

1 Identifying patterns differing 2 between high-dimensional datasets 3 with generalized contrastive PCA

4 **Eliezyer Fermino de Oliveira¹, Pranjal Garg², Jens Hjerling-Leffler³, Renata**
5 **Batista-Brito^{1,4,5}, Lucas Sjulson^{1,4,*}**

*For correspondence:
luke@sjulsonlab.org (LS)

6 ¹Dominick P. Purpura Department of Neuroscience, Albert Einstein College of Medicine,
7 Bronx, NY; ²All India Institute of Medical Sciences, Rishikesh, India; ³Department of
8 Medical Biochemistry and Biophysics, Karolinska Institutet, Stockholm SE-17177,
9 Sweden; ⁴Department of Psychiatry and Behavioral Sciences, Albert Einstein College of
10 Medicine, Bronx, NY; ⁵Department of Genetics, Albert Einstein College of Medicine,
11 Bronx, NY

13 **Abstract** High-dimensional data have become ubiquitous in the biological sciences, and it is
14 often desirable to compare two datasets collected under different experimental conditions to
15 extract low-dimensional patterns enriched in one condition. However, traditional dimensionality
16 reduction techniques cannot accomplish this because they operate on only one dataset.
17 Contrastive principal component analysis (cPCA) has been proposed to address this problem, but
18 it has seen little adoption because it requires tuning a hyperparameter resulting in multiple
19 solutions, with no way of knowing which is correct. Moreover, cPCA uses foreground and
20 background conditions that are treated differently, making it ill-suited to compare two
21 experimental conditions symmetrically. Here we describe the development of generalized
22 contrastive PCA (gcPCA), a flexible hyperparameter-free approach that solves these problems.
23 We first provide analyses explaining why cPCA requires a hyperparameter and how gcPCA avoids
24 this requirement. We then describe an open-source gcPCA toolbox containing Python and
25 MATLAB implementations of several variants of gcPCA tailored for different scenarios. Finally, we
26 demonstrate the utility of gcPCA in analyzing diverse high-dimensional biological data, revealing
27 unsupervised detection of hippocampal replay in neurophysiological recordings and
28 heterogeneity of type II diabetes in single-cell RNA sequencing data. As a fast, robust, and
29 easy-to-use comparison method, gcPCA provides a valuable resource facilitating the analysis of
30 diverse high-dimensional datasets to gain new insights into complex biological phenomena.

32 Introduction

33 Investigators in the biological sciences are increasingly collecting high-dimensional datasets that
34 are challenging to analyze, with modalities ranging from imaging to electrophysiology to single-cell
35 RNA sequencing. Dimensionality reduction algorithms such as principal components analysis (PCA)
36 and its many variants (*Pearson, 1901; Hotelling, 1933; Zou et al., 2006; Zass and Shashua, 2006; Tip-*
37 *ping and Bishop, 1999*) are used widely to help simplify these datasets and facilitate analysis. PCA
38 examines the covariance structure of the data to find dimensions that account for more variance
39 than chance; these constitute patterns that are overrepresented in the data, such as assemblies of
40 neurons whose activity fluctuates up and down together across time in a neural recording (*Chapin*

41 *and Nicoletis, 1999; Peyrache et al., 2010; Lopes-dos Santos et al., 2013; Sjulson et al., 2018*), or
42 networks of genes that are up- or down-regulated together across cells in a single-cell RNAseq
43 dataset (*Chung and Storey, 2015; Li et al., 2016*). However, in many cases, the goal is to compare
44 data collected under two different experimental conditions, which we refer to here as datasets.
45 Since PCA and other dimensionality reduction techniques operate on only one dataset, they can-
46 not take experimental conditions into account.

47 The most common approach for comparing two high-dimensional datasets is linear discrim-
48 inant analysis (*Izenman, 2008*) or its multidimensional analog, partial least squares discriminant
49 analysis (*Brereton and Lloyd, 2014*). These methods find dimensions that optimally distinguish
50 one dataset from the other, which could correspond to which neurons fire more, or which genes
51 are upregulated, in condition *A* vs. condition *B*. However, an analogous method to compare the
52 covariance structure of two datasets is not as well established. This addresses more subtle and
53 detailed questions, such as which subsets of neurons exhibit increased temporal correlations in
54 condition *A* than *B*, or which subsets of genes are more likely to be up- or downregulated together
55 in individual cells in condition *A* than *B*. Mathematically, answering these questions corresponds
56 to finding dimensions that account for more or less variance in *A* than *B*.

57 Recently, contrastive PCA (cPCA) was proposed as a method to address this problem (*Abid et al.,*
58 *2018*). Although cPCA is an important first step, it requires a hyperparameter α , which controls how
59 much covariance from the second condition to subtract from the first. The algorithm must there-
60 fore iterate over multiple choices of α with no objective criteria to determine which value of α yields
61 the correct answer. Moreover, cPCA is asymmetric, identifying the most enriched dimensions in
62 the first condition after subtracting out the second condition as background; it cannot treat the
63 two experimental conditions equally.

64 Here we propose a novel solution to these problems we call generalized contrastive PCA (gcPCA).
65 We first demonstrate the role the α hyperparameter plays in cPCA, then explain our strategy for
66 eliminating it. We then describe an open source toolbox for Python and MATLAB implementing
67 several versions of gcPCA with different objective functions that are either asymmetric or symmet-
68 ric, orthogonal or non-orthogonal, or sparse or dense, tailored to suit the specific application at
69 hand. Finally, we demonstrate the utility of gcPCA in the analysis of diverse biological datasets.

70 Results

71 The cPCA hyperparameter α compensates for bias toward high-variance dimen- 72 sions in noisy, finitely-sampled datasets

73 To explain the need for the hyperparameter α in cPCA and how we avoid it in gcPCA, we will describe
74 the objective function of each method and show how they perform in generated synthetic data. For
75 illustration purposes, we generated synthetic data for two experimental conditions containing two-
76 dimensional manifolds on a background of high-variance shared dimensions. The generated data
77 consisted of condition *A*, with a manifold (additional variance) in the 71st and 72nd dimensions
78 (ranked in order of descending variance), and condition *B*, with a manifold in the 81st and 82nd
79 dimensions (Fig. 1A). The manifold dimensions contained less total variance than most of the other
80 dimensions in the dataset, but their variance is two-fold higher in one condition relative to the other
81 (i.e., the 81st and 82nd dimensions have twice as much variance in condition *B* than condition *A*).

82 An important property of real-world biological datasets is that they are noisy and finitely sam-
83 pled. We aimed to model the finite data regime by comparing 1×10^3 samples (finite data) to
84 1×10^5 samples, which approximates infinite data. To inspect the effects of finite sampling on es-
85 timated variance, we projected the "finite" and "infinite" data onto the ground truth dimensions
86 and calculated the variance explained by each (see methods). Our results (Fig. 1B) reveal that
87 finite sampling yields noisy estimates of the true variance, with greater noise in high-variance di-
88 mensions.

