

# 1 Invariant neural subspaces 2 maintained by feedback modulation

3 Laura Bella Naumann<sup>1,2\*</sup>, Joram Keijser<sup>1</sup>, Henning Sprekeler<sup>1,2</sup>

\*For correspondence:

laurabella-naumann@bccn-berlin.  
de (LBN)

4 <sup>1</sup>Modelling of Cognitive Processes, Technical University Berlin, Germany; <sup>2</sup>Bernstein  
5 Center for Computational Neuroscience, Berlin, Germany

---

7 **Abstract** Sensory systems reliably process incoming stimuli in spite of changes in context.  
8 Most recent models accredit this context invariance to an extraction of increasingly complex  
9 sensory features in hierarchical feedforward networks. Here, we study how context-invariant  
10 representations can be established by feedback rather than feedforward processing. We show  
11 that feedforward neural networks modulated by feedback can dynamically generate invariant  
12 sensory representations. The required feedback can be implemented as a slow and spatially  
13 diffuse gain modulation. The invariance is not present on the level of individual neurons, but  
14 emerges only on the population level. Mechanistically, the feedback modulation dynamically  
15 reorients the manifold of neural activity and thereby maintains an invariant neural subspace in  
16 spite of contextual variations. Our results highlight the importance of population-level analyses  
17 for understanding the role of feedback in flexible sensory processing.

---

## 19 Introduction

20 In natural environments our senses are exposed to a colourful mix of sensory impressions. Be-  
21 haviourally relevant stimuli can appear in varying contexts, such as variations in lighting, acous-  
22 tics, stimulus position or the presence of other stimuli. Different contexts may require different  
23 responses to the same stimulus, for example when the behavioural task changes (context depen-  
24 dence). Alternatively, the same response may be required for different stimuli, for example when  
25 the sensory context changes (context invariance). Recent advances have elucidated how context-  
26 *dependent* processing can be performed by recurrent feedback in neural circuits (*Mante et al.,*  
27 *2013; Wang et al., 2018b; Dubreuil et al., 2020*). In contrast, the role of feedback mechanisms in  
28 context-*invariant* processing is not well understood.

29 In the classical view, stimuli are hierarchically processed towards a behaviourally relevant per-  
30 cept that is invariant to contextual variations. This is achieved by extracting increasingly complex  
31 features in a feedforward network (*Kriegeskorte, 2015; Zhuang et al., 2021; Yamins and DiCarlo,*  
32 *2016*). Models of such feedforward networks have been remarkably successful at learning com-  
33 plex perceptual tasks (*LeCun et al., 2015*), and they account for various features of cortical sensory  
34 representations (*DiCarlo and Cox, 2007; Kriegeskorte et al., 2008; DiCarlo et al., 2012; Hong et al.,*  
35 *2016; Cichy et al., 2016*). Yet, these models neglect feedback pathways, which are abundant in sen-  
36 sory cortex (*Felleman and Van Essen, 1991; Markov et al., 2014*) and shape sensory processing in  
37 critical ways (*Gilbert and Li, 2013*). Incorporating these feedback loops into models of sensory pro-  
38 cessing increases their flexibility and robustness (*Spoerer et al., 2017; Alamia et al., 2021; Nayebi*  
39 *et al., 2021*) and improves their fit to neural data (*Kar et al., 2019; Kietzmann et al., 2019; Nayebi*  
40 *et al., 2021*). At the neuronal level, feedback is thought to modulate rather than drive local re-  
41 sponses (*Sherman and Guillery, 1998*), for instance depending on behavioral context (*Niell and*  
42 *Stryker, 2010; Vinck et al., 2015; Kuchibhotla et al., 2017; Dipoppa et al., 2018*).

43 Here, we investigate the hypothesis that feedback modulation provides a neural mechanism  
 44 for context-invariant perception. To this end, we trained a feedback-modulated network model  
 45 to perform a context-invariant perceptual task and studied the resulting neural mechanisms. We  
 46 show that the feedback modulation does not need to be temporally or spatially precise and can be  
 47 realised by feedback-driven gain modulation in rate-based networks of excitatory and inhibitory  
 48 neurons. To solve the task, the feedback loop dynamically maintains an invariant subspace in the  
 49 population representation (*Hong et al., 2016*). This invariance is not present at the single neuron  
 50 level. Finally, we find that the feedback conveys a nonlinear representation of the context itself,  
 51 which can be hard to discern by linear decoding methods.

52 These findings corroborate that feedback-driven gain modulation of feedforward networks en-  
 53 ables context-invariant sensory processing. The underlying mechanism links single neuron mod-  
 54 ulation with its function at the population level, highlighting the importance of population-level  
 55 analyses.

## 56 Results

57 As a simple instance of a context-invariant task, we considered a dynamic version of the blind  
 58 source separation problem. The task is to recover unknown sensory sources, such as voices at a  
 59 cocktail party (*McDermott, 2009*), from sensory stimuli that are an unknown mixture of the sources.  
 60 In contrast to the classical blind source separation problem, the mixture can change in time, for  
 61 example, when the speakers move around, thus providing a time-varying sensory context. Because  
 62 the task requires a dynamic inference of the context, it cannot be solved by feedforward networks  
 63 (*Figure 1–Figure Supplement 1*) or standard blind source separation algorithms (e.g., independent  
 64 component analysis; *Bell and Sejnowski, 1995; Hyvärinen and Oja, 2000*). We hypothesised that  
 65 this dynamic task can be solved by a feedforward network that is subject to modulation from a  
 66 feedback signal. In our model the feedback signal is provided by a modulatory system that receives  
 67 both the sensory stimuli and the network output (*Figure 1a*).

### 68 Dynamic blind source separation by modulation of feedforward weights

Before we gradually take this to the neural level, we illustrate the proposed mechanism in a simple  
 example, in which the modulatory system provides a time-varying multiplicative modulation of a  
 linear two-layer network (see Methods and Models). For illustration, we used compositions of sines  
 with different frequencies as source signals ( $s$ , *Figure 1b*, top). These sources were linearly mixed to  
 generate the sensory stimuli ( $x$ ) that the network received as input;  $x = A_t s$  (*Figure 1a,b*). The linear  
 mixture ( $A_t$ ) changed over time, akin to varying the location of sound sources in a room (*Figure 1a*).  
 These locations provided a time-varying sensory context that changed on a slower timescale than  
 the sources themselves. The feedforward network had to recover the sources from the mixed  
 sensory stimuli. To achieve this, we trained the modulator to dynamically adjust the weights of the  
 feedforward network ( $W_0$ ) such that the network output ( $y$ ) matches the sources:

