1  Reviewers remark our method is intuitive and correct, and opens new directions in sparse clustering, while R1 raised
2  a concern about the extent of our contributions. Most comments mention the paper is well written while providing
3  constructive suggestions to further improve the presentation. We thank the referees for their time and feedback, and
4  provide detailed responses below:

5  R1:
6  We thank you for commenting the paper is well-written and for finding a typo. We hope you might reconsider the
7  novelty of our work if you find these responses to sufficiently address your concerns:
8  − You suggest better baselines for comparison, citing Power $k$-means [37] and matrix + tensor factorization. Our
9  manuscript *does* cite [37], and does compare with this method—see Table 3. Because power $k$-means is not designed
10  for the high-dimensional setting, it performed nearly identically with $k$-means in the simulated experiments where we
11  know the ground truth is sparse, and was omitted due to redundancy. We are happy to include its performance in those
12  comparisons as well, which will convey the same trends.
13  − Your intuition for low-rank matrix factorization is spot on, but low-rank factorization plus the constraint that one
14  seeks "hard" label assignments becomes **equivalent** to $k$-means, which we consider here. See for instance "$k$-means
15  Clustering Is Matrix Factorization" (Bauckhage 2015). Note that the number of clusters $k$ is analogous to the rank
16  $k$ of the low-rank factor; it remains nontrivial to perform feature selection jointly (i.e. simultaneously seek a sparse
17  number $s$ of informative features) as we do in the proposed method. In light of your comment, we will emphasize the
18  connection to matrix factorization in the revision. We do not compare to tensor factorization as all data we consider are
19  vector-valued and not matrix-valued.
20  − We respectfully disagree that the contribution is incremental, as the ranking-based feature selection is a marked
21  departure from the existing efforts which largely either rely on generic dimension reduction as a pre-processing step, or
22  penalization via norm-shrinkage. Instead, the proposed method allows for the exact desired number of nonzero features
23  to be specified as input, and yields a scalable approach that is appropriate in high-dimensional settings yet comparable
24  to Lloyd's algorithm in terms of simplicity and speed. As a result, our algorithm looks and functions quite differently
25  than past work on sparse $k$-means, which we have reviewed and compared to our method.

26  R2:
27  We thank you for your detailed comments and careful reading of the paper. We will further elaborate on the choice of
28  within-cluster sum of squares score as suggested during the revision.

29  R3:
30  We agree that the benchmarks and simulation details can be better described and will provide complete details in
31  the revision. You raise a good point about further competing methods– regarding COSA, we were unable to find the
32  authors' implementation and had implemented our own version whose runtime and performance was far worse than the
33  proposed method. We should also note that COSA was designed for feature weighing rather than selection, and does
34  not typically result in sparse solutions, though we will attempt to add a fair, detailed comparison in the revision. We
35  also note your suggestion regarding filter methods that focus only on feature selection and will include a comparison in
36  the final version.

37  R4:
38  You are absolutely correct that "interpretable sparsity" is overloaded here. We will clarify in the revision its twofold
39  meaning: first as you've mentioned, we can inform or "control" the sparsity level via parameter $s$. In this sense the
40  parameter $s$ is directly interpretable compared to parameters such as $\lambda$ in existing $\ell_1$ approaches. Second, our ranking
41  method selects features among the original dimensions, thus allowing them to retain their original interpretation. For
42  instance, in our mouse protein study it is important that the top ranked features identified by SKFR identify the most
43  relevant genes, as the original features correspond to expression levels along a high-dimensional space of candidate
44  genes. In this sense the *dimension reduction* is interpretable, in contrast to generic dimension reduction such as PCA,
45  where the axes (principal components) in the projected space lose their interpretation as genes. We will improve the
46  exposition to emphasize this, as well as provide further detail in a comparison table with sparse $k$-means focusing on
47  selection in the Supplement of the final draft, and thank you for these constructive comments.