89 To understand the practical consequences of this, it is helpful to start by reviewing traditional

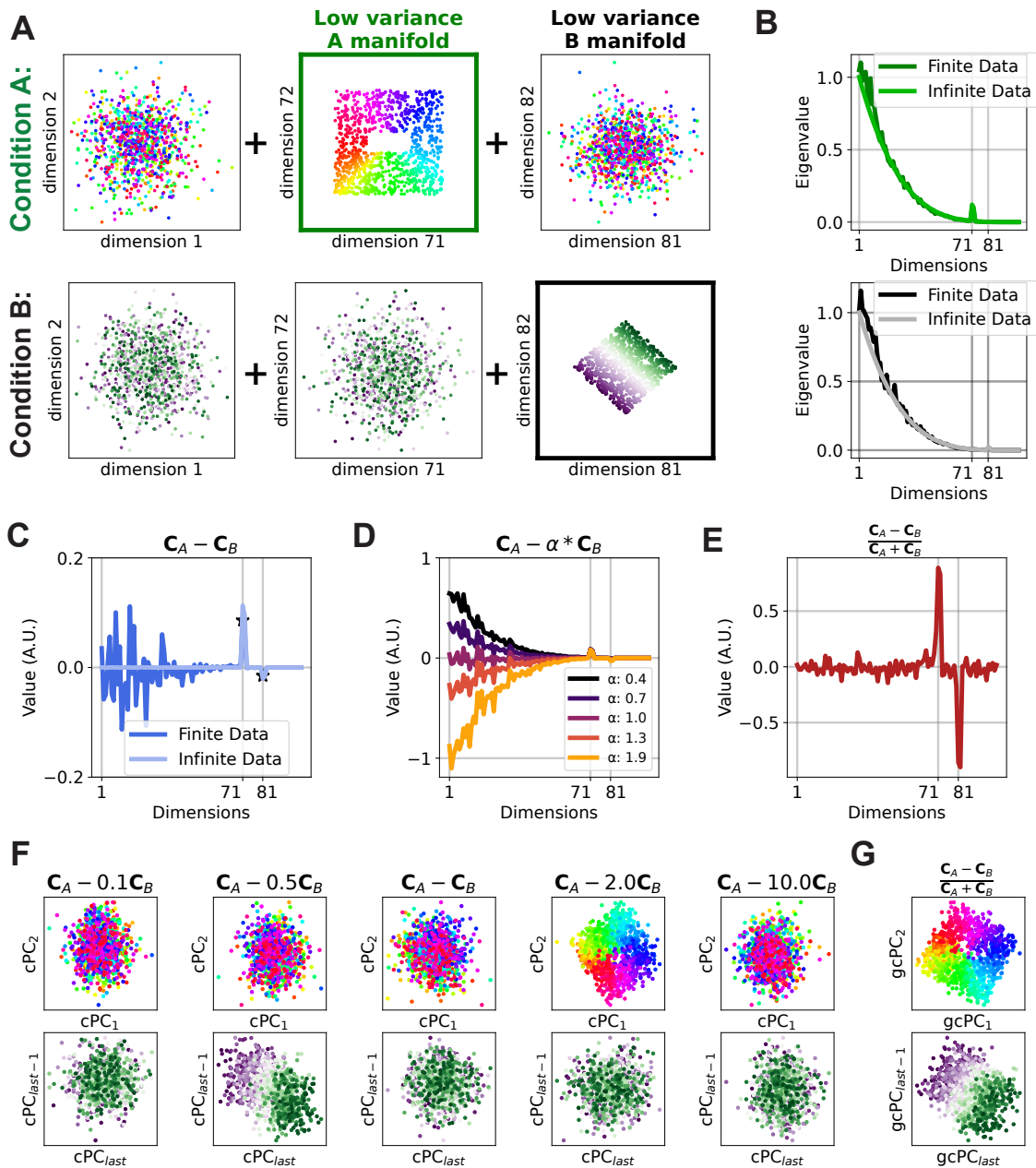


Figure 1. Generalized contrastive PCA avoids cPCA's need for the hyperparameter α in noisy, finitely-sampled data. **A** We generated two conditions in noisy synthetic data that each contain low-variance manifolds that are not present in the other. These manifolds have lower overall variance than many other dimensions and are not trivially discoverable. **B** Eigenvalue spectra for each condition estimated from finite (dark line) or infinite (light line) data. Note the sampling error in the finite data case. **C** With infinite data, eigendecomposition of $(C_A - C_B)$ suffices to extract the correct answers (dimensions 71-72 and 81-82, light lines). However, with finite data, these peaks are smaller than the sampling error in high-variance dimensions, creating a bias toward high-variance dimensions being selected. **D** cPCA uses the hyperparameter α to adjust how much influence C_B has on C_A . As α increases, the bias toward high-variance dimensions decreases until it becomes negative with $\alpha > 1$, eventually exposing the differences in lower-variance dimensions. Importantly, there is no way to know which value of α yields the correct solution. **E** Using gcPCA, the dimensions most changed in each condition are identified correctly, even with finitely-sampled data. **F** cPCA with the optimal choice of α does not extract the correct dimensions in **B**. **G** gcPCA identifies the enriched dimensions in each condition and correctly returns the low-variance manifolds. Because gcPCA is symmetric, it extracts the correct dimensions in both **A** and **B**.

90 PCA. With PCA, the principal components of a data matrix \mathbf{D} , of size $n \times p$ (samples \times features), are
91 the dimensions explaining the most variance. These can be identified by estimating the covariance
92 matrix $\mathbf{C} = \mathbf{D}^T \mathbf{D} / (n - 1)$, then solving the following quadratic optimization problem in equation 1:

$$\arg \max_{\mathbf{x} : \mathbf{x}^T \mathbf{x} = 1} \mathbf{x}^T \mathbf{C} \mathbf{x} \quad (1)$$

93
94 This problem can be solved by eigendecomposition of \mathbf{C} , yielding the matrix of eigenvectors \mathbf{X}
95 known as principal components (PCs).

96 To extend this to two datasets, a logical strategy involves formulating an objective function to
97 describe the difference in variances between the two conditions, enabling us to extract dimensions
98 that show the greatest increase in variance in A relative to B . We now have two covariance matrices,
99 \mathbf{C}_A and \mathbf{C}_B , and the contrastive PCs (cPCs) are the vectors that maximize the objective function 2:

$$\arg \max_{\mathbf{x} : \mathbf{x}^T \mathbf{x} = 1} \mathbf{x}^T (\mathbf{C}_A - \mathbf{C}_B) \mathbf{x} \quad (2)$$

100
101 Analogous to traditional PCA, this problem can be solved by eigendecomposition of $(\mathbf{C}_A - \mathbf{C}_B)$. This
102 yields cPCs that account for more variance in either condition A or B , corresponding to the eigen-
103 vectors with the largest or smallest eigenvalues, respectively. With infinite data, this approach
104 correctly finds the dimensions enriched in conditions A and B (Fig. 1C, light line), but with finite
105 data, it fails to do so because the sampling error in higher-variance dimensions is larger than the
106 true signal in the lower-variance dimensions (Fig. 1C, dark line). In other words, it has a system-
107 atic bias toward high-variance dimensions. To compensate for this effect, cPCA *Abid et al. (2018)*
108 introduces the hyperparameter α , changing the following objective function to 3:

$$\arg \max_{\mathbf{x} : \mathbf{x}^T \mathbf{x} = 1} \mathbf{x}^T (\mathbf{C}_A - \alpha \mathbf{C}_B) \mathbf{x} \quad (3)$$

109
110 As discussed in *Abid et al. (2018)*, α represents a trade-off determining the extent to which \mathbf{C}_B influ-
111 ences the identification of enriched vectors in \mathbf{C}_A . In our synthetic data, we can visually appreciate
112 the effect of different values of α in the resulting cPCA objective function value (Fig. 1D). In effect,
113 α tunes the amount of bias toward high-variance dimensions in the cPCA calculation, with $\alpha < 1$
114 biasing toward high-variance dimensions and $\alpha > 1$ biasing against them. In this case, $\alpha = 2$ yields
115 the correct solution that dimensions 71-72 are enriched in A (Fig. 1F), but other values of α yield
116 equally plausible, but incorrect, solutions. Importantly, we can only determine which solution is
117 correct because we knew the answer in advance, which is not typically the case for experimental
118 data. Further, negative values in cPCA are generally interpreted as dimensions enriched in condi-
119 tion B (*Abid et al., 2018; Boileau et al., 2020*), but our simulation shows that values of α larger than
120 1 bias the highest-variance dimensions to be negative (Fig. 1D). This creates the illusion that these
121 dimensions are enriched in condition B , even though the correct answer is that only dimensions
122 81-82 are enriched in B . Similar to the situation of finding dimensions enriched in A , the results
123 depend on the choice of α , with no way to determine which solution is correct (Fig. 1F). Importantly,
124 the range of α 's yielding the correct solution can be incredibly narrow, as in Fig. S1, where $\alpha = 2.6$
125 yields the correct solution, but 2.2 or 3.0 do not.

126 **gcPCA avoids hyperparameters by including a normalization factor**

127 Our goal for gcPCA was to eliminate the need for hyperparameters and provide unique, correct so-
128 lutions. To mitigate the bias toward high-variance dimensions, we introduce a normalization factor,
129 such as the total variance in both conditions, which can be calculated by summing the covariance
130 matrices $(\mathbf{C}_A + \mathbf{C}_B)$. The objective function then becomes (eq. 8):

$$\arg \max_{\mathbf{x} : \mathbf{x}^T \mathbf{x} = 1} \frac{\mathbf{x}^T (\mathbf{C}_A - \mathbf{C}_B) \mathbf{x}}{\mathbf{x}^T (\mathbf{C}_A + \mathbf{C}_B) \mathbf{x}} \quad (4)$$

131

132 Solving this problem is slightly more complicated than with PCA or cPCA, but ultimately it can also
133 be reduced to a computationally efficient eigenvalue problem (see Methods). The resulting general-
134 ized contrastive PCs (gcPCs) maximize relative, rather than absolute, changes in variance between
135 conditions A and B . This effectively handles the bias toward high-variance dimensions and suc-
136 cessfully extracts the ground truth dimensions in our synthetic data, even with finite sampling (Fig.
137 1E, G), and even when the range of acceptable α is narrow (Fig. S1).

138 This creates two minor complications to be aware of: first, unlike PCA or cPCA, gcPCs are not or-
139 thogonal by default. If orthogonality is important for a particular application, we have implemented
140 versions of gcPCA with an orthogonality constraint (see Methods). Second, the normalization fac-
141 tor can create numerical instability if the data are rank-deficient, yielding dimensions with zero or
142 near-zero variance in the denominator. Our implementation of gcPCA prevents this by detecting
143 and excluding those dimensions before performing the calculation (see Methods).

144 **The open-source gcPCA toolbox contains multiple gcPCA variants enabling optimal** 145 **handling of diverse use cases**

146 We developed an open-source gcPCA toolbox with implementations in Python and MATLAB of sev-
147 eral different variants of gcPCA. This toolbox is freely available at:

148 https://github.com/SjulsonLab/generalized_contrastive_PCA. Here we will present the different vari-
149 ants of gcPCA and their use cases.

150 gcPCA v1.0: traditional cPCA

151 For version 1.0, we include an implementation of the original cPCA algorithm that finds cPCs max-
152 imizing the objective function in Eqn. 3.

153 gcPCA v2.0: gcPCA maximizing A/B

154 Here we include an implementation that finds gcPCs maximizing the ratio of variance in A to B :

$$\arg \max_{\mathbf{x} : \mathbf{x}^T \mathbf{x} = 1} \frac{\mathbf{x}^T \mathbf{C}_A \mathbf{x}}{\mathbf{x}^T \mathbf{C}_B \mathbf{x}} \quad (5)$$

155

156 Like cPCA, this method is asymmetrical, meaning it is suitable for situations in which A is a fore-
157 ground condition and B is a background condition; in other words, A is presumed to be equal to B
158 with a low-dimensional pattern added, and the goal is to extract that pattern. The resulting eigen-
159 values are the ratio of the variance a given gcPC accounts for in A to the variance it accounts for in
160 B . Thus, they fall in the range $[0, \infty)$, with gcPCs enriched in A having eigenvalues > 1 .

161 gcPCA v3.0: gcPCA maximizing (A-B)/B

162 The second method developed is also asymmetrical but based on a relative change:

$$\arg \max_{\mathbf{x} : \mathbf{x}^T \mathbf{x} = 1} \frac{\mathbf{x}^T (\mathbf{C}_A - \mathbf{C}_B) \mathbf{x}}{\mathbf{x}^T (\mathbf{C}_B) \mathbf{x}} \quad (6)$$

163

164 This method is closely-related to v2.0 and is suitable for scenarios in which the investigator wishes
165 to define the gcPCs based on a relative change to a background condition (i.e., finding a 30% in-
166 crease in the neural activity in condition A relative to B). The eigenvalues returned are in the range
167 $[-1, \infty)$, with gcPCs enriched in A having eigenvalues > 0 .

168 gcPCA v4.0: gcPCA maximizing $(A-B)/(A+B)$
 169 The last of the three methods is based on a relative change:

$$\arg \max_{\mathbf{x} : \mathbf{x}^T \mathbf{x} = 1} \frac{\mathbf{x}^T (\mathbf{C}_A - \mathbf{C}_B) \mathbf{x}}{\mathbf{x}^T (\mathbf{C}_A + \mathbf{C}_B) \mathbf{x}} \quad (7)$$

170

171 This method is symmetrical, treating conditions A and B equally, and is appropriate for contrasting
 172 conditions in which B is a distinct condition and not merely a background to be removed, for exam-
 173 ple comparing neural data in sleep vs. wakefulness. The eigenvalues are in the range $[-1, 1]$ and
 174 are easily interpretable as a traditional index of the form $(A - B)/(A + B)$: 1 means that a gcPCA only
 175 accounts for variance in A , -1 means it only accounts for variance in B , and 0 means it accounts for
 176 equal variance in both. This method is fully symmetrical in the sense that switching A and B will
 177 yield the same gcPCs with the signs of the eigenvalues reversed.

