LORE: JOINTLY LEARNING THE INTRINSIC DIMENSIONALITY AND RELATIVE SIMILARITY STRUCTURE FROM ORDINAL DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

Learning the intrinsic dimensionality of subjective perceptual spaces such as taste, smell, or aesthetics from ordinal data is a challenging problem. We introduce LORE (Low Rank Ordinal Embedding), a scalable framework that jointly learns both the intrinsic dimensionality and an ordinal embedding from noisy triplet comparisons of the form, "Is A more similar to B than C?". Unlike existing methods that require the embedding dimension to be set apriori, LORE regularizes the solution using the nonconvex Schatten-p quasi norm, enabling automatic joint recovery of both the ordinal embedding and its dimensionality. We optimize this joint objective via an iteratively reweighted algorithm and establish convergence guarantees. Extensive experiments on synthetic datasets, simulated perceptual spaces, and real world crowdsourced ordinal judgements show that LORE learns compact, interpretable and highly accurate low dimensional embeddings that recover the latent geometry of subjective percepts. By simultaneously inferring both the intrinsic dimensionality and ordinal embeddings, LORE enables more interpretable and data efficient perceptual modeling in psychophysics and opens new directions for scalable discovery of low dimensional structure from ordinal data in machine learning.

1 Introduction

Learning subjective percepts (SPs), such as taste, smell, or aesthetic preference, poses unique challenges for machine learning. Traditional approaches rely on absolute queries that presuppose known perceptual axes. For example, a taste study might ask participants to rate stimuli on a 1-5 Likert scale (Likert, 1932) for "sweetness" or "bitterness". Such methods suffer from two critical flaws: (1) inconsistency, as respondents interpret scales differently (e.g., one person's "moderately sweet" is another's "very sweet") (Stewart et al., 2005), and (2) predefined conceptual frameworks that limit discovery by forcing ratings on predefined axes. Consequently, researchers risk missing latent dimensions (e.g., a "metallic" undertone in coffee) that participants lack vocabulary to describe.

In contrast, relative queries circumvent these issues by capturing perceptual relationships directly. For example, a triplet comparison like "Is coffee A more similar to coffee B or coffee C in taste?" allows participants to express nuanced judgments without relying on language or preset scales. Such relative comparisons are therefore particularly well suited for discovering the latent dimensions that organize subjective perceptual spaces.

Relative Similarity or Ordinal Embedding methods (OE) leverage these relative judgements to learn a multidimensional representation. However, all existing OE approaches require the user to specify the embedding dimension in advance (Agarwal et al., 2007; Jain et al., 2016; Tamuz et al., 2011; Terada and Luxburg, 2014; Van Der Maaten and Weinberger, 2012), with little guidance to the "true" complexity of the perceptual space. In practice, this can lead to unnecessarily high dimensional embeddings, concealing the actual structure. For instance, an OE may perfectly satisfy all triplet constraints in a 10-dimensional space, even if the underlying percept is only 2-dimensional.

Scientific discovery demands parsimony, a principle formalized as Occam's razor (Bishop and Nasrabadi, 2006). For the taste example, a 2D embedding is preferable to a 10D alternative: it is easier to interpret, less computationally intensive, and more useful for downstream analyses. In

078

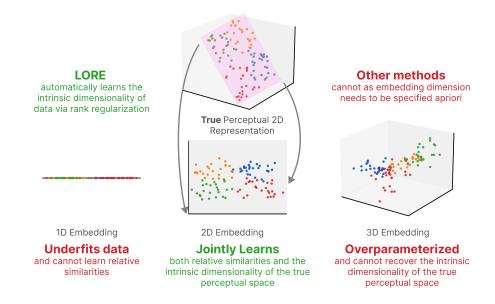


Figure 1: LORE jointly learns both the intrinsic dimensionality and relative similarities by balancing dimensionality with similarity constraints.: Other methods require the embedding dimension to be chosen in advance, making them less data driven and often suboptimal.

practical terms, a 10D taste embedding might fragment "sweetness" into several axes, complicating flavor design or neurological interpretation. Yet, most OE approaches, despite high triplet accuracy, produce overly complex models that mask the true structure of the latent percept.

To address this gap, We introduce LORE, a new ordinal embedding algorithm that jointly learns both the embedding and the intrinsic dimensionality, instead of needing to specify the dimension apriori. LORE regularizes using the nonconvex Schatten-p quasi-norm, explicitly balancing triplet accuracy with representation compactness and is optimized via an iteratively reweighted algorithm, with guarantees of convergence to stationary points. Our main contributions are:

- 1. LORE, a novel ordinal embedding algorithm that recovers latent representations that match the intrinsic dimensionality of human perceptual similarity data. LORE, jointly infers both the embedding and its dimensionality by regularizing with the nonconvex Schatten-p quasi-norm. By balancing triplet accuracy and rank regularization we can infer a compact yet accurate representation. We optimize the resulting objective using an iteratively reweighted schattent quasi norm algorithm, and provide convergence guarantees of the OE to stationary points.
- 2. LORE reliably uncovers the intrinsic dimensionality of data through an extensive evaluation where the dimensionality is known apriori. We first extensively test our algorithm on data with various dimensionality, noise levels and number of queries and demonstrate it outperforms existing methods by far in estimating intrinsic dimensionality with close to optimal performance in triplet accuracy. Secondly, we conduct a simulated perceptual experiment to model taste using an LLM as the ground truth perceptual space we try to model. See Figure 1 for a high level summary of our results.
- 3. LORE outperforms numerous state of the art methods on the large crowd sourced datasets and learns semantically interpretable axes. We find that LORE achieves a lower rank representation compared to all baselines while achieving comparable triplet accuracy on three separate crowdsourced datasets (Ellis et al., 2002; Kleindessner and Von Luxburg, 2017; Wilber et al., 2014) and learns axes which are semantically interpretable.

We anticipate LORE will be a valuable tool for mapping subtle subjective phenomena to interpretable low-dimensional spaces across psychology, neuroscience, and social science. By removing the need to hand tune embedding dimension, LORE enables data-driven discovery of subjective percepts.

2 Related Work

Ordinal Embeddings as Tools for Psychophysical Scaling: Psychophysics aims to discover the quantitative mappings that humans use to connect external stimuli to inner perceptual experiences. Psychological percepts are usually studied via *relative judgements* as they are less prone to individual biases, scale interpretation and memory limitations than absolute judgements as humans do not perceive stimuli in isolation (Stewart et al., 2005). Given the constraints of data collection, a core challenge in psychophysics is reconstructing perceptual spaces from a few human similarity judgments. OEs address this challenge; they are both query efficient and capable of reconstructing multidimensional perceptual spaces. For example, (Filip et al., 2024) derived a tactile-visual embedding for wood textures, identifying roughness and gloss as perceptually orthogonal dimensions. Moreover, active learning approaches like (Canal et al., 2020) have demonstrated how query efficiency in data collection can be further improved.

Metric Learning/Contrastive Learning are distinct from OEs: Learning from relative comparisons has been used in metric learning (Suárez-Díaz et al., 2018) and contrastive learning (Chen et al., 2020). Metric learning aims to learn a metric space from the data while contrastive learning separates similar and dissimilar datapoints. Metric Learning typically combines relative judgements and explicit representations (say images). The goal is to learn a distance metric from both sources of information. This additional representation, absent in OEs, changes the optimization problem and prevents direct transfer of metric learning approaches. Contrastive learning seeks to group similar datapoints together and push dissimilar ones apart with the presence of explicit additional supervised information which do not exist for OEs. Therefore, while metric learning and contrastive learning methods are similar in learning from relative information, they cannot be directly applied to OEs.

