\sqrt{L}$ (red curve) i.i.d. Gaussian signals can be demixed from perfect mixed invariant moments with (3.5). CPU time in seconds. Above red dots, recovery is hopeless because of an information theoretic argument. All but first plot in $\log_{10}$ scale. $L$ ranges from 2 to 100, $K$ from 1 to 10.

as well. For each value of noise level $\sigma$ on a logarithmic grid, signals are generated 6 times and two methods are run: (a) our algorithm based on (3.5), where we run the method with two random initial guesses and return the best result (based on attained cost value, which is available in practice); and (b) Expectation–Maximization (EM)—see below. We find that both methods are resilient to high noise, with EM producing more accurate estimators but our method being orders of magnitude faster for large noise.

**EM** is a classic heuristic to address estimation problems with latent variables [35]. In a nutshell (see [10] for details in the homogeneous case), assuming a current estimate $\tilde{x}_1, \ldots, \tilde{x}_K$ is correct, EM computes $w_{j,r,k}$: the probability that measurement $y_j$ comes from signal $\tilde{x}_k$ with shift $r$. Concretely, under the Gaussian noise model, these probabilities are given by

$$w_{j,r,k} \propto \exp\left(-\frac{\|R_r \tilde{x}_k - y_j\|_2^2}{2\sigma^2}\right), \qquad (4.5)$$

and global scaling is fixed using $\sum_{r,k} w_{j,r,k} = 1$ for all $j$. Then, signal estimators are updated as follows: for each $k$,

$$\tilde{x}_k \leftarrow \frac{\sum_{j,r} w_{j,r,k} R_r^{-1} y_j}{\frac{\sigma^2}{\sigma_0^2} + \sum_{j,r} w_{j,r,k}}, \qquad (4.6)$$

where $\sigma_0 > 0$ comes from a prior on $x_k$'s coming from a distribution $\mathcal{N}(0, \sigma_0^2 I)$, where we pick $\sigma_0^2 = 1$ to match the true generating model, as a hint to EM. (Setting $\sigma_0^2 = 10^9$ to make the prior almost uninformative does not change the results much.) The probabilities and the updated estimators can be computed in $O(NKL \log L)$ flops using FFTs, parallelizable over $N$ and $K$—we use Matlab's built-in parallelization of FFTs over $N$. We iterate until two subsequent iterates differ by less than $K \cdot 10^{-5}$ in metric (4.2). (Results are robust against this choice.)

As Figure 4.3 reveals, the number of EM iterations grows with the noise level (for $\sigma = 10^{-1}$, as little as 3 iterations suffice, while it saturates at our limit of 10 000 for $\sigma = 10^1$). Attempts to reduce
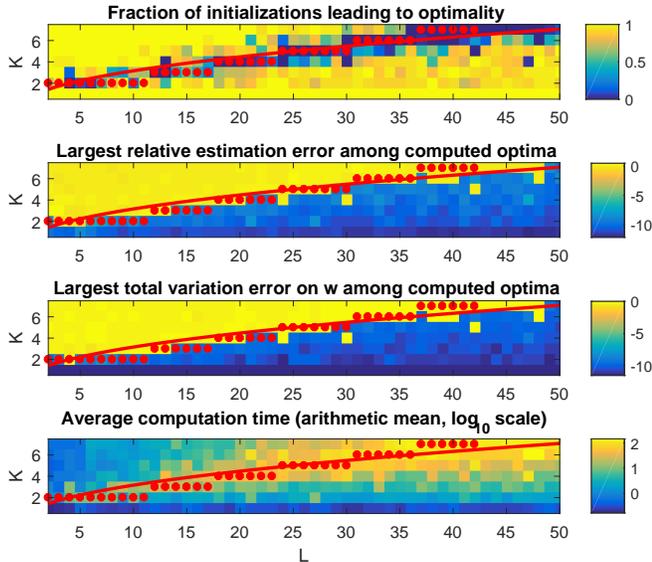
**Fig. 4.2**. Experiment 2 ($N \to \infty$) is the same as Experiment 1, except mixing probabilities $w$ are now random and unknown. For the middle plots, yellow pixels inside the blue area correspond to setups where the smallest entry in $w$ is small, which is challenging as the corresponding signal is under-represented. All but first plot in $\log_{10}$ scale. $L$ ranges from 2 to 50, $K$ from 1 to 7.

the strong effect of a large $N$ on the complexity of EM using batch iterations were not successful for this experiment. EM attains the most accurate estimators. One failure mode of EM is when the sum in the denominator of (4.6) (almost) vanishes for some signal: one of the estimators $\tilde{x}_k$ is given (almost) none of the observations, and this situation endures through iterations. Strangely, this occurs at high SNR, visible in Figure 4.3.