178 gcPCA v2.1, v3.1, and v4.1: Orthogonal gcPCA

179 Unlike PCs or cPCs, gcPCs are not orthogonal by default. Because orthogonality may be important
 180 for some applications, we also include versions of gcPCA with an orthogonality constraint (see
 181 Methods). gcPCA v2.1 is the orthogonal version of 2.0, and v3.1 is the orthogonal version of v3.0,
 182 and v4.1 is the orthogonal version of v4.0.

183 Sparse vs. dense gcPCA

184 In high-dimensional datasets, it is often desirable to perform feature selection for easy interpre-
 185 tation of the results. With that in mind, we have also developed sparse versions of gcPCA. This
 186 method works by first finding the gcPCs based on the objective function selected, then performing
 187 feature selection using an $L1$ lasso penalty (see Methods). All non-orthogonal versions of gcPCA
 188 can be run in either sparse or dense mode, but sparsification cannot be used with the orthogonality
 189 constraint (Zou *et al.*, 2006).

Table 1. gcPCA variants in the gcPCA toolbox

gcPCA version	Symmetric	Orthogonal	Sparse solution	Note
v1	✗	✓	✓	Equivalent to cPCA
v2	✗	✗	✓	Objective function 5
v2.1	✗	✓	✗	Objective function 5
v3	✗	✗	✓	Objective function 6
v3.1	✗	✓	✗	Objective function 6
v4	✓	✗	✓	Objective function 7
v4.1	✓	✓	✗	Objective function 7

190 **Successful extraction of facial expression features with gcPCA**

191 To illustrate the utility of gcPCA with real-world datasets, we first use the Chicago Face Dataset (Ma
 192 *et al.*, 2015), which contains faces with different facial expressions. Here we used happy and angry
 193 expressions as condition A and neutral faces as condition B (Fig. 2A). This dataset is useful for
 194 two reasons: 1) the categorical separation of facial expressions allows an easy evaluation of the
 195 contrastive methods, and 2) the dimensions can be visually inspected for features that are being
 196 discovered by the method.

197 We first applied cPCA to the two conditions and used its automatic α selection algorithm to
 198 pick two different α values. The automatic α selection algorithm developed by *Abid et al. (2018)*
 199 finds representative α 's that yield different cPC embeddings so that the investigator can choose
 200 the appropriate α value. The algorithm is explained fully in the original paper (see supplementary

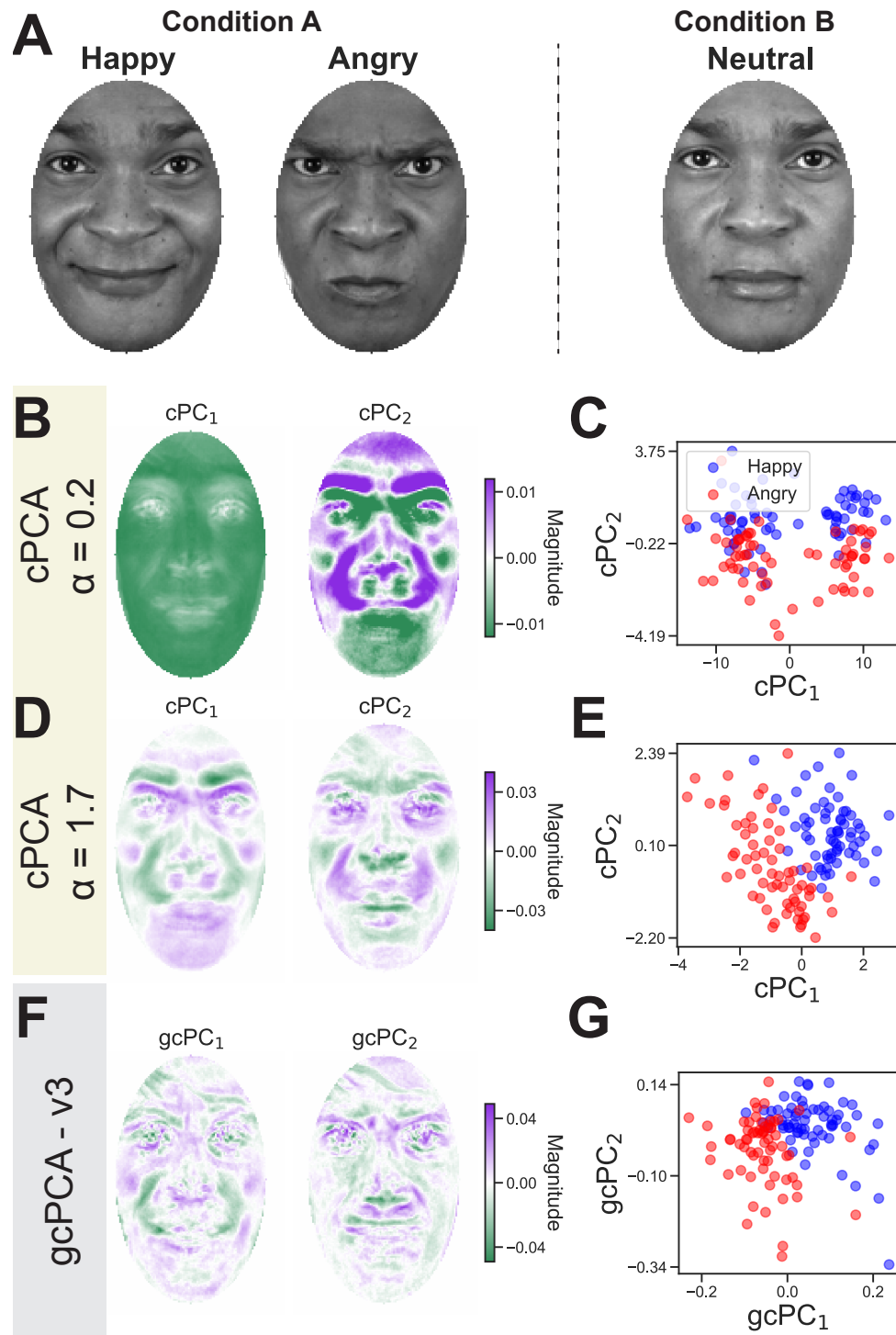


Figure 2. gcPCA correctly extracts contrastive facial features in the Chicago Face Dataset. **A** For this test, contrastive methods were applied to a set of happy and angry face images (Condition *A*) versus neutral face images from the same subjects (Condition *B*). Facial expression changes along the happy-angry axis were therefore the low-dimensional pattern that is enriched in condition *A*. **B** The first α identified by cPCA's algorithm, $\alpha = 0.2$, yields loadings similar to the first PCA dimensions, *i.e.* eigenfaces. **C** Projecting the faces onto the cPCs reveals clusters unrelated to the happy and angry facial expressions in condition *A*, indicating an incorrect solution. **D** cPCA with $\alpha = 1.7$ correctly reveals features associated with the expected facial expressions in condition *A*, such as furrowed eyebrows and the region around the mouth and nose. **E** Projecting the faces onto these cPCs reveals the separation of happy and angry faces along the first cPC. Importantly, it was only possible to determine this answer was correct because we knew the labels in advance. **F** Dimensions identified by gcPCA correctly reflect features related to the facial expressions in condition *A*. **G** Projecting the faces onto the first two gcPCs also reveals the separation of happy and angry faces along the first gcPC.

201 methods - algorithm 2 *Abid et al. (2018)*). Briefly, it calculates cPCA for an array of different α 's
202 (default: 40 different α values ranging from 0.01 to 1000 spaced on a log scale), defining a subspace
203 with the top k cPCs (default: 2 dimensions), and calculating an affinity matrix between subspaces
204 of different α 's. This affinity matrix is then clustered to find p clusters, and the medoid of each
205 cluster is a candidate α .

206 The first value returned is $\alpha = 0.2$, and we can see that the first two cPC loadings resemble
207 "eigenfaces" (*Turk and Pentland, 1991*), the largest principal components of facial images (Fig. 2B).
208 This suggests that this α is too small, leading cPCA to extract the highest-variance components in
209 condition *A*. Looking at the individual faces projected on these cPCs, two clusters can be identified
210 that separate cleanly along cPC1 (Fig. 2C). However, this solution is incorrect because these clusters
211 do not reflect the two facial expressions comprising condition *A*. Instead, they represent skin color,
212 which should account for equal variance in both datasets because images of the same subjects
213 are present in both conditions. For the second α value returned ($\alpha = 1.7$), the cPC loadings exhibit
214 features specific to condition *A* (Fig. 2D), and data projected onto these cPCs recovers the different
215 expressions (Fig. 2E). It is important to note that if we did not have the class labels *a priori*, we
216 could easily believe $\alpha = 0.2$ was the correct answer because it produces better clustering than for
217 $\alpha = 1.7$ (Fig. 2C,E). Using gcPCA v3.0 (asymmetric, non-orthogonal) also reveals features specific to
218 condition *A* (Fig. 2F), and the two expressions in the dataset can be distinguished by their projection
219 onto gcPC₁ (Fig. 2G), without the requirement of fine-tuning a hyperparameter.

