

Unsupervised discovery of temporal sequences in high-dimensional datasets, with applications to neuroscience

Emily L. Mackevicius^{1†}, Andrew H. Bahle^{1†}, Alex H. Williams², Shijie Gu^{1,3}, Natalia I. Denissenko¹, Mark S. Goldman^{4*}, Michale S. Fee^{1*}

*For correspondence:

fee@mit.edu (MSF);
msgoldman@ucdavis.edu
(MSG)

†These authors contributed
equally to this work

¹McGovern Institute for Brain Research, Department of Brain and Cognitive Sciences, MIT; ²Stanford Neurosciences Program; ³School of Life Sciences and Technology, ShanghaiTech University; ⁴Center for Neuroscience, Physiology and Behavior, UC Davis

1

Abstract

2 The ability to identify interpretable, low-dimensional features that capture the dynamics
3 of large-scale neural recordings is a major challenge in neuroscience. Repeated temporal
4 patterns (sequences) are not succinctly captured by traditional dimensionality reduction
5 techniques, so neural data is often aligned to behavioral task references. We describe a
6 task-independent, unsupervised method, which we call seqNMF, that provides a
7 framework for extracting sequences from high-dimensional datasets, and assessing the
8 significance in held-out data. We test seqNMF on simulated datasets under a variety of
9 noise conditions, and also on several neural datasets. In a hippocampal dataset, seqNMF
10 identifies neural sequences that match those calculated manually by reference to
11 behavioral events. In a songbird dataset, seqNMF discovers abnormal motor sequences
12 in birds that lack stereotyped songs. Thus, by identifying temporal structure directly from
13 neural data, seqNMF enables dissection of complex neural circuits in the absence of
14 reliable temporal references from stimuli or behavioral outputs.
15

16

Introduction

17 The ability to detect and analyze temporal sequences embedded in a complex sensory
18 stream is an essential cognitive function, and as such is a necessary capability of neuronal
19 circuits in the brain [13, 28, 5, 26], as well as artificial intelligence systems [14, 57]. The
20 detection and characterization of temporal structure in signals is also useful for the
21 analysis of many forms of physical and biological data. In neuroscience, recent advances
22 in technology for electrophysiological and optical measurements of neural activity have
23 enabled the simultaneous recording of hundreds or thousands of neurons [9, 33, 52, 29],
24 in which neuronal dynamics are often structured in sparse sequences [23, 24, 44, 19, 45].
25

26 Such sequences can be identified by averaging across multiple trials, but only in cases
27 where an animal receives a temporally precise sensory stimulus, or generates a sufficiently
28 stereotyped motor output.

29 However, it could be useful to extract sequences on a moment-to-moment basis
30 (without averaging), for example to study internal neuronal dynamics in the brain during
31 learning, sleep, or diseased states. In these applications, it is not possible to use external
32 timing references, and sequences must be extracted directly from the neuronal data. A
33 traditional unsupervised approach for directly extracting structure in neuronal data is
34 dimensionality reduction. Intuitively, sequences may be thought of as low dimensional,
35 and yet dimensionality reduction techniques such as PCA and NMF do not work for
36 sequences, because those methods only model synchronous patterns of activity.

37 Alternative approaches that search for repeating neural patterns require surprisingly
38 challenging statistical analysis [7, 41, 49]. While progress has been made in analyzing non-
39 synchronous sequential patterns using statistical models that capture cross-correlations
40 between pairs of neurons [51, 21, 53, 59, 22], such methods may not have statistical power
41 to scale to patterns that include many (more than a few dozen) neurons, may require long
42 periods ($\geq 10^5$ timebins) of stationary data, and may have challenges in dealing with (non-
43 sequential) background activity. For a review highlighting features and limitations of these
44 methods see [49]. Here we took an alternative, matrix factorization-based, approach that
45 aims to extract sequences. We reasoned that this approach would complement existing
46 methods by providing a more holistic and potentially simpler description of neural firing
47 dynamics.

48 One promising method for the unsupervised detection of temporal patterns is convo-
49 lutional non-negative matrix factorization (convNMF) [56, 55] (Figure 1), which has been
50 applied to the analysis of audio signals such as speech [43, 55, 61], as well as neural
51 signals [47]. ConvNMF identifies exemplar patterns (factors) in conjunction with the times
52 and amplitudes of pattern occurrences. This strategy eliminates the need to average
53 activity aligned to any external behavioral references. While convNMF produces excellent
54 reconstructions of the data, it does not automatically produce the minimal number of
55 factors required. Indeed, if the number of factors in the convNMF model is greater than
56 the true number of sequences, the algorithm returns overly complex and redundant
57 factorizations. These redundant factorizations are different each time the algorithm is run,
58 producing inconsistent results [47]. Notably, there is nothing in the convNMF algorithm
59 that favors the minimal factorization, as would be favored by the principle of ‘Occam’s
60 Razor’.

61 Here we describe a modification of the convNMF algorithm that suppresses redundant
62 factors, biasing the results toward factorizations with a minimal number of factors. This is
63 achieved by adding a penalty term to the convNMF cost function. Unlike other common
64 approaches such as sparsity regularization [65, 43, 50, 47] that constrain the make-up
65 of each factor, our regularization penalizes the correlations between factors that result
66 from redundant factorizations. We build on earlier applications of soft-orthogonality
67 constraints to NMF [10] to capture the types of temporally offset correlations that may
68 occur in the convolutional case.

69 Our algorithm, which we call seqNMF, produces minimal and consistent factorizations
70 in synthetic data, including under a variety of noise conditions, with high similarity to
71 ground-truth sequences. We further tested seqNMF on hippocampal spiking data in

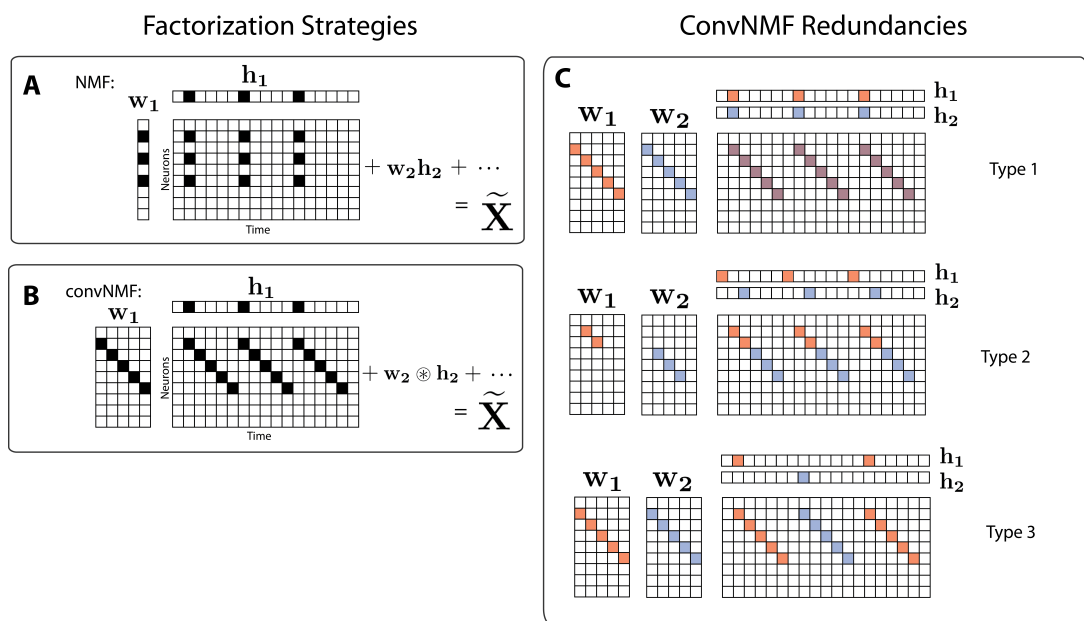


Figure 1. Introduction to convNMF factorization

(A) NMF (non-negative matrix factorization) approximates a dataset containing N neurons at T timepoints as a sum of K rank-one matrices. Each matrix is generated as the outer product of two nonnegative vectors: \mathbf{w}_k of length N , which stores a neural ensemble, and \mathbf{h}_k of length T , which holds the times at which the neural ensemble is active, and the relative amplitudes of this activity. (B) Convolutional NMF also approximates an $N \times T$ dataset as a sum of K matrices. Each matrix is generated as the convolution of two components: a non-negative matrix \mathbf{w}_k of dimension $N \times L$ that stores a sequential pattern of the N neurons at L lags, and a vector of temporal loadings, \mathbf{h}_k , which holds the times at which each factor pattern is active in the data, and the relative amplitudes of this activity. (C) Three types of inefficiencies are present in unregularized convNMF: Type 1 in which two factors are used to reconstruct the same instance of a sequence, Type 2 in which two factors reconstruct a sequence in a piecewise manner, and Type 3 in which two factors are used to reconstruct different instances of the same sequence. For each case, the factors (\mathbf{W} and \mathbf{H}) are shown, as well as the reconstruction ($\tilde{\mathbf{X}} = \mathbf{W} \otimes \mathbf{H} = \mathbf{w}_1 \otimes \mathbf{h}_1 + \mathbf{w}_2 \otimes \mathbf{h}_2$).

72 which neural sequences have previously been described. Finally, we use seqNMF to
 73 extract sequences in a functional calcium imaging dataset recorded in vocal/motor cortex
 74 of untutored songbirds that sing pathologically variable songs. We found that repeat-
 75 able neural sequences are activated in an atypical and overlapping fashion, suggesting
 76 potential neural mechanisms for this pathological song variability.

77 Results

78 Matrix factorization framework for unsupervised discovery of fea- 79 tures in neural data

80 Matrix factorization underlies many well known unsupervised learning algorithms [60]
 81 with applications to neuroscience [15], including principal component analysis (PCA) [46],
 82 non-negative matrix factorization (NMF) [34], dictionary learning, and k-means clustering.
 83 We start with a data matrix, \mathbf{X} , containing the activity of N neurons at T times. If the
 84 neurons exhibit a single repeated pattern of synchronous activity, the entire data matrix
 85 can be reconstructed using a column vector \mathbf{w} representing the neural pattern, and a row
 86 vector \mathbf{h} representing the times and amplitudes at which that pattern occurs (temporal
 87 loadings). In this case, the data matrix \mathbf{X} is mathematically reconstructed as the outer

88 product of \mathbf{w} and \mathbf{h}). If multiple patterns are present in the data, then each pattern can be
 89 reconstructed by a separate outer product, where the reconstructions are summed to
 90 approximate the entire data matrix (Figure 1A) as follows:

$$\mathbf{X}_{nt} \approx \tilde{\mathbf{X}}_{nt} = \sum_{k=1}^K \mathbf{W}_{nk} \mathbf{H}_{kt} = (\mathbf{WH})_{nt} \quad (1)$$

91 Where \mathbf{X}_{nt} is the $(nt)^{th}$ element of matrix \mathbf{X} . Here, in order to store K different patterns,
 92 \mathbf{W} is a $N \times K$ matrix containing the K exemplar patterns, and \mathbf{H} is a $K \times T$ matrix containing
 93 the K timecourses:

$$\mathbf{W} = \begin{bmatrix} | & | & & | \\ \mathbf{w}_1 & \mathbf{w}_2 & \cdots & \mathbf{w}_K \\ | & | & & | \end{bmatrix} \quad \mathbf{H} = \begin{bmatrix} - & \mathbf{h}_1 & - \\ - & \mathbf{h}_2 & - \\ \vdots & & \\ - & \mathbf{h}_K & - \end{bmatrix} \quad (2)$$

94 Given a data matrix with unknown patterns, the goal of these unsupervised learning
 95 algorithms is to discover a small set of patterns (\mathbf{W}) and a corresponding set of temporal
 96 loading vectors (\mathbf{H}) that approximate the data. In the case that the number of patterns
 97 (K) is sufficiently small (less than N and T), this corresponds to a dimensionality reduction,
 98 whereby the data is expressed in more compact form. NMF additionally requires that
 99 \mathbf{W} and \mathbf{H} must contain only non-negative numbers. The discovery of unknown factors
 100 is often accomplished by minimizing the following cost function, which measures (using
 101 the Frobenius norm, $\|\mathbf{M}\|_F = \sqrt{\sum_{ij} \mathbf{M}_{ij}^2}$) the element-by-element sum of all squared errors
 102 between a reconstruction $\tilde{\mathbf{X}} = \mathbf{WH}$ and the original data matrix \mathbf{X} :

$$(\mathbf{W}^*, \mathbf{H}^*) = \arg \min_{\mathbf{W}, \mathbf{H}} \|\tilde{\mathbf{X}} - \mathbf{X}\|_F^2 \quad (3)$$

103 The factors \mathbf{W}^* and \mathbf{H}^* that minimize this cost function produce an optimal recon-
 104 struction $\tilde{\mathbf{X}}^* = \mathbf{W}^* \mathbf{H}^*$. While this general strategy works well for extracting synchronous
 105 activity, it is unsuitable for discovering temporally extended patterns—first, because each
 106 element in a sequence must be represented by a different factor, and second, because
 107 NMF assumes that the columns of the data matrix are independent ‘samples’ of the
 108 data, so permutations in time have no effect on the factorization of a given dataset. It is
 109 therefore necessary to adopt a different strategy for temporally extended features.

110 *Convolutional non-negative matrix factorization (convNMF)*

111 Convolutional NMF (convNMF) [56, 55] extends NMF to provide a framework for extracting
 112 temporal patterns and sequences from data. While classical NMF represents each pattern
 113 as a single vector (Figure 1A), convNMF explicitly represents an exemplar pattern of neural
 114 activity over a brief period of time; the pattern is stored as an $N \times L$ matrix, where each
 115 column (indexed by $\ell = 1$ to L) indicates the activity of neurons at different timelags
 116 within the pattern (Figure 1B, where we call this matrix pattern \mathbf{w}_1 by analogy with NMF).
 117 The times at which this pattern/sequence occurs are stored using timeseries vector \mathbf{h}_1 ,
 118 as for NMF. The reconstruction is produced by convolving the $N \times L$ pattern with the
 119 timeseries \mathbf{h}_1 (Figure 1B).

120 If the dataset contains multiple patterns, each pattern is captured by a different $N \times L$
 121 matrix and a different associated timeseries vector \mathbf{h} . A collection of K different patterns

Table 1. Notation for convolutional matrix factorization

Shift operator

The operator $\overset{\ell \rightarrow}{\mathbf{H}}$ shifts a matrix \mathbf{H} in the \rightarrow direction by ℓ timebins:

$$(\mathbf{H})_{\cdot,t} \overset{\ell \rightarrow}{=} \mathbf{H}_{\cdot,(t-\ell)} \text{ and likewise } (\mathbf{H})_{\cdot,t} \overset{\ell \leftarrow}{=} \mathbf{H}_{\cdot,(t+\ell)}$$

where \cdot indicates all elements along the respective matrix dimension.

The shift operator inserts zeros when $(t - \ell) < 0$ or $(t + \ell) > T$

Tensor convolution operator

Convolutional matrix factorization reconstructs a data matrix \mathbf{X}

using a $N \times K \times L$ tensor \mathbf{W} and a $K \times T$ matrix \mathbf{H} :

$$\tilde{\mathbf{X}} = \mathbf{W} \otimes \mathbf{H} = \sum_{\ell} \mathbf{W}_{\cdot,\ell} \overset{\ell \rightarrow}{\mathbf{H}}$$

Note that each neuron n is reconstructed as the sum of k convolutions:

$$\tilde{\mathbf{X}}_{nt} = \sum_k \sum_{\ell} \mathbf{W}_{nk\ell} \mathbf{H}_{k(t-\ell)} \equiv (\mathbf{W} \otimes \mathbf{H})_{nt}$$

Transpose tensor convolution operator

The following quantity is useful in several contexts:

$$\mathbf{W} \otimes^{\top} \mathbf{X} = \sum_{\ell} (\mathbf{W}_{\cdot,\ell})^{\top} \overset{\ell \leftarrow}{\mathbf{X}}$$

Note that each element $(\mathbf{W} \otimes^{\top} \mathbf{X})_{kt} = \sum_{\ell} (\mathbf{W}_{\cdot,k\ell})^{\top} \mathbf{X}_{\cdot,(t+\ell)} = \sum_{\ell} \sum_n \mathbf{W}_{nk\ell} \mathbf{X}_{n(t+\ell)}$ measures

the overlap (correlation) of factor k with the data at time t

CNMF reconstruction

$$\mathbf{X} \approx \tilde{\mathbf{X}} = \sum_k \mathbf{W}_{\cdot,k} \otimes \mathbf{H}_k = \mathbf{W} \otimes \mathbf{H}$$

Note that NMF is a special case of CNMF, where $L = 1$

L1 norm excluding diagonal elements

For any $K \times K$ matrix \mathbf{C} , $\|\mathbf{C}\|_{1,i \neq j} \equiv \sum_k \sum_{j \neq k} \mathbf{C}_{jk}$

Special matrices

$\mathbf{1}$ is a $K \times K$ matrix of ones

\mathbf{I} is the $K \times K$ identity matrix

\mathbf{S} is a $T \times T$ smoothing matrix: $\mathbf{S}_{ij} = 1$ when $|i - j| < L$ and otherwise $\mathbf{S}_{ij} = 0$

122 can be compiled together into an $N \times K \times L$ array (also known as a tensor) \mathbf{W} and a
 123 corresponding $K \times T$ timeseries matrix \mathbf{H} . Analogous to NMF, convNMF generates a
 124 reconstruction of the data as a sum of K convolutions between each neural activity
 125 pattern (\mathbf{W}), and its corresponding temporal loadings (\mathbf{H}):

$$\mathbf{X}_{nt} \approx \tilde{\mathbf{X}}_{nt} = \sum_{k=1}^K \sum_{\ell=0}^{L-1} \mathbf{W}_{nk\ell} \mathbf{H}_{k(t-\ell)} \equiv (\mathbf{W} \circledast \mathbf{H})_{nt} \quad (4)$$

126 where the tensor/matrix convolution operator \circledast (notation summary, Table 1) reduces to
 127 matrix multiplication in the $L = 1$ case, which is equivalent to standard NMF. The quality
 128 of this reconstruction can be measured using the same cost function shown in Equation
 129 3, and \mathbf{W} and \mathbf{H} may be found iteratively using similar multiplicative gradient descent
 130 updates to standard NMF [34, 56, 55].

131 While convNMF can perform extremely well at reconstructing sequential structure,
 132 it can be challenging to use when the number of sequences in the data is not known
 133 [47]. In this case, a reasonable strategy would be to choose K at least as large as the
 134 number of sequences that one might expect in the data. However, if the number of
 135 sequences is than the actual number of sequences, convNMF will identify more significant
 136 factors than are minimally required. This is because each sequence in the data may
 137 be approximated equally well by a single sequential pattern or by a linear combination
 138 of multiple partial patterns. A related problem is that running convNMF from different
 139 random initial conditions produces inconsistent results, finding different combinations
 140 of partial patterns on each run [47]. These inconsistency errors fall into three main
 141 categories (Figure 1C):

- 142 • *Type 1*: Two or more factors are used to reconstruct the same instances of a se-
 143 quence.
- 144 • *Type 2*: Two or more factors are used to reconstruct temporally different parts of
 145 the same sequence, for instance the first half and the second half.
- 146 • *Type 3*: Identical factors are used to reconstruct different instances of a sequence.

147 Together, these inconsistency errors manifest as strong correlations between different
 148 redundant factors, as seen in the similarity of their temporal loadings (\mathbf{H}) and/or their
 149 exemplar activity patterns (\mathbf{W}).

150 *SeqNMF: A constrained convolutional non-negative matrix factorization*

151 Regularization is a common technique in optimization that allows the incorporation of
 152 constraints or additional information with the goal of improving generalization perfor-
 153 mance or simplifying solutions to resolve degeneracies [25]. To reduce the occurrence of
 154 redundant factors (and inconsistent factorizations) in convNMF, we sought a principled
 155 way of penalizing the correlations between factors by introducing a penalty term, \mathcal{R} , into
 156 the convNMF cost function of the following form:

$$(\mathbf{W}^*, \mathbf{H}^*) = \arg \min_{\mathbf{W}, \mathbf{H}} \left(\|\tilde{\mathbf{X}} - \mathbf{X}\|_F^2 + \mathcal{R} \right) \quad (5)$$

157 In this section, we will motivate a novel cost function that effectively minimizes the number
 158 of factors by penalizing spatial and temporal correlations between different factors. We
 159 will build up the full cost function by addressing, one at a time, the types of correlations
 160 generated by each type of error.

161 Regularization has previously been used in NMF to address the problem of duplicated
 162 factors, which, similar to Type 1 errors above, present as correlations between the \mathbf{H} 's [10].
 163 Such correlations are measured by computing the correlation matrix $\mathbf{H}\mathbf{H}^T$, which contains
 164 the correlations between the temporal loadings of every pair of factors. The regularization
 165 may be implemented using the penalty term $\mathcal{R} = \lambda \|\mathbf{H}\mathbf{H}^T\|_{1, i \neq j}$, where the norm $\|\cdot\|_{1, i \neq j}$
 166 sums the absolute value of every matrix entry except those along the diagonal (notation
 167 summary, Table 1) so that correlations between different factors are penalized, while
 168 the requisite correlation of each factor with itself is not. Thus, during the minimization
 169 process, similar factors compete, and a larger amplitude factor drives down the \mathbf{H} of a
 170 correlated smaller factor. The parameter λ controls the magnitude of the penalty term \mathcal{R} .

171 In convNMF, a penalty term based on $\mathbf{H}\mathbf{H}^T$ yields an effective method to prevent errors
 172 of Type 1, because it penalizes the associated zero lag correlations. However, it does
 173 not prevent errors of the other types, which exhibit different types of correlations. For
 174 example Type 2 errors result in correlated temporal loadings that have a small temporal
 175 offset and thus are not detected by $\mathbf{H}\mathbf{H}^T$. One simple way to address this problem is
 176 to smooth the \mathbf{H} 's in the penalty term with a square window of length $2L - 1$ using
 177 the smoothing matrix \mathbf{S} ($S_{ij} = 1$ when $|i - j| < L$ and otherwise $S_{ij} = 0$). The resulting
 178 penalty, $\mathcal{R} = \lambda \|\mathbf{H}\mathbf{S}\mathbf{H}^T\|$, allows factors with small temporal offsets to compete, effectively
 179 preventing errors of Type 1 and 2.