Intrinsic Dimensionality recovery is critical for psychophysics: A core goal in psychophysics is to recover the latent internal representations that individuals use to perceive psychophysical stimuli. Each representation is composed of two important characteristics: how well the representation recovers the ordinal relationships between the percepts and the *intrinsic rank or dimensionality* of the representation obtained. While OEs are able to maintain ordinal consistency (Vankadara et al., 2023), they are unable to identify the intrinsic dimensionality as we show in this paper. This is a key limitation of OEs that reduces their utility for psychophysical analysis. (Künstle et al., 2022) addressed this problem by modelling it as a multiple hypothesis test with separate embeddings trained for each candidate dimension and triplet accuracies used to estimate the true intrinsic rank. This approach, however, has two main limitations:

- Hypothesis dependence: It requires predefining plausible dimensionalities, risking model
 misspecification and reducing statistical power if the true dimensionality exceeds hypothesized bounds.
- 2. Lack of Scalability: Training multiple embeddings for each hypothesized rank is computationally expensive and quickly becomes prohibitive for a greater number of percepts. This is especially problematic for active querying where efficiency is critically important.

Building on these limitations, we propose a method to **jointly infer both dimensionality and multidimensional representations** via a novel OE method, eliminating the need for explicit hypothesis enumeration. For psychophysics, this enables recovery of perceptual geometry without prior assumptions on dimensionality. For machine learning, it offers a scalable approach to uncovering low dimensional structure directly from ordinal data.

3 Background on Ordinal Embeddings

The ordinal embedding problem seeks to learn an embedding matrix $Z \in \mathbb{R}^{N \times d'}$ from triplet judgements from the true perceptual space lying in an unknown $P \in \mathbb{R}^{N \times d}$ where $d \ll N$ is the intrinsic dimensionality or the intrinsic rank of the perceptual space. Z is learned indirectly via noisy similarity triplet comparisons where the anchor percept a is more similar or closer in the perceptual space to percept i than percept j into an embedding space of dimension d'. Specifically, this is denoted by $(a,i,j)=t\in T$ where $d(P_{a,:},P_{i,:})< d(P_{a,:},P_{j,:})$ where $P_{a,:},P_{i,:},P_{j,:}$ are the rows indexed by percepts a,i,j respectively in P and d(.,.) is the Euclidean distance between the unknown percepts. A central challenge is that intrinsic rank is unknown and the embedding dimension is set heuristically.

163

166

167168

170171

172

173

174175176

177

178

179

181

182

183

184

185

186

187

188

189

190

191

192

193 194

195

196 197

198 199

200

203

204

205

211

213

214

Optimizes Over Recovers Rank Method **Scalability Good Triplet Accuracy GNMDS** Gram Matrix X CKL Gram Matrix X **FORTE** Gram Matrix X t-STE Embedding X **SOE** Embedding **OENN** Embedding X LORE (ours) Embedding

Table 1: Characterization of Different Ordinal Embedding Algorithms

Though this framework is relatively simple, solving OEs efficiently can be challenging as the OE problem is NP-Hard (Bower et al., 2018), most loss functions are nonconvex and efficient learning demands at least $\mathcal{O}(Nd \log N)$ actively sampled triplets (Jain et al., 2016). As a result, the choice of optimization framework is crucial to obtaining a good OE.

Gram matrix approaches optimize a positive semidefinite matrix $G = ZZ^T \in \mathbb{R}^{N \times N}$ that capture the pairwise differences. While theoretically appealing because they are agnostic to the embedding dimension during optimization, they require enforcing PSD constraints that are not scalable for large N. Early methods like Generalized Non Metric Multi Dimensional Scaling (GNMDS) (Agarwal et al., 2007) and probabilistic models like Crowd Kernel Learning (CKL) suffer from limited accuracy or poor scalability. Fast Ordinal Triplet Embedding (FORTE) accelerates this with a kernelized nonconvex triplet loss optimized by efficient Projected Gradient Descent (PGD) and line search.

Direct embedding approaches optimize Z which leads to faster gradient updates that scale with the smaller $\mathcal{O}(Nd')$ versus $\mathcal{O}(N^2)$ with Gram matrix approaches. Examples include t-distributed Stochastic Triplet Embedding (t-STE) (Van Der Maaten and Weinberger, 2012) and Soft Ordinal Embedding (SOE) (Terada and Luxburg, 2014) with the latter widely used for its efficiency and high accuracy. A deep learning variant, Ordinal Embedding Neural Network (OENN) (Vankadara et al., 2023) underperforms likely due to the limited supervisory signal in purely ordinal data.

However, a shared fundamental limitation of all existing methods is the inability to recover the intrinsic rank d, which risks overparameterizing the true perceptual latent space.

4 Methods

We introduce a scalable ordinal embedding (OE) framework that jointly learns both the embedding and the intrinsic rank of the perceptual space. To ensure computational efficiency on large datasets (large T and N) we directly optimize the embedding Z instead of the Gram matrix G. The key insight is that we want the learning algorithm to adaptively the select embedding dimensionality as needed to fit the percepts well but not use any more extra space than necessary. Therefore, a natural approach is to penalize the rank of the learned embedding via regularization.

As SOE has the best properties of all the OEs that optimize over Z, we extend it with regularization. As the rank constraint is NP-Hard and non-convex (Fazel et al., 2001) a common approach is to regularize with the nuclear norm instead where $\|Z\|_* = \sum_{i=1}^{\min} {\{N,d'\} \atop i=1} \sigma_i(Z)$, where $\sigma_i(Z)$ is the ith singular value (Candes and Recht, 2008; Fazel et al., 2001). The objective then becomes:

$$\min_{\boldsymbol{Z}} \ \Psi(\boldsymbol{Z}) = \sum_{(a,i,j) \in T} \max \bigl\{ 0, 1 + d \bigl(\boldsymbol{Z}_{a,:}, \boldsymbol{Z}_{i,:} \bigr) - d \bigl(\boldsymbol{Z}_{a,:}, \boldsymbol{Z}_{j,:} \bigr) \bigr\} + \lambda \ \|\boldsymbol{Z}\|_*.$$

Though the nuclear norm is convex and relatively easy to optimize, it uniformly shrinks all of singular values (Negahban and Wainwright, 2011; Zhang, 2010). Recent theoretical and empirical evidence indicates that the nonconvex Schatten-p quasi norm $\|Z\|_p^p = \sum_{i=1}^{\min{\{N,d\}}} \sigma_i(Z)^p = \sum_{i=1}^{\min{\{N,d\}}} g[\sigma_i(Z)]$ for 0 , recovers the intrinsic rank for low rank recovery problems better than the nuclear norm can (Lu et al., 2014; Marjanovic and Solo, 2012). The Schatten-<math>p quasi norm

generalizes the nuclear norm by penalizing larger singular values less severely which is shown to aid in intrinsic rank recovery. We leverage this property and for the first time, to our knowledge, integrate the Schatten quasi-norm into a scalable ordinal embedding framework, allowing implicit perceptual rank discovery as seen below in

$$\min_{\boldsymbol{Z}} \ \Psi(\boldsymbol{Z}) = \sum_{(a,i,j) \in T} \max \bigl\{ 0, 1 + d \bigl(\boldsymbol{Z}_{a,:}, \boldsymbol{Z}_{i,:} \bigr) - d \bigl(\boldsymbol{Z}_{a,:}, \boldsymbol{Z}_{j,:} \bigr) \bigr\} + \lambda \ \|\boldsymbol{Z}\|_p^p.$$

Though incorporating the Schatten Quasi-Norm improves rank recovery properties, it also introduces additional nonconvexity into the regularizer that makes optimization more challenging. To overcome the inherent non-differentiability of the ordinal loss and the complexity of nonconvex regularization, we smooth the hinge triplet loss with the softplus function (Dugas et al., 2001). This transformation makes the objective differentiable except where the embedding collapses ($Z_{a,:} = Z_{i,:}$ or $Z_{a,:} = Z_{j,:}$). However, collapses can be avoided with wide initializations of Z. This smoothing enables provable convergence and is empirically essential, as it mitigates zero gradient plateaus to facilitate training on large datasets. Then the objective function is defined as

$$\min_{\boldsymbol{Z}} \ \Psi(\boldsymbol{Z}) = \sum_{(a,i,j) \in T} \log \left(1 + \exp \left(1 + d \left(\boldsymbol{Z}_{a,:}, \boldsymbol{Z}_{i,:}\right) - d \left(\boldsymbol{Z}_{a,:}, \boldsymbol{Z}_{j,:}\right)\right)\right)) + \sum_{i=1}^{\min\{N,d'\}} \sigma_i(\boldsymbol{Z})^p.$$

Despite smoothing the ordinal loss, our objective remains highly nonconvex due to the Schatten-p quasi-norm regularization, which makes reliable optimization difficult. Standard gradient methods often get stuck in poor local minima or fail to converge. To overcome this, we use an iteratively reweighted algorithm inspired by (Sun et al., 2017). At each step, the algorithm minimizes a weighted surrogate of the original objective, leading to steady improvement even in complex landscapes. As established in Theorem 1, this procedure is guaranteed to converge to a stationary point, ensuring robust and reliable learning.