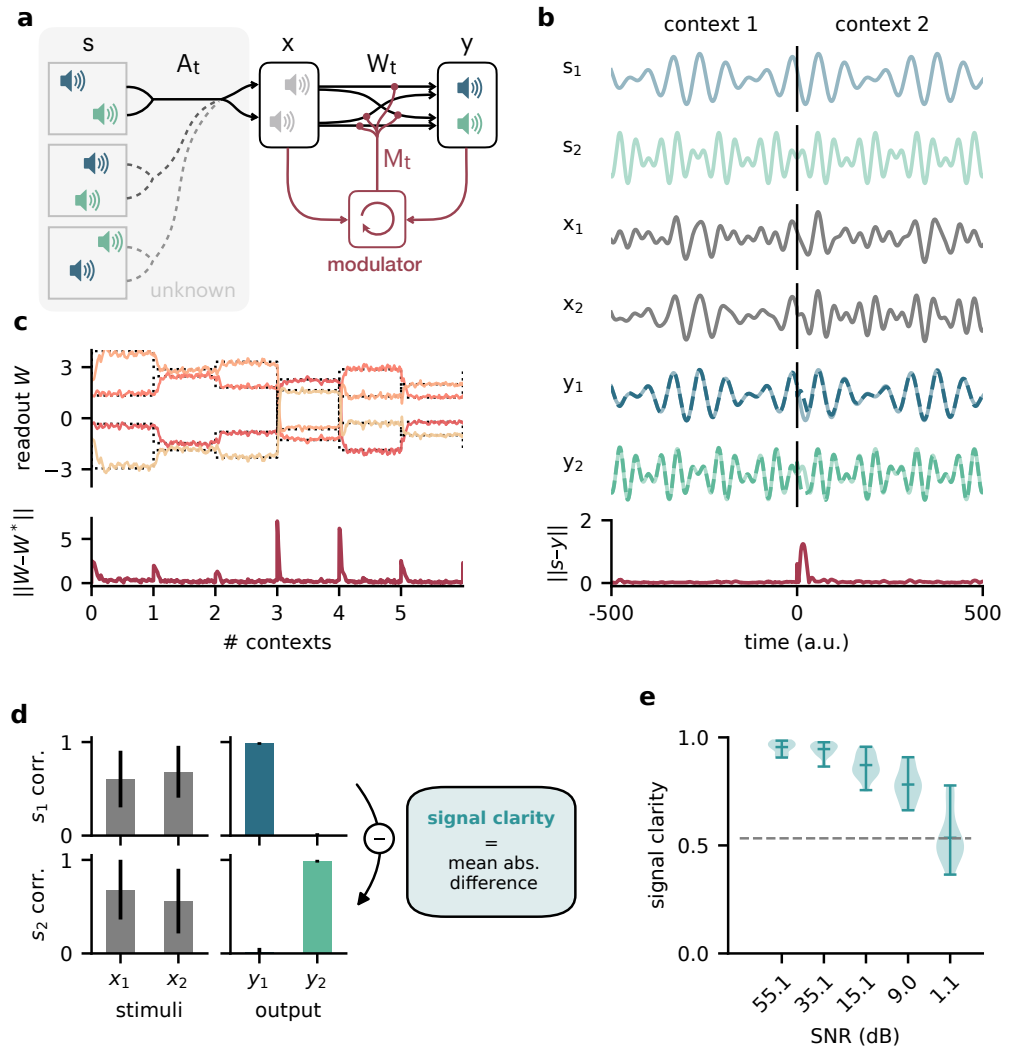
$$y = W_t x = (M_t \odot W_0) x$$

$$M_t = \text{modulator}(\text{history of } x, y).$$

69 Because the modulation requires a dynamic inference of the context, the modulator is a recurrent  
 70 neural network. The modulator was trained using supervised learning. Afterwards, its weights  
 71 were fixed and it no longer had access to the target sources (see Methods and Models, *Figure 8*).  
 72 The modulator therefore had to use its recurrent dynamics to determine the appropriate modula-  
 73 tory feedback for the time-varying context, based on the sensory stimuli and the network output.  
 74 Put differently, the modulator had to learn an internal model of the sensory data and the contexts,  
 75 and use it to establish the desired context invariance in the output.

76 After learning, the modulated network disentangled the sources, even when the context changed  
 77 (*Figure 1b, Figure 1–Figure Supplement 1a,b*). Context changes produced a transient error in the

78 network's output, but it quickly resumed matching the sources (**Figure 1b**, bottom). The transient  
 79 errors occur, because the modulator needs time to infer the new context from the time-varying



**Figure 1.** Dynamic blind source separation by modulation of feedforward connections.

**a.** Schematic of the feedforward network model receiving feedback modulation from a modulator (a recurrent network). **b.** Top: Sources ( $s_{1,2}$ ), sensory stimuli ( $x_{1,2}$ ) and network output ( $y_{1,2}$ ) for two different source locations (contexts). Bottom: Deviation of output from the sources. **c.** Top: Modulated readout weights across 6 contexts (source locations); dotted lines indicate the true weights of the inverted mixing matrix. Bottom: Deviation of readout from target weights. **d.** Correlation between the sources and the sensory stimuli (left), the network outputs (center), and calculation of the *signal clarity* (right). Errorbars indicate standard deviation across 20 contexts. **e.** Violin plot of the signal clarity for different noise levels in the sensory stimuli across 20 different contexts.

**Figure 1-Figure supplement 1.** The dynamic blind source separation task cannot be solved with a feedforward network.

**Figure 1-Figure supplement 2.** Robustness of the feedback-driven modulation mechanism.

**Figure 1-Figure supplement 3.** Model performance for two different sets of source signals.

**Figure 1-Figure supplement 4.** Model performance for three source signals.

**Figure 1-Figure supplement 5.** The modulated network model generalises across frequencies.

**Figure 1-Figure supplement 6.** The modulator learns a model of the sources and contexts, and infers the current context from the stimuli.

80 inputs, before it can provide the appropriate feedback signal to the feedforward network (*Figure 1–*  
81 *Figure Supplement 6a*, cf. *Figure 1–Figure Supplement 1g-i*). The modulated feedforward weights  
82 inverted the linear mixture of sources by switching on the same timescale (*Figure 1c*).

83 To quantify how well the sources were separated, we measured the correlation coefficient of  
84 the outputs with each source over several contexts. Consistent with a clean separation, we found  
85 that each of the two outputs strongly correlated with only one of the sources. In contrast, the sen-  
86 sory stimuli showed a positive average correlation for both sources, as expected given the positive  
87 linear mixture (*Figure 1d*, left). We determined the *signal clarity* as the absolute difference between  
88 the correlation with the first compared to the second source, averaged over the two outputs, nor-  
89 malised by the sum of the correlations (*Figure 1d*, right; see *Methods and Models*). The signal  
90 clarity thus determines the degree of signal separation, where a value close to 1 indicates a clean  
91 separation as in *Figure 1d*. Note that the signal clarity of the sensory stimuli is around 0.5 and can  
92 be used as a reference.

93 We next probed the network’s robustness by adding noise to the sensory stimuli. We found that  
94 the signal clarity gradually decreased with increasing noise levels, but only degraded to chance per-  
95 formance when the signal-to-noise ratio was close to 1 (1.1 dB, *Figure 1e*, *Figure 1–Figure Supple-*  
96 *ment 2e*). The network performance did not depend on the specific source signals (*Figure 1–Figure*  
97 *Supplement 3*) or the number of sources (*Figure 1–Figure Supplement 4*), as long as it had seen  
98 them during training. Yet, because the network had to learn an internal model of the task, we  
99 expected a limited degree of generalisation to new situations. Indeed, the network was able to  
100 interpolate between source frequencies seen during training (*Figure 1–Figure Supplement 5*), but  
101 failed on sources and contexts that were qualitatively different (*Figure 1–Figure Supplement 6b-d*).  
102 The specific computations performed by the modulator are therefore idiosyncratic to the prob-  
103 lem at hand. Hence, we did not investigate the internal dynamics of the modulator in detail, but  
104 concentrated on its effect on the feedforward network.

105 Since feedback-driven modulation enables flexible context-invariant processing in a simple ab-  
106 stract model, we wondered how this mechanism might be implemented at the neural level. For  
107 example, how does feedback-driven modulation function when feedback signals are slow and im-  
108 precise? And how does the modulation affect population activity? In the following, we will gradually  
109 increase the model complexity to account for biological constraints and pinpoint the population-  
110 level mechanisms of feedback-mediated invariance.

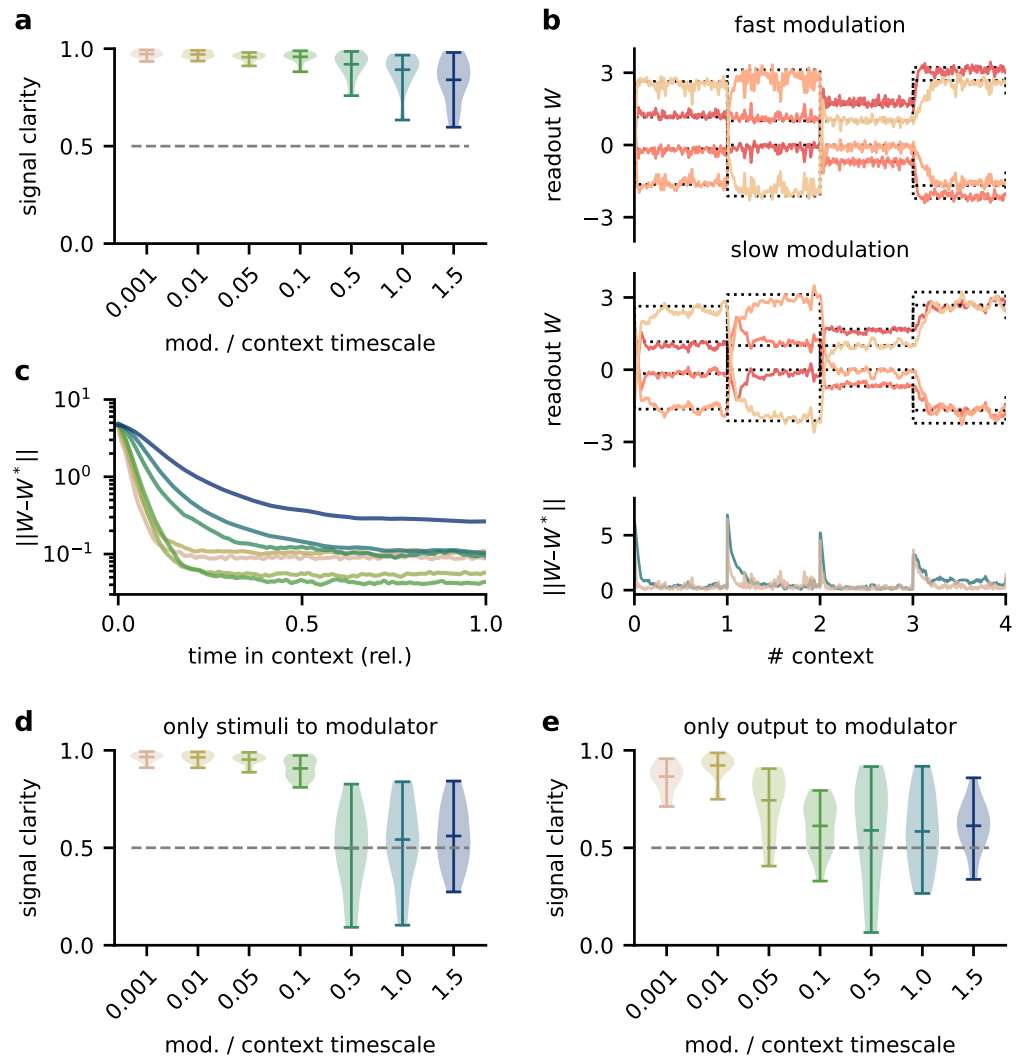
### 111 **Invariance can be established by slow feedback modulation**