About parallelization: our method uses 30 cores to compute features $M_1, M_2, M_3$ in parallel over $N$—this is the bottleneck—then uses a single core to optimize from two random initial guesses sequentially (these could be done in parallel.) EM uses about 16 cores (number chosen by Matlab) to compute FFTs in parallel over $N$. Code: https://github.com/NicolasBoumal/HeterogeneousMRA.

## 5. CONCLUSIONS

We explored non-convex optimization for heterogeneous MRA, based on invariant features of order up to three. Key properties of our method are that (a) it never seeks to estimate the shifts or classes of the observations (which may be impossible to estimate at low SNR), and (b) once invariant features are computed (this takes linear time and is parallelizable in $N$), the complexity of the method no longer depends on the number of measurements $N$, which may need to be large for low SNR. Unfortunately, we are not able to theoretically explain the performance of the second stage at the moment.

In numerical experiments not shown here, with mixed moments known exactly (infinite data regime), we noticed that if we overestimate $K$ and allow the algorithm to optimize for $w$, then the algorithm still recovers the true signals, and assigns probability close to 0 to the extra estimated signals. If we underestimate $K$, then the algorithm tends to better estimate those signals which have larger probability $w[k]$. We intend to explore this further.
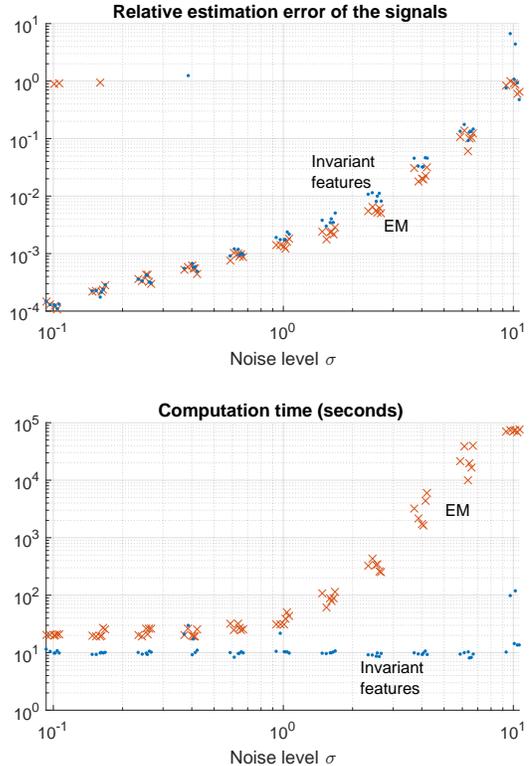


**Fig. 4.3**. Experiment 3 ($L = 50$, $K = 2$, $N = 10^6$) shows both our algorithm (blue dots) and EM (red crosses) are resilient to high noise levels. EM is more accurate and our method is orders of magnitude faster. Jitter is added on the $\sigma$ axis for visualization. At $\sigma = 10$, EM attains its iteration limit of $10^4$.

A key motivation for this work is the problem of heterogeneity in cryo-EM and in XFEL, both imaging techniques in structural biology. Heterogeneous MRA acts as a simplified model for those applications. Seminal work by Kam [27], who showed how moment-based approaches for cryo-EM without heterogeneity can work, suggests our findings here may translate to handle heterogeneity in cryo-EM and XFEL.

# Acknowledgments

## 6. REFERENCES

[1] R. Diamond, "On the multiple simultaneous superposition of molecular structures by rigid body transformations," *Protein Science*, vol. 1, no. 10, pp. 1279–1287, 1992.