Theorem (LORE converges to a stationary point) The sequence of OEs generated by the LORE algorithm $\{Z^k\}_{k=1,2,3,...}$ converges. i.e.

$$\sum_{k=1}^{+\infty} \lVert Z^{k+1} - Z^k \rVert_F < +\infty$$

Proof Sketch: We use the general framework for nonconvex Schatten Quasi-Norm optimization as seen in (Sun et al., 2017) but crucially, check the specific conditions for the LORE objective. The full proof is in Appendix A.

Our convergence guarantee is significant because, for ordinal embedding problems, stationary points are widely believed to be nearly as good as global optima in objective value. This is supported empirically (Vankadara et al., 2023) and theoretically. (Bower et al., 2018) proved that for certain OE settings with d=2, all local optima are global. Moreover, when sufficient triplet data is available, sub-optimal local minima are rarely observed. Building on these insights, we expect that our method will also recover high quality embeddings in realistic settings. Our experimental results confirm that LORE learns high accuracy ordinal embeddings, even with the inherent nonconvexity of the objective.

We implement the optimization using an efficient iteratively reweighted algorithm, seen in Algorithm 1, that updates the embedding and regularization at each step. In the typical regime where the embedding dimension is much smaller than the number of items and triplets, each iteration requires $\mathcal{O}(d'(T+Nd'))$ operations, making LORE scalable to large datasets. Additional implementation specifics are in Appendix B.

In summary, our methodological contributions are: (1) formulating a new ordinal embedding approach that jointly learns ordinal embeddins and intrinsic rank using Schatten quasi-norm regularization; (2) establishing an efficient optimization strategy based on iteratively reweighted minimization tailored for this nonconvex objective, along with convergence guarantees and (3) providing a scalable algorithm suitable for large scale perceptual similarity data.

5 Results

Algorithm 1: Learning LORE

270

271272

273

274

276

277 278

279

280

281

282

283

284

285

286

287

288

290

291

292

293

294

295

296

297

298

299 300

301 302

303

304

305

306 307

308 309

310

311

312

313

314

315 316

317

319

320

321

322

323

```
procedure LORE (\mathbf{Z}^0 \in \mathbb{R}^{\{N \times d'\}}, T, \lambda, \mu, \text{tol})
2:
            prev objs \leftarrow [\infty]
3:
            for k = 0, 1, 2, .... do
                \sigma \leftarrow \text{Singular Values}\left(oldsymbol{Z}^k
ight)
4:
                5:
                if |\text{curr_obj} - \text{prev_objs}[-1]| < \text{tol then}

⊳ Convergence check

6:
               egin{aligned} U, S, V^T &\leftarrow 	ext{SVD}\left(oldsymbol{Z}^k - rac{1}{\mu} 
abla_{oldsymbol{Z}^k} f(oldsymbol{Z}^k)
ight) \\ S^k &\leftarrow S - rac{p}{\mu} \sigma^{p-1} \end{aligned}
7:
8:
9:
                oldsymbol{S}^k \leftarrow 	ext{sorted} \left( oldsymbol{S}^k ig[ oldsymbol{S}^k > 0 ig], 	ext{ descending} 
ight) \ oldsymbol{Z}^{k+1} \leftarrow oldsymbol{U} oldsymbol{S}^k oldsymbol{V}^T
10:
11:
12:
                prev_objs[k] \leftarrow curr_obj
                if \|Z^{k+1} - Z^0\|_{\infty} < \text{tol then}
13:
                                                                                                       > Check if close to stationary point
14:
            return Z^{k+1}
15:
```

In this section, we present five pieces of empirical evidence in support of our method's claims. First, we outline our experimental setup, including baseline methods and the generative process for producing ordinal embedding (OE) tasks. Next, we show how to select regularization levels for LORE. We then benchmark LORE against standard baselines across key metrics, followed by a comparison on proxy large language model (LLM) generated perceptual spaces. Finally, we assess performance on real, crowdsourced triplet data involving human judgments and see that LORE's learned axes have semantic meaning. Collectively, these results demonstrate that LORE is uniquely effective at jointly learning high quality ordinal embeddings with intrinsic rank recovery, a property no other existing OE method can.

5.1 Setup

We benchmark LORE primarily against (1) SOE and (2) FORTE. These methods represent the best performing direct and Gram matrix OE approaches, respectively, as established by prior work (Vankadara et al., 2023). For our last experiment on crowdsourced data, which is computationally less demanding, we also compare against t-STE and CKL.

Our core evaluation criteria are:

- <u>Test Triplet Accuracy</u>: The proportion of held-out triplets correctly satisfied by the learned embedding and the primary metric in the OE literature.
- <u>Measured Rank</u>: The effective rank of the learned embedding, as a measure of intrinsic dimensionality recovery.

An ideal ordinal embedding should achieve high test triplet accuracy while maintaining a measured rank close to the true intrinsic rank of the underlying perceptual space.

We systematically vary four factors: fraction of queries, intrinsic rank, number of percepts, and noise level in the generative model. For synthetic experiments, we generate perceptual spaces of specified size and rank, sample noisy triplets to mimic human responses, and fit each method before evaluating on test triplets. Our synthetic data model generates a random perceptual space of specified rank and number of percepts, followed by sampling triplets with replacement to simulate human queries. We then use a standard approach (Canal et al., 2020; Vankadara et al., 2023) to model response uncertainty by sampling Gaussian noise independently and adding to each triplet distance. The resulting triplet data is then used to fit all OE algorithms, which are evaluated on held-out test triplets for both accuracy and measured rank. Unless otherwise stated, all experiments use query_fraction = 0.1, p = 0.5, d = 5, N = 50, noise = 0.1, 30 trials, and embedding dimension d' = 15.

5.2 LORE HAS WIDE AND STABLE REGULARIZATION SETTINGS

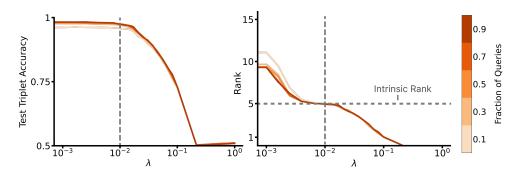


Figure 2: LORE has high test triplet accuracy and intrinsic rank recovery. (Left) Mean test triplet accuracy vs λ for LORE as Fraction of Queries varies. (Right) Mean measured rank vs λ for LORE as Fraction of Queries varies.

Our first set of experiments explores whether LORE admits a regularization regime that yields both high test triplet accuracy and reliable intrinsic rank recovery. As shown in Figure 2, across a broad range of the regularization parameter ($\lambda \approx 0.01$), LORE achieves nearly perfect test triplet accuracy and accurate intrinsic rank recovery, even as the fraction of queried triplets varies. Further results in Appendix C show that LORE performs similarly with varying noise and number of percepts. These pieces of evidence confirm that this high performance in the same robust regularization range persists with different noise levels, numbers of percepts and intrinsic rank. Thus, LORE is robust in hyperparameter selection for various dataset conditions.