112 Among the many modulatory mechanisms, even the faster ones are believed to operate on timescales  
113 of hundreds of milliseconds (*Bang et al., 2020; Molyneaux and Hasselmo, 2002*), raising the ques-  
114 tion if feedback-driven modulation is sufficiently fast to compensate for dynamic changes in envi-  
115 ronmental context.

116 To investigate how the timescale of modulation affects the performance in the dynamic blind  
117 source separation task, we trained network models, in which the modulatory feedback had an  
118 intrinsic timescale that forced it to be slow. We found that the signal clarity degraded only when this  
119 timescale was on the same order of magnitude as the timescale of contextual changes (*Figure 2a*).  
120 Note that timescales in this model are relative, and could be arbitrarily rescaled. While slower  
121 feedback modulation produced a larger initial error (*Figure 2b,c*), it also reduced the fluctuations  
122 in the readout weights such that they more closely follow the optimal weights (*Figure 2b*). This  
123 speed-accuracy trade-off explains the lower and more variable signal clarity for slow modulation  
124 (*Figure 2a*), because the signal clarity was measured over the whole duration of a context and the  
125 transient onset error dominated over the reduced fluctuations.

126 To determine architectural constraints on the modulatory system, we asked how these results  
127 depended on the input it received. So far, the modulatory system received the feedforward net-  
128 work’s inputs (the sensory stimuli) and its outputs (the inferred sources, see *Figure 1a*), but are  
129 both of these necessary to solve the task? We found that when the modulatory system only re-

130 ceived the sensory stimuli, the model could still learn the task, though it was more sensitive to  
 131 slow modulation (**Figure 2d**, Supp. **Figure 2–Figure Supplement 1**). When the modulatory system  
 132 had to rely on the network output alone, task performance was impaired even for fast modulation  
 133 (**Figure 2e**, **Figure 2–Figure Supplement 1**). Thus, while the modulatory system is more robust to  
 134 slow modulation when it receives the network output, the output is not sufficient to solve the task.  
 135 Taken together, these results show that the biological timescale of modulatory mechanisms  
 136 does not pose a problem for flexible feedback-driven processing, as long as the feedback modulation  
 137 changes on a faster timescale than variations in the context. In fact, slow modulation can



**Figure 2.** The network model is not sensitive to slow feedback modulation.

**a.** Signal clarity in the network output for varying timescales of modulation relative to the intervals at which the source locations change. **b.** Modulated readout weights across 4 source locations (contexts) for fast (top) and slow (center) feedback modulation; dotted lines indicate the optimal weights (the inverse of the mixing matrix). Bottom: deviation of the readout weights from the optimal weights for fast and slow modulation. Colours correspond to the relative timescales in (a). Fast and slow timescales are 0.001 and 1, respectively. **c.** Mean deviation of readout from optimal weights within contexts; averaged over 20 contexts. Colours code for timescale of modulation (see (a)). **d. & e.** Same as (a) but for models in which the modulatory system only received the sensory stimuli  $x$  or the network output  $y$ , respectively.

**Figure 2–Figure supplement 1.** Robustness to slow feedback modulation depends on the inputs to the modulatory system.

138 increase processing accuracy by averaging out fluctuations in the feedback signal. Nevertheless,  
139 slow modulation likely requires the modulatory system to receive both the input and output of the  
140 sensory system it modulates.

### 141 **Invariance can be established by spatially diffuse feedback modulation**

142 Neuromodulators are classically believed to diffusely affect large areas of the brain. Furthermore,  
143 signals in the brain are processed by populations of neurons. We wondered if the proposed mod-  
144 ulation mechanism is consistent with such biological constraints. We therefore extended the net-  
145 work model such that the sensory stimuli are projected to a population of 100 neurons. A fixed  
146 linear readout of this population determined the network output. The neurons in the population  
147 received spatially diffuse modulatory feedback (*Figure 3a*) such that the feedback modulation af-  
148 fected neighbouring neurons similarly. We here assume that all synaptic weights to a neuron re-  
149 ceive the same modulation, such that the feedback performs a gain modulation of neural activ-  
150 ity (*Ferguson and Cardin, 2020*). The spatial specificity of the modulation was determined by the  
151 number of distinct feedback signals and their spatial spread (*Figure 3b, Figure 3–Figure Supple-*  
152 *ment 1a*).

153 This population-based model with less specific feedback modulation could still solve the dy-  
154 namic blind source separation task. The diffuse feedback modulation switched when the context  
155 changed, but was roughly constant within contexts (*Figure 3c*), as in the simple model. The effec-  
156 tive weight from the stimuli to the network output also inverted the linear mixture of the sources  
157 (*Figure 3–Figure Supplement 1d*, cf. *Figure 1c*).

158 We found that only a few distinct feedback signals were needed for a clean separation of the  
159 sources across contexts (*Figure 3d*). Moreover, the feedback could have a spatially broad effect on  
160 the modulated population without degrading the signal clarity (*Figure 3e, Figure 3–Figure Supple-*  
161 *ment 1*), consistent with the low dimensionality of the context.

162 We conclude that, in our model, neuromodulation does not need to be spatially precise to en-  
163 able flexible processing. Given that the suggested feedback-driven modulation mechanism works  
164 for slow and diffuse feedback signals, it could in principle be realised by neuromodulatory path-  
165 ways present in the brain.

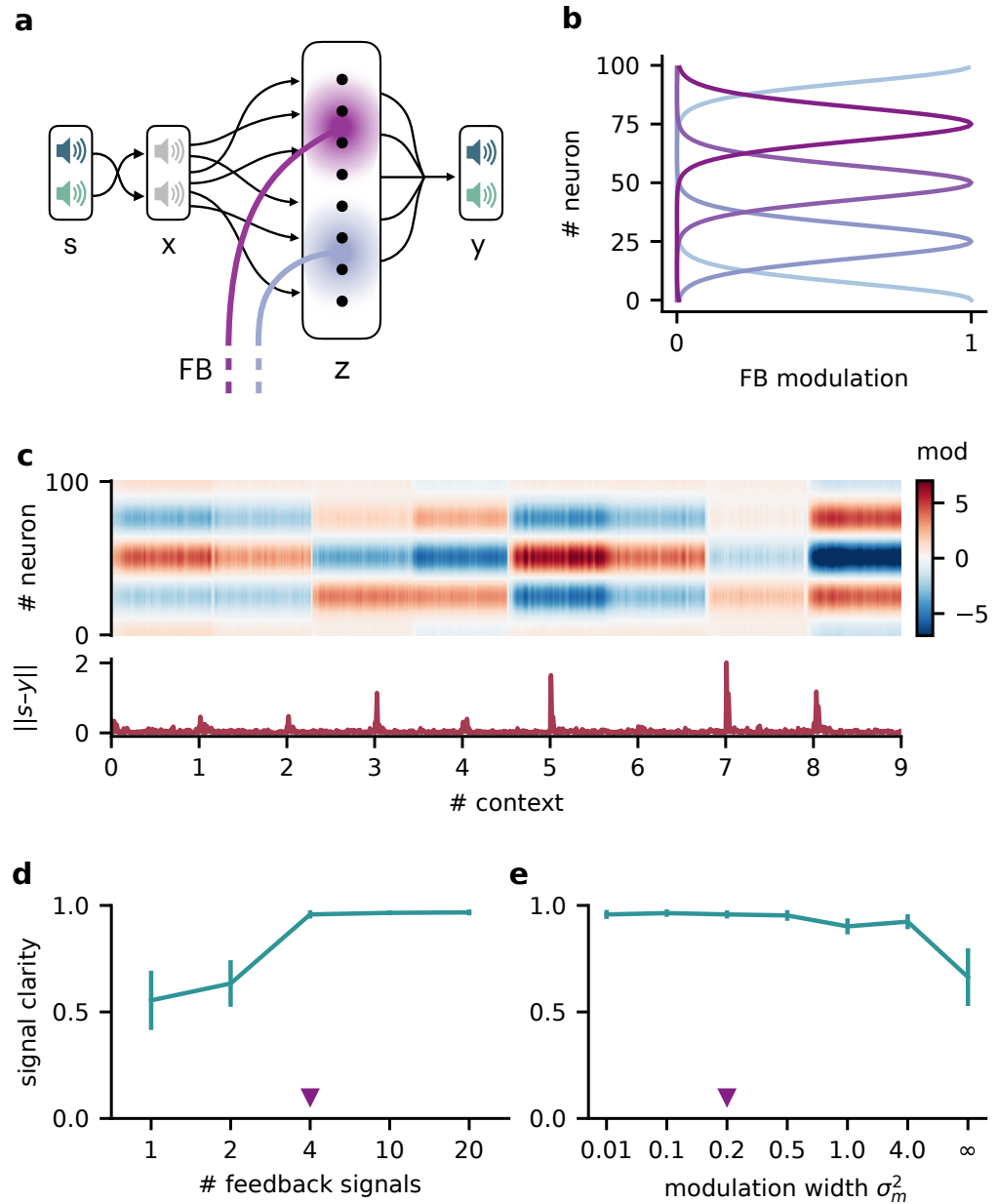
### 166 **Invariance emerges at the population level**