180 Unfortunately this penalty does not prevent errors of Type 3, in which redundant
 181 factors with highly similar patterns in \mathbf{W} are used to explain different instances of the
 182 same sequence. Such factors have temporal loadings that are segregated in time, and
 183 thus have low correlations, to which the cost term $\|\mathbf{H}\mathbf{S}\mathbf{H}^T\|$ is insensitive. One way to
 184 resolve errors of Type 3 might be to include an additional cost term that penalizes the
 185 similarity of the factor patterns in \mathbf{W} . A challenge with this approach is that, in the
 186 convNMF framework, there is no constraint on temporal translations of the sequence
 187 within \mathbf{W} . For example, if two redundant factors contain identical sequences that are
 188 simply offset by one time bin (in the L dimension), then these patterns would have zero
 189 correlation. Such offsets might be accounted for by smoothing the \mathbf{W} matrices in time
 190 before computing the correlation (Table 3), analogous to $\|\mathbf{H}\mathbf{S}\mathbf{H}^T\|$. The general approach
 191 of adding an additional cost term for \mathbf{W} correlations has the disadvantage that it requires
 192 setting an extra parameter, namely the λ associated with this cost.

193 Thus, we chose an alternative approach to resolve errors of Type 3 that simultaneously
 194 detects correlations in \mathbf{W} and \mathbf{H} using a single correlation cost term. We note that, for
 195 Type 3 errors, redundant \mathbf{W} patterns have a high degree of overlap with the data at the
 196 same times, even though their temporal loadings are segregated at different times. To
 197 introduce competition between these factors, we first compute, for each pattern in \mathbf{W} its
 198 overlap with the data at each time t . This quantity is captured in symbolic form by $\mathbf{W} \circledast \mathbf{X}$
 199 (See Table 1). We then compute the pairwise correlation between the temporal loading of
 200 each factor and the overlap of every other factor with the data. The correlation cost sums
 201 up these correlations across all pairs of factors, implemented as follows:

$$\mathcal{R} = \lambda \|\mathbf{W} \circledast \mathbf{X}\mathbf{S}\mathbf{H}^T\|_{1, i \neq j} \quad (6)$$

202 When incorporated into the update rules, this causes any factor that has a high overlap
 203 with the data to suppress the temporal loadings (\mathbf{H}) of any other factors active at that
 204 time. Thus, factors compete to explain each feature of the data, favoring solutions that

205 use a minimal set of factors to give a good reconstruction. We refer to this minimal set as
206 an efficient factorization. The resulting global cost function is:

$$(\mathbf{W}^*, \mathbf{H}^*) = \arg \min_{\mathbf{W}, \mathbf{H}} \left(\|\tilde{\mathbf{X}} - \mathbf{X}\|_F^2 + \lambda \|(\mathbf{W} \circledast \mathbf{X}) \mathbf{S} \mathbf{H}^\top\|_{1, i \neq j} \right) \quad (7)$$

207 The update rules for \mathbf{W} and \mathbf{H} are based on the derivatives of this global cost function,
208 leading to a simple modification of the standard multiplicative update rules used for NMF
209 and convNMF [34, 56, 55] (Table 3). Note that the addition of this correlation cost term
210 does not formally constitute regularization, because it also includes a contribution from
211 the data matrix \mathbf{X} , rather than just the model variables \mathbf{W} and \mathbf{H} .

212 Below, we test the performance of this penalty based on correlations between factors.
213 We will later consider different approaches to adding penalties to the convNMF cost
214 function, including an L1 norm penalty. We will also examine a parameter sweep of
215 the number of factors (K), as well as additional penalties to bias the tradeoff between
216 temporal or pattern correlations.

217 **Testing the performance of seqNMF on simulated sequences**

218 To compare the performance of seqNMF to unregularized convNMF, we simulated neural
219 sequences of a sort commonly encountered in neuronal data (Figure 2A). The simulated
220 data were used to test several aspects of the seqNMF algorithm: convergence, consis-
221 tency of factorizations, the ability of the algorithm to discover the correct number of
222 sequences in the data, and robustness to noise. As an initial pass, simulated datasets
223 were constructed by placing three ground-truth sequences at random non-overlapping
224 times. Each sequence ensemble consisted of 10 neurons evenly spaced throughout
225 a duration of 30 timesteps. The resulting data matrix had a total duration of 15000
226 timesteps and contained on average 60 ± 6 instances of each of the three sequences.
227 The seqNMF algorithm was run for 1000 iterations and reliably converged to a stable
228 asymptotic value of root-mean-squared-error (RMSE) (Figure 2B). RMSE reached to within
229 10% of the asymptotic value within 100 iterations.

230 *Consistency of seqNMF factorization*

231 We set out to determine if seqNMF exhibits the desirable property of consistency—namely
232 whether it returns similar sequences each time it is run on the same dataset using different
233 random initializations of \mathbf{W} and \mathbf{H} . Consistency was assessed as the extent to which there
234 is a good one-to-one match between factors across different runs (Methods 10). Due to
235 the inefficiencies outlined in Figure 1C, with K larger than the true number of sequences,
236 convNMF yielded low consistency scores typically ranging from 0.2 to 0.4 on a scale from
237 zero to one (Figure 2C, orange). In contrast, seqNMF factorizations were nearly identical
238 across different fits of noiseless data, producing consistency scores that were always
239 higher than any we measured for convNMF, and typically (>80% of the time) higher than
240 0.99 (Figure 2C, gray). Both convNMF and seqNMF had near-perfect reconstruction error
241 for all combinations of K and L that exceed the number and duration of sequences in
242 the data (not shown). However, convNMF exhibited low consistency scores, a problem
243 that was further exacerbated for larger values of K . In contrast, seqNMF exhibited high
244 consistency scores across a wide range of values of both K and L .

245 We also tested the consistency of seqNMF factorizations for the case in which a
246 population of neurons is active in multiple different sequences. Such neurons that are

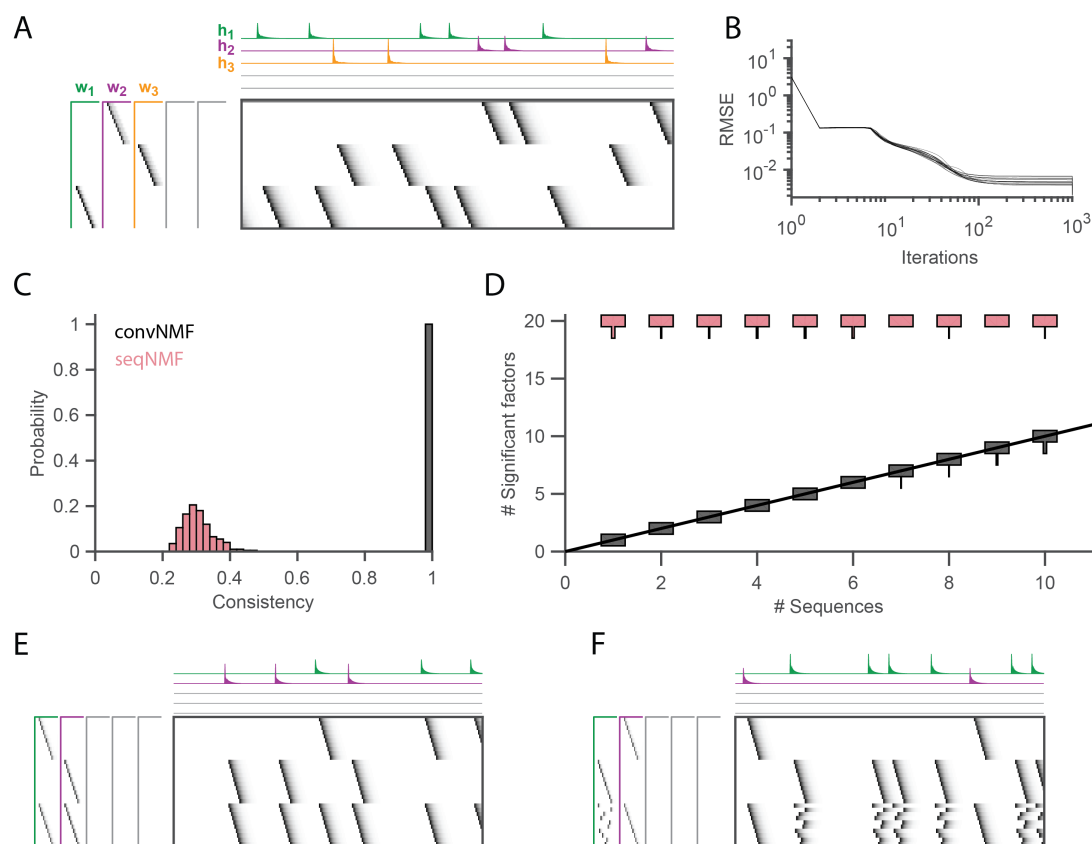


Figure 2. Testing seqNMF on simulated data

(A) A simulated dataset with three sequences and a seqNMF factorization ($K = 20$, $L = 50$, $\lambda = 0.003$). Each significant seqNMF factor is shown in a different color. At left are the exemplar patterns (\mathbf{W}) and on top are the timecourses (\mathbf{H}). (B) RMSE as a function of seqNMF iteration number. SeqNMF was run on a simulated dataset with three sequences and 15000 timebins (≈ 60 instances of each sequence). Twenty independent runs of seqNMF are shown. On this dataset, seqNMF converges to within 10% of the asymptotic error value within ≈ 100 iterations. (C) SeqNMF is more consistent than unregularized convNMF across 400 independent fits ($K = 20$, $L = 50$, $\lambda = 0.003$). (D) Plot showing the number of statistically significant factors vs. true number of simulated sequences for seqNMF and convNMF, for data containing between 1 and 10 sequences. Shown for each case is a vertical histogram representing the number of significant factors over 20 runs ($K = 20$, $L = 50$, $\lambda = 0.003$). (E) A seqNMF factorization of two simulated neural sequences with shared neurons that participate at the same latency in both sequences. (F) A seqNMF factorization of two simulated neural sequences with shared neurons that participate at different latencies in each sequence.

247 shared across different sequences have been observed in several neuronal datasets
248 [44, 45, 24]. For one test, we constructed two sequences in which shared neurons were
249 active at a common pattern of latencies in both sequences; in another test, shared
250 neurons were active in a different pattern of latencies in each sequence. In both tests,
251 seqNMF achieved near-perfect reconstruction error, and consistency was similar to the
252 case with no shared neurons (Figure 2E, F).

253 *Validating the statistical significance of extracted sequences*

254 To assess statistical significance, one can apply seqNMF to a subset of the data and
255 measure whether the extracted sequences appear in held-out data substantially more
256 than sequences drawn from a null model. We measured the appearance of sequences in
257 held-out data by measuring their overlap with held-out data, $\mathbf{W} \otimes \mathbf{X}^T$. The overlap is high
258 at timepoints at which the sequence occurs (relative to other timepoints). For a sequence
259 that matches ground truth in synthetic data, this distribution of overlap values exhibits a
260 heavy tail, indicating the presence of large outliers that correspond to times where the
261 extracted sequence appears in held out data. In contrast, a candidate sequence that does
262 not reliably occur in the held-out data produces a distribution of overlaps that appears
263 more symmetric (Figure S1).

264 While there are many ways of detecting outliers and quantifying “heavy-tailedness”
265 of a distribution, we use the skewness (the third central moment) as a simple measure.
266 In particular, we generate null distributions by circularly shifting the pattern matrices
267 \mathbf{W} along the time lag dimension (see Methods 10) and compare the skewness of these
268 distributions to the skewness of the distribution produced by the unshifted \mathbf{W} .

269 Runs of seqNMF on simulated and real data have revealed that the algorithm produces
270 two types of factors that can be immediately ruled out as candidate sequences: 1)
271 empty factors with zero amplitude in all neurons at all lags and 2) factors that have
272 amplitude in only one neuron. The latter case occurs often in datasets where one neuron
273 is substantially more active than other neurons, and thus accounts for a large amount
274 of variance in the data. SeqNMF also occasionally generates factors that appear to
275 capture one moment in the test data, especially in short datasets, where this can account
276 for a substantial fraction of the data variance. Such sequences are easily identified as
277 non-significant when tested on held-out data using the skewness test.

278 Note that if λ is set too small, seqNMF will produce multiple redundant factors to
279 explain one sequence in the data. In this case, each redundant candidate sequence will
280 pass the significance test outlined here. We will address below a procedure for choosing
281 λ and methods for determining the number of sequences.

282 *SeqNMF extracts the correct number of sequences in noise-free synthetic data*

283 A successful factorization should contain the same number of significant factors as
284 exist sequences in the data, at least in datasets for which the number of sequences
285 is unambiguous. To compare the ability of seqNMF and convNMF to recover the true
286 number of patterns in a dataset, we generated simulated noise-free data containing
287 between 1 and 10 different sequences. We then ran many independent fits of these data,
288 using both seqNMF and convNMF, and measured the number of significant factors. We
289 found that convNMF overestimates the number of sequences in the data, returning K
290 significant factors on nearly every run. In contrast, seqNMF tends to return a number of

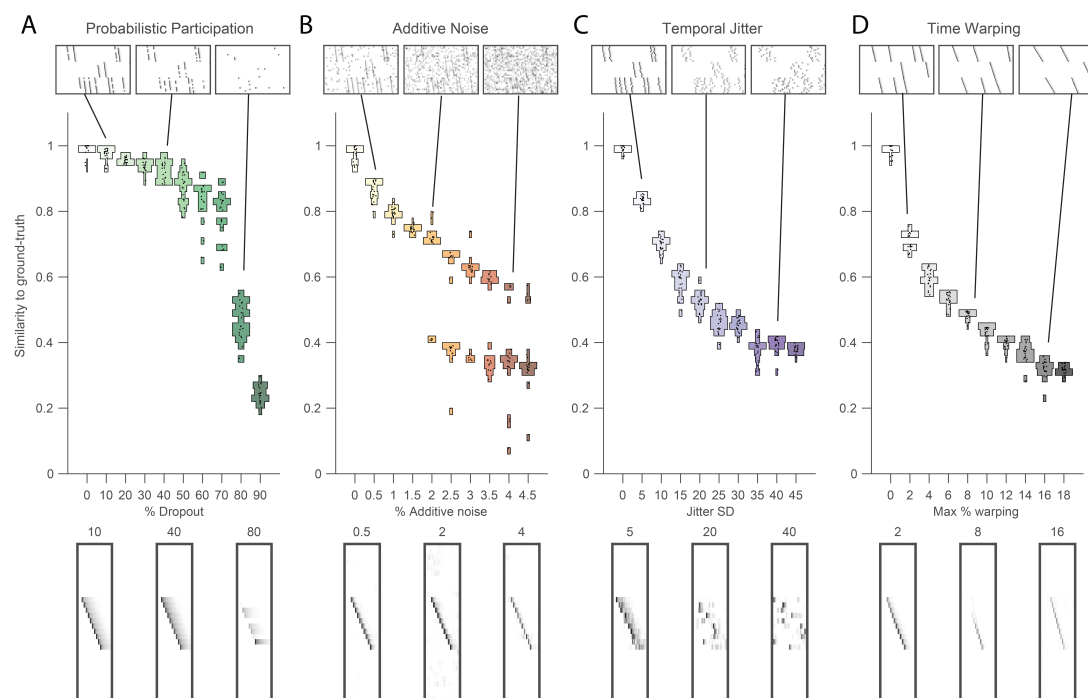


Figure 3. Testing seqNMF performance on sequences contaminated with noise

Performance of seqNMF was tested under 4 different noise conditions: **(A)** probabilistic participation, **(B)** additive noise, **(C)** timing jitter, and **(D)** sequence warping. For each noise type, we show: (top) examples of synthetic data at 3 different noise levels; (middle) similarity of seqNMF factors to ground-truth factors across a range of noise levels, showing 20 fits for each noise level; and (bottom) example of one of the W 's extracted at 3 different noise levels (same noise levels as data shown in top row). SeqNMF was run with $K = 20$, $L = 50$. In these examples, the algorithm was run with λ chosen using the automated procedure outlined in Figure 4. For results at different values of λ , see Figure S2.

291 significant factors that closely matches the actual number of sequences (Figure 2D).

292 *Robustness to noisy data*

293 SeqNMF was able to correctly extract sequences even in data corrupted by noise of types
 294 commonly found in neural data. We consider four common types of noise: participation
 295 noise, in which individual neurons participate probabilistically in instances of a sequence;
 296 additive noise, in which neuronal events occur randomly outside of normal sequence
 297 patterns; temporal jitter, in which the timing of individual neurons is shifted relative to
 298 their typical time in a sequence; and finally, temporal warping, in which each instance
 299 of the sequence occurs at a different randomly selected speed. To test the robustness
 300 of seqNMF to each of these noise conditions, we factorized data containing three neural
 301 sequences at a variety of noise levels. The value of λ was chosen using methods
 302 described in the next section. SeqNMF proved relatively robust to all four noise types,
 303 as measured by quantifying the similarity between seqNMF factors and ground-truth
 304 sequences (Methods section 10, Figure 3). For low noise conditions, seqNMF produced
 305 factors that were highly similar to ground-truth; this similarity gracefully declined as noise
 306 increased. Visualization of the extracted factors revealed that they tend to qualitatively
 307 match ground-truth sequences even in the presence of high noise (Figure 3). Together,
 308 these findings suggest that seqNMF is suitable for extracting sequence patterns from
 309 neural data with realistic forms of noise.

310 We also tested the performance of seqNMF as a function of dataset size. To do so,
311 we generated data of different sizes containing different numbers of instances of the
312 underlying ground-truth sequences, ranging from 1 to 20. For intermediate levels of
313 additive noise, we found that 3 examples of each sequence were sufficient for seqNMF
314 to correctly extract factors with similarity scores within 10% of asymptotic performance
315 (Figure S3).

316 *Method for choosing an appropriate value of λ*

317 Here we present procedures for guiding the choice of λ in seqNMF that address two
318 goals of regularization: to simplify the solution space of ill-posed problems and to reduce
319 overfitting. The choice of λ controls a trade-off between reconstruction accuracy and the
320 efficiency/consistency of the resulting factorizations (Figure 4). The goal is to reconstruct
321 only the repeating temporal patterns in the data and to do so with an efficient, maximally
322 uncorrelated set of factors. We will first describe a procedure that balances a measure
323 of correlation between factors with reconstruction error. We then describe a procedure
324 based on cross-validation in held-out data. Both of these procedures are validated under
325 a variety of noise conditions using simulated data for which the ground truth factors are
326 known.

327 In the first procedure we measure the effect of λ on both reconstruction error and
328 correlation cost in synthetic datasets containing three sequences (Figure 4). For any
329 given factorization, the reconstruction error is $\|\mathbf{W} \otimes \mathbf{H} - \mathbf{X}\|_F^2$, and the efficiency may be
330 estimated using the correlation cost term ($\|(\mathbf{W} \otimes \mathbf{X})\mathbf{S}\mathbf{H}^T\|_{1,i \neq j}$). SeqNMF was run with many
331 random initializations over a range of λ spanning six orders of magnitude. For small λ , the
332 behavior of seqNMF approaches that of convNMF, producing a large number of redundant
333 factors with high correlation cost. In the regime of small λ , correlation cost saturates at
334 a large value and reconstruction error saturates at a minimum value (Figure 4A). At the
335 opposite extreme, in the limit of large λ , seqNMF returns a single significant factor with
336 zero correlation cost because all other factors have been suppressed to zero amplitude.
337 In this limit, the single factor is unable to reconstruct multi-sequence data, resulting
338 in large reconstruction error. Between these extremes, there exists a region in which
339 increasing λ produces a rapidly increasing reconstruction error and a rapidly decreasing
340 correlation cost. Following the intuition that the optimal choice of λ for seqNMF would
341 lie in this cross-over region where the costs are balanced, we set out to quantitatively
342 identify, for known synthetic sequences, the optimal λ at which seqNMF has the highest
343 probability of recovering the correct number of significant factors, and at which these
344 factors most closely match the ground truth sequences.

345 The following procedure was implemented: For a given dataset, seqNMF is run several
346 times at a range of values of λ , and saturating values of reconstruction cost and correlation
347 cost are recorded (at the largest and smallest values of λ). Costs are normalized to vary
348 between 0 and 1, and the value of λ at which the reconstruction and correlation cost
349 curves intersect is determined (Figure 4B). This intersection point, λ_0 , then serves as a
350 precise reference by which to determine the correct choice of λ . We then separately
351 calibrated the reference λ_0 to the λ 's that performed well in synthetic datasets, with and
352 without noise, for which the ground-truth is known. This analysis revealed that values of λ
353 between $2\lambda_0$ and $5\lambda_0$ performed well across different noise types and levels (Figure 4B,C).
354 For additive noise, performance was better when λ was chosen to be near λ_0 , while with

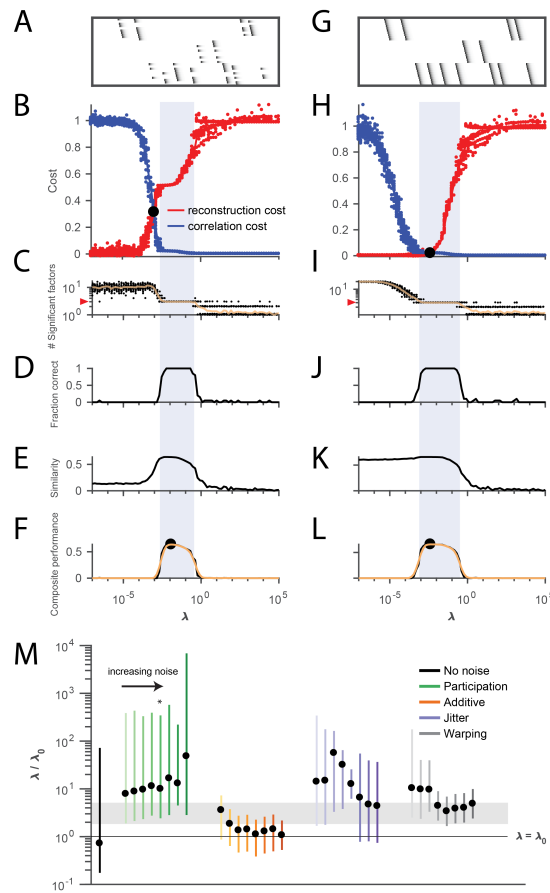


Figure 4. Procedure for choosing λ for a new dataset based on finding a balance between reconstruction cost and correlation cost

(A) Simulated data containing three sequences in the presence of participation noise (50% participation probability). This noise condition is used for the tests in (B-F). (B) Normalized reconstruction cost ($\|\tilde{\mathbf{X}} - \mathbf{X}\|_F^2$) and correlation cost ($\|(\mathbf{W} \otimes \mathbf{X})\mathbf{S}\mathbf{H}^T\|_{1, i \neq j}$) as a function of λ for these data. The cross-over point λ_0 is marked with a black circle. (C) The number of significant factors obtained from 20 fits of these data as a function of λ (mean number plotted in orange). The correct number of factors (three) is marked by a red triangle. (D) The fraction of fits returning the correct number of significant factors as a function of λ . (E) Similarity of the top three factors to ground-truth (noiseless) factors as a function of λ . (F) Composite performance measured as the product of the curves shown in (D) and (E), (smoothed curve plotted in orange with a circle marking the peak). Shaded region indicates the range of λ that works well (\pm half height of composite performance). (G) Simulated data containing three noiseless sequences. (H-L) Same as (B-F) but for the noiseless data. (M) Summary plot showing the range of values of λ (vertical bars), relative to the cross-over point λ_0 , that work well for each noise condition (\pm half height points of composite performance; note that this curve is a product of two other curves, and thus narrower, giving a conservative estimate of the range of effective λ s). Circles indicate the value of λ at the peak composite performance. For each noise type, results for the all non-zero noise levels from Figure 3 are shown (increasing color saturation at high noise levels; Green, participation: 90, 80, 70, 60, 50, 40, 30, and 20%; Orange, additive noise 0.5, 1, 2, 2.5, 3, 3.5, and 4%; Purple, jitter: SD of the distribution of random jitter: 5, 10, 15, 20, 25, 30, 35, 40, and 45 timesteps; Grey, timewarp: 13, 16, 20, 26, 33, 40, 46, and 53%). The noise type and level in panels (A-F) is indicated by *. Gray band indicates a range between $2\lambda_0$ and $5\lambda_0$. Values of λ in this range tended to perform well across the different noise conditions. In real data, it may be useful to explore a wider range of λ .