[2] D. L. Theobald and P. A. Steindel, "Optimal simultaneous superpositioning of multiple structures with missing data," *Bioinformatics*, vol. 28, no. 15, pp. 1972–1979, 2012.

[3] W. Park, C. R. Midgett, D. R. Madden, and G. S. Chirikjian, "A stochastic kinematic model of class averaging in single-particle electron microscopy," *The International journal of robotics research*, vol. 30, no. 6, pp. 730–754, 2011.

[4] W. Park and G. S. Chirikjian, "An assembly automation approach to alignment of noncircular projections in electron microscopy," *IEEE Transactions on Automation Science and Engineering*, vol. 11, no. 3, pp. 668–679, 2014.

[5] J. P. Zwart, R. van der Heiden, S. Gelsema, and F. Groen, "Fast translation invariant classification of HRR range profiles in a zero phase representation," *IEE Proceedings-Radar, Sonar and Navigation*, vol. 150, no. 6, pp. 411–418, 2003.

[6] R. Gil-Pita, M. Rosa-Zurera, P. Jarabo-Amores, and F. López-Ferreras, "Using multilayer perceptrons to align high range resolution radar signals," in *International Conference on Artificial Neural Networks*, pp. 911–916, Springer, 2005.

[7] I. L. Dryden and K. V. Mardia, *Statistical shape analysis*, vol. 4. J. Wiley Chichester, 1998.

[8] H. Foroosh, J. B. Zerubia, and M. Berthod, "Extension of phase correlation to subpixel registration," *IEEE transactions on image processing*, vol. 11, no. 3, pp. 188–200, 2002.

[9] D. Robinson, S. Farsiu, and P. Milanfar, "Optimal registration of aliased images using variable projection with applications to super-resolution," *The Computer Journal*, vol. 52, no. 1, pp. 31–42, 2009.

[10] T. Bendory, N. Boumal, C. Ma, Z. Zhao, and A. Singer, "Bispectrum inversion with application to multireference alignment," *IEEE Transactions on Signal Processing*, vol. 66, no. 4, pp. 1037–1050, 2018.

[11] A. Bandeira, P. Rigollet, and J. Weed, "Optimal rates of estimation for multi-reference alignment," *arXiv preprint arXiv:1702.08546*, 2017.

[12] A. Perry, J. Weed, A. Bandeira, P. Rigollet, and A. Singer, "The sample complexity of multi-reference alignment," *arXiv preprint arXiv:1707.00943*, 2017.

[13] A. S. Bandeira, M. Charikar, A. Singer, and A. Zhu, "Multireference alignment using semidefinite programming," in *Proceedings of the 5th conference on Innovations in theoretical computer science*, pp. 459–470, ACM, 2014.

[14] E. Abbe, J. Pereira, and A. Singer, "Sample complexity of the boolean multireference alignment problem," to appear in *The IEEE International Symposium on Information Theory*, 2017.

[15] A. Bartesaghi, A. Merk, S. Banerjee, D. Matthies, X. Wu, J. L. Milne, and S. Subramaniam, "2.2 Å resolution cryo-EM structure of $\beta$-galactosidase in complex with a cell-permeant inhibitor," *Science*, vol. 348, no. 6239, pp. 1147–1151, 2015.

[16] B. W. McNeil and N. R. Thompson, "X-ray free-electron lasers," *Nature photonics*, vol. 4, no. 12, pp. 814–821, 2010.

[17] J. Frank, *Three-dimensional electron microscopy of macromolecular assemblies: visualization of biological molecules in their native state*. Oxford University Press, 2006.

[18] E. Nogales, "The development of cryo-EM into a mainstream structural biology technique," *Nature Methods*, vol. 13, no. 1, pp. 24–27, 2016.

[19] P. Schwander, R. Fung, and A. Ourmazd, "Conformations of macromolecules and their complexes from heterogeneous datasets," *Phil. Trans. R. Soc. B*, vol. 369, no. 1647, p. 20130567, 2014.

[20] J. Frank and A. Ourmazd, "Continuous changes in structure mapped by manifold embedding of single-particle data in cryo-EM," *Methods*, vol. 100, pp. 61–67, 2016.