167 Having established that slow and spatially diffuse feedback modulation enables context-invariant  
168 processing, we next investigated the underlying mechanisms at the single neuron and population  
169 level. Given that the readout of the population activity was fixed, it is not clear how the context-  
170 dependent modulation of single neurons could give rise to a context-independent network output.  
171 One possible explanation is that some of the neurons are context-invariant and are ex-  
172 ploited by the readout. However, a first inspection of neural activity indicated that single neurons  
173 are strongly modulated by context (*Figure 4a*). To quantify this, we determined the signal clarity for  
174 each neuron at each stage of the feedforward network, averaged across contexts (*Figure 4b*). As  
175 expected, the signal clarity was low for the sensory stimuli. Intriguingly, the same was true for all  
176 neurons of the modulated neural population, indicating no clean separation of the sources at the  
177 level of single neurons. Although most neurons had a high signal clarity in some of the contexts,  
178 there was no group of neurons that consistently represented one or the other source (*Figure 4c*).  
179 Furthermore, the average signal clarity of the neurons did not correlate with their contribution to  
180 the readout (*Figure 4d*). Since single neuron responses were not invariant, context invariance must  
181 arise at the population level.

182 To confirm this, we asked how well the sources could be decoded at different stages of the  
183 feedforward network. We trained a single linear decoder of the sources on one set of contexts  
184 and tested its generalisation to novel contexts. We found that the decoding performance was  
185 poor for the sensory stimuli (*Figure 4e*), indicating that these did not contain a context-invariant

186 representation. In contrast, the sources could be decoded with high accuracy from the modulated  
 187 population.

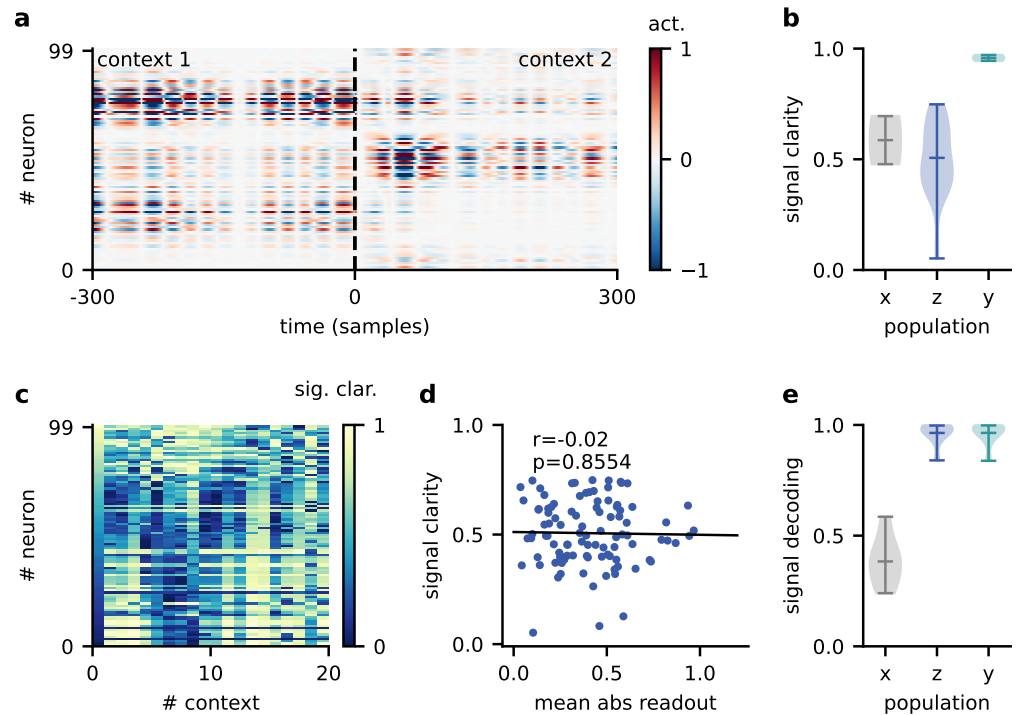
188 This demonstrates that while individual neurons were not invariant, the population activity con-  
 189 tained a context-invariant subspace. In fact, the population had to contain an invariant subspace,  
 190 because the fixed linear readout of the population was able to extract the sources across contexts.



**Figure 3.** Feedback modulation in the model can be spatially diffuse.

**a.** Schematic of the feedforward network with a population that receives diffuse feedback-driven modulation.  
**b.** Spatial spread of the modulation mediated by 4 modulatory feedback signals with a width of 0.2. **c.** Top: Per neuron modulation during 8 different contexts. Bottom: Corresponding deviation of the network output from sources. **d.** Mean signal clarity across 20 contexts for different numbers of feedback signals; modulation width is 0.2. Error bars indicate standard deviation. Purple triangle indicates default parameters used in (c). **e.** Same as (d) but for different modulation widths; number of feedback signals is 4. The modulation width " $\infty$ " corresponds to uniform modulation across the population.

**Figure 3-Figure supplement 1.** Robustness to the spatial scale of feedback modulation.



**Figure 4.** Invariance emerges at the population level. **a.** Population activity in two contexts. **b.** Violin plot of the signal clarity in the sensory stimuli ( $x$ ), neural population ( $z$ ), and network output ( $y$ ), computed across 20 different contexts. **c.** Signal clarity of single neurons in the modulated population for different contexts. **d.** Correlation between average signal clarity over contexts and magnitude of neurons' readout weight. Corresponding Pearson  $r$  and  $p$ -value are indicated in the panel. **e.** Violin plot of the linear decoding performance of the sources from different stages of the feedforward network, computed across 20 contexts. The decoder was trained on a different set of 20 contexts.