355 other noise types, performance was better at higher λ s ($\approx 5\lambda_0$). For all of the data shown
356 in Figure 3, we chose $\lambda = 2\lambda_0$. Figure S2 shows how choosing $\lambda = \lambda_0$ for additive noise
357 and $\lambda = 5\lambda_0$ for the other noise types yields slightly improved performance. Note that the
358 procedure for choosing λ does not need to be run on every dataset analyzed, rather, only
359 when seqNMF is applied to a new type of data for which a reasonable range of λ is not
360 already known. Similar ranges of λ appeared to work for datasets with different numbers
361 of ground-truth sequences—for the datasets used in Figure 2D, a range of λ between
362 0.001 and 0.01 returned the correct number of sequences at least 90% of the time for
363 datasets containing between 1 and 10 sequences (Figure S4). Furthermore, this method
364 for choosing λ also works on datasets containing sequences with shared neurons (Figure
365 S5).

366 Our second method for choosing λ directly tests generalization error by randomly
367 holding out a small subset of elements in the data matrix [64, 6] (Figure S6). This held-out
368 set is only used to test the performance of seqNMF, but is not used for fitting. At high
369 values of λ , seqNMF extracts only one factor, which exhibits similar reconstruction error on
370 training data and held-out test data. At low values of λ , seqNMF extracts a large number
371 of factors, yielding better reconstruction error on the training data, but the performance
372 of these factors on held-out data is often far worse, corresponding to overfitting. At
373 intermediate values of λ , within the optimal range described above, there was often a
374 minimum in the reconstruction error on held-out data (test error). This corresponds to
375 the classical approach for choosing regularization strength using cross-validation. In some
376 datasets, the minimum in test error can be subtle or nonexistent, so we instead identify
377 the λ corresponding to the rapid divergence between training error and test error (Figure
378 S6C). In many of our test datasets, this divergence point agrees with the ground-truth and
379 with the procedure described above based on the crossover between correlation cost
380 and reconstruction cost. One caution in using the cross-validation method to choose an
381 optimal λ is that it fails on synthetic datasets that have zero or very low noise (because
382 of a lack of overfitting), as well as in datasets with temporal warping. More broadly,
383 difficulties using cross-validation to choose λ may reflect that the primary function of the
384 seqNMF penalty term is to reduce factor correlations and redundancies, not to minimize
385 over-fitting.

386 *Can we choose K rather than choosing λ ?*

387 A goal of the seqNMF correlation cost term is to limit the factorization to a small number
388 of non-redundant factors. An alternative approach may be to directly constrain the
389 number of factors (K) in the convNMF algorithm without regularization. If the number
390 of underlying sequences in the data is unambiguous and is precisely known, as for the
391 simulated datasets described above, then this approach works well, yielding factorizations
392 close to ground truth sequences. We have found that the number of underlying sequences
393 can sometimes be estimated by running convNMF for all reasonable values of K and
394 selecting the value that yields the best cross-validated performance on held-out data. This
395 method works reasonably well for simulated datasets with participation noise, additive
396 noise, or temporal jitter over a range of noise levels that might be expected in real neural
397 data. In some cases, there is a clear minimum in the test error at the correct K. In other
398 cases there is a distinguishing feature such as a kink or a plateau in the test error as a
399 function of K that could potentially be used to estimate the correct number of sequences

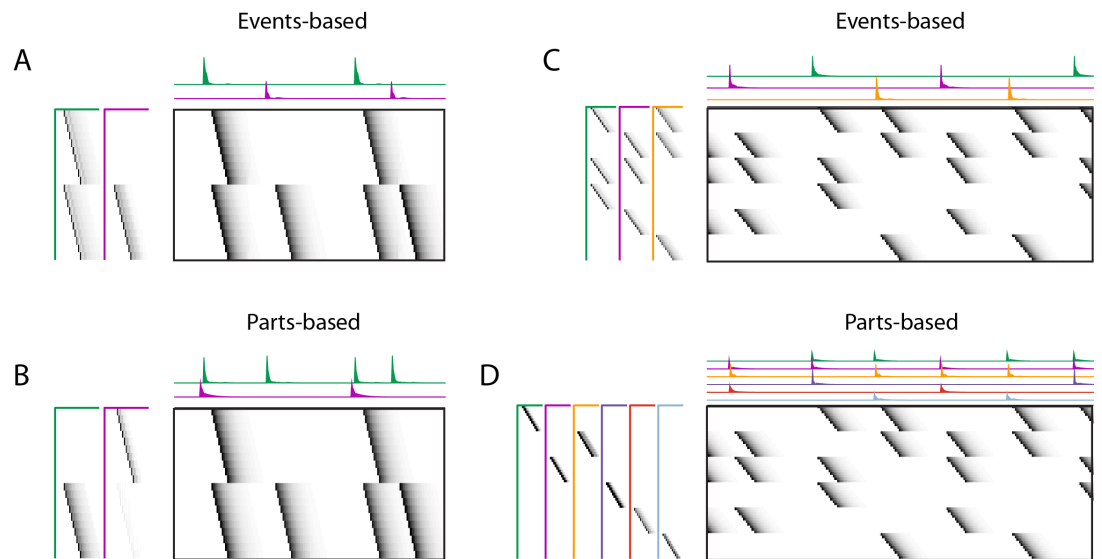


Figure 5. Using penalties to bias towards events-based and parts-based factorizations
 Datasets that have neurons shared between multiple sequences can be factorized in different ways, emphasizing discrete temporal events (events-based) or component neuronal ensembles (parts-based), by penalizing correlations in \mathbf{H} or \mathbf{W} respectively. (Left) A dataset with two different ensembles of neurons that participate in two different types of events, with (A) events-based and (B) parts-based factorizations. (Right) A dataset with six different ensembles of neurons that participate in three different types of events, with (C) events-based and (D) parts-based factorizations.

400 (Figure S7). Notably, this method fails to identify the number of underlying sequences in
 401 the case of temporal warping—an issue to which we will return in the next section.

402 *Strategies for dealing with ambiguous cases*

403 In some datasets, there is not a unique answer for the desired factorization of sequences.
 404 A common example of such ambiguity arises when neurons are shared between different
 405 sequences, as is shown in Figure 5A and B. In this case, there are two ensembles of
 406 neurons (1 and 2), that participate in two different types of events. In one event type,
 407 ensemble 1 is active alone, while in the other event type, ensemble 1 is coactive with
 408 ensemble 2. There are two different reasonable factorizations of these data. In one
 409 factorization the two different ensembles are separated into two different factors, while
 410 in the other factorization the two different event types are separated into two different
 411 factors. We refer to these as ‘parts-based’ and ‘events-based’ respectively. Note that
 412 these different factorizations may correspond to different intuitions about underlying
 413 mechanisms. ‘Parts-based’ factorizations will be particularly useful for clustering neurons
 414 into ensembles, and ‘events-based’ factorizations will be particularly useful for correlating
 415 neural events with behavior.

416 We have found that seqNMF and convNMF can produce either type of factorization,
 417 depending on initial conditions and the structure of shared neurons in the data. It
 418 may therefore be useful to explicitly control the tendency to produce these different
 419 factorizations by the addition of penalties on either \mathbf{W} or \mathbf{H} correlations. Note that in the
 420 ‘events-based’ factorization, the \mathbf{H} s are orthogonal (uncorrelated) while the \mathbf{W} s have high
 421 overlap; in the ‘parts-based’ factorization, the \mathbf{W} s are orthogonal while the \mathbf{H} s are strongly
 422 correlated. Note that these correlations in \mathbf{W} or \mathbf{H} are unavoidable in the presence of

423 shared neurons and the presence of such correlations does not indicate a redundant
424 factorization. Update rules to implement penalties on correlations in \mathbf{W} or \mathbf{H} are provided
425 in Table 3 with derivations in Appendix 1. Figure S9 shows examples of using these
426 penalties on the songbird dataset described in Figure 7.

427 Another type of ambiguity arises from the presence of systematic variations in the
428 amplitude or timing of neuronal participation in a sequence. A notable example of this
429 is data with temporal warping. In the case of high λ , seqNMF extracts a single factor
430 for the underlying ground truth sequence. In contrast, at lower λ seqNMF extracts
431 multiple factors for the underlying ground truth sequence, corresponding to slower and
432 faster variations of the sequence, effectively tiling the space of warped sequences at a
433 finer granularity depending on the strength of the penalty (λ). Note that each of these
434 factorizations corresponds to a reasonable interpretation, in the context of seqNMF,
435 for the same underlying timewarping process. Different neural datasets may require
436 estimating warping with different degrees of precision, depending on the behavior being
437 studied, leading to different reasonable choices of λ .

438 Another case requiring a choice between different reasonable levels of λ occurs when
439 a sequence exhibits two variants in which, for example, two subensembles of neurons
440 participate with different amplitudes in different instances of the sequence. Depending
441 on the desired level of granularity, controlled by the choice of λ , this dataset could be
442 factorized either as a single sequence or as two sequences. Any example in which a
443 sequence has multiple close variants, either in the timing or activity of different neurons,
444 can lead to this type of ambiguity. Depending on what type of factorization is desired,
445 a different value of λ might be preferable. In real datasets, it can be useful to explore
446 the factorization for different values of λ between λ_0 and $10\lambda_0$. There may often be a
447 range of λ that give rise to different reasonable factorizations. Note that high λ risks
448 missing sequences, especially sequences that occur rarely or include only a small number
449 of neurons, and low λ may give rise to redundant factors.

450 *Addition of a sparsity penalty to seqNMF or convNMF*

451 Sparsity regularization is a widely used strategy for achieving more interpretable and
452 generalizable results across a variety of algorithms and datasets [65], including convNMF
453 [43, 50]. In some of our datasets, we found it useful to include $L1$ regularization for
454 sparsity. The multiplicative update rules in the presence of $L1$ regularization are included
455 in Table 3, and as part of our code package. Sparsity on the matrices \mathbf{W} and \mathbf{H} may be
456 particularly useful in cases when sequences are repeated rhythmically (Figure S8). For
457 example, the addition of a sparsity regularizer on the \mathbf{W} update will bias the \mathbf{W} exemplars
458 to include only a single repetition of the repeated sequence, while the addition of a
459 sparsity regularizer on \mathbf{H} will bias the \mathbf{W} exemplars to include multiple repetitions of the
460 repeated sequence. This gives one fine control over how much structure in the signal
461 to pack into \mathbf{W} versus \mathbf{H} . Like the ambiguities described above, these are both valid
462 interpretations of the data, but each may be more useful in different contexts.

463 **Application of seqNMF to hippocampal sequences**

464 To test the ability of seqNMF to discover patterns in electrophysiological data, we analyzed
465 spiking activity in datasets of simultaneously recorded hippocampal neurons acquired in
466 the Buszaki lab and available from a public repository (<https://crcns.org/data-sets/hc>) [2, 1].

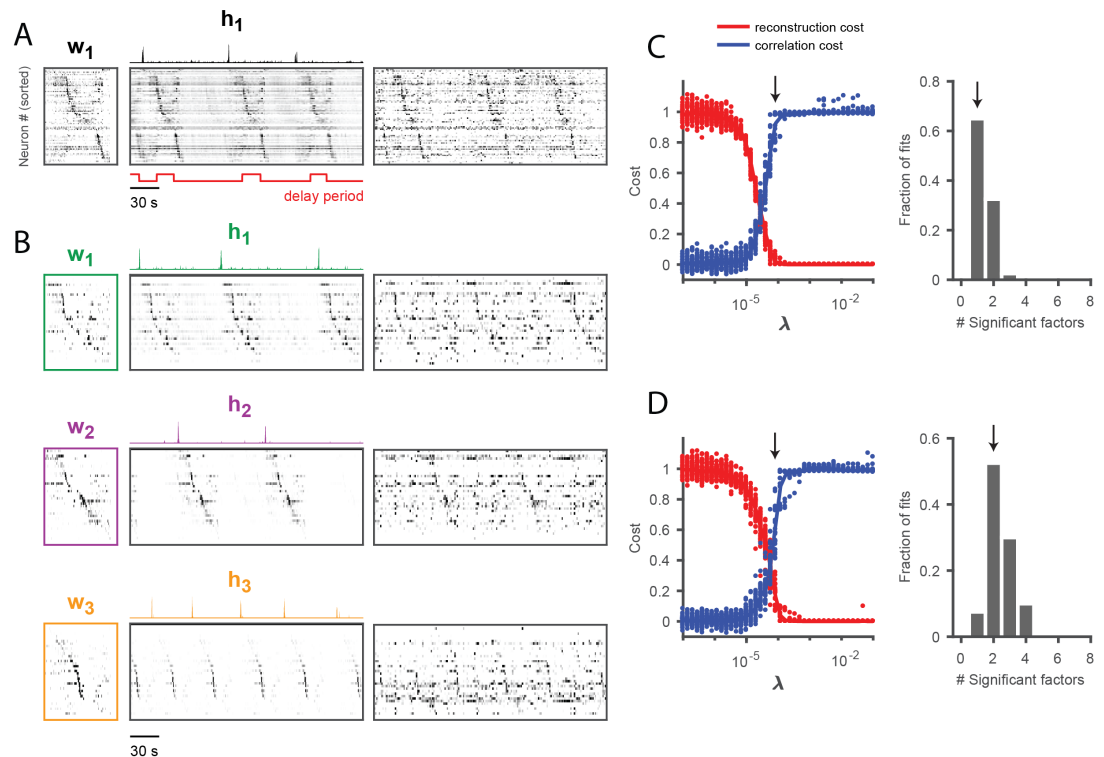


Figure 6. Application of seqNMF to extract hippocampal sequences from two rats **(A)** Firing rates of 110 neurons recorded in the hippocampus of Rat 1 during an alternating left-right task with a delay period [2], as well as the single significant extracted seqNMF factor. Neurons are sorted according to the latency of their peak activation within the factor. The red line shows the onset and offset of the forced delay periods, during which the animal ran on a treadmill **(B)** Firing rates of 43 hippocampal neurons recorded in Rat 2 during the same task [1]. Neurons are sorted according to the latency of their peak activation within each of the three significant extracted sequences. Both seqNMF reconstruction of each factor (left) and raw data (right) are shown. The first two factors correspond to left and right trials, and the third corresponds to running along the stem of the maze. **(C)** (Left) Reconstruction (red) and correlation (blue) costs as a function of λ for Rat 1. Arrow indicates $\lambda = 8 \times 10^{-5}$, used for seqNMF factorization shown in (A). (Right) Histogram of the number of significant factors across 30 runs of seqNMF. **(D)** Same as in (C) but for Rat 2. Arrow indicates $\lambda = 8 \times 10^{-5}$ used for factorization shown in (B).

467 The data were acquired in two rats as part of published studies describing sequences in
468 the hippocampus [45, 16]. In these experiments, rats were trained to alternate between
469 left and right turns in a T-maze to earn a water reward. Between alternations, the rats ran
470 on a running wheel during an imposed delay period lasting either 10 or 20 seconds. By
471 averaging spiking activity during the delay period, the authors reported long temporal
472 sequences of neural activity spanning the delay. In some rats, the same sequence
473 occurred on left and right trials, while in other rats, different sequences were active in the
474 delay period during the different trial types.

475 Without reference to the behavioral landmarks, seqNMF was able to extract sequences
476 in both datasets. The automated method described above was used to choose λ (Figure
477 6). In Rat 1, with $\lambda = 2\lambda_0$, most runs of seqNMF extracted a single significant factor,
478 corresponding to a sequence active throughout the running wheel delay period and
479 immediately after, when the rat runs up the stem of the maze (Figure 6B). Some runs
480 of seqNMF extracted two factors, splitting the delay period sequence and the maze
481 stem sequence; this is a reasonable interpretation of the data, and likely results from
482 variability in the relative timing of running wheel and maze stem traversal. At somewhat
483 lower values of λ , seqNMF more often split these sequences into two factors. At even
484 lower values of λ , seqNMF produced more significant factors. Such higher granularity
485 factorizations may correspond to real variants of the sequences, as they generalize to
486 held-out data (Figure S7J).

487 In Rat 2, at a λ of $1.5\lambda_0$, three significant factors were typically identified (Figure 6C).
488 The first two correspond to distinct sequences active for the duration of the delay period
489 on alternating trials. The third sequence was active immediately following each of the
490 alternating sequences, corresponding to the time at which the animal exits the wheel
491 and runs up the stem of the maze. Taken together, these results suggest that seqNMF
492 can detect multiple neural sequences without the use of any behavioral landmarks.
493 Having validated this functionality in both simulated data and previously published neural
494 sequences, we then applied seqNMF to find structure in a novel dataset, in which the
495 ground truth is unknown and difficult to ascertain using previous methods.

496 **Application of seqNMF to abnormal sequence development in avian** 497 **motor cortex**

498 We applied seqNMF to analyze new functional imaging data recorded in songbird HVC
499 during singing. Normal adult birds sing a highly stereotyped song, making it possible to
500 detect sequences by averaging neural activity aligned to the song. Using this approach, it
501 has been shown that HVC neurons generate precisely timed sequences that tile each song
502 syllable [23, 48, 37]. In contrast to adult birds, young birds sing highly variable babbling
503 vocalizations, known as subsong, for which HVC is not necessary [3]. The emergence of
504 sequences in HVC occurs gradually over development, as the song matures from subsong
505 to adult song [44].

506 Songbirds learn their song by imitation and must hear a tutor to develop normal adult
507 vocalizations. Birds isolated from a tutor sing highly variable and abnormal songs as
508 adults [18]. Such 'isolate' birds provide an opportunity to study how the absence of normal
509 auditory experience leads to pathological vocal/motor development. However, the high
510 variability of pathological 'isolate' song makes it difficult to identify neural sequences
511 using the standard approach of aligning neural activity to vocal output.

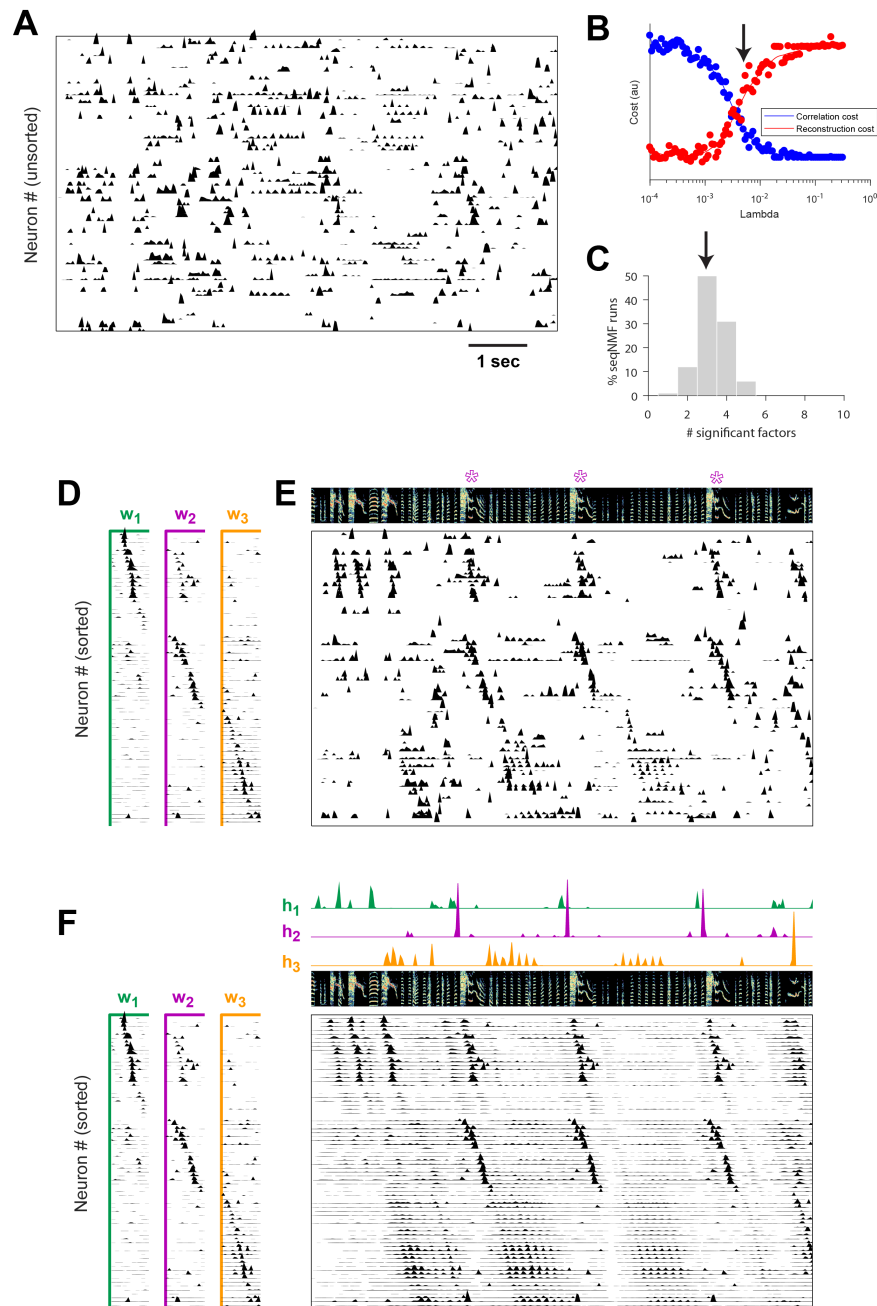


Figure 7. SeqNMF applied to calcium imaging data from a singing isolate bird reveals abnormal sequence deployment

(A) Functional calcium signals recorded from 75 neurons, unsorted, in a singing isolate bird. **(B)** Reconstruction and correlation cost as a function of λ . The arrow at $\lambda = 0.005$ indicates the value selected for the rest of the analysis. **(C)** Number of significant factors for 100 runs of seqNMF with $K = 10$, $\lambda = 0.005$. Arrow indicates 3 is the most common number of significant factors. **(D)** SeqNMF factor exemplars (W 's), Neurons are grouped according to the factor in which they have peak activation, and within each group neurons are sorted by the latency of their peak activation within the factor **(E)** The same data shown in (A), after sorting neurons by their latency within each factor as in (D). A spectrogram of the bird's song is shown at top, with a purple "*" denoting syllable variants correlated with w_2 . **(F)** Same as (E), but showing reconstructed data rather than calcium signals. Shown at top are the temporal loadings (**H**) of each factor.