5.3 LORE OUTPERFORMS BASELINES IN RANK RECOVERY; MATCHES IN TRIPLET ACCURACY

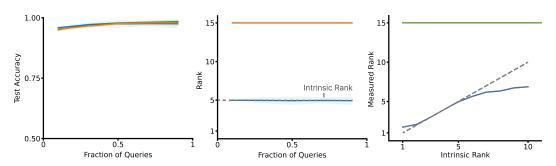


Figure 3: Only LORE can recover the intrinsic rank while maintaining comparable test triplet accuracy: (Left) Mean test triplet accuracy vs fraction of queries used. (Center) Mean measured rank vs fraction of queries used. (Right) Mean Measured Rank vs Intrinsic Rank. The gray dotted line indicates the ideal case where the measured rank is equal to the intrinsic rank.

Figure 3 (left, center) demonstrates that LORE uniquely recovers the true intrinsic rank of the embedding across all tested query fractions, while baseline methods consistently default to the maximum allowed dimension. Importantly, LORE matches the test triplet accuracy of the best baseline across all conditions, achieving low rank solutions without sacrificing predictive performance.

Additionally, Figure 3 (right) shows that as the true intrinsic rank increases, only LORE tracks this change, whereas all other methods ignore the underlying complexity and fail to adapt. While some loss in rank recovery is observed at higher true ranks (expected due to fixed number of triplets and the curse of dimensionality (Bishop and Nasrabadi, 2006)), LORE consistently outperforms competitors in recovering reduced the intrinsic rank. Further results in Appendix D confirm that LORE maintains an advantage across different noise and percept counts. Thus, LORE is the only method to reliably recover both accurate ordinal embeddings and the intrinsic rank. These results highlight the practical value of LORE for applications where discovering latent structure is critical.

5.4 LORE RECOVERS INTRINSIC RANK IN A SIMULATED PERCEPTUAL EXPERIMENT

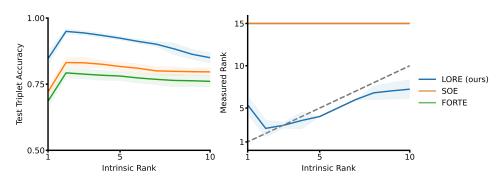


Figure 4: LORE outperforms baselines for both test triplet accuracy and intrinsic rank for a simulated LLM perceptual experiment. (Left) Mean test triplet accuracy vs intrinsic rank. (Right) Mean measured rank vs intrinsic rank. The gray dotted line is the ideal case where the measured rank is equal to the intrinsic rank. Shaded Areas indicate ± 2 Standard Deviations.

Human perceptual experiments are a key application of LORE, but for real datasets, the true intrinsic rank is unknown. To address this, we leverage recent findings that large language models (LLMs) encode human-aligned perceptual information across domains such as taste, pitch, and timbre (Marjieh et al., 2024). We therefore use an LLM embedding as a realistic proxy of the true perceptual space. In our experiment, we obtain SBERT embeddings (Reimers and Gurevych, 2019) for 50 randomly chosen foods, then restrict their dimensionality by applying truncated SVD (ranks 1-10). From this space, we generate noisy triplet comparisons (sampling 5% of the total, 30 repetitions per configuration, with noise 0.1) to mimic the data limited regime typical in human experiments.

Figure 4 summarizes the results. Even in this highly undersampled setting, LORE closely tracks the true intrinsic rank across all tested values, while all baselines default to the ambient dimension. Furthermore, LORE significantly outperforms baselines in test triplet accuracy, demonstrating robust ordinal embedding recovery even with noise, small amounts of data and realistic semantic structure.

5.5 LORE LEARNS LOW RANK AND ACCURATE REPRESENTATIONS ON CROWDSOURCED DATA

To test LORE in practical, noisy settings with unknown intrinsic rank, we evaluate it alongside baselines on three representative crowdsourced human similarity datasets covering food images (Wilber et al., 2014), musical artists (Ellis et al., 2002), and car images (Kleindessner and Von Luxburg, 2017). These datasets differ in size, query semantics, and noise, reflecting the variety and challenges of real world ordinal data with results in Table 2. Further dataset details are in Appendix F.

All methods are trained on a random sample comprising of 90% of the total triplets with added Gaussian noise and the same embedding dimension. Across all datasets, LORE has comparable test triplet accuracy to existing methods, but uniquely yields a substantially lower rank (for e.g., Food-100: LORE gets rank 3.3 vs. 15 for the others), suggesting an intrinsic low rank structure. For musicians, low triplet counts relative to the number of percepts limit rank recovery, and for cars, extreme noise is reflected in low accuracy for all methods, but LORE's embeddings remain significantly more compact. These results show that only LORE recovers low rank structure from real data without sacrificing significant accuracy, enabling practical perceptual modeling.

5.6 LORE'S LEARNED AXES ARE SEMANTICALLY INTERPRETABLE

Figure 5 shows that, without semantic supervision, LORE's first three axes for Food-100, the same embedding used for Table 2, each align with interpretable food properties: from sweet to savory (Axis 1), dense to light (Axis 2), and carb-rich to protein/vegetable (Axis 3). The last axis is slightly less coherent, as is expected for axes linked to smaller singular values. These results demonstrate that LORE actually recovers semantically meaningful latent dimensions while recovering a low rank embedding. Consequently, the axes are interpretable and this property is invaluable for scientific discovery where the subjective percept is not well understood.

Table 2: Comparison of OEs on Real Life Ordinal Datasets

Method	Food-100			Musicians			Cars		
Metric ± Std	Test Acc.	Rank	Time (s)	Test Acc.	Rank	Time (s)	Test Acc.	Rank	Time (s)
LORE (Ours)	82.45 ± 0.27	$\begin{array}{c} 3.3 \pm \\ 0.47 \end{array}$	6.64 ± 3.90	75.63 ± 0.94	27.8 ± 0.55	$13.82 \pm \\ 9.72$	52.12 ± 1.22	$\begin{array}{c} 3\pm \\ 0.45 \end{array}$	4.45 ± 1.62
SOE	82.34 ± 0.32	$\begin{array}{c} 15 \pm \\ 0.00 \end{array}$	27.09 ± 1.38	${81.41 \pm \atop 0.93}$	$\begin{array}{c} 30 \pm \\ 0.0 \end{array}$	${28.45 \pm \atop 2.20}$	53.17 ± 1.42	$15.0\pm\\0.0$	5.53 ± 1.22
FORTE	81.73 ± 0.46	$\begin{array}{c} 15 \pm \\ 0.00 \end{array}$	6.34 ± 0.52	69.94 ± 1.61	$\begin{array}{c} 30 \pm \\ 0.0 \end{array}$	8.63 ± 2.79	52.91 ± 0.84	$15.0\pm\\0.0$	0.85 ± 0.18
t-STE	82.79 ± 0.24	$\begin{array}{c} 15 \pm \\ 0.00 \end{array}$	$^{40.93\pm}_{20.14}$	$79.49 \pm \\ 1.52$	$\begin{array}{c} 30 \pm \\ 0.0 \end{array}$	${98.97 \pm \atop 81.26}$	53.70 ± 1.15	$15.0\pm\\0.0$	15.13 ± 4.29
CKL	82.75 ± 0.20	$\begin{array}{c} 15 \pm \\ 0.00 \end{array}$	18.41 ± 7.89	$78.05 \pm \\ 0.96$	$\begin{array}{c} 30 \pm \\ 0.0 \end{array}$	24.3 ± 10.51	54.06 ± 1.19	$15.0\pm\\0.0$	4.85 ± 0.39

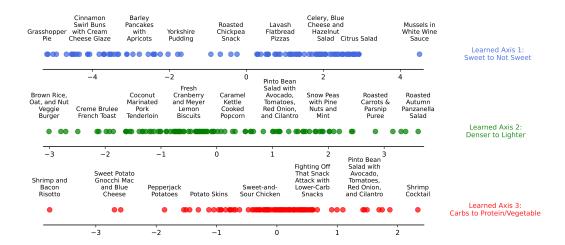


Figure 5: **LORE's learned axes are semantically interpretable**: Food groups as axis value varies for the first three learned axes of the LORE embedding learned on the Food-100 dataset. (Same embedding as one learned for Table 2).