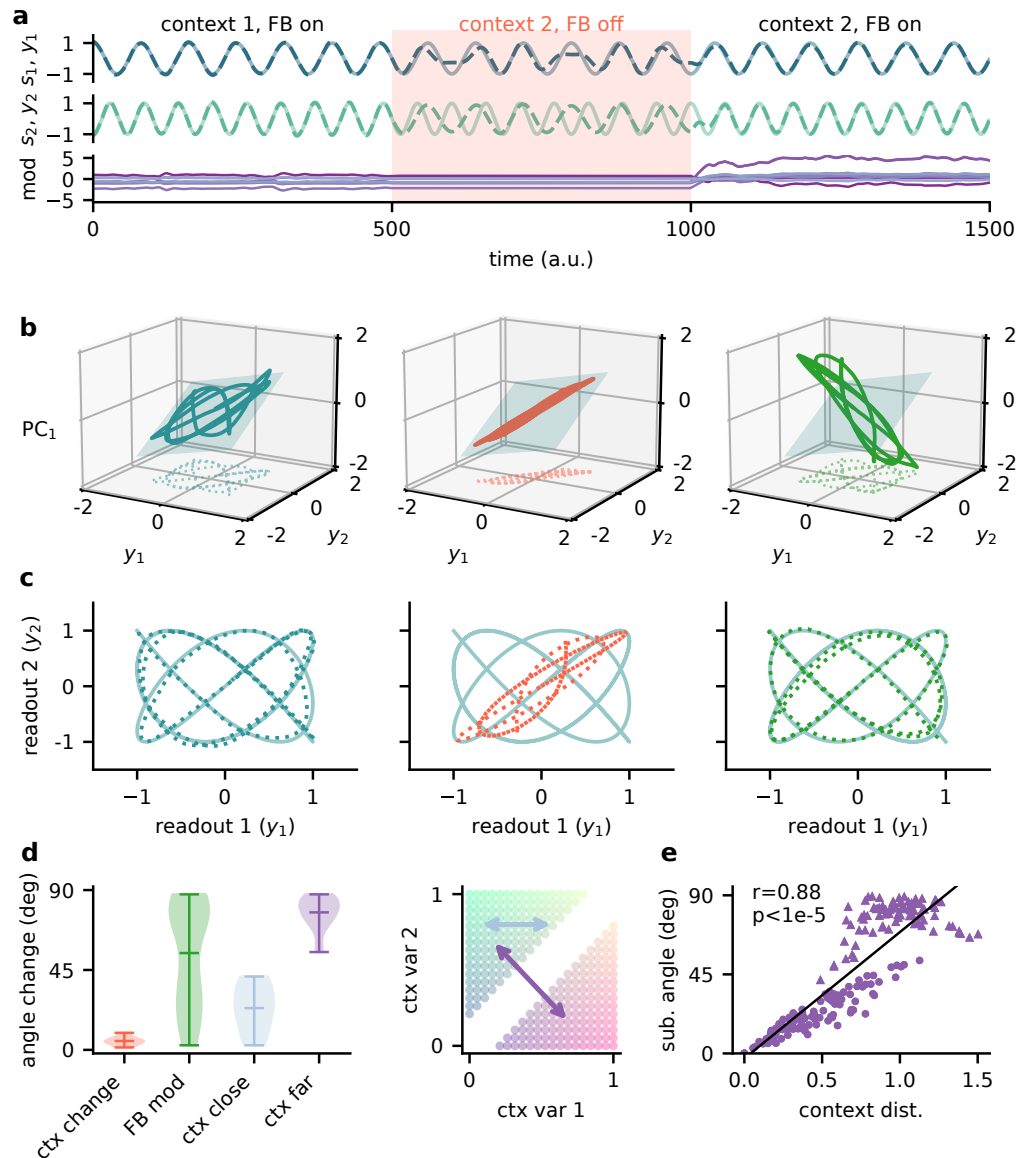
191 However, the linear decoding approach shows that this subspace can be revealed from the popu-  
 192 lation activity itself with only a few contexts and no knowledge of how the neural representation  
 193 is used downstream. The same approach could therefore be used to reveal context-invariant sub-  
 194 spaces in neural data from population recordings. Note, that the learned readout and the decoder  
 195 obtained from population activity are not necessarily identical, due to the high dimensionality of  
 196 the population activity compared to the sources.

### 197 **Feedback re-oriens the population representation**

198 The question remains how exactly the context-invariant subspace is maintained by feedback mod-  
 199 ulation. In contrast to a pure feedforward model of invariant perception (*Kriegeskorte, 2015*;  
 200 *Yamins and DiCarlo, 2016*), feedback-mediated invariance requires time to establish after contex-  
 201 tual changes. Experimentally, hallmarks of this adaptive process should be visible when comparing  
 202 the population representations immediately after a change and at a later point in time. Our model  
 203 allows to cleanly separate the early and the late representation by freezing the feedback signals  
 204 in the initial period after a contextual change (*Figure 5a*), thereby disentangling the effects of feed-  
 205 back and context on population activity.

206 The simulated experiment consisted of three stages: First, the feedback was intact for a particu-  
 207 lar context and the network outputs closely tracked the sources. Second, the context was changed  
 208 but the feedback modulation was frozen at the same value as before. As expected, this produced  
 209 deviations of the output from the sources. Third, for the same context the feedback modulation  
 210 was turned back on, which reinstated the source signals in the output. In this experiment, we used





**Figure 5.** Feedback re-orientates the population representation.

**a.** Network output (top) and feedback modulation (bottom) for two contexts. The feedback modulation is frozen for the initial period after the context changes. **b.** Population activity in the space of the two readout axes and the first principal component. Projection onto the readout is indicated in the bottom (see (c)). The signal representation is shown for different phases of the experiment. Left: context 1 with intact feedback, center: context 2 with frozen feedback, right: context 2 with intact feedback. The blue plane spans the population activity subspace in context 1 (left). **c.** Same as (b), but projected onto the readout space (dotted lines in (b)). The light blue trace corresponds to the sources. **d.** Left: Change in subspace orientation across 40 repetitions of the experiment, measured by the angle between the original subspace and the subspace for context changes (ctx change), feedback modulation (FB mod) and feedback modulation for similar contexts (ctx close) or dissimilar contexts (ctx far). Right: two-dimensional context space, defined by the coefficients in the mixing matrix. Arrows indicate similar (light blue) and dissimilar contexts (purple). **e.** Distance between pairs of contexts versus the angle between population activity subspaces for these contexts. Circles indicate similar contexts (from the same side of the diagonal, see (d)) and triangles dissimilar contexts (from different sides of the diagonal). Pearson  $r$  and  $p$ -value indicated in the panel.

**Figure 5-Figure supplement 1.** Principal component analysis captures the low-dimensional population subspaces and the subspace re-orientation with feedback.

211 pure sines as signals for visualisation purposes (*Figure 5a,c*). To visualise the population activity  
212 in the three stages of the experiment, we considered the space of the two readout dimensions  
213 and the first principal component (*Figure 5b*). We chose this space rather than, e.g., the first three  
214 principal components (*Figure 5–Figure Supplement 1*), because it provides an intuitive illustration  
215 of the invariant subspace.

216 Because the sources were two-dimensional, the population activity followed a pattern within  
217 a two-dimensional subspace (*Figure 5b, left; Figure 5–Figure Supplement 1a*). For intact feedback,  
218 this population activity matched the sources when projected onto the readout (*Figure 5c, left*).  
219 Changing the context while freezing the feedback rotated and stretched this representation within  
220 the same subspace, such that the readout did not match the sources (*Figure 5b & c, center*). Would  
221 turning the feedback modulation back on simply reverse this transformation to re-establish an in-  
222 variant subspace? We found that this was not the case. Instead, the feedback rotated the represen-  
223 tation out of the old subspace (*Figure 5b, right*), thereby re-orienting it into the invariant readout  
224 (*Figure 5c, right*).

225 To quantify the transformation of the population representation, we repeated this experiment  
226 multiple times and determined the angle between the neural subspaces. Consistent with the il-  
227 lustration in *Figure 5b*, changing the context did not change the subspace orientation, whereas  
228 unfreezing the feedback caused a consistent re-orientation (*Figure 5d*). The magnitude of this sub-  
229 space re-orientation depended on the similarity of the old and new context. Similar contexts gen-  
230 erally evoked population activity with similar subspace orientations (*Figure 5d,e*). This highlights  
231 that there is a consistent mapping between contexts and the resulting low-dimensional population  
232 activity.

233 In summary, the role of feedback-driven modulation in our model is to re-orient the population  
234 representation in response to changing contexts such that an invariant subspace is preserved.