512 Using seqNMF, we were able to identify repeating neural sequences in isolate song-
513 birds (Figure 7A). At the chosen λ (Figure 7B), seqNMF typically extracts three significant
514 sequences (Figure 7C). Similarly, our masked cross-validation test indicated good convNMF
515 performance at $K = 3$, with over-fitting starting at $K = 4$ (Figure S7I). The extracted
516 sequences include sequences deployed during syllables of abnormally long and variable
517 durations (Figure 7D-F).

518 In addition, the extracted sequences exhibit properties not observed in normal adult
519 birds. We see an example of two distinct sequences that sometimes, but not always,
520 co-occur (Figure 7). We observe that a short sequence occurs alone on some syllable
521 renditions, while on other syllable renditions, a second longer sequence is generated
522 simultaneously. This probabilistic overlap of different sequences is highly atypical in normal
523 adult birds [23, 36, 48, 37]. Furthermore, this pattern of neural activity is associated
524 with abnormal variations in syllable structure—in this case resulting in a longer variant
525 of the syllable when both sequences co-occur. This acoustic variation is a characteristic
526 pathology of isolate song [18]. Thus, even though we observe HVC generating some
527 sequences in the absence of a tutor, it appears that these sequences are deployed in a
528 highly abnormal fashion.

529 **Application of seqNMF to a behavioral dataset: song spectrograms**

530 Although we have focused on the application of seqNMF to neural activity data, this
531 method naturally extends to other types of high-dimensional datasets, including behavioral
532 data with applications to neuroscience. The neural mechanisms underlying song
533 production and learning in songbirds is an area of active research. However, the identification
534 and labeling of song syllables in acoustic recordings is challenging, particularly in
535 young birds where song syllables are highly variable. Because automatic segmentation
536 and clustering often fail, song syllables are still routinely labelled by hand [44]. We tested
537 whether seqNMF, applied to a spectrographic representation of zebra finch vocalizations,
538 is able to extract meaningful features in behavioral data. SeqNMF correctly identified
539 repeated acoustic patterns in juvenile songs, placing each distinct syllable type into a
540 different factor (Figure 8). The resulting classifications agree with previously published
541 hand-labeled syllable types [44]. A similar approach could be applied to other behavioral
542 data, for example movement data or human speech, and could facilitate the study of
543 neural mechanisms underlying even earlier and more variable stages of learning. Indeed,
544 convNMF was originally developed for application to spectrograms [56]; notably it has
545 been suggested that auditory cortex may use similar computations to represent and
546 parse natural song statistics [40].

547 **Discussion**

548 As neuroscientists strive to record larger datasets, there is a need for rigorous tools to
549 reveal underlying structure in high-dimensional data [20, 54, 11, 8]. In particular, sequential
550 structure is increasingly regarded as a fundamental property of neuronal circuits
551 [23, 24, 44, 45], but standardized statistical approaches for extracting such structure have
552 not been widely adopted or agreed upon.

553 Here, we explored a simple matrix factorization-based approach to identify neural
554 sequences [47]. The convNMF model elegantly captures sequential structure in an un-
555 supervised manner [56, 55]. However, in datasets where the number of sequences is

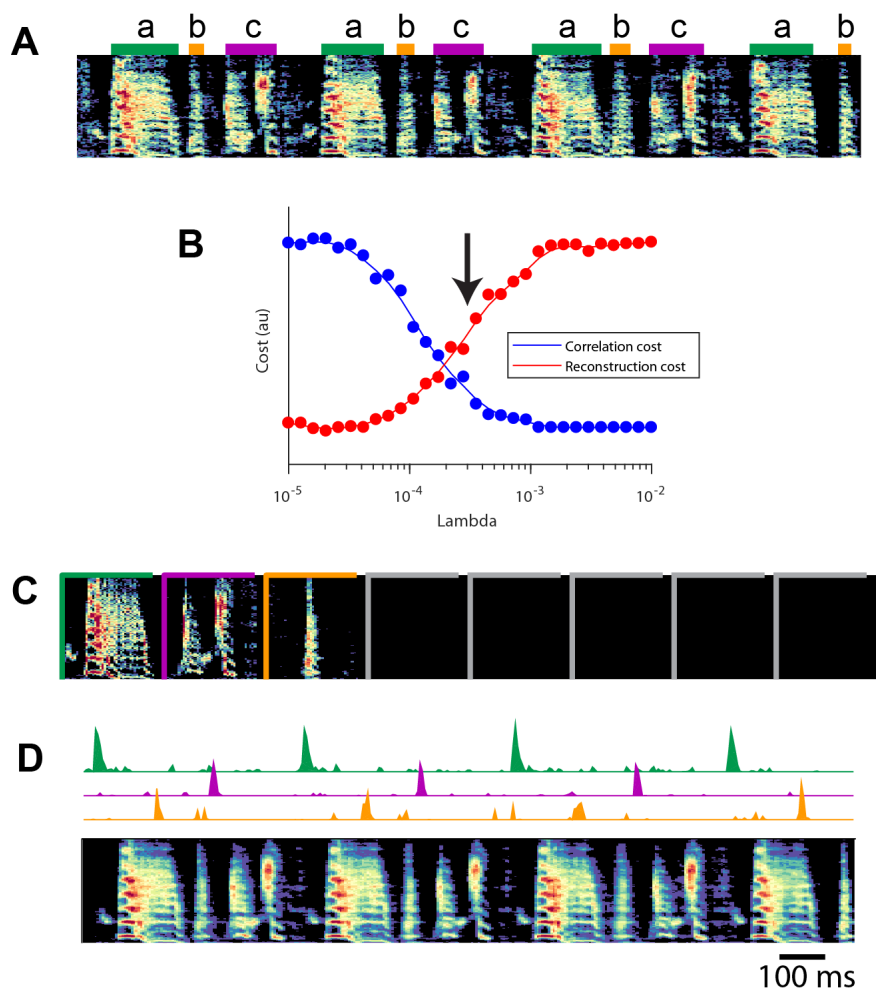


Figure 8. SeqNMF applied to song spectrograms

(A) Spectrogram of juvenile song, with hand-labeled syllable types [44]. (B) Reconstruction cost and correlation cost for these data as a function of λ . Arrow denotes $\lambda = 0.0003$, which was used to run seqNMF (C) SeqNMF W 's for this song, fit with $K = 8$, $L = 200ms$, $\lambda = 0.0003$. Note that there are three non-empty factors, corresponding to the three hand-labeled syllables a, b, and c. (D) SeqNMF H 's (for the three non-empty factors) and seqNMF reconstruction of the song shown in (A) using these factors.

556 not known, convNMF may return redundant, inefficient, or inconsistent factorizations. In
557 order to resolve these challenges, we introduced a new regularization (penalty) term to
558 encourage the model to identify sparse and non-redundant sequential firing patterns.
559 Furthermore, we carefully explored the robustness of this method to noise and devel-
560 oped procedures for choosing hyperparameters (K and λ) based on cross-validation and
561 assessing the significance of identified sequences based on shuffled null distributions.
562 Our results show that seqNMF is highly robust to many forms of noise. For example, even
563 when (synthetic) neurons participate probabilistically in sequences at a rate of 50%, the
564 model typically identifies factors with greater than 80% similarity to the ground truth (Fig-
565 ure 3A). Additionally, seqNMF performs well even with limited data, successfully extracting
566 sequences that only appear a handful of times in a noisy data stream (Figure S3)

567 Prior investigations of neural sequences have relied on manual alignment of neural
568 activity to behavioral events, such as animal position for the case of hippocampal and
569 cortical sequences [24, 45], or syllable onset for the case of songbird vocalizations [23].
570 This approach is not ideally suited for the case of highly variable behaviors, such as in
571 early learning and development [44]. For example, the analysis of neural activity in singing
572 juvenile birds has been challenging because of the difficulty in identifying distinct syllable
573 types on which to perform the temporal alignment. This problem would also apply to
574 isolate song birds because of the pathologically variable nature of their vocalizations. By
575 applying seqNMF, we were able to identify neural sequences without reference to song
576 syllables, enabling future work into the neural basis of singing in isolate birds.

577 As in many data analysis scenarios, a variety of statistical approaches may be brought
578 to bear on finding sequences in neural data. A classic method is to construct cross-
579 correlogram plots, showing spike time correlations between pairs of neurons at various
580 time lags. However, other forms of spike rate covariation, such as trial-to-trial gain mod-
581 ulation, can produce spurious peaks in this measure [7]; recent work has developed
582 statistical corrections for these effects [51]. After significant pairwise correlations are
583 identified, one can heuristically piece together pairs of neurons with significant interac-
584 tions into a sequence. This bottom-up approach may be better than seqNMF at detecting
585 sequences involving small numbers of neurons if such microsequences contribute only
586 a small amount of variance in the overall dataset. On the other hand, this bottom-up
587 approach may fail to identify long sequences with high participation noise or jitter in each
588 neuron [49]. One can think of seqNMF as a complementary top-down approach, which
589 performs very well in the high-noise regime since it learns a template sequence at the
590 level of the full population that is robust to noise at the level of individual units.

591 Statistical models with a dynamical component, such as Hidden Markov Models
592 (HMMs) [38], linear dynamical systems [30], and models with switching dynamics [35],
593 can also capture sequential firing patterns. These methods will typically require many
594 hidden states or latent dimensions to capture sequences, similar to PCA and NMF which
595 require many components to recover sequences. However, since dynamical models are
596 much more constrained than PCA or NMF, they can yield more interpretable results. For
597 example, visualizing the transition matrix of an HMM can provide insight into the order
598 in which hidden states of the model are visited, mapping onto different sequences that
599 manifest in population activity [38]. One advantage of this approach is that it can model
600 sequences that occasionally end prematurely, while seqNMF will always produce the full
601 sequence. On the other hand, this pattern completion property makes seqNMF robust

602 to participation noise and jitter. In contrast, a standard HMM must pass through each
603 hidden state to model a sequence, and therefore will have trouble whenever one of these
604 hidden states is skipped. Thus, we expect HMMs (or related models) and seqNMF to
605 exhibit complementary strengths and weaknesses.

606 Another contribution of our work is a natural framework in which to bias factorizations
607 towards parts-based versus events-based solutions. While existing computational work
608 has focused on neural sequences that do not have ensembles of shared neurons, such
609 shared populations have been observed during song learning [44], demonstrating that
610 neural sequences in real biological data can substantially overlap. Such shared sequences
611 can lead to different reasonable factorizations of the data that may correspond to dif-
612 ferent interpretations of underlying mechanisms. For example, we found that neural
613 sequences in HVC of isolated songbirds are well-described by both parts- or events-based
614 factorizations (figure S9), each of which could correspond to a different biophysical model
615 of sequence generation. This capacity for a combinatorial description of overlapping
616 sequences distinguishes convNMF and seqNMF from clustering methods [22, 39] and
617 methods based on hypothesis testing [49, 51], which seek to identify full snapshots of
618 repeated population firing patterns rather than parts- or events-based representations.
619 Another difference between these methods and seqNMF, particularly when using an
620 events-based factorization, is its ability to model different amplitudes in the sequences by
621 changing the magnitude of the event loadings in \mathbf{H} .

622 More generally, a key strength of seqNMF is that it can be easily tuned to the require-
623 ments and goals of a particular analysis. In addition to changing between a parts- and
624 events-based factorization, one can tune the overall sparsity in the model by classic L1
625 regularization. Future work could incorporate outlier detection into the objective function
626 as has been done in other matrix factorization models [42]. One could also incorporate
627 additional parameters to model changes in neural sequences across trials or days during
628 development or learning of a new behavior, similar to extensions of PCA and NMF to
629 multi-trial data [63]. Thus, adding convolutional structure to factorization-based mod-
630 els of neural data represents a rich opportunity for future developments in statistical
631 methodology.

632 Despite limiting ourselves to a relatively simple model for the purposes of this paper,
633 we extracted biological insights that were difficult to achieve by other methods in practical
634 experimental datasets. Overall, seqNMF can extract neural sequences from large-scale
635 population recordings without reference to stereotyped behavior or rigid sensory stimuli,
636 enabling the dissection of neural circuit activity during rich and variable animal behaviors.

637 Acknowledgements

638 This work was supported by a grant from the Simons Collaboration for the Global Brain,
639 the National Institutes of Health (NIH) [grant number R01 DC009183] and the G. Harold
640 & Leila Y. Mathers Charitable Foundation. ELM received support through the NDSEG
641 Fellowship program. AHB received support through NIH training grant 5T32EB019940-
642 03. MSG received support from the NIH [grant number U19NS104648]. AHW received
643 support from the U.S. Department of Energy Computational Science Graduate Fellowship
644 (CSGF) program. Thanks to Pengcheng Zhou for advice on his CNMF_E calcium data
645 cell extraction algorithm. Thanks to Wiktor Młynarski for helpful convNMF discussions.
646 Thanks to Michael Stetner, Galen Lynch, Nhat Le, Dezhe Jin, Edward Nieh, Adam Charles

647 and Jane Van Velden for comments on the manuscript and on our code package. Special
648 thanks to the 2017 Methods in Computational Neuroscience course [supported by NIH
649 grant R25 MH062204 and Simons Foundation] at the Woods Hole Marine Biology Lab,
650 where this collaboration was started.

651 **Author contributions**

652 ELM, AHB, AHW, MSG and MSF conceived the project, based on previous discussions of
653 MSG, MSF and ELM. ELM, AHB and MSF designed and tested the seqNMF regularizers,
654 the method for validating the significance of sequences in a held-out dataset, and the
655 method for choosing λ . ELM, AHB, AHW, and MSF designed and tested the method for
656 measuring RMSE on a masked test set. ELM and AHB wrote the algorithm and demo code.
657 ELM and NID collected the imaging data in singing birds. ELM and SG analyzed imaging
658 data. All authors contributed to writing the manuscript.

659 **Methods and Materials**

660 **Table of key resources**

661 Key resources, and references for how to access them, are listed in Table 2.

662 **Contact for resource sharing**

663 Further requests should be directed to Michale Fee (fee@mit.edu).

664 **Software and data availability**

665 Our seqNMF MATLAB code is publicly available as a github repository, along with some of
666 our data for demonstration:

667 <https://github.com/FeeLab/seqNMF>

668 The repository includes the seqNMF function, as well as helper functions for selecting
669 λ , testing the significance of factors, plotting, and other functions. It also includes a
670 demo script that goes through an example of how to select λ for a new dataset, test for
671 significance of factors, plot the seqNMF factorization, switch between parts-based and
672 events-based factorizations, and calculate cross-validated performance on a masked test
673 set.

674 We plan to post more of our data publicly on the CRCNS data-sharing platform.

675 **Generating simulated data**

676 We simulated neural sequences containing between 1 and 10 distinct neural sequences
677 in the presence of various noise conditions. Each neural sequence was made up of 10
678 consecutively active neurons, each separated by three timebins. The binary activity matrix
679 was convolved with an exponential kernel ($\tau = 10$ timebins) to resemble neural calcium
680 imaging activity.

681 **SeqNMF algorithm details**

682 Our algorithm for seqNMF (convNMF with additional regularization to promote efficient
683 factorizations) is a direct extension of the multiplicative update convNMF algorithm [56],
684 and draws on previous work regularizing NMF to encourage factor orthogonality [10].

Table 2. Key resources

Software/algorithm	Source	Link to code
seqNMF	This paper	https://github.com/FeeLab/seqNMF
convNMF	[56, 55]	https://github.com/colinvaz/nmf-too
Sparse convNMF	[43, 50]	https://github.com/colinvaz/nmf-too
Soft orthogonal NMF	[10]	
Other NMF extensions	[12]	
NMF	[34]	
CNMF_E (cell extraction)	[66]	https://github.com/zhoup/CNMF_E
MATLAB	MathWorks	www.mathworks.com
Dataset	Source	Link to data
HVC, Isolate songbird	This paper	To be uploaded to CRCNS after publ
Hippocampus, running wheel task (rat 1)	[2]	/hc-5
hline Hippocampus, running wheel task (rat 2)	[1]	https://crcns.org/data-sets/hc/hc-3
Other	Source	Link
Zebra finches (<i>Taeniopygia guttata</i>)	MIT animal facility	
AAV9.CAG.GCaMP6f.WPRE.SV40	[9]	https://pennvectorcore.med.upenn.
Miniature microscope	Inscopix	https://www.inscopix.com/nvista

Table 3. Regularized NMF and convNMF: cost functions and algorithms

NMF

$$\mathcal{L} = \frac{1}{2} \|\tilde{\mathbf{X}} - \mathbf{X}\|_2^2 + \mathcal{R}$$

$$\tilde{\mathbf{X}} = \mathbf{W}\mathbf{H}$$

$$\mathbf{W} \leftarrow \mathbf{W} \times \frac{\mathbf{X}\mathbf{H}^\top}{\tilde{\mathbf{X}}\mathbf{H}^\top + \frac{d\mathcal{R}}{d\mathbf{W}}}$$

$$\mathbf{H} \leftarrow \mathbf{H} \times \frac{\mathbf{W}^\top \mathbf{X}}{\mathbf{W}^\top \tilde{\mathbf{X}} + \frac{d\mathcal{R}}{d\mathbf{H}}}$$

convNMF

$$\mathcal{L} = \frac{1}{2} \|\tilde{\mathbf{X}} - \mathbf{X}\|_2^2 + \mathcal{R}$$

$$\tilde{\mathbf{X}} = \mathbf{W} \circledast \mathbf{H}$$

$$\mathbf{W}_{\cdot\ell} \leftarrow \mathbf{W}_{\cdot\ell} \times \frac{\mathbf{X} \overset{\ell \rightarrow \top}{\mathbf{H}}}{\tilde{\mathbf{X}} \overset{\ell \rightarrow \top}{\mathbf{H}} + \frac{d\mathcal{R}}{d\mathbf{W}_{\cdot\ell}}}$$

$$\mathbf{H} \leftarrow \mathbf{H} \times \frac{\mathbf{W}^\top \mathbf{X}}{\mathbf{W} \circledast \tilde{\mathbf{X}} + \frac{d\mathcal{R}}{d\mathbf{H}}}$$

L1 regularization for \mathbf{H} ($L1$ for \mathbf{W} is analogous)

$$\mathcal{R} = \lambda \|\mathbf{H}\|_1$$

$$\frac{d\mathcal{R}}{d\mathbf{W}_{\cdot\ell}} = 0$$

$$\frac{d\mathcal{R}}{d\mathbf{H}} = \lambda \mathbf{1}$$

Soft orthogonality for \mathbf{H}

$$\mathcal{R} = \frac{\lambda}{2} \|\mathbf{H}\mathbf{H}^\top\|_{1, i \neq j}$$

$$\frac{d\mathcal{R}}{d\mathbf{W}_{\cdot\ell}} = 0$$

$$\frac{d\mathcal{R}}{d\mathbf{H}} = \lambda(\mathbf{1} - \mathbf{I})\mathbf{H}$$

Smoothed soft orthogonality for \mathbf{H} (favors 'events-based')

$$\mathcal{R} = \frac{\lambda}{2} \|\mathbf{H}\mathbf{S}\mathbf{H}^\top\|_{1, i \neq j}$$

$$\frac{d\mathcal{R}}{d\mathbf{W}_{\cdot\ell}} = 0$$

$$\frac{d\mathcal{R}}{d\mathbf{H}} = \lambda(\mathbf{1} - \mathbf{I})\mathbf{H}\mathbf{S}$$

Smoothed soft orthogonality for \mathbf{W} (favors 'parts-based')

$$\mathcal{R} = \frac{\lambda}{2} \|\mathbf{W}_{flat}^\top \mathbf{W}_{flat}\|_{1, i \neq j}$$

$$\frac{d\mathcal{R}}{d\mathbf{W}_{\cdot\ell}} = \lambda \mathbf{W}_{flat}(\mathbf{1} - \mathbf{I})$$

$$\text{where } (\mathbf{W}_{flat})_{nk} = \sum_{\ell} \mathbf{W}_{nk\ell}$$

$$\frac{d\mathcal{R}}{d\mathbf{H}} = 0$$

Smoothed cross-factor orthogonality (main seqNMF \mathcal{R})

$$\mathcal{R} = \lambda \|\mathbf{W} \circledast \mathbf{X}\mathbf{S}\mathbf{H}^\top\|_{1, i \neq j}$$

$$\frac{d\mathcal{R}}{d\mathbf{W}_{\cdot\ell}} = \lambda \overset{\leftarrow \ell}{\mathbf{X}} \mathbf{S} \mathbf{H}^\top (\mathbf{1} - \mathbf{I})$$

$$\frac{d\mathcal{R}}{d\mathbf{H}} = \lambda(\mathbf{1} - \mathbf{I})\mathbf{W} \circledast \mathbf{X}\mathbf{S}$$

685 The uniqueness and consistency of traditional NMF has been better studied than
 686 convNMF, but in special cases, NMF has a unique solution comprised of sparse, ‘parts-
 687 based’ features that can be consistently identified by known algorithms [17, 4]. However,
 688 this ideal scenario does not hold in many practical settings. In these cases, NMF is
 689 sensitive to initialization, resulting in potentially inconsistent features. This problem can
 690 be addressed by introducing additional constraints or regularization terms that encourage
 691 the model to extract particular, e.g. sparse or approximately orthogonal, features [27, 31].
 692 Both theoretical work and empirical observations suggest that these modifications result
 693 in more consistently identified features [58, 31].