6 Discussion

In this work, we introduced LORE, a framework for jointly learning the intrinsic rank and the true perceptual latent space via an ordinal embedding. Our results show that LORE consistently recovers low-dimensional representations, with ranks that closely match ground truth while maintaining competitive test triplet accuracy. On real crowdsourced data, LORE also uncovers interpretable axes aligned with meaningful semantic concepts, making subjective perceptual spaces easier to analyze.

One limitation of this work is the absence of theoretical guarantees for exact rank recovery or optimal embeddings. Our method empirically performs well, but its theoretical underpinnings remain an open question. We also note that LORE's optimization is only guaranteed to reach stationary points, not global minima.

Future directions include developing theoretical guarantees and optimization refinements to improve reliability, as well as exploring active learning to collect perceptual data more efficiently. Finally, we hope this work inspires further applied and theoretical advances, expanding the use of LORE for uncovering the structure of perceptual spaces across a range of domains.

7 REPRODUCIBILITY STATEMENT

We have taken several steps to ensure the reproducibility of our work. All code and detailed documentation for our artificial human experiments and crowdsourced human experiments are included in the supplemental material, with full experimental configurations provided in Appendix H and Appendix I.

For the synthetic experiments and baseline comparisons, which require large scale parallelization and substantial computational resources, we do not provide raw code. Instead, we describe in detail the procedures and parameter settings necessary to reproduce them in Appendix G.

Our main theoretical result, establishing convergence of LORE to a local optimum, includes a complete proof with all required assumptions in Appendix A. To support practical use, we provide a demo (in the supplemental material) showing how LORE can be applied to new datasets, along with additional implementation details in Appendix B.

Upon acceptance, we will release the full code and demo on GitHub and integrate the implementation of LORE into cblearn (Künstle and Luxburg, 2024), a Python package for ordinal embeddings and comparison-based machine learning. We believe these efforts will make our work fully reproducible and accessible to the community.

REFERENCES

- Agarwal, S., Wills, J., Cayton, L., Lanckriet, G., Kriegman, D., and Belongie, S. Generalized non-metric multidimensional scaling. *Artificial Intelligence and Statistics*, 11–18, 2007.
- Bishop, C. M., and Nasrabadi, N. M. *Pattern recognition and machine learning* (Vol. 4, Issue 4). Springer, 2006.
- Bolte, J., Daniilidis, A., Ley, O., and Mazet, L. Characterizations of Łojasiewicz inequalities: subgradient flows, talweg, convexity. *Transactions of the American Mathematical Society*, 362(6), 3319–3363, 2010.
- Bower, A., Jain, L., and Balzano, L. The landscape of non-convex quadratic feasibility. 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 3974–3978, 2018.
- Canal, G., Fenu, S., and Rozell, C. Active ordinal querying for tuplewise similarity learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(4), 3332–3340, 2020.
- Candes, E. J., and Recht, B. *Exact Matrix Completion via Convex Optimization* (Issue arXiv:805.4471). arXiv, 2008. https://doi.org/10.48550/arXiv.0805.4471
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. *International Conference on Machine Learning*, 1597–1607, 2020.
- Dugas, C., Bengio, Y., Dupont, F., Marcotte, H., Vincent, P., Garcia, C., and Maillet, D. Incorporating Second-Order Functional Knowledge for Better Option Pricing. *Advances in Neural Information Processing Systems (Neurips 2000)*, 13, 2001.
- Ellis, D. P., Whitman, B., Berenzweig, A., and Lawrence, S. *The quest for ground truth in musical artist similarity*, 2002.
- Fazel, M., Hindi, H., and Boyd, S. P. A rank minimization heuristic with application to minimum order system approximation. *Proceedings of the 2001 American Control Conference.(Cat. No. 01ch37148)*, 6, 4734–4739, 2001.
- Filip, J., Lukavskỳ, J., Děchtěrenko, F., Schmidt, F., and Fleming, R. W. Perceptual dimensions of wood materials. *Journal of Vision*, 24(5), 12, 2024.
 - Jain, L., Jamieson, K. G., and Nowak, R. Finite sample prediction and recovery bounds for ordinal embedding. *Advances in Neural Information Processing Systems*, 29, 2016.

565

567

569

570

571

574

575

577

578

579

580

581 582

583

584

- Kleindessner, M., and Von Luxburg, U. Lens depth function and k-relative neighborhood graph: versatile tools for ordinal data analysis. *Journal of Machine Learning Research*, 18(58), 1–52, 2017.
- Künstle, D.-E., and Luxburg, U. von. cblearn: Comparison-based Machine Learning in Python. Journal of Open Source Software, 9(98), 6139, 2024.
- Künstle, D.-E., Luxburg, U. von, and Wichmann, F. A. Estimating the perceived dimension of psychophysical stimuli using triplet accuracy and hypothesis testing. *Journal of Vision*, 22(13), 5, 2022.
 - Likert, R. A technique for the measurement of attitudes. *Archives of Psychology*, 1932.
- Lu, C., Tang, J., Yan, S., and Lin, Z. Generalized Nonconvex Nonsmooth Low-Rank Minimization.

 2014 IEEE Conference on Computer Vision and Pattern Recognition, 4130–4137, 2014. https://doi.org/10.1109/CVPR.2014.526
- Marjanovic, G., and Solo, V. On *l_q* optimization and matrix completion. *IEEE Transactions on Signal Processing*, *60*(11), 5714–5724, 2012.
 - Marjieh, R., Sucholutsky, I., Rijn, P. van, Jacoby, N., and Griffiths, T. L. Large language models predict human sensory judgments across six modalities. *Scientific Reports*, 14(1), 21445, 2024.
 - Negahban, S., and Wainwright, M. J. Estimation of (near) low-rank matrices with noise and highdimensional scaling, 2011.
- Reimers, N., and Gurevych, I. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, 2019. https://arxiv.org/abs/1908.10084
 - Stewart, N., Brown, G. D., and Chater, N. Absolute identification by relative judgment. *Psychological Review*, 112(4), 881, 2005.
 - Sun, T., Jiang, H., and Cheng, L. Convergence of proximal iteratively reweighted nuclear norm algorithm for image processing. *IEEE Transactions on Image Processing*, 26(12), 5632–5644, 2017.
 - Suárez-Díaz, J. L., García, S., and Herrera, F. A tutorial on distance metric learning: Mathematical foundations, algorithms, experimental analysis, prospects and challenges (with appendices on mathematical background and detailed algorithms explanation). *Arxiv Preprint Arxiv:1812.05944*, 2018.
 - Tamuz, O., Liu, C., Belongie, S., Shamir, O., and Kalai, A. T. *Adaptively Learning the Crowd Kernel* (Issue arXiv:1105.1033). arXiv, 2011. https://doi.org/10.48550/arXiv.1105.1033
 - Terada, Y., and Luxburg, U. Local ordinal embedding. *International Conference on Machine Learning*, 847–855, 2014.
 - Van Der Maaten, L., and Weinberger, K. Stochastic triplet embedding. 2012 IEEE International Workshop on Machine Learning for Signal Processing, 1–6, 2012.
 - Vankadara, L. C., Lohaus, M., Haghiri, S., Wahab, F. U., and Von Luxburg, U. Insights into ordinal embedding algorithms: A systematic evaluation. *Journal of Machine Learning Research*, 24(191), 1–83, 2023.
- Wang, H., Wang, Y., and Yang, X. Efficient Active Manifold Identification via Accelerated Iteratively
 Reweighted Nuclear Norm Minimization. *Journal of Machine Learning Research*, 25(319), 1–44, 2024.
- Wilber, M., Kwak, I., and Belongie, S. Cost-effective hits for relative similarity comparisons. *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, 2, 227–233, 2014.
- Zhang, C.-H. Nearly unbiased variable selection under minimax concave penalty, 2010.