220 **Applying gcPCA to neurophysiological recordings reveals hippocampal replay** 221 **without *a priori* knowledge of replay content**

222 A key application for gcPCA is neuronal recordings, which frequently contain hundreds of iso-
223 lated single units (*de Oliveira et al., 2022; Jun et al., 2017*). A well-studied neurophysiological
224 phenomenon is hippocampal replay, in which hippocampal neurons encoding spatial trajectories
225 traversed during a behavioral task "replay" the same activity patterns in post-task periods (*Wil-*
226 *son and McNaughton, 1994*). It is important to note that spatial location is a continuous variable,
227 so we are not expecting gcPCA or cPCA to find dimensions that cluster the activity into different
228 groups, as with the facial expression data. Instead, we are hoping to see replay of the firing pat-
229 terns that encode the linear track the animal recently explored (*Wilson and McNaughton, 1994;*
230 *Foster, 2017*). For our analysis, we used a previously published dataset recorded from hippocam-
231 pal CA1 *Girardeau et al. (2017)* where rats learn the location of an aversive air puff on a linear track.
232 The air puff is only delivered when the rat is running in one direction, called the danger run, and
233 the other direction is the safe run (Fig. 3A). Using previously-established methods (*Kudrimoti et al.,*
234 *1999*), *Girardeau et al. (2017)* found that hippocampal neurons exhibit reactivation of the task rep-
235 resentation in post-task activity when compared to pre-task activity. We tested whether cPCA or
236 gcPCA could extract hippocampal replay directly from neuronal activity by contrasting post-task
237 activity (condition *A*) with pre-task activity (condition *B*) (Fig. 3B). For this, we used cPCA and gcPCA
238 to extract cPCs or gcPCs from the pre- and post-task data, then projected the during-task data onto
239 the cPCs/gcPCs to test whether any spatial structure was discernible.

240 Using cPCA, it was not straightforward to detect replay. When we requested the cPCA algorithm
241 evaluate all the dimensions ($k = 48$), the α values identified by the automatic selection did not reveal
242 any obvious task-related spatial structure (Fig. 3C - left column, α values 0.44 and 701.70). When we
243 requested a smaller set of dimensions ($k = 2$), the automatic α selection returned several different
244 values, with one of them ($\alpha = 1.43$) revealing spatial structure in the task data (Fig. 3C - right column).
245 This reveals that although the only hyperparameter for cPCA is α in theory, the automatic alpha
246 selection algorithm depends on k , the number of components requested, constituting a second *de*
247 *facto* hyperparameter.

248 Applying gcPCA v4 (symmetrical, non-orthogonal), we readily recovered replay of task-related
249 spatial structure (Fig. 3D, top). Importantly, traditional replay analyses require knowledge of the
250 firing patterns during the task to test whether they were overrepresented in post-task sleep (*Kud-*

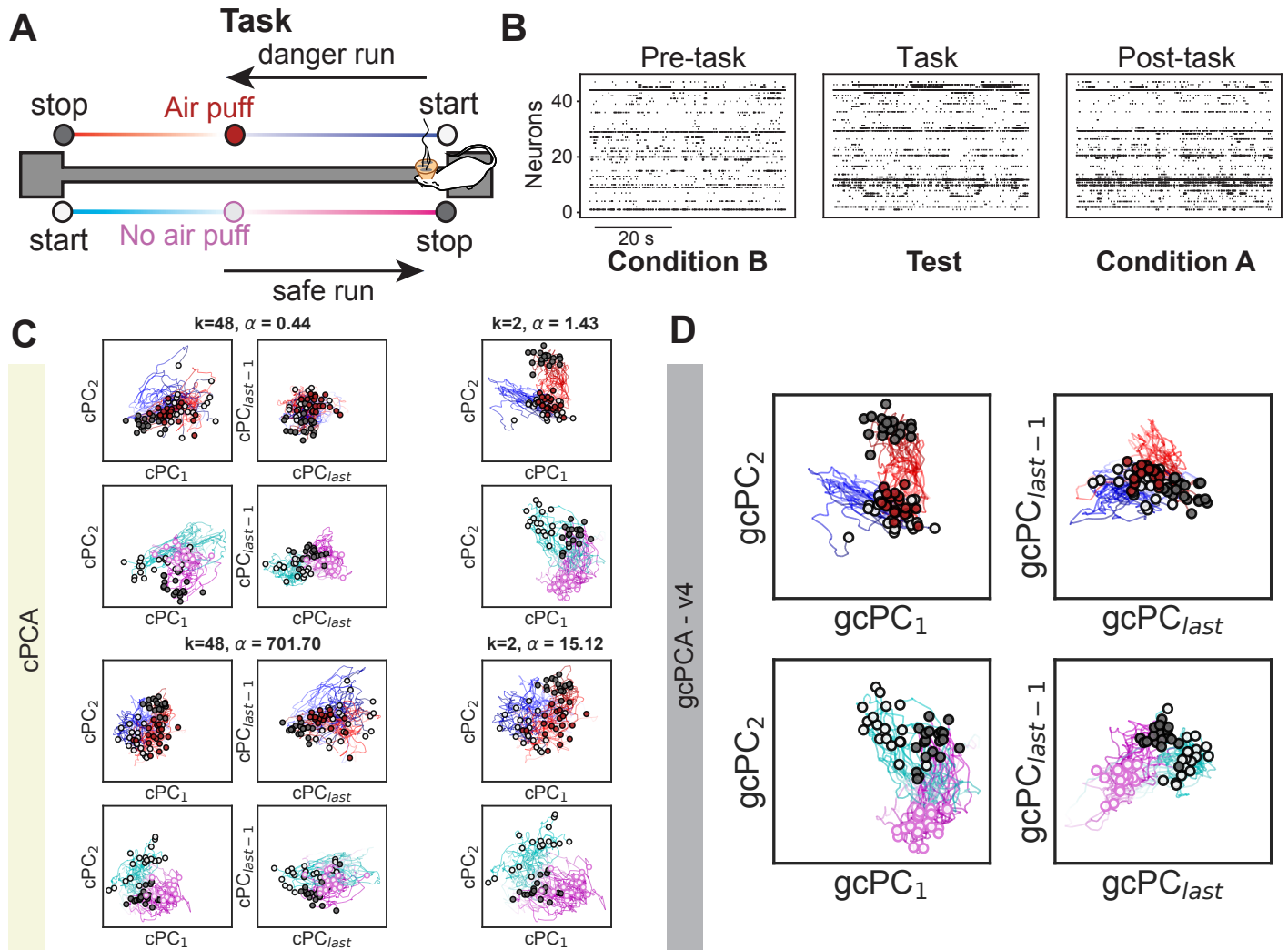


Figure 3. gcPCA v4.0 correctly identifies hippocampal replay in neurophysiological data without prior knowledge of replay content. A In *Girardeau et al. (2017)*, rats were trained to traverse a linear track where one direction has an aversive air puff and the other does not. Rats learned the location of the air puff and which direction was dangerous or safe. **B** *Girardeau et al. (2017)* recorded hippocampal neurons in pre- and post-task periods and used activity recorded during the task as a template to determine that that activity was "replayed" during post-task. We reanalyzed their published data using post-task activity as condition *A* and pre-task activity as condition *B* to identify the dimensions most enriched in post-task activity without taking task-related activity into account. **C** We first performed cPCA, then projected task-related neuronal data onto the cPCs to determine whether task-related data was enriched in the post-task period. The automatic α selection algorithm from cPCA returns various α values depending on the number of cPCs requested (parameter k). *Left Column* Representative α values returned with $k=48$ cPCs. cPC_{1-2} are the dimensions most enriched in post-task, and cPC_{last} and cPC_{last-1} are the most enriched in pre-task. No discernible spatial structure is identified, indicating that replay was not detected. *Right Column* With $k=2$, the α values returned are of different magnitudes, and one of them ($\alpha = 1.43$) reveals spatial structure related to the task, indicative of hippocampal replay. **D** gcPCA readily identifies the replay of the spatial task structure (left) with no parameter search.

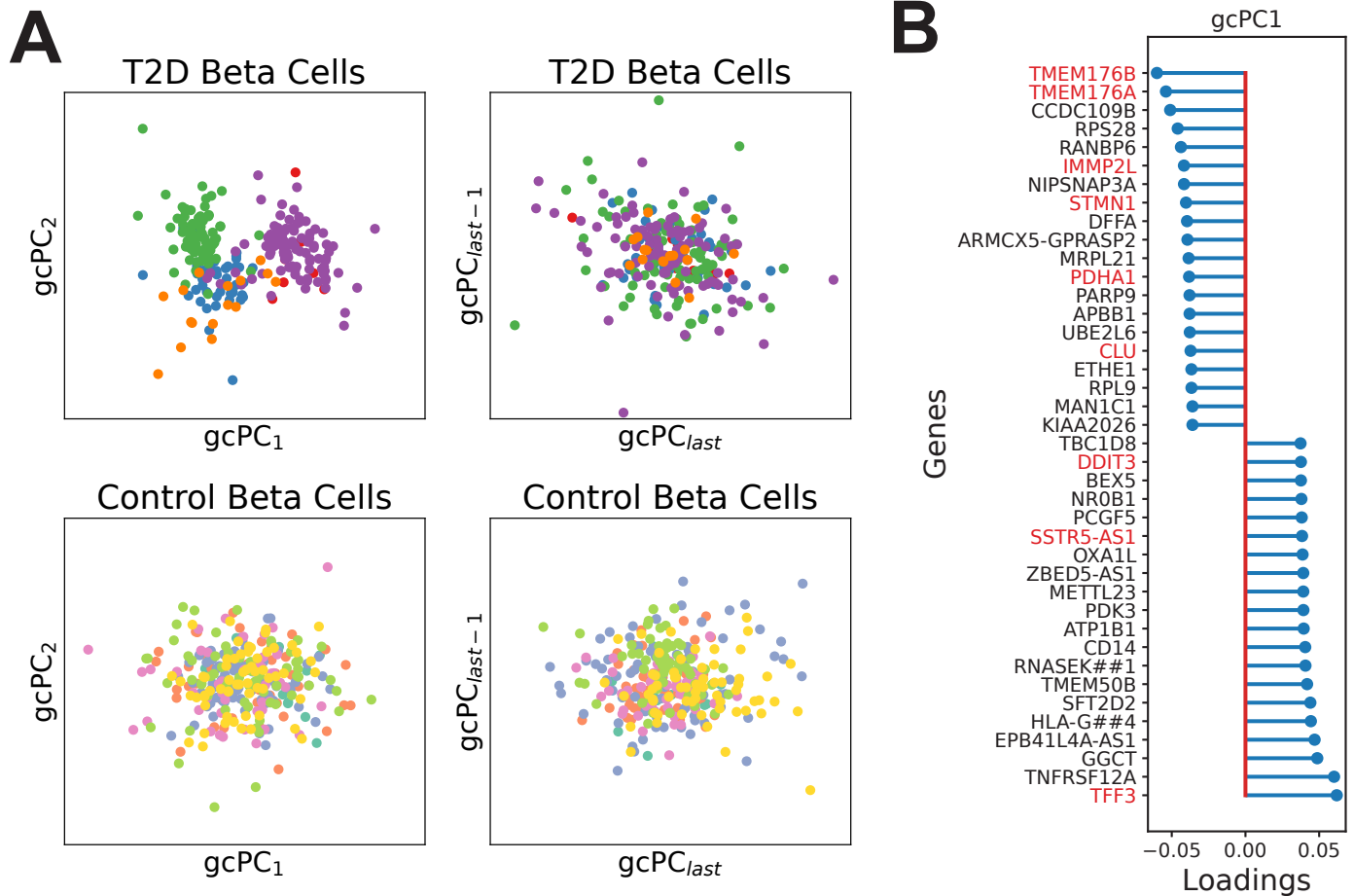


Figure 4. gcPCA v4.0 reveals possible disease heterogeneity in type II diabetes. **A** *Martínez-López et al. (2023)* performed scRNA-seq on isolated pancreatic cells from type II diabetes (T2D) patients and healthy controls. We used gcPCA v4.0 to compare the two conditions and found clustering of beta cells by donor identity (each color is a donor) in gcPC_{1,2} (top left panel). Such clustering was not found in the last gcPCs (top right panel) or in the control donors (bottom left and right panels) **B** The list of 40 genes with the highest loadings on gcPC₁ includes several previously linked to T2D (red). Notably, the top two hits, *TMEM176A* and *TMEM176B*, were shown by *Martínez-López et al. (2023)* to play functional roles in beta cell function.