694 For seqNMF, we added to the convNMF cost function a term that promotes competition
 695 between overlapping factors, resulting in the following cost function:

$$(\tilde{\mathbf{W}}, \tilde{\mathbf{H}}) = \arg \min_{\mathbf{W}, \mathbf{H}} \left(\|\tilde{\mathbf{X}} - \mathbf{X}\|_F^2 + \lambda \|\mathbf{W} \circledast \mathbf{X}\mathbf{S}\mathbf{H}^\top\|_{1, i \neq j} \right) \quad (8)$$

696 We derived the following multiplicative update rules for \mathbf{W} and \mathbf{H} (Appendix 1):

$$\mathbf{W}_{\cdot \cdot \ell} \leftarrow \mathbf{W}_{\cdot \cdot \ell} \times \frac{\mathbf{X} \left(\overset{\ell \rightarrow}{\mathbf{H}} \right)^\top}{\tilde{\mathbf{X}} \left(\overset{\ell \rightarrow}{\mathbf{H}} \right)^\top + \lambda \overset{\leftarrow \ell}{\mathbf{X}} \mathbf{S} \mathbf{H}^\top (\mathbf{1} - \mathbf{I})} \quad (9)$$

$$\mathbf{H} \leftarrow \mathbf{H} \times \frac{\mathbf{W}^\top \mathbf{X}}{\mathbf{W}^\top \tilde{\mathbf{X}} + \lambda (\mathbf{1} - \mathbf{I}) (\mathbf{W}^\top \mathbf{X} \mathbf{S})} \quad (10)$$

697 Where the division and \times are element-wise. The operator $\overset{\ell \rightarrow}{(\cdot)}$ shifts a matrix in the \rightarrow
 698 direction by ℓ timebins, i.e. a delay by ℓ timebins, and $\overset{\leftarrow \ell}{(\cdot)}$ shifts a matrix in the \leftarrow direction
 699 by ℓ timebins (notation summary, Table 1). Note that multiplication with the $K \times K$ matrix
 700 $(\mathbf{1} - \mathbf{I})$ effectively implements factor competition because it places in the k th row a sum
 701 across all other factors. These update rules are derived in Appendix 1 by taking the
 702 derivative of the cost function in Equation 8.

703 In addition to the multiplicative updates outlined in Table 3, we also renormalize so
 704 rows of \mathbf{H} have unit norm; shift factors to be centered in time such that the center of
 705 mass of each \mathbf{W} pattern occurs in the middle; and in the final iteration run one additional
 706 step of unregularized convNMF to prioritize the cost of reconstruction error over the
 707 regularization (Algorithm 1). This final step is done to correct a minor suppression in
 708 the amplitude of some peaks in \mathbf{H} that may occur within $2L$ timebins of neighboring
 709 sequences.

710 *Calculating consistency*

711 The consistency between two factorizations measures the extent to which it is possible to
 712 create a one-to-one match between factors in factorization A and factors in factorization
 713 B . Specifically, given two factorizations $(\mathbf{W}^A, \mathbf{H}^A)$ and $(\mathbf{W}^B, \mathbf{H}^B)$ respectively, consistency
 714 is measured with the following procedure:

- 715 1. For each factor number k , compute the part of the reconstruction explained by this
 716 factor in each reconstruction, $\tilde{\mathbf{X}}_k^A = \mathbf{W}_{\cdot k}^A \circledast \mathbf{H}_k^A$ and $\tilde{\mathbf{X}}_k^B = \mathbf{W}_{\cdot k}^B \circledast \mathbf{H}_k^B$.

Algorithm 1: SeqNMF

Input: Data matrix \mathbf{X} , number of factors K , factor duration L , regularization strength λ

Output: Factor exemplars \mathbf{W} , and factor timecourses \mathbf{H}

- 1 Initialize \mathbf{W} and \mathbf{H} randomly
 - 2 Iter = 1
 - 3 **while** (Iter < maxIter) & (Δ cost > tolerance) **do**
 - 4 Update \mathbf{H} using multiplicative update from Table 3
 - 5 Shift \mathbf{W} and \mathbf{H} to center \mathbf{W} 's in time
 - 6 Renormalize \mathbf{W} and \mathbf{H} so rows of \mathbf{H} have unit norm
 - 7 Update \mathbf{W} using multiplicative update from Table 3
 - 8 Iter = Iter+1
 - 9 Do one final unregularized convNMF update of \mathbf{W} and \mathbf{H}
 - 10 **return**
-

- 717 2. Reshape $\tilde{\mathbf{X}}_k^A$ and $\tilde{\mathbf{X}}_k^B$ into vectors containing all the elements of each matrix res-
718 spectively, then compute \mathbf{C} , a $K \times K$ correlation matrix where C_{ij} is the correlation
719 between the vectorized $\tilde{\mathbf{X}}_i^A$ and $\tilde{\mathbf{X}}_j^B$
- 720 3. Permute the factors greedily so factor 1 is the best matched pair of factors, factor 2
721 is the best matched pair of the remaining factors, etc. The quality of the match is
722 measured by the correlation between the reconstructions computed using just each
723 factor individually.
- 724 4. Measure consistency as the ratio of the power (sum of squared matrix elements)
725 contained on the diagonal of the permuted \mathbf{C} matrix to the total power in \mathbf{C}

726 Thus, two factorizations are perfectly consistent when there exists a permutation of factor
727 numbers for which there is a one-to-one match between what parts of the reconstruction
728 are explained by each factor.

729 *Testing the significance of each factor on held-out data*

730 In order to test whether a factor is significantly present in held-out data, we measure
731 the distribution across timebins of the overlaps of the factor with the held-out data, and
732 compare the skewness of this distribution to the null case (Figure S1). Overlap with the
733 data is measured as $\mathbf{W} \otimes \mathbf{X}$, so this quantity will be high at timepoints when the sequence
734 occurs, producing a distribution of $\mathbf{W} \otimes \mathbf{X}$ with high skew. In contrast, a distribution of
735 overlaps exhibiting low skew indicates a sequence is not present in the data, since there
736 are few timepoints of particularly high overlap. We estimate what skew levels would
737 appear by chance by constructing null factors where temporal relationships between
738 neurons have been eliminated. To create such null factors, we start from the real factors
739 then circularly shift the timecourse of each neuron by a random amount between 0 and
740 L . We measure the skew of the overlap distributions for each null factor, and ask whether
741 the skew we measured for the real factor is significant at p-value α , that is, if it exceeds
742 the $((1 - \frac{\alpha}{K}) \times 100)^{th}$ percentile of the null skews. Note the required Bonferroni correction
743 for K comparisons when testing K factors.

744 *Choosing appropriate parameters for a new dataset*

745 Choice of appropriate parameters (λ , K and L) will depend on the data type (sequence
746 length, number, and density; amount of noise; etc.).

747 In practice, we find that results are relatively robust to choice of parameters. When K
748 or L is set larger than necessary, seqNMF tends to simply leave the unnecessary factors
749 or time bins empty. For λ , the goal is to find the ‘sweet spot’ (Figure 4) to explain as much
750 data as possible while still producing sensible factorizations, that is, minimally correlated
751 factors, with low values of $\|(\mathbf{W} \otimes \mathbf{X})\mathbf{S}\mathbf{H}^T\|_{1,i \neq j}$. Our software package includes demo code
752 for determining the best parameters for a new type of data, using the following strategy:

- 753 1. Start with K slightly larger than the number of sequences anticipated in the data
- 754 2. Start with L slightly longer than the maximum expected factor length
- 755 3. Run seqNMF for a range of λ 's, and for each λ measure the reconstruction error
756 $(\|\mathbf{X} - \mathbf{W} \otimes \mathbf{H}\|_F^2)$ and the factor competition regularization term $(\|(\mathbf{W} \otimes \mathbf{X})\mathbf{S}\mathbf{H}^T\|_{1,i \neq j})$
- 757 4. Choose a λ slightly above the crossover point λ_0
- 758 5. Decrease K if desired, as otherwise some factors will be consistently empty
- 759 6. Decrease L if desired, as otherwise some time bins will consistently be empty

760 In some applications, achieving the desired accuracy may depend on choosing a λ
761 that allows some inconsistency. It is possible to deal with this remaining inconsistency
762 by comparing factors produced by different random initializations, and only considering
763 factors that arise from several different initializations, a strategy that has been previously
764 applied to standard convNMF on neural data [47].

765 During validation of our procedure for choosing λ , we compared factorizations to
766 ground truth sequences as shown in Figure 4. To find the optimal lambda we used the
767 product of two curves. The first curve was obtained by calculating the fraction of fits in
768 which the true number of sequences was recovered as a function of λ . The second curve
769 was obtained by calculating similarity to ground truth as a function of λ . Similarity to
770 ground truth is measured as the consistency the factorization and the noiseless sequences
771 used to generate the data. The product of these two curves was smoothed using a three-
772 sample boxcar sliding window, and the width was found as the values of λ on either side
773 of the peak value that correspond most closely to the half-max points of the curve.

774 *Measuring performance on noisy data by comparing seqNMF sequences to* 775 *ground-truth sequences*

776 We wanted to measure the ability of seqNMF to recover ground-truth sequences even
777 when the sequences are obstructed by noise. Our noisy data consisted of two ground-
778 truth sequences, obstructed by a variety of noise types. We first took the top seqNMF
779 factor, and made a reconstruction with only this factor. We then measured the correlation
780 between this reconstruction and reconstructions generated from each of the ground-
781 truth factors, and chose the best match. Next, we measured the correlation between the
782 remaining ground-truth reconstruction and the second seqNMF factor. The mean of these
783 two correlations was used as a measure of similarity between the seqNMF factorization
784 and the ground-truth (noiseless) sequences.

785 *Testing generalization of factorization to randomly held-out (masked) data*
 786 *entries*

787 The data matrix \mathbf{X} was divided into training data and test data by randomly selecting 5 or
 788 10% of matrix entries to hold out. Specifically, the objective function (equation 5, in the
 789 Results section) was modified to:

$$\arg \min_{\mathbf{W}, \mathbf{H}} \|\mathbf{M} \times (\mathbf{W} \otimes \mathbf{H} - \mathbf{X})\|_F^2 + \mathcal{R} \quad (11)$$

790 where \times indicates elementwise multiplication (Hadamard product) and \mathbf{M} is a binary
 791 matrix with 5 or 10% of the entries randomly selected to be zero (held-out test set) and
 792 the remaining 95 or 90% set to one (training set). To search for a solution, we reformulate
 793 this optimization problem as:

$$\begin{aligned} \arg \min_{\mathbf{W}, \mathbf{H}, \mathbf{Z}} \quad & \|\mathbf{W} \otimes \mathbf{H} - \mathbf{Z}\|_F^2 + \mathcal{R} \\ \text{subject to} \quad & \mathbf{M} \times \mathbf{Z} = \mathbf{M} \times \mathbf{X} \end{aligned} \quad (12)$$

794 where we have introduced a new optimization variable \mathbf{Z} , which can be thought of as
 795 a surrogate dataset that is equal to the ground truth data only on the training set. The
 796 goal is now to minimize the difference between the model estimate, $\tilde{\mathbf{X}} = \mathbf{W} \otimes \mathbf{H}$, and the
 797 surrogate, \mathbf{Z} , while constraining \mathbf{Z} to equal \mathbf{X} at unmasked elements (where $m_{ij} = 1$) and
 798 allowing \mathbf{Z} to be freely chosen at masked elements (where $m_{ij} = 0$). Clearly, at masked
 799 elements, the best choice is to make \mathbf{Z} equal to the current model estimate $\tilde{\mathbf{X}}$ as this
 800 minimizes the cost function without violating the constraint. This leads to the following
 801 update rules which are applied cyclically to update \mathbf{Z} , \mathbf{W} , and \mathbf{H} .

$$\mathbf{Z}_{nt} \leftarrow \begin{cases} \mathbf{X}_{nt} & \text{if } \mathbf{M}_{nt} = 1 \\ (\mathbf{W} \otimes \mathbf{H})_{nt} & \text{if } \mathbf{M}_{nt} = 0 \end{cases} \quad (13)$$

$$\mathbf{W}_{..\ell} \leftarrow \mathbf{W}_{..\ell} \times \frac{\mathbf{Z} \left(\begin{smallmatrix} \ell \rightarrow \\ \mathbf{H} \end{smallmatrix} \right)^\top}{\tilde{\mathbf{X}} \left(\begin{smallmatrix} \ell \rightarrow \\ \mathbf{H} \end{smallmatrix} \right)^\top + \lambda \mathbf{Z} \mathbf{S} \mathbf{H}^\top (\mathbf{I} - \mathbf{I})} \quad (14)$$

$$\mathbf{H} \leftarrow \mathbf{H} \times \frac{\mathbf{W}^\top \mathbf{Z}}{\mathbf{W}^\top \tilde{\mathbf{X}} + \lambda (\mathbf{I} - \mathbf{I}) (\mathbf{W}^\top \mathbf{Z})} \quad (15)$$

802 The measure used for testing generalization performance was RMSE. For the testing
 803 phase, RMSE was computed from the difference between $\tilde{\mathbf{X}}$ and the data matrix \mathbf{X} only
 804 for held-out entries.

805 *Algorithm speed*

806 In practice, our algorithm converges rapidly: fewer than 100 iterations on a typical 150
 807 neuron by 10,000 time point data matrix. Typically, 100 iterations on such data take
 808 less than 30 seconds on a standard PC. However, applications to much larger datasets
 809 may require faster performance. In these cases, we recommend running seqNMF on
 810 smaller subsets of the dataset, perhaps by incorporating seqNMF regularization into an
 811 online version of convNMF [62], and/or parallelizing the algorithm by running it on shorter

812 datasets and merging/recombining factors that are common across these shorter runs
813 (finding common factors by e.g. [47]).

814 *Notes on data preprocessing*

815 While seqNMF is generally quite robust, proper preprocessing of the data can be important
816 to obtaining reasonable factorizations on real neural data. A key principle is that, in
817 minimizing the reconstruction error, seqNMF is most strongly influenced by parts of
818 the data that exhibit high variance. This can be problematic if the regions of interest in
819 the data have relatively low amplitude. For example, high firing rate neurons may be
820 prioritized over those with lower firing rate. As an alternative to subtracting the mean
821 firing rate of each neuron, which would introduce negative values, neurons could be
822 normalized divisively or by subtracting off a NMF reconstruction fit in method that forces
823 a non-negative residual [32]. Additionally, variations in behavioral state may lead to
824 seqNMF factorizations that prioritize regions of the data with high variance and neglect
825 other regions. It may be possible to mitigate these effects by normalizing data, or by
826 restricting analysis to particular subsets of the data, either by time or by neuron.

827 **Hippocampus data**

828 The hippocampal data we used was collected in the Buzsaki lab [2, 1], and is publicly
829 available on the Collaborative Research in Computational Neuroscience (CRCNS) Data
830 sharing website. The dataset we refer to as 'Rat 1' is in the [hc-5](#) dataset, and the dataset
831 we refer to as 'Rat 2' is in the [hc-3](#) and dataset. Before running seqNMF, we processed
832 the data by convolving the raw spike trains with a gaussian kernel of standard deviation
833 100ms.

834 **Animal care and use**

835 We used male zebra finches (*Taeniopygia guttata*) from the MIT zebra finch breeding facility
836 (Cambridge, MA). Animal care and experiments were carried out in accordance with NIH
837 guidelines, and reviewed and approved by the Massachusetts Institute of Technology
838 Committee on Animal Care (protocol 0715-071-18).

839 In order to prevent exposure to a tutor song, birds were foster-raised by female birds,
840 which do not sing, starting on or before post-hatch day 15. For experiments, birds were
841 housed singly in custom-made sound isolation chambers.

842 **Calcium imaging**

843 The calcium indicator GCaMP6f was expressed in HVC by intercranial injection of the viral
844 vector AAV9.CAG.GCaMP6f.WPRE.SV40 [9] into HVC. In the same surgery, a cranial window
845 was made using a GRIN (gradient index) lens (1mm diameter, 4mm length, Inscopix).
846 After at least one week, in order to allow for sufficient viral expression, recordings were
847 made using the Inscopix nVista miniature fluorescent microscope.

848 Neuronal activity traces were extracted from raw fluorescence movies using the
849 CNMF_E algorithm, a constrained non-negative matrix factorization algorithm specialized
850 for microendoscope data by including a local background model to remove activity from
851 out-of-focus cells [66].

852 We performed several preprocessing steps before applying seqNMF to functional
853 calcium traces extracted by CNMF_E. First, we estimated burst times from the raw traces

854 by deconvolving the traces using an AR-2 process. The deconvolution parameters (time
855 constants and noise floor) were estimated for each neuron using the CNMF_E code
856 package [66]. Some neurons exhibited larger peaks than others, likely due to different
857 expression levels of the calcium indicator. Since seqNMF would prioritize the neurons
858 with the most power, we renormalized by dividing the signal from each neuron by the
859 sum of the maximum value of that row and the 95th percentile of the signal across all
860 neurons. In this way, neurons with larger peaks were given some priority, but not much
861 more than that of neurons with weaker signals.

862 References

- 863 [1] Multiple single unit recordings from different rat hippocampal and entorhinal regions while
864 the animals were performing multiple behavioral tasks. CRCNSorg. 2013; [https://crcns.org/
865 data-sets/hc/hc-5/about-hc-5](https://crcns.org/data-sets/hc/hc-5/about-hc-5), doi: <http://dx.doi.org/10.6080/K09G5JRZ>.
- 866 [2] Simultaneous extracellular recordings from left and right hippocampal areas CA1 and
867 right entorhinal cortex from a rat performing a left / right alternation task and
868 other behaviors. CRCNSorg. 2015; <https://crcns.org/data-sets/hc/hc-5/about-hc-5>, doi:
869 <http://dx.doi.org/10.6080/K0KS6PHF>.
- 870 [3] **Aronov D**, Andalman AS, Fee MS. A specialized forebrain circuit for vocal babbling in the
871 juvenile songbird. *Science* (New York, NY). 2008 may; 320(5876):630–4. [http://www.ncbi.nlm.
872 nih.gov/pubmed/18451295](http://www.ncbi.nlm.nih.gov/pubmed/18451295), doi: [10.1126/science.1155140](https://doi.org/10.1126/science.1155140).
- 873 [4] **Arora S**, Ge R, Kannan R, Moitra A. Computing a Nonnegative Matrix Factorization – Provably.
874 ArXiv e-prints. 2011 nov; .
- 875 [5] **Bapi RS**, Pammi VSC, Miyapuram KP, Ahmed. Investigation of sequence processing: A cognitive
876 and computational neuroscience perspective. *Current Science*. 2005; 89(10):1690–1698. [http:
877 //www.jstor.org/stable/24111208](http://www.jstor.org/stable/24111208).
- 878 [6] **Bro R**, Kjeldahl K, Smilde AK, Kiers HAL. Cross-validation of component models: A critical
879 look at current methods. *Analytical and Bioanalytical Chemistry*. 2008 Mar; 390(5):1241–1251.
880 <https://doi.org/10.1007/s00216-007-1790-1>, doi: [10.1007/s00216-007-1790-1](https://doi.org/10.1007/s00216-007-1790-1).
- 881 [7] **Brody CD**. Correlations Without Synchrony. *Neural Computation*. 1999; 11(7):1537–1551.
882 <https://doi.org/10.1162/089976699300016133>, doi: [10.1162/089976699300016133](https://doi.org/10.1162/089976699300016133).
- 883 [8] **Bzdok D**, Yeo BTT. Inference in the age of big data: Future perspectives on neuroscience.
884 *NeuroImage*. 2017; 155(Supplement C):549–564. [http://www.sciencedirect.com/science/article/
885 pii/S1053811917303816](http://www.sciencedirect.com/science/article/pii/S1053811917303816), doi: <https://doi.org/10.1016/j.neuroimage.2017.04.061>.
- 886 [9] **Chen TW**, Wardill TJ, Sun Y, Pulver SR, Renninger SL, Baohan A, Schreiter ER, Kerr
887 RA, Orger MB, Jayaraman V, Looger LL, Svoboda K, Kim DS. Ultrasensitive flu-
888 orescent proteins for imaging neuronal activity. *Nature*. 2013 jul; 499(7458):295–
889 300. <http://www.ncbi.nlm.nih.gov/pubmed/23868258>[http://www.pubmedcentral.nih.gov/
890 articlerender.fcgi?artid=PMC3777791](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3777791), doi: [10.1038/nature12354](https://doi.org/10.1038/nature12354).
- 891 [10] **Chen Z**, Cichocki A. Nonnegative matrix factorization with temporal smoothness and/or spatial
892 decorrelation constraints. In: *Laboratory for Advanced Brain Signal Processing, RIKEN, Tech.
893 Rep*; 2005. .
- 894 [11] **Churchland AK**, Abbott LF. Conceptual and technical advances define a key moment for
895 theoretical neuroscience. *Nature Neuroscience*. 2016 feb; 19(3):348–349. [http://www.nature.
896 com/doi/10.1038/nn.4255](http://www.nature.com/doi/10.1038/nn.4255), doi: [10.1038/nn.4255](https://doi.org/10.1038/nn.4255).