A Proof for Theorem 1

Theorem (LORE converges to a stationary point) The sequence of OEs generated by the LORE algorithm $\{Z^k\}_{k=1,2,3,...}$ converges. i.e.

$$\sum_{k=1}^{+\infty} \lVert Z^{k+1} - Z^k
Vert_F < +\infty$$

Proof: We use the general framework for nonconvex Schatten Quasi-Norm optimization as seen in (Sun et al., 2017) but check the specific conditions for the LORE objective.

Let us split up the objective as follows.

$$\begin{split} \min_{\boldsymbol{Z}} \ \Psi(\boldsymbol{Z}) &= \sum_{(a,i,j) \in T} \log \left(1 + \exp \left(1 + d \left(\boldsymbol{Z}_{a,:}, \boldsymbol{Z}_{i,:} \right) - d \left(\boldsymbol{Z}_{a,:}, \boldsymbol{Z}_{j,:} \right) \right) \right) + \lambda \ \| \boldsymbol{Z} \|_p^p \\ &= \underbrace{\sum_{(a,i,j) \in T} \log \left(1 + \exp \left(1 + d \left(\boldsymbol{Z}_{a,:}, \boldsymbol{Z}_{i,:} \right) - d \left(\boldsymbol{Z}_{a,:}, \boldsymbol{Z}_{j,:} \right) \right) \right)}_{f(\boldsymbol{Z})} \ + \sum_{i=1}^{\min\{N,d'\}} \lambda g[\sigma_i(\boldsymbol{Z})]. \end{split}$$

There are four assumptions we need to satisfy to apply the general result from (Sun et al., 2017).

A1. *f* is differentiable and has a Lipschitz gradient:

This stems from smoothing the triplet loss with the softplus function. The composition makes f differentiable everywhere except at degenerate collapse points (which do not arise with practical initializations). The log-sum-exp structure ensures (locally) Lipschitz gradients (Chen et al., 2020).

A2. *g* is concave, nondecreasing, and Lipschitz:

$$q(x) = x^p$$
 for $x > 0$ has these properties.

A3. Ψ is coercive:

For our objective, suppose $\|Z\|_F \to \infty$. Then the sum of squared singular values diverges, so at least one $\sigma(Z) \to \infty$. As g(Z) is the sum of all $\sigma(Z)^p$, and p > 0, we therefore have $g(Z) \to \infty$ and thus $\Psi(Z) \to \infty$.

A4. Ψ has the Kurdyka Lojasciewicz (KL) property:

As established in (Bolte et al., 2010), sums of o-minimal (definable) functions, such as our loss and regularizer, possess the KL property.

With all required assumptions satisfied, Theorem 1 of (Sun et al., 2017) applies and guarantees:

$$\sum_{k=1}^{\infty} \lVert Z^{k+1} - Z^k \rVert_F < \infty.$$

Therefore, the LORE algorithm converges to a stationary point.

B IMPLEMENTATION DETAILS FOR LORE

The optimization algorithm used for LORE is an adaptation of the original algorithm from (Sun et al., 2017).

The function takes the initialized embedding Z^0 , the regularization parameter λ , the Lipschitz constant of $\nabla f(.)$, μ , and the tolerance for convergence tol. The exact Lipschitz constant of f(.) is not known but was be empirically estimated to be strictly greater than 0.013 by the Power Iteration method. The algorithm initializes the ordinal embedding and iteratively updates it by minimizing the smoothed ordinal loss plus Schatten-p regularization. Each step performs a proximal gradient update and singular value thresholding, repeating until convergence. Based on prior literature in the Schatten-p quasi-norm optimization literature (Lu et al., 2014; Sun et al., 2017; Wang et al., 2024), we fix p=0.5 as it has been shown to have good empirical results.

If we consider the most standard operational setting, i.e. $d' < N \ll T$, then the time complexity of each iteration is $\mathcal{O}(d'(T+Nd'))$. The dominating terms here are the number of percepts and the number of triplets used with the most intensive operation is in line 8 where the gradient of f is calculated and a singular value decomposition is subsequently performed. As a result, LORE is scalable for higher N and d'. One does need to be careful as T could scale with $\mathcal{O}(N^3)$ if many triplets are chosen which could slow down each iteration.

C ADDITIONAL PLOTS FOR REGULARIZATION OF LORE

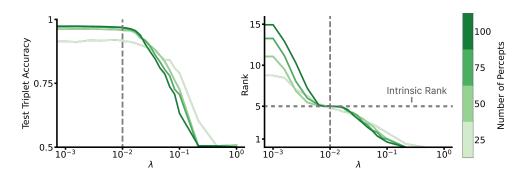


Figure 6: LORE has high test triplet accuracy and intrinsic rank recovery. (Left) Mean test triplet accuracy vs λ for LORE as number of percepts varies. (Right) Mean measured rank vs λ for LORE as number of percepts varies.

Figure 6 shows the test triplet accuracy and intrinsic rank recovery of LORE as the number of percepts varies. We see that with greater number of percepts rank recovery stays roughly constant whereas the test triplet accuracy increases significantly from 25-50 percepts. Baseline parameters are intrinsic rank = 5, fraction of queries = 0.1, noise = 0.1.

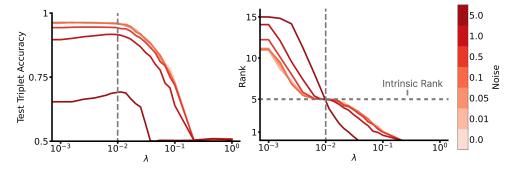


Figure 7: LORE has high test triplet accuracy and intrinsic rank recovery. (Left) Mean test triplet accuracy vs λ for LORE as noise varies. (Right) Mean measured rank vs λ for LORE as noise varies.

Figure 7 shows the test triplet accuracy and intrinsic rank recovery of LORE as the noise varies. We see that with greater noise rank recovery and test triplet accuracy both decrease. There is a dramatic drop in test triplet accuracy from 1 to 5 noise. Baseline parameters are intrinsic rank = 5, number of percepts = 50, fraction of queries = 0.1.

These results, together with Figure 2, show that LORE is quite robust to the various knobs that can be tuned for OE algorithms.

D ADDITIONAL PLOTS COMPARING LORE TO BASELINES

72.7

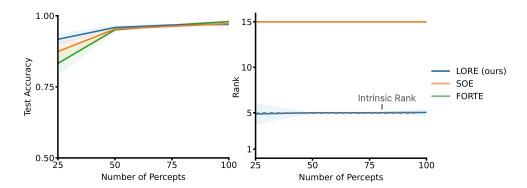


Figure 8: Only LORE can recover the intrinsic rank while maintaining comparable test triplet accuracy. (Left) Mean test triplet accuracy vs number of percepts used for LORE and the baselines. (Right) Mean measured rank vs number of percepts used for LORE and the baselines.

Figure 8 shows the test triplet accuracy and intrinsic rank recovery of LORE and the baselines as the number of percepts varies. We see that with greater number of percepts rank recovery stays roughly constant, though spread decreases, for LORE from 25-50 percepts. Baselines again cannot recover the intrinsic rank at all. Test triplet accuracy increases from 25-50 percepts for all OE algorithms. Baseline parameters are intrinsic rank = 5, fraction of queries = 0.1, noise = 0.1.