251 *rimoti et al., 1999; Wilson and McNaughton, 1994; Foster, 2017*). gcPCA was able to extract signa-
252 tures of replay without prior knowledge of the task-related activity. This may prove useful in many
253 analogous situations in which the experimenter does not have prior knowledge of the pattern they
254 are searching for.

255 We also took advantage of the symmetric nature of gcPCA and investigated the patterns en-
256 riched in pre-task activity, $gcPC_{last}$ and $gcPC_{last-1}$. As expected, they did not exhibit task-related spa-
257 tial structure (Fig. 3D, top), suggesting they contain replay of environments other than the linear
258 track.

259 **Applying gcPCA to scRNA-seq data prioritizes disease genes and reveals disease** 260 **heterogeneity in type II diabetes**

261 Another key application of gcPCA is high-dimensional omics datasets, which often reflect a com-
262 parison of two conditions (*e.g.* disease vs. healthy). For this analysis, we used published single-cell
263 RNA sequencing data from pancreatic beta cells taken from healthy controls or patients with type
264 II diabetes (T2D) (*Martínez-López et al. (2023)*). We used T2D beta cells as condition *A* and healthy
265 control beta cells as condition *B* to test whether gcPCA 4.0 could identify groups of genes that
266 vary more among T2D patients or controls. $gcPC_1$ and $gcPC_2$ therefore represent axes along which
267 cells from patients vary more and the $gcPC_{last}$ and $gcPC_{last-1}$ represent axes along which control
268 cells vary more. We found that T2D and control data had similar levels of variability (Fig. 4A), but
269 in T2D patients this variability exhibited clear donor-based clustering (Fig. 4A, top left) that was
270 not observed in controls (Fig. 4A, bottom right). Several of the genes with the highest loadings
271 on $gcPC_1$ have been previously implicated in T2D (*IMMP2L (Diabetes Genetics Initiative of Broad*
272 *Institute of Harvard and MIT, Lund University, and Novartis Institutes of BioMedical Research*
273 *et al., 2007; Greenwald et al., 2019)*, *STMN1 (Horn et al., 2016)*, *PDHA1 (Srinivasan et al., 2010)*, *CLU*
274 *(Kim et al., 2001, 2006)*, *DDIT3 (Li et al., 2022; Yong et al., 2021)*, *SSTR5-AS1 (Jian and Felsenfeld,*
275 *2018)*, *TFF3 (Fueger et al., 2008)*, Fig. 4B, red), notably the top two hits, which were the *TMEM176A/B*
276 genes that *Martínez-López et al. (2023)* demonstrated are functionally important for T2D-related
277 beta-cell function. The fact that gcPCA identifies known T2D-related genes and that clustering is
278 specific to cases suggests that gcPCA is revealing disease heterogeneity (*Ahlqvist et al., 2020*).

279 **Discussion**

280 Discovering low-dimensional patterns that vary between conditions in high-dimensional datasets
281 is a crucial analysis in many research contexts. Here we present gcPCA, a method that achieves this
282 by examining the covariance structure of the datasets to find dimensions that exhibit the largest
283 relative changes in variance between conditions. This work builds on the pioneering insights in
284 the development of cPCA *Abid et al. (2018)* but solves cPCA's key problem, the requirement for
285 the hyperparameter α . Here we showed that the function of α is to compensate for bias toward
286 high-variance dimensions in noisy, finitely-sampled data. Further, we showed how this can be cir-
287 cumvented by introducing a normalization factor. Previous work *Abid et al. (2018); Boileau et al.*
288 *(2020)* has focused on developing and improving methods to find appropriate choices for α , but
289 with gcPCA we chose instead to eliminate the α hyperparameter entirely. Importantly, the advan-
290 tage of our approach is not merely that it is computationally cheaper than scanning a range of
291 α 's; it is that in most real-world cases there is no way to know whether a given choice of α yields a
292 correct solution.

293 We wish to address a common point of confusion by reiterating how gcPCA differs from LDA or
294 PLS. These are methods that find patterns optimally distinguishing two datasets, but gcPCA finds
295 patterns that exhibit more within-dataset variability in one dataset than another. As a fictitious ex-
296 ample, LDA might find a height/weight dimension distinguishing a university rugby team from the
297 general population because rugby players are taller and heavier. In contrast, gcPCA would be more
298 likely to find an age/education-level dimension because those features exhibit more variability in

299 the general population than in a university team. Despite the simplicity of this analysis, it reveals
300 interesting phenomena in high-dimensional biological data such as hippocampal replay (Fig. 3) or
301 transcriptomic heterogeneity in disease states (Fig. 4).

302 gcPCA has a few caveats: first, unlike ordinary PCs or cPCs, gcPCs are not orthogonal by default.
303 Our toolbox includes versions of gcPCA with an orthogonality constraint (v2.1, v3.1, and v4.1), which
304 comes at increased computational cost because a new eigendecomposition must be performed
305 for each gcPC. Second, the normalization factor introduces the possibility of numerical instability
306 if the denominator matrix is rank-deficient, meaning it has dimensions with zero (or near-zero)
307 variance that create a "divide-by-zero" situation. However, the implementation of gcPCA in the
308 toolbox automatically excludes these dimensions if they exist. Finally, there may be situations in
309 which cPCA's α could be a feature, rather than a bug, if the investigator has prior knowledge that
310 the patterns of interest will lie in high- or low-variance dimensions. Choosing an appropriate α
311 could then intentionally bias the analysis in favor of the results of interest. In such cases, it would
312 be relatively straightforward to extend gcPCA by adding a parameter that accomplishes a similar
313 result by adjusting the eigenspectrum of the denominator matrix in the objective function, e.g.
314 $(\mathbf{C}_A - \mathbf{C}_B)$. There are also other extensions that could be added, such as contrasting more than two
315 conditions or incorporating nonlinearity, which can be relevant to specific data problems. However,
316 we leave the development of such tools for future efforts.

317 The biological sciences are currently undergoing an explosion of technologies that produce
318 high-dimensional datasets, including novel forms of microscopy and neuroimaging, high-speed
319 video tracking, Neuropixels recordings, -omics approaches with single-cell resolution, and many
320 others. In addition to the analyses of electrophysiological recordings or single-cell RNA sequencing
321 data demonstrated here, gcPCA could be applied to any of these experimental modalities. We thus
322 anticipate that the open-source gcPCA toolbox will provide a valuable resource facilitating a broad
323 range of biological investigations that require contrasting two experimental conditions.

324 Methods

325 Generalized contrastive PCA

326 Our motivation for the following method stems from eliminating the necessity of the free param-
327 eter α in the contrastive PCA method. To accomplish this, we introduce a normalization factor
328 to mitigate the bias toward high-variance dimensions. We will summarize the process of calculat-
329 ing the gcPCs using gcPCA v4.0 as an example, but v2 and v3 are analogous. gcPCA v4.0 has the
330 following objective function, as shown in equation (eq. 8):

$$\arg \max_{\mathbf{x} : \mathbf{x}^T \mathbf{x} = 1} \frac{\mathbf{x}^T (\mathbf{C}_A - \mathbf{C}_B) \mathbf{x}}{\mathbf{x}^T (\mathbf{C}_A + \mathbf{C}_B) \mathbf{x}} \quad (8)$$

331
332 A potential problem of this objective function is the denominator creating numerical instability if
333 there are vectors that have eigenvalues approaching zero in the denominator covariance matrix.
334 To address this, we consider only the principal components (J) of that matrix that have non-zero
335 eigenvalues. In this case, the matrix J is composed of the principal components of the row-wise
336 concatenated datasets A and B that have non-zero eigenvalues. We then substitute for \mathbf{x} using
337 $\mathbf{x} = \mathbf{J}\boldsymbol{\gamma}$, yielding:

$$\arg \max_{\boldsymbol{\gamma} : \boldsymbol{\gamma}^T \boldsymbol{\gamma} = 1} \frac{\boldsymbol{\gamma}^T \mathbf{J}^T (\mathbf{C}_A - \mathbf{C}_B) \mathbf{J} \boldsymbol{\gamma}}{\boldsymbol{\gamma}^T \mathbf{J}^T (\mathbf{C}_A + \mathbf{C}_B) \mathbf{J} \boldsymbol{\gamma}} \quad (9)$$

338
339 The matrix $(\mathbf{J}^T (\mathbf{C}_A + \mathbf{C}_B) \mathbf{J})$ in the denominator is guaranteed to be positive definite, allowing us to
340 find a symmetric matrix \mathbf{M} that is its square root, yielding equation (eq. 10):

$$\arg \max_{\gamma : \gamma^T \gamma = 1} \frac{\gamma^T \mathbf{J}^T (\mathbf{C}_A - \mathbf{C}_B) \mathbf{J} \gamma}{\gamma^T \mathbf{M}^T \mathbf{M} \gamma} \quad (10)$$

341

342 Let $\mathbf{y} = \mathbf{M}\gamma$, yielding:

$$\arg \max_{\mathbf{y} : \mathbf{y}^T \mathbf{y} = 1} \frac{\mathbf{y}^T \mathbf{M}^{-1} \mathbf{J}^T (\mathbf{C}_A - \mathbf{C}_B) \mathbf{J} \mathbf{M}^{-1} \mathbf{y}}{\mathbf{y}^T \mathbf{y}} \quad (11)$$

343

344 This optimization problem can be solved with the eigendecomposition of the numerator matrix.

345 The vectors in \mathbf{X} are then calculated using equation (eq. 12):

$$\mathbf{X} = \mathbf{J} \mathbf{M}^{-1} \mathbf{y} \quad (12)$$

346

347 The column vectors in \mathbf{X} are referred to as gcPCs, following the term cPCs used in *Abid et al. (2018)*.

348 The other versions of gcPCA can be solved in the exact same way, by replacing $\mathbf{C}_A - \mathbf{C}_B$ with \mathbf{C}_A (v2)

349 and replacing $\mathbf{C}_A + \mathbf{C}_B$ with \mathbf{C}_B (v2 and v3).