- 897 [12] **Cichocki A.** Nonnegative Matrix and Tensor Factorizations : Applications to Ex-
898 ploratory Multi-way Data Analysis and Blind Source Separation. Wiley; 2009.
899 [http://libproxy.mit.edu/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=](http://libproxy.mit.edu/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=nlebk&AN=287301&site=ehost-live&scope=site)
900 [nlebk&AN=287301&site=ehost-live&scope=site](http://libproxy.mit.edu/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=nlebk&AN=287301&site=ehost-live&scope=site).
- 901 [13] **Clegg BA**, Digirolamo GJ, Keele SW. Sequence learning. Trends in cognitive sciences. 1998
902 aug; 2(8):275–81. <http://www.ncbi.nlm.nih.gov/pubmed/21227209>, doi: 10.1016/S1364-
903 6613(98)01202-9.
- 904 [14] **Cui Y**, Ahmad S, Hawkins J. Continuous Online Sequence Learning with an Unsupervised
905 Neural Network Model. Neural Computation. 2016; 28(11):2474–2504. [https://doi.org/10.1162/](https://doi.org/10.1162/NECO_a_00893)
906 [NECO_a_00893](https://doi.org/10.1162/NECO_a_00893), doi: 10.1162/NECO_a_00893.
- 907 [15] **Cunningham JP**, Yu BM. Dimensionality reduction for large-scale neural recordings. Nature
908 Neuroscience. 2014 nov; 17(11):1500–1509. <http://www.nature.com/articles/nn.3776>, doi:
909 [10.1038/nn.3776](http://www.nature.com/articles/nn.3776).
- 910 [16] **Diba K**, Buzsáki G. Hippocampal Network Dynamics Constrain the Time Lag between Pyramidal
911 Cells across Modified Environments. Journal of Neuroscience. 2008; 28(50):13448–13456.
912 <http://www.jneurosci.org/content/28/50/13448>, doi: 10.1523/JNEUROSCI.3824-08.2008.
- 913 [17] **Donoho D**, Stodden V. When Does Non-Negative Matrix Factorization Give a Correct Decom-
914 position into Parts? In: Thrun S, Saul LK, Schölkopf B, editors. *Advances in Neural Information*
915 *Processing Systems 16* MIT Press; 2004.p. 1141–1148. [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/2463-when-does-non-negative-matrix-factorization-give-a-correct-decomposition-into-parts.pdf)
916 [2463-when-does-non-negative-matrix-factorization-give-a-correct-decomposition-into-parts.](http://papers.nips.cc/paper/2463-when-does-non-negative-matrix-factorization-give-a-correct-decomposition-into-parts.pdf)
917 [pdf](http://papers.nips.cc/paper/2463-when-does-non-negative-matrix-factorization-give-a-correct-decomposition-into-parts.pdf).
- 918 [18] **Fehér O**, Wang H, Saar S, Mitra PP, Tchernichovski O. De novo establishment of wild-type song
919 culture in the zebra finch. Nature. 2009 may; 459(7246):564–568. [http://www.pubmedcentral.](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2693086&tool=pmcentrez&rendertype=abstract)
920 [nih.gov/articlerender.fcgi?artid=2693086&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2693086&tool=pmcentrez&rendertype=abstract), doi:
921 [10.1038/nature07994](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2693086&tool=pmcentrez&rendertype=abstract).
- 922 [19] **Fujisawa S**, Amarasingham A, Harrison M, Buzsáki G. Behavior-dependent short-term as-
923 sembly dynamics in the medial prefrontal cortex. Nature Neuroscience. 2008; 11(7):823–833.
924 <https://www.nature.com/articles/nn.2134>, doi: 10.1038/nn.2134.
- 925 [20] **Gao P**, Ganguli S. On Simplicity and Complexity in the Brave New World of Large-Scale
926 Neuroscience. ArXiv e-prints. 2015 mar; .
- 927 [21] **Gerstein GL**, Williams ER, Diesmann M, Gründ S, Trengove C. Detecting synfire
928 chains in parallel spike data. Journal of Neuroscience Methods. 2012; 206:54–64. doi:
929 [10.1016/j.jneumeth.2012.02.003](https://doi.org/10.1016/j.jneumeth.2012.02.003).
- 930 [22] **Grossberger L**, Battaglia FP, Vinck M. Unsupervised clustering of temporal patterns in high-
931 dimensional neuronal ensembles using a novel dissimilarity measure. bioRxiv. 2018; [https://](https://www.biorxiv.org/content/early/2018/04/30/252791)
932 www.biorxiv.org/content/early/2018/04/30/252791, doi: 10.1101/252791.
- 933 [23] **Hahnloser RHR**, Kozhevnikov AA, Fee MS. An ultra-sparse code underlies the generation of
934 neural sequences in a songbird. Nature. 2002 sep; 419(6902):65–70. [http://www.ncbi.nlm.nih.](http://www.ncbi.nlm.nih.gov/pubmed/12214232)
935 [gov/pubmed/12214232](http://www.ncbi.nlm.nih.gov/pubmed/12214232), doi: 10.1038/nature00974.
- 936 [24] **Harvey CD**, Coen P, Tank DW. Choice-specific sequences in parietal cortex during a
937 virtual-navigation decision task. Nature. 2012 mar; 484(7392):62–68. [http://www.ncbi.](http://www.ncbi.nlm.nih.gov/pubmed/22419153)
938 [nlm.nih.gov/pubmed/22419153](http://www.ncbi.nlm.nih.gov/pubmed/22419153)[http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3321074)
939 [PMC3321074](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3321074)<http://www.nature.com/doi/10.1038/nature10918>, doi: 10.1038/na-
940 [ture10918](http://www.nature.com/doi/10.1038/nature10918).

- 941 [25] **Hastie T**, Tibshirani R, Friedman JHJH. The elements of statistical learning : data mining,
942 inference, and prediction. Springer; 2009.
- 943 [26] **Hawkins J**, Ahmad S. Why Neurons Have Thousands of Synapses, a Theory of Sequence
944 Memory in Neocortex. *Frontiers in Neural Circuits*. 2016; 10:23. [https://www.frontiersin.org/](https://www.frontiersin.org/article/10.3389/fncir.2016.00023)
945 [article/10.3389/fncir.2016.00023](https://www.frontiersin.org/article/10.3389/fncir.2016.00023), doi: 10.3389/fncir.2016.00023.
- 946 [27] **Huang K**, Sidiropoulos ND, Swami A. Non-Negative Matrix Factorization Revisited: Uniqueness
947 and Algorithm for Symmetric Decomposition. *IEEE Transactions on Signal Processing*. 2014
948 jan; 62(1):211–224. doi: 10.1109/TSP.2013.2285514.
- 949 [28] **Janata P**, Grafton ST. Swinging in the brain: shared neural substrates for behaviors related to
950 sequencing and music. *Nature Neuroscience*. 2003 jul; 6(7):682–687. [http://www.nature.com/](http://www.nature.com/articles/nn1081)
951 [articles/nn1081](http://www.nature.com/articles/nn1081), doi: 10.1038/nn1081.
- 952 [29] **Jun JJ**, Steinmetz NA, Siegle JH, Denman DJ, Bauza M, Barbarits B, Lee AK, Anastassiou CA,
953 Andrei A, Aydın Ç, Barbic M, Blanche TJ, Bonin V, Couto J, Dutta B, Gratiy SL, Gutnisky DA,
954 Häusser M, Karsh B, Ledochowitsch P, et al. Fully integrated silicon probes for high-density
955 recording of neural activity. *Nature*. 2017 nov; 551(7679):232–236. [http://www.nature.com/](http://www.nature.com/doi/10.1038/nature24636)
956 [doi/10.1038/nature24636](http://www.nature.com/doi/10.1038/nature24636), doi: 10.1038/nature24636.
- 957 [30] **Kao JC**, Nuyujukian P, Ryu SI, Churchland MM, Cunningham JP, Shenoy KV. Single-trial dynamics
958 of motor cortex and their applications to brain-machine interfaces. *Nature Communications*.
959 2015 07; 6:7759 EP –. <http://dx.doi.org/10.1038/ncomms8759>.
- 960 [31] **Kim J**, Park H. Sparse Nonnegative Matrix Factorization for Clustering. In: ; 2008. .
- 961 [32] **Kim M**, Smaragdis P. Efficient model selection for speech enhancement using a deflation
962 method for Nonnegative Matrix Factorization. In: 2014 IEEE Global Conference on Signal
963 and Information Processing (GlobalSIP) IEEE; 2014. p. 537–541. [http://ieeexplore.ieee.org/](http://ieeexplore.ieee.org/document/7032175/)
964 [document/7032175/](http://ieeexplore.ieee.org/document/7032175/), doi: 10.1109/GlobalSIP.2014.7032175.
- 965 [33] **Kim TH**, Zhang Y, Lecoq J, Jung JC, Li J, Zeng H, Niell CM, Schnitzer MJ. Long-Term
966 Optical Access to an Estimated One Million Neurons in the Live Mouse Cortex.
967 *Cell reports*. 2016 dec; 17(12):3385–3394. [http://www.ncbi.nlm.nih.gov/pubmed/](http://www.ncbi.nlm.nih.gov/pubmed/28009304)
968 [28009304](http://www.ncbi.nlm.nih.gov/pubmed/28009304)<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC5459490>, doi:
969 [10.1016/j.celrep.2016.12.004](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC5459490).
- 970 [34] **Lee DD**, Seung HS. Learning the parts of objects by non-negative matrix factorization. *Nature*.
971 1999; 401(6755):788–791.
- 972 [35] **Linderman S**, Johnson M, Miller A, Adams R, Blei D, Paninski L. Bayesian Learning and Inference
973 in Recurrent Switching Linear Dynamical Systems. In: Singh A, Zhu J, editors. Proceedings of the
974 20th International Conference on Artificial Intelligence and Statistics, vol. 54 of Proceedings
975 of Machine Learning Research Fort Lauderdale, FL, USA: PMLR; 2017. p. 914–922. [http:](http://proceedings.mlr.press/v54/linderman17a.html)
976 [//proceedings.mlr.press/v54/linderman17a.html](http://proceedings.mlr.press/v54/linderman17a.html).
- 977 [36] **Long MA**, Jin DZ, Fee MS. Support for a synaptic chain model of neuronal sequence generation.
978 *Nature*. 2010 nov; 468(7322):394–399. [http://www.pubmedcentral.nih.gov/articlerender.fcgi?](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2998755&tool=pmcentrez&rendertype=abstract)
979 [artid=2998755&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2998755&tool=pmcentrez&rendertype=abstract), doi: 10.1038/nature09514.
- 980 [37] **Lynch G**, Okubo T, Hanuschkin A, Hahnloser RR, Fee M. Rhythmic Continuous-Time
981 Coding in the Songbird Analog of Vocal Motor Cortex. *Neuron*. 2016 may; 90(4):877–
982 892. <http://www.ncbi.nlm.nih.gov/pubmed/27196977>[http://linkinghub.elsevier.com/retrieve/](http://linkinghub.elsevier.com/retrieve/pii/S0896627316301088)
983 [pii/S0896627316301088](http://linkinghub.elsevier.com/retrieve/pii/S0896627316301088), doi: 10.1016/j.neuron.2016.04.021.

- 984 [38] **Maboudi K**, Ackermann E, Pfeiffer BE, Foster DJ, Diba K, Kemere C. Uncovering temporal
985 structure in hippocampal output patterns. *bioRxiv*. 2018; [https://www.biorxiv.org/content/](https://www.biorxiv.org/content/early/2018/01/04/242594)
986 [early/2018/01/04/242594](https://www.biorxiv.org/content/early/2018/01/04/242594), doi: 10.1101/242594.
- 987 [39] **van der Meij R**, Voytek B. Uncovering Neuronal Networks Defined by Consistent between-
988 Neuron Spike Timing from Neuronal Spike Recordings. *eNeuro*. 2018; [http://www.eneuro.org/](http://www.eneuro.org/content/early/2018/05/08/ENEURO.0379-17.2018)
989 [content/early/2018/05/08/ENEURO.0379-17.2018](http://www.eneuro.org/content/early/2018/05/08/ENEURO.0379-17.2018), doi: 10.1523/ENEURO.0379-17.2018.
- 990 [40] **Młynarski W**, McDermott JH. Learning Midlevel Auditory Codes from Natural Sound Statistics.
991 *Neural Computation*. 2018 mar; 30(3):631–669. [https://www.mitpressjournals.org/doi/abs/10.](https://www.mitpressjournals.org/doi/abs/10.1162/neco.2018.03.01048)
992 [1162/neco.2018.03.01048](https://www.mitpressjournals.org/doi/abs/10.1162/neco.2018.03.01048), doi: 10.1162/neco_a_01048.
- 993 [41] **Mokeichev A**, Okun M, Barak O, Katz Y, Ben-Shahar O, Lampl I. Stochastic Emergence of
994 Repeating Cortical Motifs in Spontaneous Membrane Potential Fluctuations In Vivo. *Neuron*.
995 2007; 53(3):413 – 425. <http://www.sciencedirect.com/science/article/pii/S0896627307000372>,
996 doi: <https://doi.org/10.1016/j.neuron.2007.01.017>.
- 997 [42] **Netrapalli P**, U N N, Sanghavi S, Anandkumar A, Jain P. Non-convex Robust PCA. In:
998 Ghahramani Z, Welling M, Cortes C, Lawrence ND, Weinberger KQ, editors. *Advances in*
999 *Neural Information Processing Systems 27* Curran Associates, Inc.; 2014.p. 1107–1115. [http://](http://papers.nips.cc/paper/5430-non-convex-robust-pca.pdf)
1000 papers.nips.cc/paper/5430-non-convex-robust-pca.pdf.
- 1001 [43] **O’Grady PD**, Pearlmutter BA. Convolutional Non-Negative Matrix Factorisation with a Sparseness
1002 Constraint. In: *2006 16th IEEE Signal Processing Society Workshop on Machine Learning for*
1003 *Signal Processing*; 2006. p. 427–432. doi: 10.1109/MLSP.2006.275588.
- 1004 [44] **Okubo TS**, Mackevicius EL, Payne HL, Lynch GF, Fee MS. Growth and splitting of neural
1005 sequences in songbird vocal development. *Nature*. 2015 nov; 528(7582):352–357. [http://www.](http://www.ncbi.nlm.nih.gov/pubmed/26618871)
1006 [ncbi.nlm.nih.gov/pubmed/26618871](http://www.ncbi.nlm.nih.gov/pubmed/26618871)[http://www.pubmedcentral.nih.gov/articlerender.fcgi?](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4957523)
1007 [artid=PMC4957523](http://www.nature.com/doi/10.1038/nature15741)<http://www.nature.com/doi/10.1038/nature15741>, doi: 10.1038/nature15741.
1008
- 1009 [45] **Pastalkova E**, Itskov V, Amarasingham A, Buzsáki G. Internally Generated Cell Assembly
1010 Sequences in the Rat Hippocampus. *Science*. 2008; 321(5894):1322–1327. [http://science.](http://science.sciencemag.org/content/321/5894/1322)
1011 [sciencemag.org/content/321/5894/1322](http://science.sciencemag.org/content/321/5894/1322), doi: 10.1126/science.1159775.
- 1012 [46] **Pearson K**. LIII. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*. 1901
1013 nov; 2(11):559–572. <https://www.tandfonline.com/doi/full/10.1080/14786440109462720>, doi:
1014 10.1080/14786440109462720.
1015
- 1016 [47] **Peter S**, Kirschbaum E, Both M, Campbell L, Harvey B, Heins C, Durstewitz D, Diego F, Ham-
1017 precht FA. Sparse convolutional coding for neuronal assembly detection. In: Guyon I, Luxburg
1018 UV, Bengio S, Wallach H, Fergus R, Vishwanathan S, Garnett R, editors. *Advances in Neural*
1019 *Information Processing Systems 30* Curran Associates, Inc.; 2017.p. 3675–3685. [http://papers.](http://papers.nips.cc/paper/6958-sparse-convolutional-coding-for-neuronal-assembly-detection.pdf)
1020 [nips.cc/paper/6958-sparse-convolutional-coding-for-neuronal-assembly-detection.pdf](http://papers.nips.cc/paper/6958-sparse-convolutional-coding-for-neuronal-assembly-detection.pdf).
- 1021 [48] **Picardo M**, Merel J, Katlowitz K, Vallentin D, Okobi D, Benezra S, Clary R, Pnevmatikakis E,
1022 Paninski L, Long M. Population-Level Representation of a Temporal Sequence Underlying
1023 Song Production in the Zebra Finch. *Neuron*. 2016 may; 90(4):866–876. [http://www.](http://www.ncbi.nlm.nih.gov/pubmed/27196976)
1024 [ncbi.nlm.nih.gov/pubmed/27196976](http://www.ncbi.nlm.nih.gov/pubmed/27196976)[http://www.pubmedcentral.nih.gov/articlerender.fcgi?](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4941616)
1025 [artid=PMC4941616](http://linkinghub.elsevier.com/retrieve/pii/S0896627316001094)<http://linkinghub.elsevier.com/retrieve/pii/S0896627316001094>, doi:
1026 10.1016/j.neuron.2016.02.016.

- 1027 [49] **Quaglio P**, Rostami V, Torre E, Grün S. Methods for identification of spike patterns
1028 in massively parallel spike trains. *Biological Cybernetics*. 2018 apr; 112(1-2):57–80.
1029 <http://www.ncbi.nlm.nih.gov/pubmed/29651582><http://www.pubmedcentral.nih.gov/>
1030 [articlerender.fcgi?artid=PMC5908877http://link.springer.com/10.1007/s00422-018-0755-0](http://link.springer.com/10.1007/s00422-018-0755-0),
1031 doi: 10.1007/s00422-018-0755-0.
- 1032 [50] **Ramanarayanan V**, Goldstein L, Narayanan SS. Spatio-temporal articulatory movement
1033 primitives during speech production: Extraction, interpretation, and validation. *The Journal of*
1034 *the Acoustical Society of America*. 2013; 134(2):1378–1394. <https://doi.org/10.1121/1.4812765>,
1035 doi: 10.1121/1.4812765.
- 1036 [51] **Russo E**, Durstewitz D. Cell assemblies at multiple time scales with arbitrary lag constellations.
1037 *eLife*. 2017; 6:e19428. <https://doi.org/10.7554/eLife.19428>, doi: 10.7554/eLife.19428.
- 1038 [52] **Scholvin J**, Kinney JP, Bernstein JG, Moore-Kochlacs C, Kopell N, Fonstad CG, Boyden
1039 ES. Close-Packed Silicon Microelectrodes for Scalable Spatially Oversampled Neural
1040 Recording. *IEEE Transactions on Biomedical Engineering*. 2016 jan; 63(1):120–130. doi:
1041 [10.1109/TBME.2015.2406113](https://doi.org/10.1109/TBME.2015.2406113).
- 1042 [53] **Schrader S**, Grün S, Diesmann M, Gerstein GL. Detecting Synfire Chain Activity Using
1043 Massively Parallel Spike Train Recording. *Journal of Neurophysiology*. 2008 oct; 100(4):2165–
1044 2176. <http://www.ncbi.nlm.nih.gov/pubmed/18632888><http://www.pubmedcentral.nih.gov/>
1045 [articlerender.fcgi?artid=PMC2576207http://www.physiology.org/doi/10.1152/jn.01245.2007](http://www.physiology.org/doi/10.1152/jn.01245.2007),
1046 doi: 10.1152/jn.01245.2007.
- 1047 [54] **Sejnowski TJ**, Churchland PS, Movshon JA. Putting big data to good use in neuroscience.
1048 *Nature neuroscience*. 2014 nov; 17(11):1440–1. [http://www.ncbi.nlm.nih.gov/pubmed/](http://www.ncbi.nlm.nih.gov/pubmed/25349909)
1049 [25349909http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4224030](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4224030), doi:
1050 [10.1038/nn.3839](https://doi.org/10.1038/nn.3839).
- 1051 [55] **Smaragdis P**. Convolutional Speech Bases and Their Application to Supervised Speech Separation.
1052 *IEEE Transactions on Audio, Speech, and Language Processing*. 2007 jan; 15(1):1–12. doi:
1053 [10.1109/TASL.2006.876726](https://doi.org/10.1109/TASL.2006.876726).
- 1054 [56] **Smaragdis P**. In: Punttonet CG, Prieto A, editors. *Non-negative Matrix Factor Deconvolution;*
1055 *Extraction of Multiple Sound Sources from Monophonic Inputs* Berlin, Heidelberg: Springer
1056 Berlin Heidelberg; 2004. p. 494–499. https://doi.org/10.1007/978-3-540-30110-3_63, doi:
1057 [10.1007/978-3-540-30110-3_63](https://doi.org/10.1007/978-3-540-30110-3_63).
- 1058 [57] **Sutskever I**, Vinyals O, Le QV. Sequence to Sequence Learning with Neural Networks. In:
1059 Ghahramani Z, Welling M, Cortes C, Lawrence ND, Weinberger KQ, editors. *Advances in Neural*
1060 *Information Processing Systems 27* Curran Associates, Inc.; 2014.p. 3104–3112. [http://papers.](http://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural-networks.pdf)
1061 [nips.cc/paper/5346-sequence-to-sequence-learning-with-neural-networks.pdf](http://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural-networks.pdf).
- 1062 [58] **Theis FJ**, Stadlthanner K, Tanaka T. First results on uniqueness of sparse non-negative matrix
1063 factorization. In: *2005 13th European Signal Processing Conference*; 2005. p. 1–4.
- 1064 [59] **Torre E**, Canova C, Denker M, Gerstein G, Helias M, Grün S. ASSET: Analysis of Se-
1065 quences of Synchronous Events in Massively Parallel Spike Trains. *PLOS Computational*
1066 *Biology*. 2016 jul; 12(7):e1004939. <http://www.ncbi.nlm.nih.gov/pubmed/27420734><http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4946788>[http://dx.plos.org/10.](http://dx.plos.org/10.1371/journal.pcbi.1004939)
1067 [1371/journal.pcbi.1004939](http://dx.plos.org/10.1371/journal.pcbi.1004939), doi: 10.1371/journal.pcbi.1004939.
1068
- 1069 [60] **Udell M**, Horn C, Zadeh R, Boyd S. Generalized Low Rank Models. *Foundations and Trends in*
1070 *Machine Learning*. 2016; 9(1). <http://dx.doi.org/10.1561/22000000055>.