[21] E. Katsevich, A. Katsevich, and A. Singer, "Covariance matrix estimation for the cryo-EM heterogeneity problem," *SIAM journal on imaging sciences*, vol. 8, no. 1, pp. 126–185, 2015.

[22] J. Andén, E. Katsevich, and A. Singer, "Covariance estimation using conjugate gradient for 3D classification in cryo-EM," in *Biomedical Imaging (ISBI), 2015 IEEE 12th International Symposium on*, pp. 200–204, IEEE, 2015.

[23] R. R. Lederman and A. Singer, "A representation theory perspective on simultaneous alignment and classification," *arXiv preprint arXiv:1607.03464*, 2016.

[24] Y. Aizenbud and Y. Shkolnisky, "A max-cut approach to heterogeneity in cryo-electron microscopy," *arXiv preprint arXiv:1609.01100*, 2016.

[25] R. R. Lederman and A. Singer, "Continuously heterogeneous hyper-objects in cryo-EM and 3-D movies of many temporal dimensions," *arXiv preprint arXiv:1704.02899*, 2017.

[26] B. M. Sadler and G. B. Giannakis, "Shift- and rotation-invariant object reconstruction using the bispectrum," *JOSA A*, vol. 9, no. 1, pp. 57–69, 1992.

[27] Z. Kam, "The reconstruction of structure from electron micrographs of randomly oriented particles," *Journal of Theoretical Biology*, vol. 82, no. 1, pp. 15–39, 1980.

[28] R. Marabini and J. M. Carazo, "On a new computationally fast image invariant based on bispectral projections," *Pattern Recognition Letters*, vol. 17, no. 9, pp. 959–967, 1996.

[29] D. Hsu and S. Kakade, "Learning mixtures of spherical Gaussians: moment methods and spectral decompositions," in *Proceedings of the 4th conference on Innovations in Theoretical Computer Science (ITCS)*, pp. 11–20, ACM, 2013.

[30] A. Anandkumar, R. Ge, D. J. Hsu, S. M. Kakade, and M. Telgarsky, "Tensor decompositions for learning latent variable models," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 2773–2832, 2014.

[31] T. Strohmer and K. Wei, "Painless breakups–efficient demixing of low rank matrices," *Journal of Fourier Analysis and Applications*, Sep 2017.

[32] C. Aguerrebere, M. Delbracio, A. Bartesaghi, and G. Sapiro, "Fundamental limits in multi-image alignment," *IEEE Transactions on Signal Processing*, vol. 64, no. 21, pp. 5707–5722, 2016.

[33] N. Boumal, B. Mishra, P.-A. Absil, and R. Sepulchre, "Manopt, a Matlab toolbox for optimization on manifolds," *Journal of Machine Learning Research*, vol. 15, pp. 1455–1459, 2014.

[34] Y. Sun, J. Gao, X. Hong, B. Mishra, and B. Yin, "Heterogeneous tensor decomposition for clustering via manifold optimization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 3, pp. 476–489, 2016.

[35] A. Dempster, N. Laird, and D. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of the royal statistical society. Series B (methodological)*, vol. 39, no. 1, pp. 1–38, 1977.

## A. WEIGHING THE COST FUNCTION

We pick the weights in the cost function of problem (3.5) as follows. Along the way, we make a number of simplifying assumptions to keep formulas straightforward. At the ground truth signal $\mathbf{x} = (x_1, \ldots, x_K) \in (\mathbb{R}^L)^K$ and mixing probabilities $w \in \Delta_K$, the error variables are as follows:

$$E_1 = M_1 - \sum_{k=1}^K w[k]\mu_{x_k}, \quad E_2 = M_2 - \sigma^2 L \mathbf{1} - \sum_{k=1}^K w[k]P_{x_k},$$

$$E_3 = M_3 - M_1 \cdot \sigma^2 L^2 A - \sum_{k=1}^K w[k]B_{x_k}.$$

By construction, they have zero mean. If all error terms were (entry-wise) independent and Gaussian (neither is true[1]), then minimizing the sum of squared errors normalized by their individual variances would yield a maximum likelihood-type estimator. This motivates us to normalize by (approximate) variances, as follows. If $y = x_v + \varepsilon$ (ignoring the shift $R_r$ since the features are invariant under it) with $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_L)$ and $v \sim w$ (independent), then by independence of the measurements,

$$E_1 = \frac{1}{N}\sum_{j=1}^N \mu_{y_j} - \sum_{k=1}^K w[k]\mu_{x_k}$$

$$= \frac{1}{N}\sum_{j=1}^N \mu_{\varepsilon_j} + \left[\frac{1}{N}\sum_{j=1}^N \mu_{x_{v_j}} - \sum_{k=1}^K w[k]\mu_{x_k}\right].$$