### 235 **The mechanism generalises to a hierarchical Dalean network**

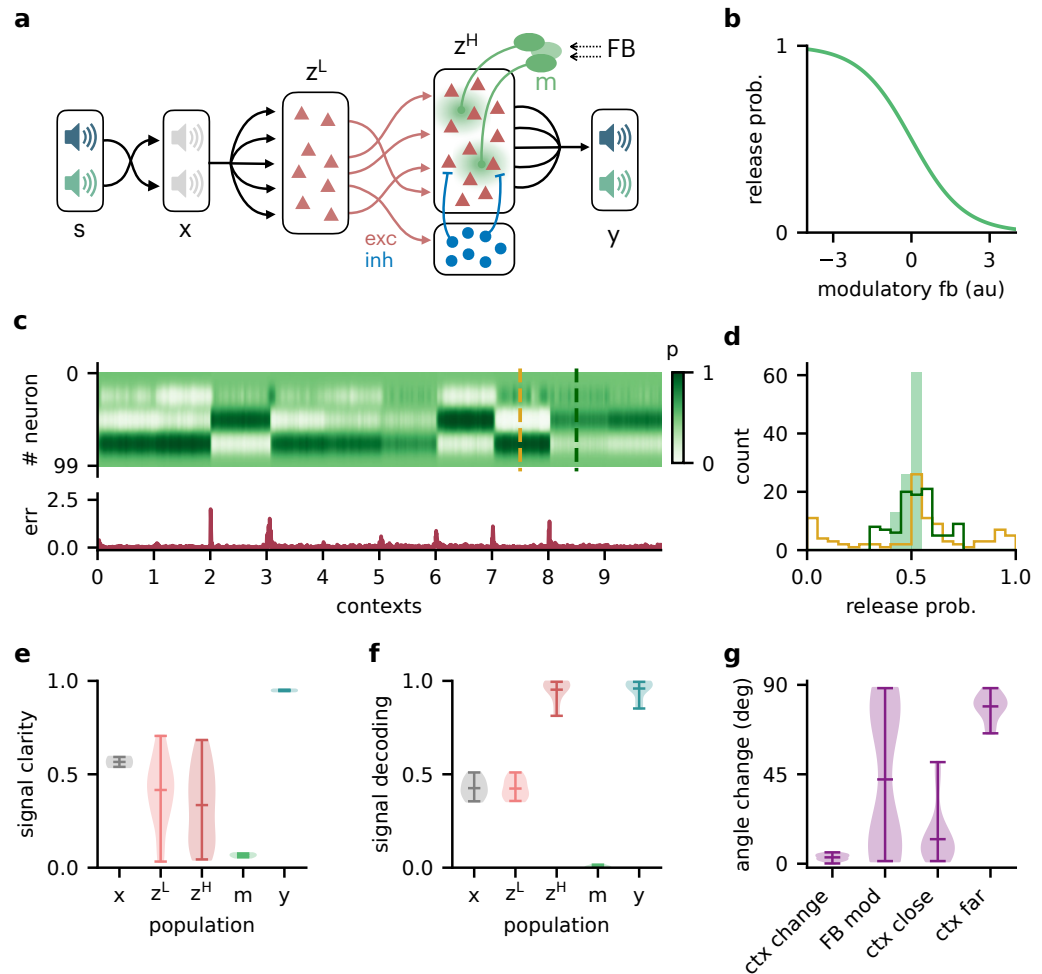
236 So far, we considered a linear network, in which neural activity could be positive and negative.  
237 Moreover, feedback modulation could switch the sign of the neurons' downstream influence, which  
238 is inconsistent with Dale's principle. We wondered if the same population-level mechanisms would  
239 operate in a Dalean network, in which feedback is implemented as a positive gain modulation. Al-  
240 though gain modulation is a broadly observed phenomenon that is attributed to a range of cellular  
241 mechanisms (*Ferguson and Cardin, 2020; Salinas and Thier, 2000*), its effect at the population level  
242 is less clear (*Shine et al., 2021*).

243 We extended the feedforward model as follows (*Figure 6a*): First, all neurons had positive firing  
244 rates. Second, we split the neural population ( $z$  in the previous model) into a "lower-level" ( $z^L$ ) and  
245 "higher-level" population ( $z^H$ ). The lower-level population served as a neural representation of the  
246 sensory stimuli, whereas the higher-level population was modulated by feedback. This allowed a  
247 direct comparison between a modulated and an unmodulated neural population. It also allowed  
248 us to include Dalean weights between the two populations. Direct projections from the lower-level  
249 to the higher-level population were excitatory. In addition, a small population of local inhibitory  
250 neurons provided feedforward inhibition to the higher-level population. Third, the modulation of  
251 the higher-level population was implemented as a local gain modulation that scaled the neural  
252 responses. As a specific realisation of gain modulation, we assumed that feedback targeted in-  
253 hibitory interneurons (e.g., in layer 1; *Abs et al., 2018; Ferguson and Cardin, 2020; Malina et al.,*  
254 *2021*) that mediate the modulation in the higher-level population (e.g., via presynaptic inhibition;  
255 *Pardi et al., 2020; Naumann and Sprekeler, 2020*). This means that stronger feedback decreased  
256 the gain of neurons (*Figure 4b*). We will refer to these modulatory interneurons as modulation  
257 units  $m$  (green units in *Figure 4a*).

258 We found that this biologically more constrained model could still learn the context-invariant  
259 processing task (*Figure 6–Figure Supplement 1a,b*). Notably, the network's performance did not  
260 depend on specifics of the model architecture, such as the target of the modulation or the number

261 of inhibitory neurons (**Figure 6–Figure Supplement 1c-e**). In analogy to the previous model, the  
 262 gain modulation of individual neurons changed with the context and thus enabled the flexible  
 263 processing required to account for varying context (**Figure 4c**). The average gain over contexts was  
 264 similar across neurons, whereas within a context the gains were broadly distributed (**Figure 4d**).

265 To verify if the task is solved by the same population-level mechanism, we repeated our pre-  
 266 vious analyses on the single neuron and population level. Indeed, all results generalised to the  
 267 Dalean network with feedback-driven gain modulation (cf. **Figure 4, Figure 5 & Figure 6**). Single  
 268 neurons in the higher- and lower-level population were not context-invariant (**Figure 6e**), but the  
 269 higher-level population contained a context-invariant subspace (**Figure 6f**). This was not the case



**Figure 6.** Feedback-driven gain modulation in a hierarchical rate network.

**a.** Schematic of the Dalean network comprising a lower- and higher-level population ( $z^L$  and  $z^H$ ), a population of local inhibitory neurons (blue) and diffuse gain modulation mediated by modulatory interneurons (green). **b.** Decrease in gain (i.e. release probability) with stronger modulatory feedback. **c.** Top: Modulation of neurons in the higher-level population for 10 different contexts. Bottom: Corresponding deviation of outputs  $y$  from sources  $s$ . **d.** Histogram of neuron-specific release probabilities averaged across 20 contexts (filled, light green) and during two different contexts (yellow & dark green, see (c)). **e.** Violin plot of signal clarity at different stages of the Dalean model: sensory stimuli ( $x$ ), lower-level ( $z^L$ ) and higher-level population ( $z^H$ ), modulatory units ( $m$ ) and network output ( $y$ ), computed across 20 contexts (cf. **Figure 4a**). **f.** Violin plot of linear decoding performance of the sources from the same stages as in (e) (cf. **Figure 4d**). **g.** Feedback modulation re-oriens the population activity (cf. **Figure 5d**).

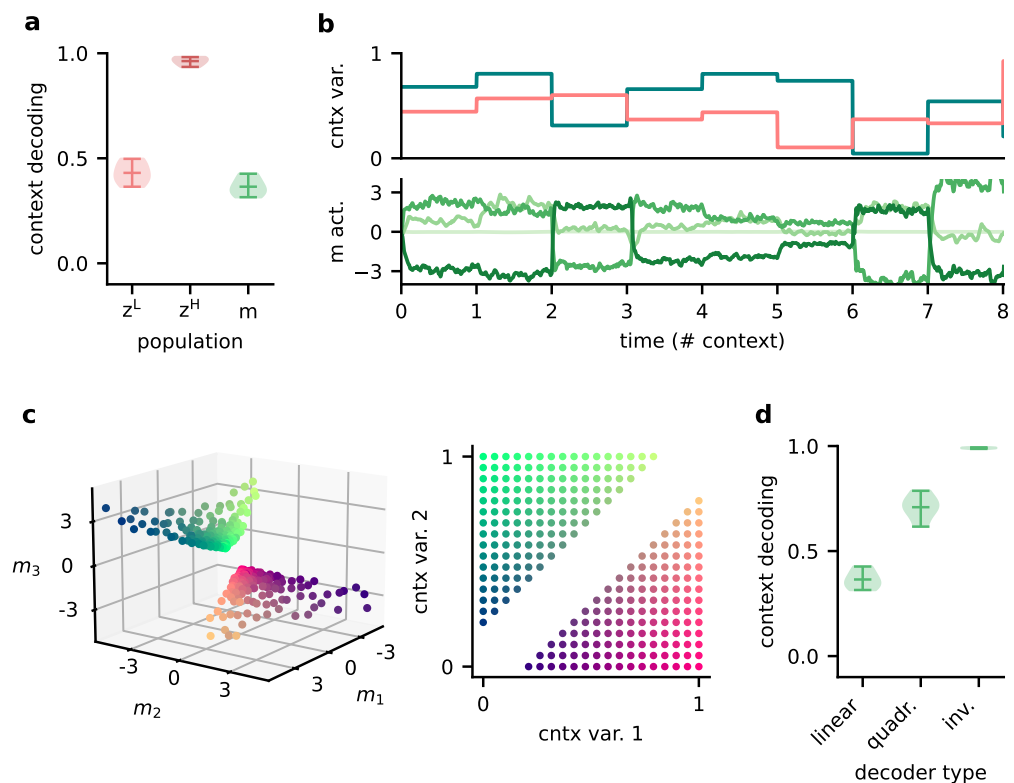
**Figure 6–Figure supplement 1.** The Dalean network can learn the dynamic blind source separation task, and the performance does not depend on specifics of the model architecture.

270 for the lower-level population, underscoring that invariant representations do not just arise from  
 271 projecting the sensory stimuli into a higher dimensional space. Instead, the invariant subspace in  
 272 the higher-level population was again maintained by the feedback modulation, which re-oriented the  
 273 the population activity in response to context changes (*Figure 6g*).