- 1071 [61] **Vaz C**, Toutios A, Narayanan S. Convex Hull Convolutional Non-negative Matrix Factorization for
1072 Uncovering Temporal Patterns in Multivariate Time-Series Data. In: *Interspeech* San Francisco,
1073 CA; 2016. p. 963–967.
- 1074 [62] **Wang D**, Vipperla R, Evans N, Zheng TF. Online Non-Negative Convolutional Pattern Learn-
1075 ing for Speech Signals. *IEEE Transactions on Signal Processing*. 2013 jan; 61(1):44–56. doi:
1076 [10.1109/TSP.2012.2222381](https://doi.org/10.1109/TSP.2012.2222381).
- 1077 [63] **Williams AH**, Kim TH, Wang F, Vyas S, Ryu SI, Shenoy KV, Schnitzer M, Kolda TG, Ganguli
1078 S. Unsupervised discovery of demixed, low-dimensional neural dynamics across multiple
1079 timescales through tensor components analysis. *Neuron*. in press; [https://www.biorxiv.org/
1080 content/early/2017/10/30/211128](https://www.biorxiv.org/content/early/2017/10/30/211128), doi: 10.1101/211128.
- 1081 [64] **Wold S**. Cross-Validatory Estimation of the Number of Components in Factor and Principal
1082 Components Models. *Technometrics*. 1978; 20(4):397–405. [https://www.tandfonline.com/doi/
1083 abs/10.1080/00401706.1978.10489693](https://www.tandfonline.com/doi/abs/10.1080/00401706.1978.10489693), doi: 10.1080/00401706.1978.10489693.
- 1084 [65] **Zhang Z**, Xu Y, Yang J, Li X, Zhang D. A survey of sparse representation: algorithms and
1085 applications. *ArXiv e-prints*. 2016 feb; .
- 1086 [66] **Zhou P**, Resendez SL, Rodriguez-Romaguera J, Jimenez JC, Neufeld SQ, Giovannucci A, Friedrich
1087 J, Pnevmatikakis EA, Stuber GD, Hen R, Kheirbek MA, Sabatini BL, Kass RE, Paninski L. Efficient
1088 and accurate extraction of in vivo calcium signals from microendoscopic video data. *eLife*.
1089 2018 feb; 7:e28728. <https://elifesciences.org/articles/28728>, doi: 10.7554/eLife.28728.

1090 **Figures**
 1091 **Supplemental Figures**

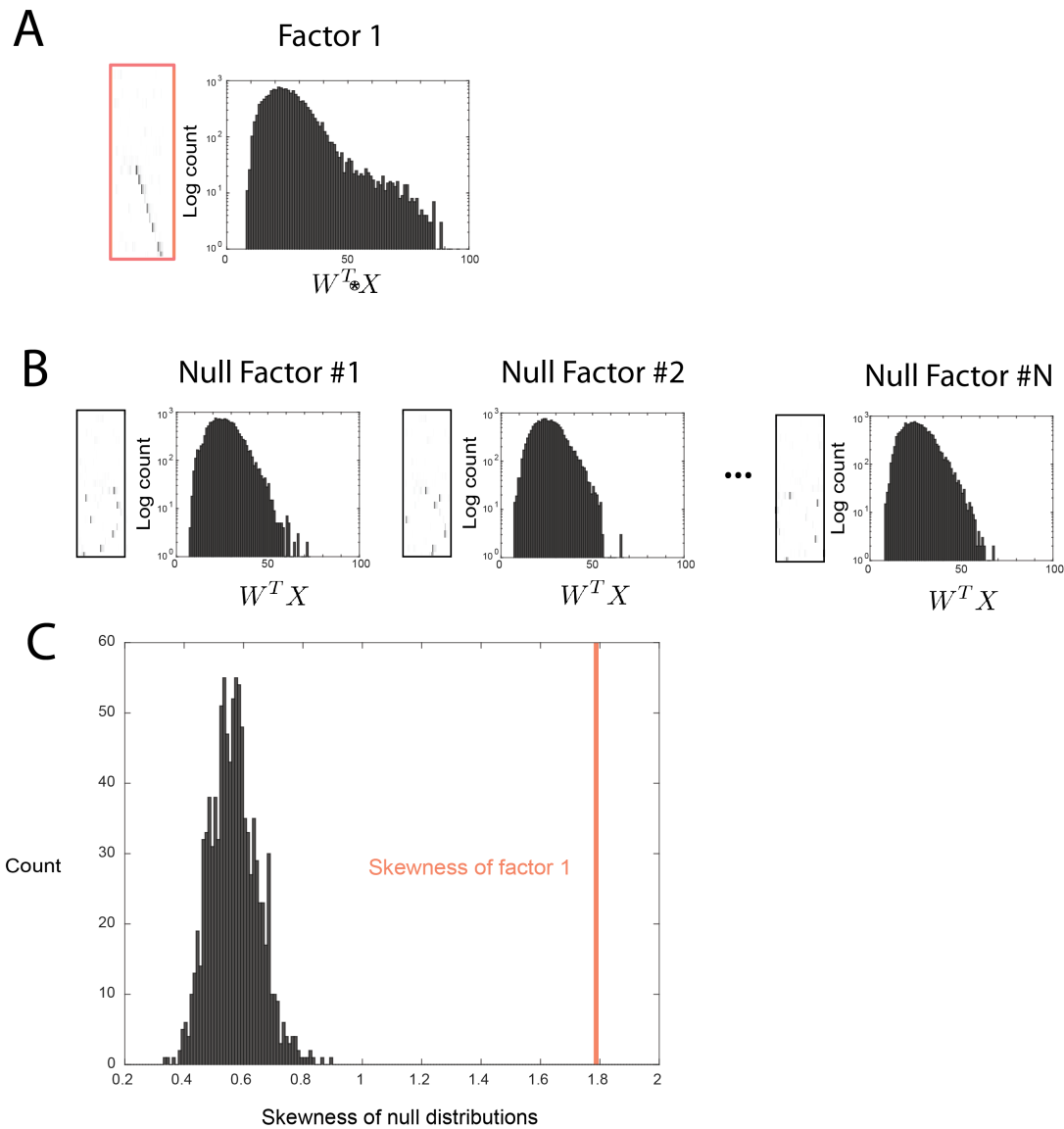


Figure S1. Outline of the procedure used to assess factor significance.

(A) The distribution of overlap values between the real factor and the held-out data. **(B)** In order to test the significance of a factor on held-out data, we constructed null (shifted) versions of the factor, and measured the distribution of overlap values ($W^T X$) between each null factor and the held-out data. **(C)** We then compared the skewness of the actual values distribution to the skewness of null distributions, and asked whether it was significantly higher than the null case.

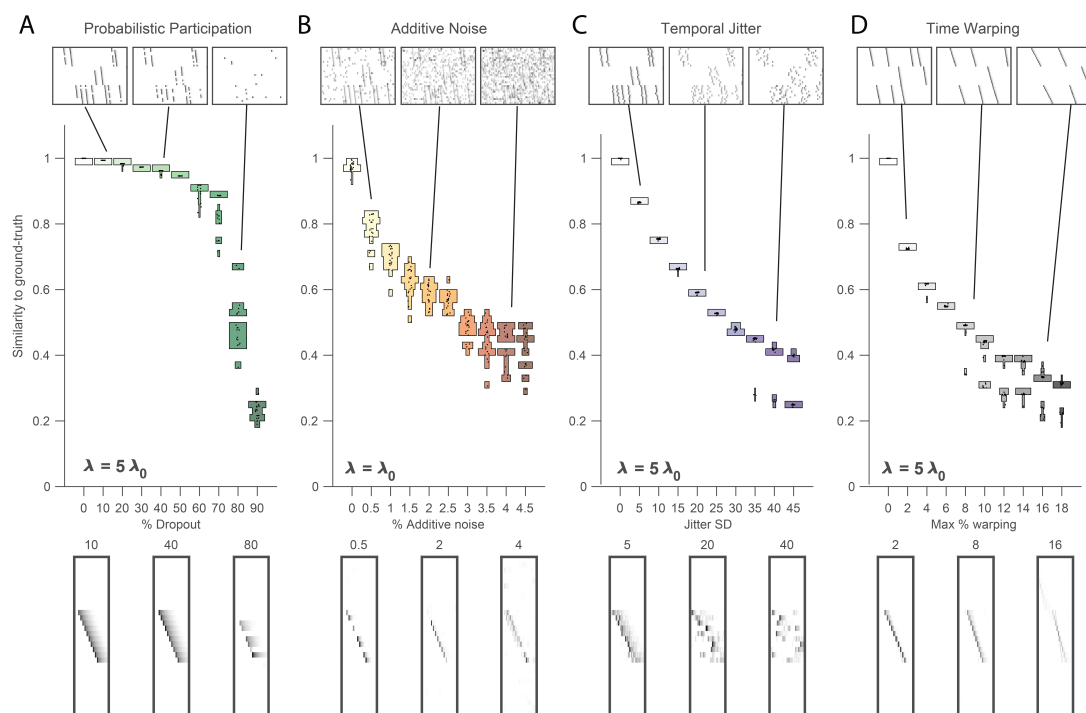


Figure S2. Robustness to noise at different values of λ

Performance of seqNMF was tested under 4 different noise conditions, at different values of λ than in Figure 3 (where $\lambda = 2\lambda_0$): **(A)** probabilistic participation, $\lambda = 5\lambda_0$, **(B)** additive noise, $\lambda = \lambda_0$ **(C)** timing jitter, $\lambda = 5\lambda_0$ and **(D)** sequence warping, $\lambda = 5\lambda_0$. For each noise type, we show: (top) examples of synthetic data at 3 different noise levels; (middle) similarity of seqNMF factors to ground-truth factors across a range of noise levels, showing 20 fits for each noise level; and (bottom) example of one of the \mathbf{W} 's extracted at 3 different noise levels (same conditions as data shown above). SeqNMF was run with $K = 20$, $L = 50$.

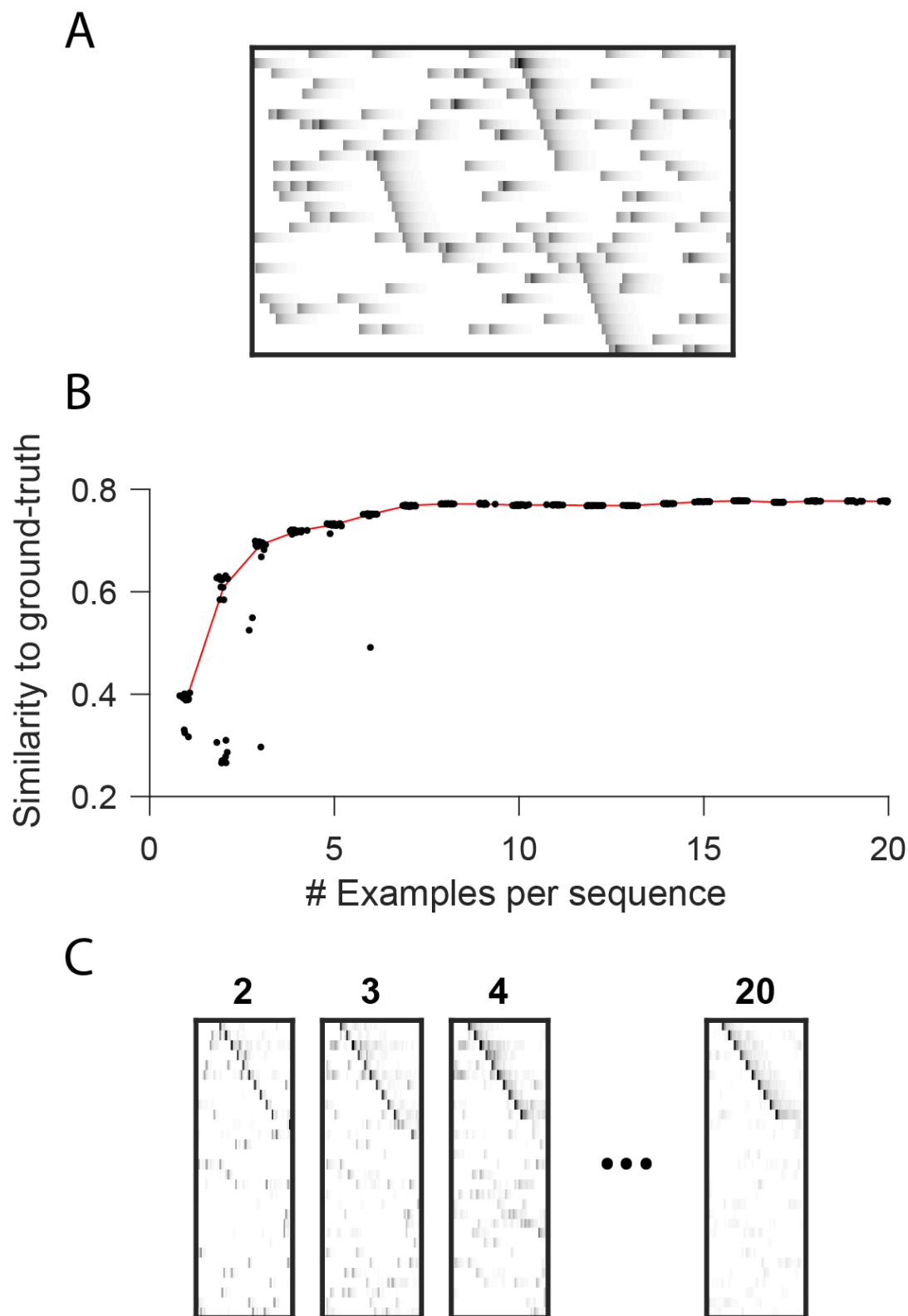


Figure S3. Effect of dataset size on seqNMF reconstruction

(A) A short (400 timestep) dataset containing one example each of three ground-truth sequences, as well as additive noise. **(B)** As a function of dataset size, similarity of extracted factors to noiseless, ground-truth factors. At each dataset size, 20 independent fits of seqNMF are shown. Median shown in red. **(C)** Example factors fit on data containing 2, 3, 4 or 20 examples of each sequence. Extracted factors were significant on held-out data compared to null (shuffled) factors even when training and test datasets each only contained only 2 examples of each sequence.

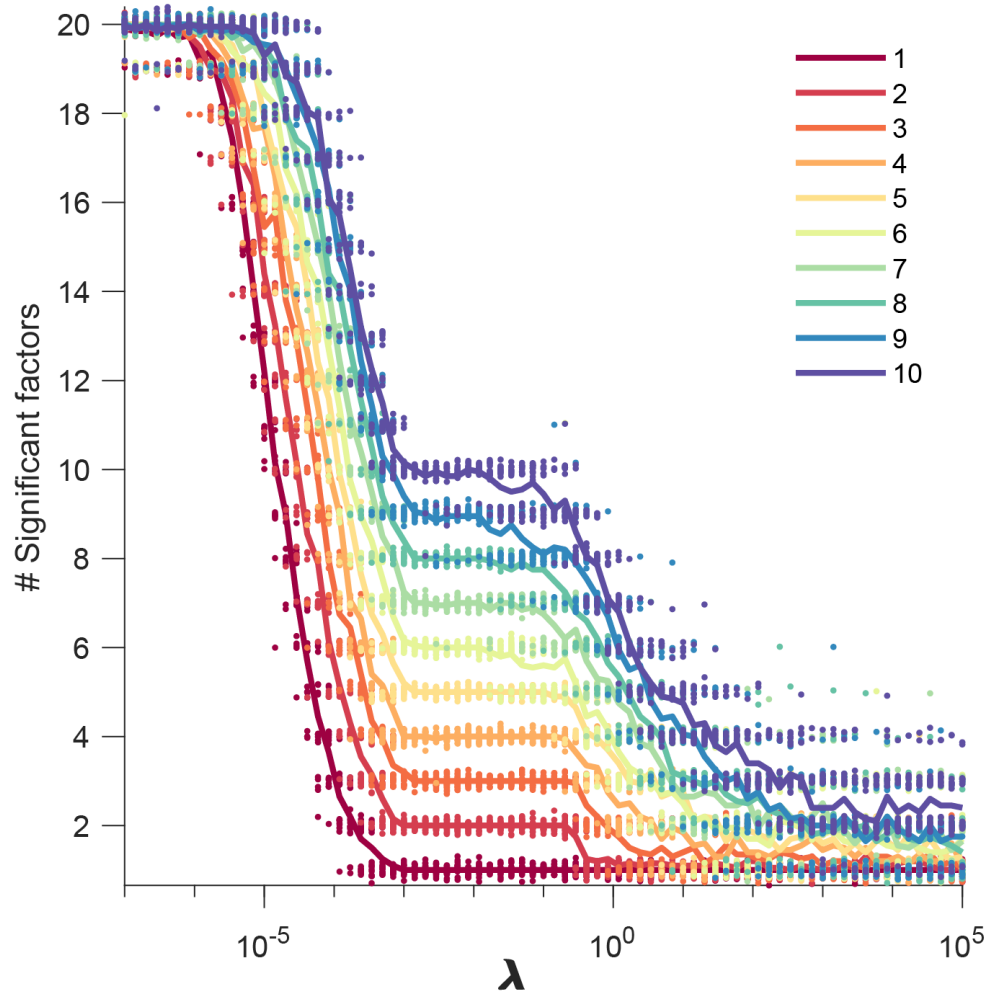


Figure S4. Number of significant factors as a function of λ for datasets containing between 1 and 10 sequences

Number of significant factors were obtained by fitting seqNMF to data containing between 1 and 10 ground truth sequences ($K = 20$, $L = 50$) for a large range of values of λ . For each dataset a λ ranging between 0.001 and 0.1 tended to return the correct number of significant sequences at least 90% of the time.

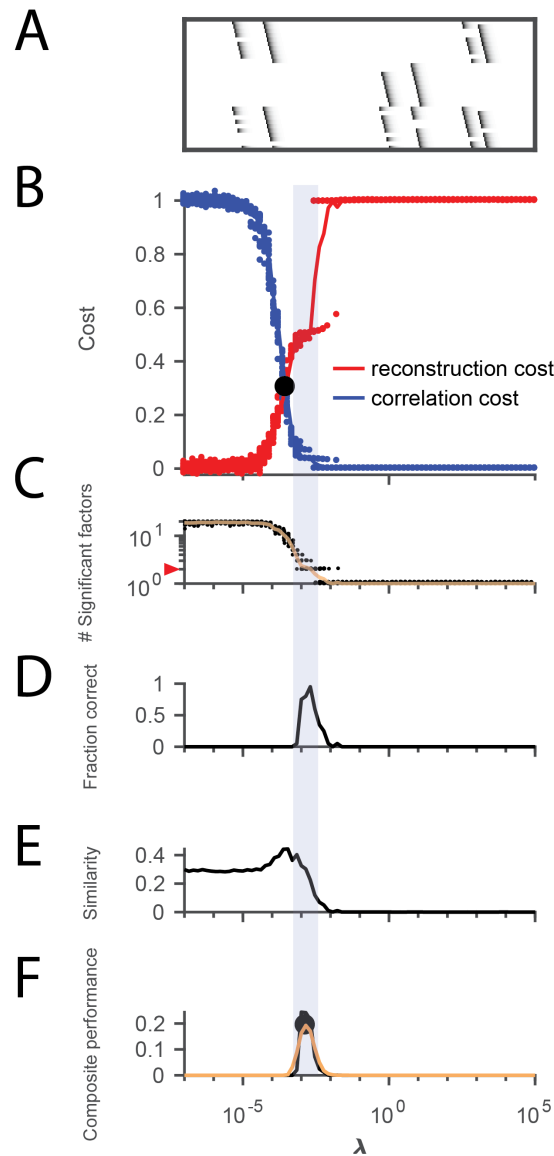


Figure S5. Procedure for choosing λ applied to data with shared neurons

(A) Simulated data containing two patterns which share 50% of their neurons, in the presence of participation noise (70% participation probability). **(B)** Normalized reconstruction cost ($\|\tilde{\mathbf{X}} - \mathbf{X}\|_F^2$) and correlation cost ($\|(\mathbf{W} \otimes \mathbf{X})\mathbf{S}\mathbf{H}^T\|_{1,i \neq j}$) as a function of λ for these data. The cross-over point λ_0 is marked with a black circle. **(C)** The number of significant factors obtained from 20 fits of these data as a function of λ (mean number plotted in orange). The correct number of factors (two) is marked by a red triangle. **(D)** The fraction of fits returning the correct number of significant factors as a function of λ . **(E)** Similarity of the top two factors to ground-truth (noiseless) factors as a function of λ . **(F)** Composite performance measured as the product of the curves shown in (D) and (E), (smoothed curve plotted in orange with a circle marking the peak). Shaded region indicates the range of λ that works well (\pm half height of composite performance). For this data set, the best performance occurs at $\lambda = 5\lambda_0$, while a range of λ between $2\lambda_0$ and $10\lambda_0$ performs well.