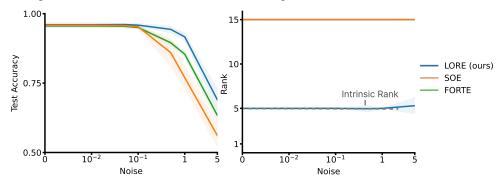


Figure 9: Only LORE can recover the intrinsic rank while maintaining comparable test triplet accuracy. (Left) Mean test triplet accuracy vs noise used for LORE and the baselines. (Right) Mean measured rank vs noise used for LORE and the baselines.

Figure 9 shows the test triplet accuracy and intrinsic rank recovery of LORE and the baselines as the noise varies. We see that with greater noise, LORE is still able to recover the intrinsic rank though spread increases with noise from 1-5. The baselines cannot recover the intrinsic rank at all. Test triplet accuracy decreases with noise for all OE algorithms though LORE still performs the best. Baseline parameters are intrinsic rank = 5, number of percepts = 50, fraction of queries = 0.1.

E EXPERIMENTAL SETUP FOR FIGURE 1

For this figure, each experiment was run for 30 runs with different random seeds. The results were averaged over the 30 runs. This was run on a server with 1 RTX3080 GPU and 128 GB of RAM.

For each seed, a separate "true" representation P of N=50 percepts and an intrinsic rank of d=5 was sampled at random with a training set of 10% of the total number of queries. 3000 held out triplets were used as the test set. $\sigma=0.1$ Gaussian noise was added to model the uncertainty of the human responses. All embedding algorithms were run with an embedding dimension (d') of 15. LORE had the regularization parameter λ set to 0.01.

Code is included in the supplemental material.

F Crowdsourced Dataset Details

Table 3: Characterization of Crowdsourced Datasets Used

There is a share the state of t									
Datasets	Number of Percepts	Number of Triplets	Triplet Type	Notes					
Food-100 (Wilber et al., 2014)	100	190,376	Compared to A, which is more similar, B or C or?	Images of foods to user. Converted data into similar- ity triplets using the python package cblearn (Künstle and Luxburg, 2024)					
Musicians (Ellis et al., 2002)	448	118,263	Compared to A, which is more similar, B or C or?	Names of musicians presented to users. Converted data into similarity triplets us- ing the python pack- age cblearn (Künstle and Luxburg, 2024)					
Cars (Kleindessner and Von Luxburg, 2017)	68	7097	Which of A, B, C is the most central?	Images of cars presented to user. Each central triplet can be converted to similarity triplets using the python package cblearn (Künstle and Luxburg, 2024)					

Of these datasets, the Cars dataset is known to be very noisy (Kleindessner and Von Luxburg, 2017; Vankadara et al., 2023). Food-100 has been used as a dataset to evaluate active querying methods (Canal et al., 2020).

G Experimental Setup for Section 5.2 and Section 5.3

This experiment was performed on a SLURM server with over 30 GPUs of varying quality and compute power. We do not include the scripts used to run those experiments as they are highly complex due to parallelism and take too long to run (over 8 days). However, a quick rundown of the experiment is given below.

A grid search over all the following parameters was performed for these experiments. The grid search was performed in parallel over 30 GPUs. Each experiment was run for 30 runs with different random seeds. The results were averaged over the 30 runs and the standard deviation was calculated.

In our experiments, the various knobs we tune are as follows.

- Number of Percepts (N): We vary it from [25, 50, 75, 100] and use 50 Percepts as a default. We do not increase the number of percepts beyond 100 as the number of queries increases combinatorially. Additionally, this is not a practical number of percepts to collect for perceptual experiments.
- True Dimension (d): We vary it from [1, 2, 3, 4, 5, 6, 7, 8, 9, 10] and use 5 as a default. We do not examine over 10 dimensions as it is not possible that many dimensions without increasing the number of percepts due to the curse of dimensionality (Bishop and Nasrabadi, 2006).
- Fraction of Queries used: we vary it from [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9] and use 0.1 as a default.
- Noise (σ^2) : we vary it from [0, 0.01, 0.05, 0.1, 0.5, 1.0, 5.0] and use 0.1 as a default.
- **Regularization** (λ): This is only for LORE but we vary it with [0, 0.001, 0.00158489, 0.00251189, 0.00398107, 0.00630957,

0.00768625, 0.00936329, 0.01, 0.01140625, 0.01389495, 0.01692667, 0.02061986, 0.02511886, 0.0305995, 0.03727594, 0.0454091, 0.05531681, 0.06738627, 0.08208914, 0.1, 0.21544347, 0.46415888, 1.] and use 0.01 as a default.

• Embedding Dimension (d'): This is the dimension of the embedding we are trying to learn. This is only for the baselines other than LORE. We vary it from [1, 2, 3, 4, 5, 6, 7, 8, 10, 12, 15] and use 15 as a default.

The metrics we measure are as follows

810

811

812

813

814

815

816

817

818

819 820

821

822

823 824

826

827

828

829

831

832

833

834

835

837

838 839

840

841 842

- Test Triplet Accuracy: The accuracy of the test triplets on the test set. This is the main metric we use to measure performance.
- Measured Rank: The rank of the embedding matrix. This is a measure of how well the
 algorithm is able to recover the intrinsic rank of the data. We measure this by taking the SVD
 of the embedding matrix and counting the number of non-zero singular values. Specifically,
 we use the rank function from the numpy library to compute the rank of the embedding
 matrix.
- Peak Signal to Noise Ratio: The PSNR is a measure of the quality of the recovered matrix. However, note that the recovered embedding matrix has to be aligned to the true percepts matrix to compute the PSNR. The specific formulation is described in Appendix J.
- Normalized Procrustes Distance: The NPD is a measure of how well the recovered matrix matches the true matrix up to rotation, scaling and translation. To perform procrustes analysis, true percepts $P \in \mathbb{R}^{\{N \times d\}}$ and the computed embedding $Z \in \mathbb{R}^{\{N \times d'\}}$ must be the same shape. Therefore, we use the same subspace alignment technique to ensure that the two matrices have the same shape. The specific formulation is described in Appendix J.

It should be noted that test triplet accuracy and measured rank are the main metrics we use to measure performance as the other metrics require knowledge of the percepts P which is not known in practice.

H Experimental Setup for Section 5.4

50 random foods were chosen from the Food-100 dataset (Wilber et al., 2014). This was run on a server with 1 RTX3080 GPU and 128 GB of RAM. The names of the specific percepts are as follows.

```
['Cinnamon Swirl Buns with Cream Cheese Glaze',
         'Shrimp and Bacon Risotto',
844
         'Shrimp Cocktail',
845
         'Homemade Cracker Jacks',
846
         'Creme Brulee French Toast',
847
         'Red Lobster Cheddar Bay Biscuits',
848
         'Apple Bacon Stuffed Sweet Potatoes',
849
         'Sweet-and-Sour Chicken',
850
         'Pumpkin-Chocolate Chunk Pancakes',
851
         'Chocolate Hazelnut Biscotti',
852
         'Eggnog Ice Cream',
853
         'Celery, Blue Cheese and Hazelnut Salad',
854
         'Shredded and Roasted Brussels Sprouts with Almonds and Parmesan',
855
         'Low-Sugar Pumpkin and Apple Crumble',
856
         'Roasted Sweet Potatoes Recipe with Double Truffle Flavor and Parmesan',
857
         'Chicken Florentine Bowtie Pasta',
858
         'White Whole Wheat Pizza Dough',
         'Chervil Mayonnaise',
860
         'Pork Tenderloin in Tomatillo Sauce',
861
         'Yellow Tomato Salad with Roasted Red Pepper, Feta, and Mint',
         'Daisy Brand Sour Cream Chocolate Cake',
863
         'Shredded Brussels Sprouts & Apples',
         'Mexican Corn Salad',
         'Potato Skins',
```