### 274 **Feedback conveys a non-linear representation of the context**

275 Since single neurons in the higher-level population were not invariant to context, the population  
 276 representation must also contain contextual information. Indeed, contextual variables could be  
 277 linearly decoded from the higher-level population activity (*Figure 7a*). In contrast, decoding the  
 278 context from the lower-level population gave much lower accuracy. This shows that the contextual  
 279 information is not just inherited from the sensory stimuli but conveyed by the feedback via the  
 280 modulatory units. We therefore expected that the modulatory units themselves would contain a  
 281 representation of the context. To our surprise, decoding accuracy on the modulatory units was  
 282 low. This seems counter-intuitive, especially since the modulatory units clearly co-varied with the  
 283 contextual variables (*Figure 7b*). To understand these seemingly conflicting results, we examined  
 284 how the context was represented in the activity of the modulation units.

285 We found that the modulation unit activity did encode the contextual variables, albeit in a non-  
 286 linear way (*Figure 7c*). The underlying reason is that the feedback modulation needs to remove  
 287 contextual variations, which requires nonlinear computations. Specifically, the blind source sepa-



**Figure 7.** Feedback conveys a non-linear representation of the context.

**a.** Linear decoding performance of the context (i.e. mixing) from the network. **b.** Context variables (e.g. source locations, top) and activity of modulatory interneurons (bottom) over contexts; one of the modulatory interneurons is silent in all contexts. **c.** Left: Activity of the three active modulatory interneurons (see b) for different contexts. The context variables are colour-coded as indicated on the right. **d.** Performance of different decoders trained to predict the context from the modulatory interneuron activity. Decoder types are a linear decoder, a decoder on a quadratic expansion and a linear decoder trained to predict the inverse of the mixing matrix.

288 ration task requires an inversion of the linear mixture of sources. Consistent with this idea, non-  
289 linear decoding approaches performed better (*Figure 7d*). In fact, the modulatory units contained  
290 a linear representation of the "inverse context" (i.e., the inverse mixing matrix, see *Methods and*  
291 *Models*).

292 In summary, the higher-level population provides a linear representation not only of the stimuli,  
293 but also of the context. In contrast, the modulatory units contained a nonlinear representation of  
294 the context, which could not be extracted by linear decoding approaches. We speculate that if  
295 contextual feedback modulation is mediated by interneurons in layer 1, they should represent the  
296 context in a nonlinear way.

## 297 Discussion

298 Accumulating evidence suggests that sensory processing is strongly modulated by top-down feed-  
299 back projections (*Gilbert and Li, 2013; Keller and Mrsic-Flogel, 2018*). Here, we demonstrate that  
300 feedback-driven gain modulation of a feedforward network could underlie stable perception in  
301 varying contexts. The feedback can be slow, spatially diffuse and low-dimensional. To elucidate  
302 how the context invariance is achieved, we performed single neuron and population analyses. We  
303 found that invariance was not evident at the single neuron level, but only emerged in a subspace of  
304 the population representation. The feedback modulation dynamically transformed the manifold  
305 of neural activity patterns such that this subspace was maintained across contexts. Our results pro-  
306 vide further support that gain modulation at the single cell level enables non-trivial computations  
307 at the population level (*Failor et al., 2021; Shine et al., 2021*).

## 308 Invariance in sensory processing

309 As an example of context-invariant sensory processing, we chose a dynamic variant of the blind  
310 source separation task. This task is commonly illustrated by a mixture of voices at a cocktail party  
311 (*Cherry, 1953; McDermott, 2009*). For auditory signals, bottom-up mechanisms of frequency segre-  
312 gation can provide a first processing step for the separation of multiple sound sources (*Bronkhorst,*  
313 *2015; McDermott, 2009*). However, separating more complex sounds requires additional active top-  
314 down processes (*Parthasarathy et al., 2020; Oberfeld and Kloeckner-Nowotny, 2016*). In our model  
315 top-down feedback guides the source separation itself, while the selection of a source would occur  
316 at a later processing stage – consistent with recent evidence for "late selection" (*Brodbeck et al.,*  
317 *2020; Yahav and Golumbic, 2021*).

318 Although blind source separation is commonly illustrated with auditory signals, the suggested  
319 mechanism of context-invariant perception is not limited to a given sensory modality. The key  
320 nature of the task is that it contains stimulus dimensions that need to be encoded (the sources)  
321 and dimensions that need to be ignored (the context). In visual object recognition, for example,  
322 the identity of visual objects needs to be encoded, while contextual variables such as size, location,  
323 orientation, or surround need to be ignored. Neural hallmarks of invariant object recognition are  
324 present at the population level (*DiCarlo and Cox, 2007; DiCarlo et al., 2012; Hong et al., 2016*), and  
325 to some extent also on the level of single neurons (*Quiroga et al., 2005*). Classically, the emergence  
326 of invariance has been attributed to the extraction of invariant features in feedforward networks  
327 (*Riesenhuber and Poggio, 1999; Wiskott and Sejnowski, 2002; DiCarlo and Cox, 2007; Kriegeskorte,*  
328 *2015*), but recent work also highlights the role of recurrence and feedback (*Gilbert and Li, 2013;*  
329 *Kar et al., 2019; Kietzmann et al., 2019; Thorat et al., 2021*). Here, we focused on the role of  
330 feedback, but clearly, feedforward and feedback processes are not mutually exclusive and likely  
331 work in concert to create invariance. Their relative contribution to invariant perception requires  
332 further studies and may depend on the invariance in question.

333 Similarly, how invariance can be learned will depend on the underlying mechanism. The feedback-  
334 driven mechanism we propose is reminiscent of meta-learning consisting of an inner and an outer  
335 loop (*Hochreiter et al., 2001; Wang et al., 2018a*). In the inner loop, the modulatory system infers  
336 the context to modulate the feedforward network accordingly. This process is unsupervised. In the

337 outer loop, the modulatory system is trained to generalise across contexts. Here, we performed  
338 this training using supervised learning, which requires the modulatory system to experience the  
339 sources in isolation (or at least obtain an error signal). Such an identification of the individual  
340 sources could, e.g., be aided by other sensory modalities (*McDermott, 2009*). However, the op-  
341 timisation of the modulatory system does not necessarily require supervised learning. It could  
342 also be guided by task demands via reinforcement learning, or by unsupervised priors such as a  
343 non-Gaussianity of the outputs.

#### 344 **Mechanisms of feedback-driven gain modulation**

345 There are different ways in which feedback can affect local processing. Here, we focused on gain  
346 modulation (*McAdams and Maunsell, 1999; Reynolds and Heeger, 2009; Vinck et al., 2015*). Neu-  
347 ronal gains can be modulated by a range of mechanisms (*Ferguson and Cardin, 2020; Shine et al.,*  
348 *2021*). In our model, the mechanism needs to satisfy a few key requirements: i) the modulation  
349 is not uniform across the population, ii) it operates on a timescale similar to that of changes in  
350 context, and iii) it is driven by a brain region that has access to the information needed to infer the  
351 context.

352 Classical neuromodulators such as acetylcholine (*Disney et al., 2007; Kawai et al., 2007*), dopamine  
353 (*Thurley et al., 2008*) or serotonin (*Azimi et al., 2020*) are signalled through specialised neuromod-  
354 ulatory pathways from subcortical nuclei (*van den Brink et al., 2019*). These neuromodulators can  
355 control the neural gain depending on behavioural states such as arousal, attention or expectation  
356 of rewards (*Ferguson and Cardin, 2020; Hasselmo and McGaughy, 2004; Bayer and Glimcher, 2005;*  
357 *Polack et al., 2013; Kuchibhotla et al., 2017*). Their effect is typically thought to be brain-wide and  
358 long-lasting, but recent advances in measurement techniques (*Sabatini and Tian, 2020; Lohani*  
359 *et al., 2020*) indicate that it could be area- or even layer-specific, and vary on sub-second time  
360 scales (*Lohani et al., 2020; Bang et al., 2020; Poorthuis et al., 2013; Pinto et al., 2013*).

361 More specific feedback projections arrive in layer 1 of the cortex, where they target the distal  
362 dendrites of pyramidal cells and inhibitory interneurons (*Douglas and Martin, 2004; Roth et al.,*  
363 *2016; Marques et al., 2018*). Dendritic input can change the gain of the neural transfer function on  
364 fast timescales (*Larkum et al., 2004; Jarvis et al., 2018*). The spatial scale of the modulation will  
365 depend on the spatial spread of the feedback projections and the dendritic arbourisation. Feed-  
366 back to layer 1 interneurons provides an alternative mechanism of local gain control. In particular,  
367 neuron-derived neurotrophic factor-expressing interneurons (NDNF) in layer 1 receive a variety  
368 of top-down feedback projections and produce GABAergic volume transmission (*Abs et al., 2018*),  
369 thereby down-regulating synaptic transmission (*Miller, 1998; Laviv et al., 2010*). This gain modu-  
370 lation can act on a timescale of hundreds of milliseconds (*Branco and Staras, 2009; Urban-Ciecko*  
371 *et al., 2015; Malina et al., 2021; Molyneaux and Hasselmo, 2002*), and, although generally consid-  
372 ered diffuse, can also be synapse type-specific (*Chittajallu et al., 2013*).