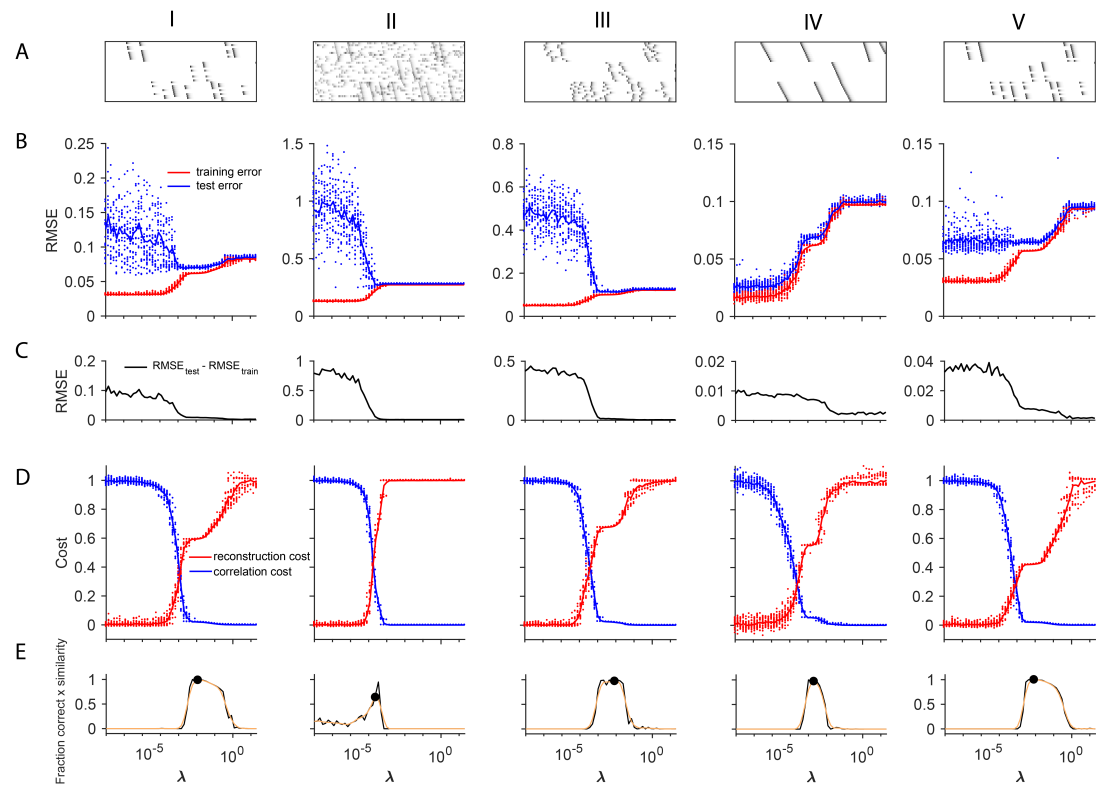


Figure S6. Using cross-validation on held-out (masked) data to choose λ

A method for choosing a reasonable value of λ based on cross validation is shown for five different noise types (each column shows a different noise type; from left to right: **(I)** participation probability, **(II)** additive noise, **(III)** jitter, **(IV)** temporal warping), and **(V)** a lower level of participation noise. For all fits, 10% of the data was held out as the test set. **(A)** Examples of each dataset. **(B)** Test error (blue) and training error (red) as a function of λ for each of the different noise conditions. **(C)** The difference between the test error and training error values shown above. **(D)** Normalized reconstruction cost ($\|\tilde{\mathbf{X}} - \mathbf{X}\|_F^2$) and correlation cost ($\|(\mathbf{W} \otimes \mathbf{X})\mathbf{S}\mathbf{H}^T\|_{1,i \neq j}$) as a function of λ for each of the different noise conditions. **(E)** Composite performance as a function of λ . Panels D and E are identical to those in Figure 4, and are included here for comparison. **(V)** These data have a lower amount of participation noise than (I). Note that in low-noise conditions, test error may not exhibit a minima within the optimal range of λ .

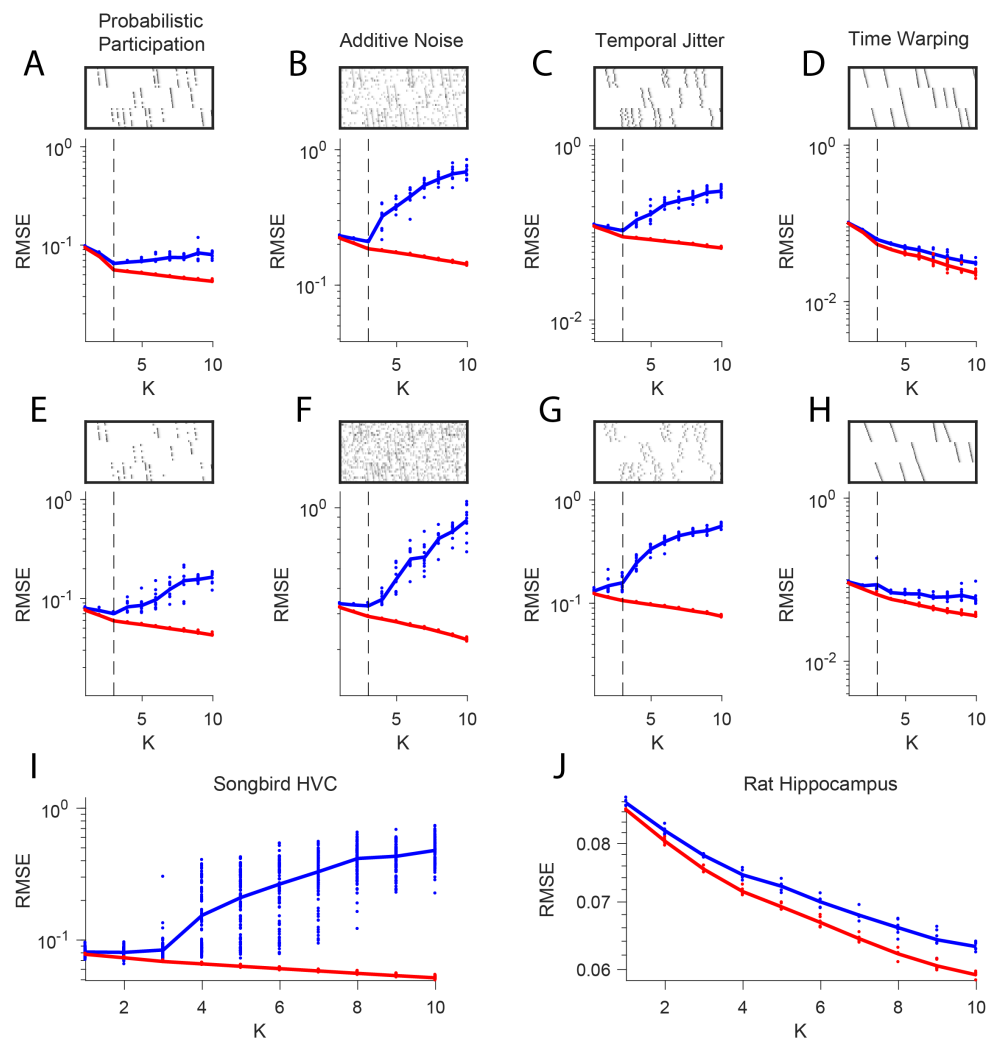


Figure S7. Estimating the number of sequences in a dataset using cross-validation on randomly masked held-out datapoints
(A) Reconstruction error (RMSE) for test data (red) and training data (blue) plotted as a function of the number of components (K) used in convNMF. Twenty independent convNMF fits are shown for each value of K . This panel shows results for 1% participation noise. For synthetic data fit fits, 10% of the data was held out as the test set. For neural data 5% of the data was held out. Other noise conditions are shown as follows: **(B)** jitter noise (10 timestep SD); **(C)** warping (13%); **(D)** higher additive noise (2.5%); **(E)** higher jitter noise (25 timestep SD); **(F)** higher warping (33%) **(G)** Reconstruction error vs. K for neuronal data collected from premotor cortex (area HVC) of a singing bird (Figure 7) and **(H)** hippocampus of a rat performing a left-right alternation task (Figure 6).

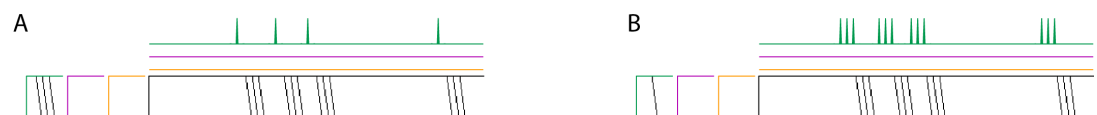


Figure S8. Biasing factorizations between sparsity in \mathbf{W} or \mathbf{H}
 Two different factorizations of the same simulated data, where a sequence is always repeated precisely three times. Both yield perfect reconstructions, and no cross-factor correlations. The factorizations differ in the amount of features placed in \mathbf{W} versus \mathbf{H} . Both use $K = 3$ and $\lambda = 0.001$. **(A)** Factorization achieved using a sparsity penalty on \mathbf{H} , with $\lambda_{L1H} = 1$. **(B)** Factorization achieved using a sparsity penalty on \mathbf{W} , with $\lambda_{L1W} = 1$.

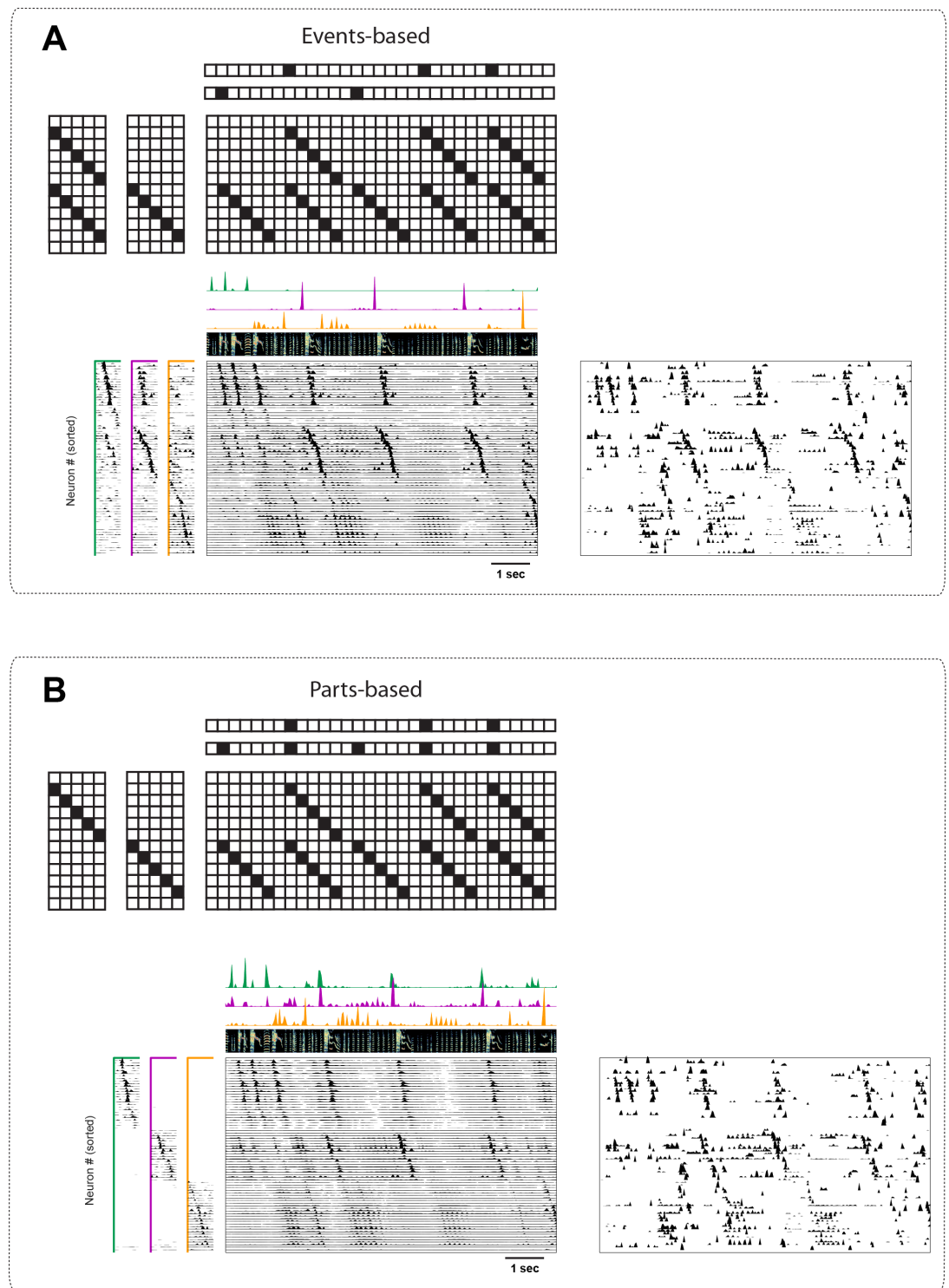


Figure S9. Events-based and parts-based factorizations of songbird data

Illustration of a trade-off between parts-based (\mathbf{W} is more strictly orthogonal) and events-based (\mathbf{H} is more strictly orthogonal) factorizations in a dataset where some neurons are shared between different sequences. The same data as in Figure 7 is factorized using smoothed soft orthogonality on \mathbf{H} (**A**, events-based), or on \mathbf{W} (**B**, parts-based). Below each motivating cartoon factorization, we show seqNMF fits (\mathbf{W} and \mathbf{H} together with the reconstruction) of the data in Figure 7. The right panels contain the raw data sorted according to these factorizations. Favoring events-based or parts-based factorizations is a matter of preference. Parts-based factorizations are particularly useful for separating neurons into ensembles. Events-based factorizations are particularly useful for identifying what neural events occur when.

Appendix 1

Deriving multiplicative update rules

Standard gradient descent methods for minimizing a cost function must be adapted when solutions are constrained to be non-negative, since gradient descent steps may result in negative values. Lee and Seung invented an elegant and widely-used algorithm for non-negative gradient descent that avoids negative values by performing multiplicative updates [34]. They derive these multiplicative updates by choosing an adaptive learning rate that makes additive terms cancel from standard gradient descent on the cost function. We will reproduce their derivation here, and detail how to extend it to the convolutional case [56] and apply several forms of regularization [43, 50, 10]. See Table 3 for a compilation of cost functions, derivatives and multiplicative updates for NMF and convNMF under several different regularization conditions.

Standard NMF

NMF performs the factorization $\mathbf{X} \approx \tilde{\mathbf{X}} = \mathbf{W}\mathbf{H}$. NMF factorizations seek to solve the following problem:

$$(\tilde{\mathbf{W}}, \tilde{\mathbf{H}}) = \arg \min_{\mathbf{W}, \mathbf{H}} \mathcal{L}(\mathbf{W}, \mathbf{H}) \quad (16)$$

$$\mathcal{L}(\mathbf{W}, \mathbf{H}) = \frac{1}{2} \|\tilde{\mathbf{X}} - \mathbf{X}\|_F^2 \quad (17)$$

$$\tilde{\mathbf{W}}, \tilde{\mathbf{H}} \geq 0 \quad (18)$$

This problem is convex in \mathbf{W} and \mathbf{H} separately, not together, so a local minimum is found by alternating \mathbf{W} and \mathbf{H} updates. Note that:

$$\frac{d}{d\mathbf{W}} \mathcal{L}(\mathbf{W}, \mathbf{H}) = \tilde{\mathbf{X}}\mathbf{H}^\top - \mathbf{X}\mathbf{H}^\top \quad (19)$$

$$\frac{d}{d\mathbf{H}} \mathcal{L}(\mathbf{W}, \mathbf{H}) = \mathbf{W}^\top \tilde{\mathbf{X}} - \mathbf{W}^\top \mathbf{X} \quad (20)$$

Thus, gradient descent steps for \mathbf{W} and \mathbf{H} are:

$$\mathbf{W} \leftarrow \mathbf{W} - \eta_{\mathbf{W}} (\tilde{\mathbf{X}}\mathbf{H}^\top - \mathbf{X}\mathbf{H}^\top) \quad (21)$$

$$\mathbf{H} \leftarrow \mathbf{H} - \eta_{\mathbf{H}} (\mathbf{W}^\top \tilde{\mathbf{X}} - \mathbf{W}^\top \mathbf{X}) \quad (22)$$

To arrive at multiplicative updates, Lee and Seung [34] set:

$$\eta_{\mathbf{W}} = \frac{\mathbf{W}}{\mathbf{W}\mathbf{H}\mathbf{H}^\top} \quad (23)$$

$$\eta_{\mathbf{H}} = \frac{\mathbf{H}}{\mathbf{W}^\top \mathbf{W}\mathbf{H}} \quad (24)$$

Thus, the gradient descent updates become multiplicative:

$$\mathbf{W} \leftarrow \mathbf{W} \times \frac{\mathbf{X}\mathbf{H}^\top}{\mathbf{W}\mathbf{H}\mathbf{H}^\top} = \mathbf{W} \times \frac{\mathbf{X}\mathbf{H}^\top}{\tilde{\mathbf{X}}\mathbf{H}^\top} \quad (25)$$

$$\mathbf{H} \leftarrow \mathbf{H} \times \frac{\mathbf{W}^\top\mathbf{X}}{\mathbf{W}^\top\mathbf{W}\mathbf{H}} = \mathbf{H} \times \frac{\mathbf{W}^\top\mathbf{X}}{\mathbf{W}^\top\tilde{\mathbf{X}}} \quad (26)$$

where the division and \times are element-wise.

Standard convNMF

Convolutional NMF factorizes data $\mathbf{X} \approx \tilde{\mathbf{X}} = \sum_{\ell} \mathbf{W}_{\cdot\cdot\ell} \overset{\ell \rightarrow}{\mathbf{H}} = \mathbf{W} \circledast \mathbf{H}$. convNMF factorizations seek to solve the following problem:

$$(\tilde{\mathbf{W}}, \tilde{\mathbf{H}}) = \arg \min_{\mathbf{W}, \mathbf{H}} \mathcal{L}(\mathbf{W}, \mathbf{H}) \quad (27)$$

$$\mathcal{L}(\mathbf{W}, \mathbf{H}) = \frac{1}{2} \|\tilde{\mathbf{X}} - \mathbf{X}\|_F^2 \quad (28)$$

$$\tilde{\mathbf{W}}, \tilde{\mathbf{H}} \geq 0 \quad (29)$$

The derivation above for standard NMF can be applied for each ℓ , yielding the following update rules for convNMF [56]:

$$\mathbf{W}_{\cdot\cdot\ell} \leftarrow \mathbf{W}_{\cdot\cdot\ell} \times \frac{\mathbf{X}\mathbf{H}^\top}{\tilde{\mathbf{X}}\mathbf{H}^\top} \quad (30)$$

$$\mathbf{H} \leftarrow \mathbf{H} \times \frac{\sum_{\ell} \mathbf{W}_{\cdot\cdot\ell}^\top \overset{\leftarrow \ell}{\mathbf{X}}}{\sum_{\ell} \mathbf{W}_{\cdot\cdot\ell}^\top \overset{\leftarrow \ell}{\tilde{\mathbf{X}}}} = \mathbf{H} \times \frac{\mathbf{W}_{\circledast}^\top \mathbf{X}}{\mathbf{W}_{\circledast}^\top \tilde{\mathbf{X}}} \quad (31)$$

Where the operator $\ell \rightarrow$ shifts a matrix in the \rightarrow direction by ℓ timebins, i.e. a delay by ℓ timebins, and $\leftarrow \ell$ shifts a matrix in the \leftarrow direction by ℓ timebins (Table 1). Note that NMF is a special case of convNMF where $L = 1$.

Incorporating regularization terms

Suppose we want to regularize by adding a new term, \mathcal{R} to the cost function:

$$(\tilde{\mathbf{W}}, \tilde{\mathbf{H}}) = \arg \min_{\mathbf{W}, \mathbf{H}} \mathcal{L}(\mathbf{W}, \mathbf{H}) \quad (32)$$

$$\mathcal{L}(\mathbf{W}, \mathbf{H}) = \frac{1}{2} \|\tilde{\mathbf{X}} - \mathbf{X}\|_F^2 + \mathcal{R} \quad (33)$$

$$\tilde{\mathbf{W}}, \tilde{\mathbf{H}} \geq 0 \quad (34)$$

Using a similar trick to Lee and Seung, we choose a $\eta_{\mathbf{W}}, \eta_{\mathbf{H}}$ to arrive at a simple multiplicative update. Below is the standard NMF case, which generalizes trivially to the convNMF case.

Note that:

$$\frac{d\mathcal{L}}{d\mathbf{W}} = \tilde{\mathbf{X}}\mathbf{H}^\top - \mathbf{X}\mathbf{H}^\top + \frac{d\mathcal{R}}{d\mathbf{W}} \quad (35)$$

$$\frac{d\mathcal{L}}{d\mathbf{H}} = \mathbf{W}^\top\tilde{\mathbf{X}} - \mathbf{W}^\top\mathbf{X} + \frac{d\mathcal{R}}{d\mathbf{H}} \quad (36)$$

We set:

$$\eta_{\mathbf{W}} = \frac{\mathbf{W}}{\tilde{\mathbf{X}}\mathbf{H}^\top + \frac{d\mathcal{R}}{d\mathbf{W}}} \quad (37)$$

$$\eta_{\mathbf{H}} = \frac{\mathbf{H}}{\mathbf{W}^\top\tilde{\mathbf{X}} + \frac{d\mathcal{R}}{d\mathbf{H}}} \quad (38)$$

Thus, the gradient descent updates become multiplicative:

$$\mathbf{W} \leftarrow \mathbf{W} - \eta_{\mathbf{W}} \frac{d\mathcal{L}}{d\mathbf{W}} = \mathbf{W} \times \frac{\mathbf{X}\mathbf{H}^\top}{\tilde{\mathbf{X}}\mathbf{H}^\top + \frac{d\mathcal{R}}{d\mathbf{W}}} \quad (39)$$

$$\mathbf{H} \leftarrow \mathbf{H} - \eta_{\mathbf{H}} \frac{d\mathcal{L}}{d\mathbf{H}} = \mathbf{H} \times \frac{\mathbf{W}^\top\mathbf{X}}{\mathbf{W}^\top\tilde{\mathbf{X}} + \frac{d\mathcal{R}}{d\mathbf{H}}} \quad (40)$$

where the division and \times are element-wise.

This framework enables flexible incorporation of different types of regularization or penalty terms into the multiplicative NMF update algorithm. This framework also extends naturally to the convolutional case. See Table 3 for examples of several regularization terms, including $L1$ sparsity [43, 50] and soft orthogonality [10], as well as the terms we introduce here to combat the types of inefficiencies and cross correlations we identified in convolutional NMF, namely, smoothed orthogonality for \mathbf{H} and \mathbf{W} , and smoothed cross-factor orthogonality, the primary seqNMF regularization term. For the seqNMF regularization term, $\lambda \|(\mathbf{W}^\top \otimes \mathbf{X})\mathbf{S}\mathbf{H}^\top\|_{1, i \neq j}$, the multiplicative update rules are:

$$\mathbf{W}_{.. \ell} \leftarrow \mathbf{W}_{.. \ell} \times \frac{\mathbf{X} \left(\begin{smallmatrix} \ell \rightarrow \\ \mathbf{H} \end{smallmatrix} \right)^\top}{\tilde{\mathbf{X}} \left(\begin{smallmatrix} \ell \rightarrow \\ \mathbf{H} \end{smallmatrix} \right)^\top + \lambda \overset{\leftarrow \ell}{\mathbf{X}} \mathbf{S} \mathbf{H}^\top (\mathbf{1} - \mathbf{I})} \quad (41)$$

$$\mathbf{H} \leftarrow \mathbf{H} \times \frac{\mathbf{W}^\top \otimes \mathbf{X}}{\mathbf{W}^\top \otimes \tilde{\mathbf{X}} + \lambda (\mathbf{1} - \mathbf{I}) (\mathbf{W}^\top \otimes \mathbf{X} \mathbf{S})} \quad (42)$$

Where the division and \times are element-wise. Note that multiplication with the $K \times K$ matrix $(\mathbf{1} - \mathbf{I})$ effectively implements factor competition because it places in the k th row a sum across all other factors.