The bracketed term is zero for homogeneous MRA. We neglect it for the heterogeneous case. Then, by independence,

$$\mathrm{Var}\{E_1\} \approx \frac{1}{N}\mathrm{Var}\{\mu_\varepsilon\} = \frac{\sigma^2}{NL}.$$

Proceeding similarly for $E_2$ we first get

$$E_2[k] = \frac{1}{N}\sum_{j=1}^N |\hat{y}_j[k]|^2 - \sigma^2 L - \sum_{k'=1}^K w[k']|\hat{x}_{k'}[k]|^2$$

$$= \frac{1}{N}\sum_{j=1}^N \left[|\hat{x}_{v_j}[k]|^2 + |\hat{\varepsilon}_j[k]|^2 + 2\Re\{\overline{\hat{x}_{v_j}[k]}\hat{\varepsilon}_j[k]\}\right]$$

$$- \sigma^2 L - \sum_{k'=1}^K w[k']|\hat{x}_{k'}[k]|^2$$

$$\approx \frac{1}{N}\sum_{j=1}^N \left[|\hat{\varepsilon}_j[k]|^2 - \sigma^2 L + 2\Re\{\overline{\hat{x}_{v_j}[k]}\hat{\varepsilon}_j[k]\}\right],$$

where again we neglected a term which vanishes exactly in the homogeneous case. The terms $|\hat{\varepsilon}_j[k]|^2$ and $\Re\{\overline{\hat{x}_{v_j}[k]}\hat{\varepsilon}_j[k]\}$ are uncorrelated because $\hat{\varepsilon}_j[k]$ is distributed isotropically in the complex plane (in particular, the phase is uniform.) Thus, we can separate the variance computation in two parts:

$$\mathrm{Var}\{E_2[k]\} \approx \frac{1}{N}\mathrm{Var}\{|\hat{\varepsilon}[k]|^2\} + \frac{1}{N}\mathrm{Var}\{2\Re\{\overline{\hat{x}_v[k]}\hat{\varepsilon}[k]\}\}.$$

We can easily understand the distribution of $\hat{\varepsilon}$. Indeed, let $F$ be the DFT matrix so that $\hat{\varepsilon} = F\varepsilon$ for $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_L)$. Noting that $FF^* = LI_L$, we get: $\mathbb{E}\{\hat{\varepsilon}\hat{\varepsilon}^*\} = F\mathbb{E}\{\varepsilon\varepsilon^*\}F^* = \sigma^2 LI_L$. Thus,

---

[1] In particular, certain entries of the power spectrum and the bispectrum are repeated.

$\hat{\varepsilon} \sim \mathbb{CN}(0, L\sigma^2 I_L)$. Let $\hat{\varepsilon}[k] = z_1 + iz_2$, with $z_1, z_2 \sim \mathcal{N}(0, \frac{\sigma^2 L}{2})$. Then,

$$\mathrm{Var}\{|\hat{\varepsilon}[k]|^2\} = \mathrm{Var}\{z_1^2 + z_2^2\} = \sigma^4 L^2.$$

On the other hand, using $\mathbb{E}\{\hat{\varepsilon}[k]^2\} = 0$ due to uniform phase again and independence of $\varepsilon$ and $v$,

$$\mathrm{Var}\{2\Re\{\overline{\hat{x}_v[k]}\hat{\varepsilon}[k]\}\} = \mathbb{E}\{(\overline{\hat{x}_v[k]}\hat{\varepsilon}[k] + \hat{x}_v[k]\overline{\hat{\varepsilon}[k]})^2\}$$

$$= 2\mathbb{E}\{|\hat{x}_v[k]|^2\}\mathbb{E}\{|\hat{\varepsilon}[k]|^2\} = 2\sigma^2 L\sum_{k'=1}^K w[k']|\hat{x}_{k'}[k]|^2.$$