350

351 In our case, the gcPCs (\mathbf{X}) are not guaranteed to be orthogonal due to the presence of the \mathbf{M}^{-1}
352 matrix. By default, gcPCA returns non-orthogonal gcPCs. If orthogonality is desired (as in gcPCA
353 v4.1), we iteratively shrink matrix \mathbf{J} to remove the subspace spanned by the gcPCs (*i.e.* columns of
354 \mathbf{X}) that have already been found. At each step, we compute the largest eigenvector \mathbf{x} of equation
355 11, then project it into the feature space with equation 12 and concatenate it column-wise into the
growing matrix \mathbf{X} . To shrink \mathbf{J} , we first regress out the \mathbf{x} from \mathbf{J} , as shown in equation 13:

$$\hat{\mathbf{J}} = \mathbf{J} - \mathbf{x} \mathbf{x}^T \mathbf{J} \quad (13)$$

356

357 We then use SVD to get the left singular vectors of $\hat{\mathbf{J}}$, and we define $\tilde{\mathbf{J}}$ as the first $n - i$ of these (on
358 the i -th iteration). $\tilde{\mathbf{J}}$ serves as an orthonormal basis for the subspace of \mathbf{J} that is orthogonal to \mathbf{X} ,
359 and we use $\tilde{\mathbf{J}}$ as the new \mathbf{J} in eq. 9 for the next iteration. This process continues until n gcPCs are
360 found, which can be the number of features in the dataset, the minimum rank of the conditions,
361 or a number specified by the user for the gcPCs to be extracted.

362 Sparse gcPCA

363 We developed an extension for sparse gcPCA using a similar approach as sparse PCA (*Zou et al.,*
364 *2006*) and sparse cPCA (*Boileau et al., 2020*). Here we will first review the sparse PCA framework,
365 then explain how we adapt it for gcPCA.

366 Sparse PCA method

367 Sparse PCA was first proposed as a reinterpretation of PCA as a regression problem. In brief, given
368 a matrix $\mathbf{X}_{n \times k}$, where the first k ordinary PCs are organized column-wise and are orthonormal, PCA
369 can be seen as minimizing the following objective:

$$\arg \min_{\mathbf{X}_{n \times k} : \mathbf{X}^T \mathbf{X} = \mathbf{I}} \|\mathbf{D} - \mathbf{D} \mathbf{X} \mathbf{X}^T\|^2 \quad (14)$$

370

371 Where \mathbf{D} is the data matrix of size features \times samples. To achieve sparse loadings in the first k PCs,
372 *Zou et al. (2006)* proposes the use of elastic net regularization, as shown in the following objective
function:

$$\arg \min_{\mathbf{Y}_{n \times k}, \mathbf{B}_{n \times k} : \mathbf{Y}^T \mathbf{Y} = \mathbf{I}} \|\mathbf{D} - \mathbf{D} \mathbf{B} \mathbf{Y}^T\|^2 + \lambda \sum_{j=1}^k \|\beta_j\|^2 + \sum_{j=1}^k \lambda_{1,j} \|\beta_j\|_1 \quad (15)$$

373 Where \mathbf{B} is the matrix containing the sparse PCs β_j , and \mathbf{Y} is a column-wise matrix that projects
 374 data from the sparse PC space to the feature space. The elastic net is the combination of the ridge
 375 penalty ($\lambda \sum_{j=1}^k \|\beta_j\|^2$) and the lasso penalty ($\sum_{j=1}^k \lambda_{1,j} \|\beta_j\|_1$). The ridge penalty is used to correct a
 376 rank-deficient matrix \mathbf{D} for numerical purposes. In *Zou et al. (2006)*, the same λ was used for all
 377 k components while using a different $\lambda_{1,j}$ for every component. To numerically solve equation 15,
 378 (*Zou et al., 2006*) use an alternating algorithm in which \mathbf{Y} is held constant as we solve for \mathbf{B} , then \mathbf{B}
 379 is held constant while we update \mathbf{Y} , and this is repeated until the algorithm converges. \mathbf{Y} is initially
 380 set to be equal to the ordinary PCs (\mathbf{X}), then we can find each column of \mathbf{B} by the following elastic
 381 net regression:

$$\hat{\beta}_j = \arg \min_{\beta_j} \|\mathbf{D}\mathbf{y}_j - \mathbf{D}\beta_j\|^2 + \lambda \|\beta_j\|^2 + \lambda_{1,j} \|\beta_j\|_1 \quad (16)$$

382 Next, \mathbf{B} is fixed, meaning the penalty terms can be ignored, and the new \mathbf{Y} is defined as:

$$\arg \min_{\mathbf{Y}^T \mathbf{Y} = \mathbf{I}_{k \times k}} \|\mathbf{D} - \mathbf{D}\mathbf{B}\mathbf{Y}^T\|^2 \quad (17)$$

383 The solution to 17 can be found by a reduced rank form of Procrustes rotation. Using SVD, we find

$$(\mathbf{D}^T \mathbf{D})\mathbf{B} = \mathbf{U}\mathbf{S}\mathbf{V}^T \quad (18)$$

384 We then set $\hat{\mathbf{Y}} = \mathbf{U}\mathbf{V}^T$. To solve eq. 16, *Zou et al. (2006)* has shown that is only necessary to know
 385 the Gram matrix $\mathbf{D}^T \mathbf{D}$. For a fixed \mathbf{Y} , finding β_j is equivalent to minimizing:

$$\begin{aligned} & \|\mathbf{D}\mathbf{y}_j - \mathbf{D}\beta_j\|^2 + \lambda \|\beta_j\|^2 + \lambda_{1,j} \|\beta_j\|_1 \\ & = (\mathbf{y}_j - \beta_j)^T \mathbf{D}^T \mathbf{D} (\mathbf{y}_j - \beta_j) + \lambda \|\beta_j\|^2 + \lambda_{1,j} \|\beta_j\|_1 \end{aligned} \quad (19)$$

386 If the covariance matrix of \mathbf{D} is known (denoted below as Σ), the term $\mathbf{D}^T \mathbf{D}$ can be replaced with
 387 Σ . For solving the eq. 16, the \mathbf{D} matrix can be replaced by $\Sigma^{\frac{1}{2}}$, which is the square root matrix of Σ .
 388 Resulting in the updated equation 20

$$\hat{\beta}_j = \arg \min_{\beta_j} \|\Sigma^{\frac{1}{2}} \mathbf{y}_j - \Sigma^{\frac{1}{2}} \beta_j\|^2 + \lambda \|\beta_j\|^2 + \lambda_{1,j} \|\beta_j\|_1 \quad (20)$$

389 Sparse gcPCA method

390 We can implement sparse gcPCA by adapting the sparse PCA method presented and following sim-
 391 ilar steps proposed in *Boileau et al. (2020)*. For sparse gcPCA, the covariance matrix Σ is replaced
 392 with the matrix Θ which reflects the appropriate objective function of the version used. Following
 393 gcPCA objective function 11, $\Theta = \mathbf{M}^{-1} (\mathbf{C}_A - \mathbf{C}_B) \mathbf{M}^{-1}$, where \mathbf{M} is the square root matrix of $\mathbf{C}_A + \mathbf{C}_B$.
 394 For version 1 (equivalent to cPCA), we instead use $\Theta = \mathbf{C}_A - \alpha \mathbf{C}_B$, and for versions 2 and 3 we
 395 change Θ to match their respective objective functions, as mentioned previously. We removed the
 396 \mathbf{J} matrix from the objective function so the sparsity is enforced in the features and not in the prin-
 397 cipal components. The components \mathbf{X} are the gcPCs identified by the ordinary gcPCA algorithm.
 398 Following the numerical solution for sparse PCA presented before, the sparse gcPCA is obtained
 399 by the following alternating algorithm until convergence:

400 **B given Y:** Each sparse gcPC (denoted here as β_j) was found according to the following elastic net
 401 solution:

$$\hat{\beta}_j = \arg \min_{\beta_j} \|\Theta^{\frac{1}{2}} \mathbf{y}_j - \Theta^{\frac{1}{2}} \beta_j\|^2 + \lambda \|\beta_j\|^2 + \lambda_{1,j} \|\beta_j\|_1 \quad (21)$$

402 Where \mathbf{y}_j is the j^{th} gcPC, and β is the sparse gcPC. The ridge penalty λ is used to fix rank-deficient
 403 matrices. To simplify our approach, we used the same λ_1 for all components instead of a different
 404 $\lambda_{1,j}$ for every j^{th} component. Therefore, the eq. 21 is reduced to a lasso regression and is solved
 405 through least angle regression, similar to *Zou et al. (2006)*.

