1 We would like to thank the reviewers for their insightful comments. Responses below:

2 **Common Points**:

3 • We want to emphasize that our work is primarily theoretical, building on [27] (published in IEEE Transaction
4 on Information Theory) by improving sample complexity bounds and more importantly eliminating the need
5 for strong assumptions. We believe this significantly enhances our understanding of this fundamental problem.

6 • Given this, we feel that experiments do not add much to the paper. Synthetic experiments will simply confirm
7 our theorems. We urge the reviewers to evaluate the paper on the basis of the theoretical results and techniques.

8 • We also note that many published NeurIPS do not contain empirical evaluations, and this is not a barrier to
9 publication or even NeurIPS awards. For example one of last year's best paper awards went to *Nearly Tight*
10 *Sample Complexity Bounds for Learning Mixtures of Gaussians via Sample Compression Schemes*, a purely
11 theoretical contribution, on a problem quite close to what we study here.

12 **Reviewer 1:**

13 • Thanks for pointing out the relevant papers. The second paper is quite relevant as it studies mixture model
14 parameter recovery in a compressive sensing framework. We will add the citations.

15 • An illustration is a great suggestion! We will add one to the final version to help explain our algorithms.

16 **Reviewer 2:**

17 • Exponential SNR is common in statistical problems involving mixtures. Examples include trace reconstruction
18 and learning binomial mixture where the dependence is known to be optimal. The exponential dependence
19 in our result arises from a reduction to learning Gaussian mixtures. This latter problem is extremely well
20 studied (and out of scope for our paper), and while we are not aware of better results for Gaussian mixtures,
21 any improvement would immediately yield improvements for our setting.

22 • As in [27], we assume $L$ is a constant. We can extract the dependence on $L$ from our proof for the final version,
23 but note that it is fairly complicated as $L$ is implicit in the definitions of the constants in Eqs (A), (B), (C).

24 • $\epsilon$-precision means that all coordinates are integer-multiples of $\epsilon$ (this assumption is also used in [27]). We can
25 add the definition to the camera ready.

26 • The set of responses recovered after using 2k rows of the Vandermonde matrix is provably unique, and it is
27 possible to recover the unknown vector uniquely using efficient decoding algorithms borrowing from literature
28 of coding theory (such as, Berlekamp-Welch Decoder, see, Arora and Barak, Computational Complexity,
29 Sec. 19.3).

30 • There are at most $L$ unique responses (fixed sensing vector, only $L$ possible choices for $\beta$), and by a coupon
31 collector phenomenon $L \log(Lk^2)$ measurements suffices to ensure there are exactly $L$, with high probability.

32 • The interval must belong to the set of positive real numbers so that we can apply Lemma 1. If the difference
33 polynomials have no roots in an interval, then by continuity, the ordering of the polynomials $f^1, f^2, \ldots, f^L$
34 does not change in that interval. Thus the ordering is consistent.

35 • We will add the definition of the minimum distance estimator to the final version.

36 • We will clarify the paragraph after Lemma 2. In short, the $\beta^j$s are on an $\epsilon$-grid, but $\langle v, \beta^j \rangle$ may not be if we
37 naively run the noiseless algorithm. (The $\epsilon$-grid property is crucial for our mixture learning algorithm.) This
38 motivates the new decoding strategy, which ensures $\langle v, \beta^j \rangle$ is on the $\epsilon$-grid.

39 **Reviewer 3:**

40 • You are right that Vandermonde matrices are ill-conditioned – this is precisely why we do not use them for the
41 noisy case. That is also the point of having different treatments for noiseless and noisy cases, with the later
42 being more involved. This is similar to classical results in compressed sensing.

43 • Due to space restrictions we cannot include the entire proof in the body. But all the details are in the appendix.

44 • The constants are simply artifacts of the proof.

45 • Note that every triple (good or bad) contains vectors with integral entries. Thus we can learn the parameters
46 $\langle v, \beta^j \rangle$ for the three vectors. Then to check if the condition holds, note that for fixed $\beta$ and our choice of
47 vectors, we have $\langle v_1, \beta \rangle + \langle v_2, \beta \rangle = \langle v_3, \beta \rangle$, so we can match up the learned parameters. If we find two such
48 identities involving a single parameter, we know that triplet is not good.