Overall,

$$\mathrm{Var}\{E_2[k]\} \approx \frac{\sigma^2 L}{N}\left(\sigma^2 L + 2\sum_{k'=1}^K w[k']P_{x_{k'}}[k]\right).$$

The sum is nothing but the mixed power spectrum, which we can estimate from the data: this could be used as weight directly. To simplify even further, assuming the power spectra of the signals are not too far from the constant $P$ (in the experiments, $x$ has i.i.d. standard entries, so the power spectrum is close to $L$), we can approximate the variance as a constant:

$$\mathrm{Var}\{E_2[k]\} \approx \frac{\sigma^2 L}{N}\left(\sigma^2 L + 2P\right).$$

We now turn to the bispectrum: each measurement $y_j$ contributes eight terms to $E_3[k, \ell]$ through $\hat{y}_j[k]\overline{\hat{y}_j[\ell]}\hat{y}_j[\ell - k]$ and $\hat{y}_j = \hat{x}_{v_j} + \hat{\varepsilon}_j[k]$ (again, ignoring the shift $R_{r_j}$ since the bispectrum is invariant under it.) In approximating the variance, we aim to identify the leading terms in $\sigma$ for low and for high SNR. As for $E_1$ and $E_2$, a term independent of $\sigma$ which vanishes exactly in the homogeneous case is neglected here. Thus, it remains to identify the terms which scale as $\sigma^2$ and $\sigma^6$. These come from:

$$\mathrm{Var}\{\hat{\varepsilon}[k]\overline{\hat{\varepsilon}[\ell]}\hat{\varepsilon}[\ell - k]\}$$

for $\sigma^6$ and from

$$\mathrm{Var}\{\hat{x}_v[k]\overline{\hat{x}_v[\ell]}\hat{\varepsilon}[\ell - k]\},$$
$$\mathrm{Var}\{\hat{x}_v[k]\overline{\hat{\varepsilon}[\ell]}\hat{x}_v[\ell - k]\},$$
$$\mathrm{Var}\{\hat{\varepsilon}[k]\overline{\hat{x}_v[\ell]}\hat{x}_v[\ell - k]\}$$

for $\sigma^2$. Aiming for a crude approximation, for the $\sigma^6$ term we consider only the case where $\varepsilon[k], \varepsilon[\ell], \varepsilon[\ell - k]$ are independent, in which case

$$\mathrm{Var}\{\hat{\varepsilon}[k]\overline{\hat{\varepsilon}[\ell]}\hat{\varepsilon}[\ell - k]\} = \mathbb{E}\{|\hat{\varepsilon}[k]|^2\}^3 = (\sigma^2 L)^3.$$

(For values of $k, \ell$ where independence does not hold, an extra constant would appear.) For the $\sigma^2$ terms, the first one expands as

$$\mathrm{Var}\{\hat{x}_v[k]\overline{\hat{x}_v[\ell]}\hat{\varepsilon}[\ell - k]\} = \mathbb{E}\{|\hat{x}_v[k]|^2|\hat{x}_v[\ell]|^2\}\mathbb{E}\{|\hat{\varepsilon}[\ell - k]|^2\}.$$

The first expectation (over $v$) can be estimated from the data: it is a mixture of fourth order moments. Unfortunately, estimating it accurately would require $O(\sigma^8)$ observations. Alternatively, for the homogeneous case, it could be approximated using the estimated power spectrum. Simpler still, as we do here, assuming the power spectra of the signals to be estimated are close to $P$, we approximate

$$\mathrm{Var}\{\hat{x}_v[k]\overline{\hat{x}_v[\ell]}\hat{\varepsilon}[\ell - k]\} \approx \sigma^2 LP^2.$$

There are three such terms, so that overall we get the approximation:

$$\mathrm{Var}\{E_3[k, \ell]\} \approx \frac{\sigma^2 L}{N}\left(\sigma^4 L^2 + 3P^2\right).$$