406 **Y given B:** Using fixed \mathbf{B} , we can find a new $\hat{\mathbf{Y}}$ with a Procrustes rotation using SVD:

$$\Theta^{\frac{1}{2}} \Theta^{\frac{1}{2}} \mathbf{B} = \mathbf{USV}^T \quad (22)$$

407 We can then determine $\hat{\mathbf{Y}} = \mathbf{UV}^T$. These steps are repeated until loadings converge. The main
 408 caveat with this approach is that for gcPCA v3 or v4, the matrix Θ can have negative eigenvalues,
 409 which prevents the calculation of the square root matrix $\Theta^{\frac{1}{2}}$. To overcome this problem we follow
 410 similar steps to *Boileau et al. (2020)*, which we briefly replicate here. In gcPCA v3 or v4, positive
 411 and negative eigenvalues have a clear interpretation: for the matrix Θ , positive eigenvalues denote
 412 vectors with larger variance in condition A , while negative eigenvalues denote vectors with larger
 413 variance in condition B . We can then perform eigendecomposition of Θ and replace any negative
 414 eigenvalue with zeros to make Θ_+ positive semi-definite:

$$\begin{aligned} \Theta &= \mathbf{V}\Lambda\mathbf{V}^T \\ \Theta_+ &= \mathbf{V}\mathbf{S}\mathbf{V}^T \end{aligned} \quad (23)$$

where $\mathbf{S}_i = \begin{cases} \Lambda_i & \text{if } \Lambda_i > 0 \\ 0 & \text{otherwise} \end{cases}$

415 We can then define the square root matrix $\Theta_+ = \Theta_+^{\frac{1}{2}} \Theta_+^{\frac{1}{2}}$. We note that this matrix is only able to
 416 solve sparse gcPCA for vectors enriched in condition A . Here we propose a solution for finding
 417 the vectors also in condition B . As mentioned previously, the sign of the eigenvalues of matrix Θ
 418 indicates whether they are more expressed in condition $A(+)$ or $B(-)$. To solve the sparse gcPCA for
 419 condition B , we turn any positive eigenvalue to zero and switch the sign of negative eigenvalues to
 420 positive:

$$\begin{aligned} \Theta &= \mathbf{V}\Lambda\mathbf{V}^T \\ \Theta_- &= \mathbf{V}\mathbf{D}\mathbf{V}^T \end{aligned} \quad (24)$$

where $\mathbf{D}_i = \begin{cases} -\Lambda_i & \text{if } \Lambda_i < 0 \\ 0 & \text{otherwise} \end{cases}$

421 Although Θ_- does not contain negative eigenvalues, it still represents the dimensions most ex-
 422 pressed in condition B . This procedure is equivalent to switching the order of conditions A and B
 423 as the eigenvalues would be flipped in sign. The sparse gcPCs are found separately for Θ_+ and Θ_-
 424 and are later concatenated for the final sparse gcPCs.

425 **Synthetic data generation**

426 In the synthetic data, we generated two conditions with 1×10^5 samples and 100 dimensions. The
 427 dimensions were sampled from a Gaussian distribution ($\mu = 0$ and $\sigma = 1$) and then orthogo-
 428 nalized using singular value decomposition and picking the left singular vectors. In each condition,
 429 we created a pattern in the samples that was to be discovered. In condition A , we took dimensions
 430 71 and 72 and drew the samples from a uniform distribution $([0, 1])$. In dimension 71 we replaced
 431 any value from 0.3 and 0.7 with a different uniform distribution $([0, 0.4])$. In dimension 72 all the
 432 values between 0.4 and 0.6 were replaced with another uniform distribution $([0, 0.4])$. This created
 433 the square with a square hole in the middle in Fig. 1 A. The values were then offset by 0.5, the sam-
 434 ples were sorted by the angle they formed in each dimension, calculated by the inverse tangent
 435 ($\tan^{-1} \frac{x_{71}}{x_{72}}$). The samples in each dimension were normalized by their l_2 -norm. In condition B ,
 436 we generated the samples of dimensions 81 and 82 from a uniform distribution $([0, 1])$, sorted the
 437 sample values based on dimension 81, and rotated both dimensions by 45 degrees. This created

438 the diamond shape seen in Fig. 1 A. The sample values were later normalized by their l_2 -norm.
439 The samples for all the other dimensions were drawn from a Gaussian distribution ($\mu = 0$ and
440 $\sigma = 1$), and then normalized by their l_2 -norm. The magnitude of each dimension was established
441 as a line with a negative slope, starting at value 10 in the 1st dimension and ending at 0.001 in
442 the 100th dimension. For condition *A*, we doubled the magnitude in dimensions 71 and 72, while
443 in condition *B* we doubled the magnitude in dimensions 81 and 82. Identifying these changes in
444 magnitude is the goal of contrastive methods. Because the samples were drawn from a normal
445 distribution, the dimensions will display correlations among them. To estimate the total variance
446 explained by each dimension, we use a QR decomposition approach described in *Zou et al. (2006)*.
447 In brief, let Z be a matrix containing scores of each dimension generated, the variance is usually
448 calculated through $\text{tr}(Z^T Z)$, where tr is the trace of the matrix. However, in correlated scores, this
449 estimate is too optimistic. Using regression projection, it is possible to find the linear relationships
450 of the dimensions and correct to find the adjusted total variance. *Zou et al. (2006)* shows that this
451 is equivalent to using QR decomposition in Z , such that $Z = QR$ where Q is orthonormal and R is
452 upper triangular, and calculating the adjusted variance as follows:

$$\text{adjusted variance}_i = \sum_{j=1}^k R_{ij}^2 \quad (25)$$

453

454 **Face dataset**

455 For the facial expression analysis, we used the Chicago Face Database *Ma et al. (2015)*. This
456 database consists of neutral and emotional expression faces, and for the analysis we used a sub-
457 set of samples that had happy and angry faces alongside neutral ones. We used only male faces
458 to reduce variability in feature positioning. The images used were cropped by an ellipse (length
459 of 75 pixels and width of 45 pixels) centered in the face to focus on the facial expression rather
460 than other features such as hairstyle or shoulders. Each sample image was flattened from a two-
461 dimensional matrix to a vector, and all the flattened samples were then concatenated, resulting in
462 a matrix of samples x features. Each feature was z-scored and normalized by its l^2 -norm. Condition
463 *A* consisted of all the samples of happy and angry facial expressions, while condition *B* samples
464 were neutral expressions.

465 **Hippocampal electrophysiology data**

466 We used a previously published hippocampal electrophysiology dataset, with the experimental
467 details listed in the original publication *Girardeau et al. (2017)*. In brief, Long-Evans rats were im-
468 planted with silicon probes in the dorsal hippocampus CA1 region (either left or right hemisphere),
469 and neuronal activity was isolated through automatic spike sorting and manually curated. Animals
470 were trained to collect water rewards at the end of a linear track, and an air puff was introduced
471 at a fixed location for every lap, in only one of the directions. Recordings consisted of task, where
472 the animal learned the air puff location, and periods of pre- and post-task activity. For testing the
473 contrastive methods, we used pre-task recordings as condition *B* and post-task recordings as con-
474 dition *A*. For our analysis, we only used neurons that had a minimum firing rate of 0.01 spikes/s
475 during the task. We binned the neural data using a bin size of 10 ms and smoothed using a rolling
476 average with a Gaussian window of size 5 bins. The data was then z-scored and normalized by the
477 norm before testing the contrastive methods. The task data was then projected on the contrastive
478 dimensions for evaluation.

479 **Pancreatic single-cell RNA sequencing**

480 For the single-cell RNA sequencing data analysis, we used a previously published dataset *Martínez-
481 López et al. (2023)*, available at GEO accession GSE153855, consisting of scRNA-seq data from
482 human pancreatic islet cells from patients with type II diabetes and healthy controls. We used

483 the annotated dataset to identify the beta cells, which were identified previously by the authors
484 **Martínez-López et al. (2023)**. For condition *A* we used the beta cells from subjects that had type II
485 diabetes, and for condition *B* we used the beta cells from healthy patients. We used the expres-
486 sion values in reads per kilobase of the gene model and million mappable reads (RPKM). The
487 values were log-transformed, and all the features were centered before the analysis. gcPCA was
488 performed using the same set of genes used in the analysis by **Martínez-López et al. (2023)**.

489 Acknowledgments

490 We would like to thank Soyoun Kim, Ruben Coen-Cagli, Ehsan Sabri, Wenzhu Mowrey, Cleiton
491 Lopes-Aguiar, and members of the Sjulson and Batista-Brito lab for valuable conversations and
492 insightful comments on the manuscript. This work was supported by funds from the National In-
493 stitute on Drug Abuse (DP1 DA051608 and R01 DA051652 to LS), as well as from the Whitehall
494 Foundation and McManus Charitable Trust.

495 Author contributions

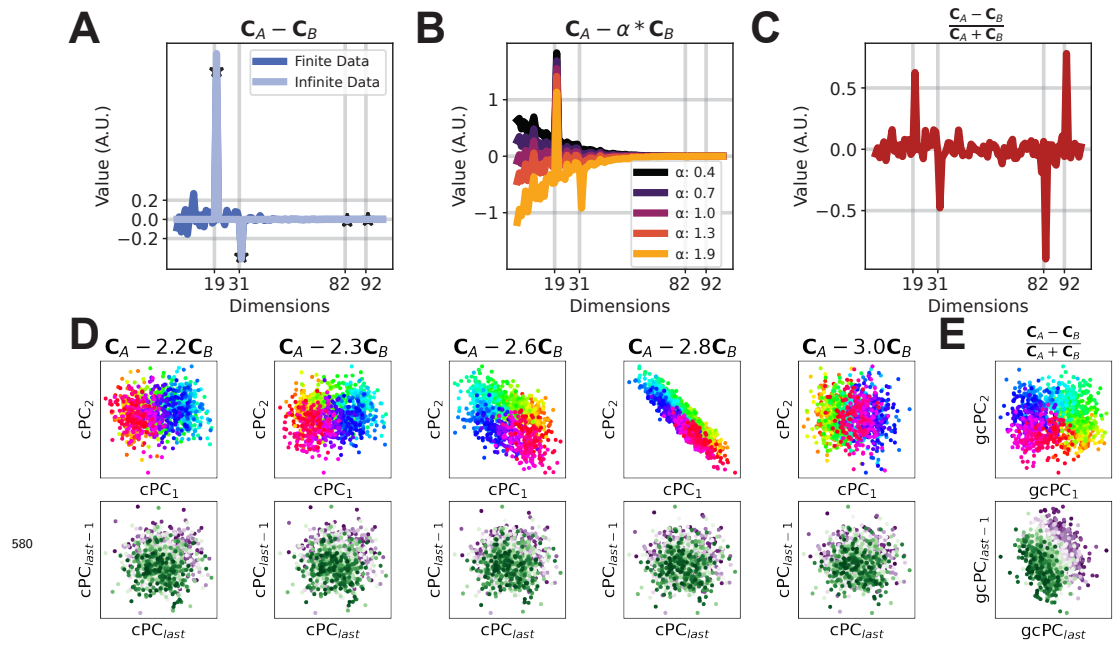
496 L.S. and E.F.O. conceived the project and developed the mathematical framework. E.F.O., P.G.,
497 and L.S. wrote the toolbox code and performed data analysis. L.S. supervised the project with
498 assistance from J.H.L. and R.B.B.