'Caramel Kettle Cooked Popcorn',

- Roasted Garlic + Veggie Tostadas',

 Bear Seared Scallops with Baby Grant
 - 'Pan Seared Scallops with Baby Greens and Citrus Mojo Vinaigrette',
- Lemon Cranberry Scones',
- Warm Butternut and Chickpea Salad with Tahini Dressing',
- Fighting Off That Snack Attack with Lower-Carb Snacks',
- 'Chicken with Forty Cloves of Garlic',
- Edna Mae's Sour Cream Pancakes',
- 'Sweet Potato Gnocchi Mac and Blue Cheese',
- 872 'Yorkshire Pudding',
- 873 'Luscious Lemon Squares',
- ₈₇₄ 'Japanese Pizza',
- 'Grilled Asparagus & Feta Salad',
- 'Grilled Corn Salad',
- 'Garlic Meatball Pasta',
- 'Roasted Autumn Panzanella Salad',
- 'Coconut Marinated Pork Tenderloin',
- 'Black Raspberry Sorbet',
- 'Mini Whole Wheat BBQ Chicken Calzones',
- 'Mussels in White Wine Sauce',
- 'Brown Rice, Oat, and Nut Veggie Burger',
- 'Dark Chocolate Cookies',
- 'Citrus Salad',

891

892

893

894 895

896

898

899

900

901 902 903

904

905 906

907

908

909

910

911

912

913

914

915

916

917

- ⁸⁸⁵ 'Roasted Carrots & Parsnip Puree',
- 'South African Cheese, Grilled Onion & Tomato Panini (Braaibroodjie)',
 - 'Pinto Bean Salad with Avocado, Tomatoes, Red Onion, and Cilantro']

These names are passed to the SBERT library (Reimers and Gurevych, 2019) with the "all-mpnet-base-v2" model to get a 768 dimensional LLM embedding. To simulate various possible intrinsic ranks, we use the truncated singular value decomposition to constrain the "true" perceptual representations of foods to intrinsic ranks 1-10. Specifically, the truncated SVD is the following.

The singular value decomposition of a matrix $P' \in \mathbb{R}^{\{N \times d\}}$ is given by

$$P' = U\Sigma V^T$$

Here $\boldsymbol{U} \in \mathbb{R}^{\{N \times N\}}$ $\boldsymbol{S} \in \mathbb{R}^{\{N \times 768\}}$ and $\boldsymbol{V} \in \mathbb{R}^{\{768 \times 768\}}$. \boldsymbol{U} and \boldsymbol{V} have orthonormal columns and $\boldsymbol{\Sigma}$ is a diagonal matrix with singular values $\sigma_1 \geq \sigma_2 \geq ... \geq \sigma_N > 0$. The intrinsic rank of the matrix is the number of non-zero singular values, which in this case is N before truncation. We can truncate the SVD to fix an intrinsic rank of d.

If $U_d \in \mathbb{R}^{N \times d}$ and $V_d \in \mathbb{R}^{768 \times d}$ are the first d columns of U and V respectively, $\Sigma_d \in \mathbb{R}^{d \times d}$ with the biggest d singular values in the diagonal entries and otherwise 0 then we can write the truncated SVD of P' to get the our "true" perceptual representation of the foods as

$$P = U_d \Sigma_d V_d^T$$

Then, we query just 5% of the total possible triplets (2940 out of a possible 58800) with 0.1 noise added to the triplet comparisons thirty times, independently, before training the various OE algorithms. The metrics we measure are as follows.

- **Test Triplet Accuracy**: The accuracy of the test triplets on the test set (fixed at 3000 queries not in the train set and chosen at random). This is the main metric we use to measure performance.
- Measured Rank: The rank of the embedding matrix. This is a measure of how well the algorithm is able to recover the intrinsic rank of the data. We measure this by taking the SVD of the embedding matrix and counting the number of non-zero singular values. Specifically, we use the rank function from the numpy library to compute the rank of the embedding matrix.
- Peak Signal to Noise Ratio: The PSNR is a measure of the quality of the recovered matrix. However, note that the recovered embedding matrix has to be aligned to the true

percepts matrix to compute the PSNR. The specific formulation is described in Appendix J. (We do not report these in the paper)

• Normalized Procrustes Distance: The NPD is a measure of how well the recovered matrix matches the true matrix up to rotation, scaling and translation. To perform procrustes analysis, true percepts $P \in \mathbb{R}^{\{N \times d\}}$ and the computed embedding $Z \in \mathbb{R}^{\{N \times d'\}}$ must be the same shape. Therefore, we use the same subspace alignment technique to ensure that the two matrices have the same shape. The specific formulation is described in Appendix J. (We do not report these in the paper)

Code for this experiment is included in the supplemental material.

I Experimental Setup for Section 5.5

This was run on a server with 1 RTX3080 GPU and 128 GB of RAM.

We learn OEs for foods from the Food-100 dataset (Wilber et al., 2014). This dataset contains 100 foods with 171,388 crowdsourced triplet comparisons. We choose a random 50 foods from the dataset and restrict our analysis to the triplets which only contain those 50 foods. The foods chosen here are the same ones from Appendix H. Then, we randomly sample 90% of the triplet data, add noise with scale 0.1, train OE algorithms and test on the remaining 10%. We repeat for a total of 30 independent train test splits.

For LORE, we set the regularization parameter, λ , to 0.01. For all OE methods, we set the number of dimensions of the OE, d', to 15.

Note that for this experiment unlike in Appendix H, we do not have access to the true percepts P and therefore cannot compute the PSNR or NPD. We only report the test triplet accuracy and measured rank.

Code for this experiment is included in the supplemental material.

J FORMULATION OF OTHER METRICS

Code for all of these implementations is included in the supplemental material.

J1. SUBSPACE ALIGNMENT

To perform procrustes analysis, true percepts $P \in \mathbb{R}^{\{N \times d\}}$ and the computed embedding $Z \in \mathbb{R}^{\{N \times d'\}}$ must be the same shape.

Specifically we compute

$$m{P_c} = m{P} - \mathbf{1}_N m{\mu}_{m{P}}^T$$
 and $m{Z_c} = m{Z} - \mathbf{1}_N m{\mu}_{m{Z}}^T$

Then, we compute the tikhonov regularized projection matrix to prevent numerical instability due to ill conditioning. We use a regularization parameter of $\eta = 1e - 3$.

$$\boldsymbol{A} = (\boldsymbol{Z}_c^T \boldsymbol{Z}_c + \eta \boldsymbol{I}_d) \boldsymbol{Z}_c^T \boldsymbol{P}_c$$

Then, we can get the aligned ordinal embedding $m{Z}_{ ext{aligned}} = m{Z}_c m{A} + \mathbf{1}_N \mu_{m{P}}^T.$

J2. NORMALIZED PROCRUSTES DISTANCE

Now that we have an aligned matrix the same shape as P, the normalized procrustes distance between the aligned embedding and the true percepts can be computed as

Normalized Procrustes Distance = $\frac{\|P - Z_{\text{aligned}}\|_F}{\|Z_c\|_F}$

J3. PEAK SIGNAL TO NOISE RATIO

The Peak Signal to Noise Ratio (PSNR) is a measure of the quality of the recovered matrix and is defined as $20\log_{10}\left(\frac{\max(\mathbf{Z}_{\text{aligned}})}{\|\mathbf{Z}_{\text{aligned}}-\mathbf{P}\|_F}\right)$ where \mathbf{P} is the true matrix.

K LLM USAGE

In this work, we leverage the use of large language models for two purposes. (1) to refine the writing by eliminating grammatical errors and improving flow. However, these were only used at the individual paragraph level rather than whole sections and (2) to discover similar papers during the literature review for the related work. Specifically, we searched for terms like "distance metric learning", "contrastive learning", "psychophysical scaling" etc.