499 References

- 500 **Abid A**, Zhang MJ, Bagaria VK, Zou J. Exploring patterns enriched in a dataset with contrastive principal compo-
501 nent analysis. *Nat Commun*. 2018 May; 9(1):2134.
- 502 **Ahlqvist E**, Prasad RB, Groop L. Subtypes of type 2 diabetes determined from clinical parameters. *Diabetes*.
503 2020 Oct; 69(10):2086–2093.
- 504 **Boileau P**, Hejazi NS, Dudoit S. Exploring high-dimensional biological data with sparse contrastive principal
505 component analysis. *Bioinformatics*. 2020 Jun; 36(11):3422–3430.
- 506 **Brereton RG**, Lloyd GR. Partial least squares discriminant analysis: taking the magic away: PLS-DA: taking the
507 magic away. *J Chemom*. 2014 Apr; 28(4):213–225.
- 508 **Chapin JK**, Nicolelis MA. Principal component analysis of neuronal ensemble activity reveals multidimensional
509 somatosensory representations. *J Neurosci Methods*. 1999 Dec; 94(1):121–140.
- 510 **Chung NC**, Storey JD. Statistical significance of variables driving systematic variation in high-dimensional data.
511 *Bioinformatics*. 2015 Feb; 31(4):545–554.
- 512 **Diabetes Genetics Initiative of Broad Institute of Harvard and MIT, Lund University, and Novartis Insti-
513 tutes of BioMedical Research**, Saxena R, Voight BF, Lyssenko V, Burt NP, de Bakker PIW, Chen H, Roix JJ,
514 Kathiresan S, Hirschhorn JN, Daly MJ, Hughes TE, Groop L, Alshuler D, Almgren P, Florez JC, Meyer J, Ardlie K,
515 Bengtsson Boström K, Isomaa B, et al. Genome-wide association analysis identifies loci for type 2 diabetes
516 and triglyceride levels. *Science*. 2007 Jun; 316(5829):1331–1336.
- 517 **Foster DJ**. Replay Comes of Age. *Annu Rev Neurosci*. 2017 Jul; 40(1):581–602.
- 518 **Fueger PT**, Schisler JC, Lu D, Babu DA, Mirmira RG, Newgard CB, Hohmeier HE. Trefoil factor 3 stimulates
519 human and rodent pancreatic islet beta-cell replication with retention of function. *Mol Endocrinol*. 2008
520 May; 22(5):1251–1259.
- 521 **Girardeau G**, Inema I, Buzsáki G. Reactivations of emotional memory in the hippocampus–amygdala system
522 during sleep. *Nat Neurosci*. 2017 Sep; 20(11):1634–1642.
- 523 **Greenwald WW**, Chiou J, Yan J, Qiu Y, Dai N, Wang A, Nariai N, Aylward A, Han JY, Kadakia N, Regue L, Okino
524 ML, Drees F, Kramer D, Vinckier N, Minichiello L, Gorkin D, Avruch J, Frazer KA, Sander M, et al. Pancreatic
525 islet chromatin accessibility and conformation reveals distal enhancer networks of type 2 diabetes risk. *Nat*
526 *Commun*. 2019 May; 10(1):2078.

- 527 **Horn S**, Kirkegaard JS, Hoelper S, Seymour PA, Rescan C, Nielsen JH, Madsen OD, Jensen JN, Krüger M, Grønborg
528 M, Ahnfelt-Rønne J. Research resource: A dual proteomic approach identifies regulated islet proteins during
529 β -cell mass expansion in vivo. *Mol Endocrinol*. 2016 Jan; 30(1):133–143.
- 530 **Hotelling H**. Analysis of a complex of statistical variables into principal components. *J Educ Psychol*. 1933 Sep;
531 24(6):417–441.
- 532 **Izenman AJ**. Linear Discriminant Analysis. In: Izenman AJ, editor. *Modern Multivariate Statistical Techniques:*
533 *Regression, Classification, and Manifold Learning* New York, NY: Springer New York; 2008.p. 237–280.
- 534 **Jian X**, Felsenfeld G. Insulin promoter in human pancreatic β cells contacts diabetes susceptibility loci and
535 regulates genes affecting insulin metabolism. *Proc Natl Acad Sci U S A*. 2018 May; 115(20):E4633–E4641.
- 536 **Jun JJ**, Steinmetz NA, Siegle JH, Denman DJ, Bauza M, Barbarits B, Lee AK, Anastassiou CA, Andrei A, Aydın
537 Ç, Barbic M, Blanche TJ, Bonin V, Couto J, Dutta B, Gratiy SL, Gutnisky DA, Häusser M, Karsh B, Ledochow-
538 itsch P, et al. Fully integrated silicon probes for high-density recording of neural activity. *Nature*. 2017 Nov;
539 551(7679):232–236.
- 540 **Kim BM**, Han YM, Shin YJ, Min BH, Park IS. Clusterin expression during regeneration of pancreatic islet cells in
541 streptozotocin-induced diabetic rats. *Diabetologia*. 2001 Dec; 44(12):2192–2202.
- 542 **Kim BM**, Kim SY, Lee S, Shin YJ, Min BH, Bendayan M, Park IS. Clusterin induces differentiation of pancreatic
543 duct cells into insulin-secreting cells. *Diabetologia*. 2006 Feb; 49(2):311–320.
- 544 **Kudrimoti HS**, Barnes CA, McNaughton BL. Reactivation of hippocampal cell assemblies: effects of behavioral
545 state, experience, and EEG dynamics. *J Neurosci*. 1999 May; 19(10):4090–4101.
- 546 **Li J**, Klughammer J, Farlik M, Penz T, Spittler A, Barbieux C, Berishvili E, Bock C, Kubicek S. Single-cell transcrip-
547 tomes reveal characteristic features of human pancreatic islet cell types. *EMBO Rep*. 2016 Feb; 17(2):178–187.
- 548 **Li J**, Inoue R, Togashi Y, Okuyama T, Satoh A, Kyohara M, Nishiyama K, Tsuno T, Miyashita D, Kin T, Shapiro
549 AMJ, Chew RSE, Teo AKK, Oyadomari S, Terauchi Y, Shirakawa J. Imeglimin ameliorates β -cell apoptosis by
550 modulating the endoplasmic reticulum homeostasis pathway. *Diabetes*. 2022 Mar; 71(3):424–439.
- 551 **Ma DS**, Correll J, Wittenbrink B. The Chicago face database: A free stimulus set of faces and norming data.
552 *Behav Res Methods*. 2015 Dec; 47(4):1122–1135.
- 553 **Martínez-López JA**, Lindqvist A, Lopez-Pascual A, Chen P, Shcherbina L, Chriett S, Skene NG, Prasad RB, Lancien
554 M, Johnson PF, Eliasson P, Louvet C, Muñoz-Manchado AB, Sandberg R, Hjerling-Leffler J, Wierup N. Single-cell
555 mRNA-regulation analysis reveals cell type-specific mechanisms of type 2 diabetes; 2023.
- 556 **de Oliveira EF**, Kim S, Qiu TS, Peyrache A, Batista-Brito R, Sjulson L. Off-manifold coding in visual cortex revealed
557 by sleep; 2022.
- 558 **Pearson K**. LIII. *On lines and planes of closest fit to systems of points in space*. The London, Edinburgh, and Dublin
559 Philosophical Magazine and Journal of Science. 1901 Nov; 2(11):559–572.
- 560 **Peyrache A**, Benchenane K, Khamassi M, Wiener SI, Battaglia FP. Principal component analysis of ensemble
561 recordings reveals cell assemblies at high temporal resolution. *J Comput Neurosci*. 2010 Aug; 29(1-2):309–
562 325.
- 563 **Lopes-dos Santos V**, Ribeiro S, Tort ABL. Detecting cell assemblies in large neuronal populations. *J Neurosci*
564 *Methods*. 2013 Nov; 220(2):149–166.
- 565 **Sjulson L**, Peyrache A, Cumpelik A, Cassataro D, Buzsáki G. Cocaine Place Conditioning Strengthens Location-
566 Specific Hippocampal Coupling to the Nucleus Accumbens. *Neuron*. 2018 Jun; 98(5):926–934.e5.
- 567 **Srinivasan M**, Choi CS, Ghoshal P, Pliss L, Pandya JD, Hill D, Cline G, Patel MS. β -Cell-specific pyruvate dehydro-
568 genase deficiency impairs glucose-stimulated insulin secretion. *Am J Physiol Endocrinol Metab*. 2010 Dec;
569 299(6):E910–7.
- 570 **Tippling ME**, Bishop CM. Probabilistic Principal Component Analysis. *J R Stat Soc Series B Stat Methodol*. 1999
571 Sep; 61(3):611–622.
- 572 **Turk M**, Pentland A. Eigenfaces for recognition. *J Cogn Neurosci*. 1991; 3(1):71–86.
- 573 **Wilson MA**, McNaughton BL. Reactivation of hippocampal ensemble memories during sleep. *Science*. 1994
574 Jul; 265(5172):676–679.

- 575 **Yong J**, Parekh VS, Reilly SM, Nayak J, Chen Z, Lebeaupin C, Jang I, Zhang J, Prakash TP, Sun H, Murray S, Guo
576 S, Ayala JE, Satin LS, Saltiel AR, Kaufman RJ. Chop/Ddit3 depletion in β cells alleviates ER stress and corrects
577 hepatic steatosis in mice. *Sci Transl Med*. 2021 Jul; 13(604):eaba9796.
- 578 **Zass R**, Shashua A. Nonnegative Sparse PCA. *Adv Neural Inf Process Syst*. 2006 Dec; p. 1561–1568.
- 579 **Zou H**, Hastie T, Tibshirani R. Sparse Principal Component Analysis. *J Comput Graph Stat*. 2006; 15(2):265–286.



580