373 The question remains where in the brain the feedback signals originate. Our model requires the  
374 responsible network to receive feedforward sensory input to infer the context. In addition, feed-  
375 back inputs from higher-level sensory areas to the modulatory system allow a better control of the  
376 modulated network state. Higher-order thalamic nuclei are ideally situated to integrate different  
377 sources of sensory inputs and top-down feedback (*Sampathkumar et al., 2021*) and mediate the  
378 resulting modulation by targeting layer 1 of lower-level sensory areas (*Purushothaman et al., 2012;*  
379 *Roth et al., 2016; Sherman, 2016*). In our task setting, the inference of the context requires the in-  
380 tegration of sensory signals over time and therefore recurrent neural processing. For this kind of  
381 task, thalamus may not be the site of contextual inference, because it lacks the required recur-  
382 rent connectivity (*Halassa and Sherman, 2019*). However, contextual inference may be performed  
383 by higher-order cortical areas, and could either be relayed back via the thalamus or transmitted  
384 directly, for example, via cortico-cortical feedback connections.

### 385 **Testable predictions**

386 Our model makes several predictions that could be tested in animals performing invariant sensory  
387 perception. Firstly, our model indicates that invariance across contexts may only be evident at the  
388 neural population level, but not on the single cell level. Probing context invariance at different  
389 hierarchical stages of sensory processing may therefore require population recordings and corre-  
390 sponding statistical analyses such as neural decoding (*Glaser et al., 2020*). Secondly, we assumed  
391 that this context invariance is mediated by feedback modulation. The extent to which context in-  
392 variance is enabled by feedback on a particular level of the sensory hierarchy could be studied  
393 by manipulating feedback connections. Since layer 1 receives a broad range of feedback inputs  
394 from different sources, this may require targeted manipulations. If no effect of feedback on con-  
395 text invariance is found, this may either indicate that feedforward mechanisms dominate or that  
396 the invariance in question is inherited from an earlier stage, in which it may well be the result of  
397 feedback modulation. Given that feedback is more pronounced in higher cortical areas (*McAdams*  
398 *and Maunsell, 1999; Pardi et al., 2020*), we expect that the contribution of feedback may play a  
399 larger role for the more complex forms of invariance further up in the sensory processing hierar-  
400 chy. Thirdly, for feedback to mediate context invariance, the feedback projections need to contain  
401 a representation of the contextual variables. Our findings suggest, however, that the detection  
402 of this representation may require a non-linear decoding method. Finally, a distinguishing fea-  
403 ture of feedback and feedforward mechanisms is that feedback mechanisms take more time. We  
404 found that immediately following a sudden contextual change, the neuronal representation initially  
405 changes within the manifold associated with the previous context. Later, the feedback reorients  
406 the manifold to reestablish the invariance on the population level. Whether these dynamics are  
407 a signature of feedback processing or also present in feedforward networks will be an interesting  
408 question for future work.

### 409 **Comparison to prior work**

410 Computational models have implicated neuronal gain modulation for a variety of functions (*Salin-*  
411 *inas and Sejnowski, 2001; Reynolds and Heeger, 2009*). Even homogeneous changes in neuronal  
412 gain can achieve interesting population effects (*Shine et al., 2021*), such as orthogonalisation of  
413 sensory responses (*Failor et al., 2021*). More heterogeneous gain modulation provides additional  
414 degrees of freedom that enables, for example, attentional modulation (*Reynolds and Heeger, 2009;*  
415 *Carandini and Heeger, 2012*), coordinate transformations (*Salinas and Thier, 2000*) and – when am-  
416 plified by recurrent dynamics – a rich repertoire of neural trajectories (*Stroud et al., 2018*). Gain  
417 modulation has also been suggested as a means to establish invariant processing (*Salinas and Ab-*  
418 *bott, 1997*), as a biological implementation of dynamic routing (*Olshausen et al., 1993*). While the  
419 modulation in these models of invariance can be interpreted as an abstract form of feedback, the  
420 resulting effects on the population level were not studied.

421 An interesting question is by which mechanisms the appropriate gain modulation is computed.  
422 In previous work, gain factors were often learned individually for each context, for example by gra-  
423 dient descent or Hebbian plasticity (*Olshausen et al., 1993; Salinas and Abbott, 1997; Stroud et al.,*  
424 *2018*), mechanisms that may be too slow to achieve invariance on a perceptual timescale (*Wiskott,*  
425 *2006*). In our model, by contrast, the modulation is dynamically controlled by a recurrent network.  
426 Once it has been trained, such a recurrent modulatory system can rapidly infer the current con-  
427 text, and provide an appropriate feedback signal on a timescale only limited by the modulatory  
428 mechanism.

### 429 **Limitations and future work**

430 In our model, we simplified many aspects of sensory processing. Using simplistic sensory stimuli  
431 – compositions of sines – allowed us to focus on the mechanisms at the population level, while  
432 avoiding the complexities of natural sensory stimuli and deep sensory hierarchies. Although we  
433 do not expect conceptual problems in generalising our results to more complex stimuli, such as

434 speech or visual stimuli, the associated computational challenges are substantial. For example,  
435 the feedback in our model was provided by a recurrent network, whose parameters were trained  
436 by back-propagating errors through the network and through time. This training process can get  
437 very challenging for large networks and long temporal dependencies (*Bengio et al., 1994; Pascanu*  
438 *et al., 2013*).

439 In our simulations we trained the whole model – the modulatory system, the sensory represen-  
440 tation and the readout. For the simplistic stimuli we used, we observed that the training process  
441 mostly concentrated on optimising the modulatory system and readout, while a random mapping  
442 of sensory stimuli to neural representations seemed largely sufficient to solve the task. For more  
443 demanding stimuli, we expect that the sensory representation the modulatory system acts upon  
444 may become more important. A well-suited representation could minimise the need for modula-  
445 tory interventions (*Finn et al., 2017*), in a coordinated interaction of feedforward and feedback.

446 To understand the effects of feedback modulation on population representations, we included  
447 biological constraints in the feedforward network and the structure of the modulatory feedback.  
448 However, we did not strive to provide a biologically plausible implementation for the computation  
449 of the appropriate feedback signals, and instead used an off-the-shelf recurrent neural network  
450 (*Hochreiter and Schmidhuber, 1997*). The question how these signals could be computed in a  
451 biologically plausible way remains for future studies. The same applies to the question how the  
452 appropriate feedback signals can be learned by local learning rules (*Lillicrap et al., 2020*) and how  
453 neural representations and modulatory systems learn to act in concert.

## 454 **Methods and Models**

455 To study how feedback-driven modulation can enable flexible sensory processing, we built models  
456 of feedforward networks that are modulated by feedback. The feedback was dynamically gener-  
457 ated by a modulatory system, which we implemented as a recurrent network. The weights of the  
458 recurrent network were trained such that the feedback modulation allowed the feedforward net-  
459 work to solve a flexible invariant processing task.

### 460 **The dynamic blind source separation task**

461 As an instance of flexible sensory processing we used a dynamic variant of blind source separation.  
462 In classical blind source separation, two or more unknown time-varying sources  $\vec{s}(t)$  need to be  
463 recovered from a set of observations (i.e. sensory stimuli)  $\vec{x}(t)$ . The sensory stimuli are composed  
464 of an unknown linear mixture of the sources such that  $\vec{x}(t) = A\vec{s}(t)$  with a fixed mixing matrix  $A$ .  
465 Recovering the sources requires to find weights  $W$  such that  $W\vec{x}(t) \approx \vec{s}(t)$ . Ideally,  $W$  is equal to the  
466 pseudo-inverse of the unknown mixing matrix  $A$ , up to permutations.

In our dynamic blind source separation task, we model variations in the stimulus context by  
changing the linear mixture over time – albeit on a slower timescale than the time-varying signals.  
Thus, the sensory stimuli are constructed as

$$\vec{x}(t) = A(t)\vec{s}(t) + \sigma_n \vec{\xi}(t) \quad , \quad (1)$$

467 where  $A(t)$  is a time-dependent mixing matrix and  $\sigma_n$  is the amplitude of additive white noise  $\vec{\xi}(t)$ .  
468 The time-dependent mixing matrix determines the current context and was varied in discrete time  
469 intervals  $n_t$ , meaning that the mixing matrix  $A(t)$  (i.e. the context) was constant for  $n_t$  samples before  
470 it changed. The goal of the dynamic blind source separation task is to recover the original signal  
471 sources  $\vec{s}$  from the sensory stimuli  $\vec{x}$  across varying contexts. Thus, the network model output  
472 needs to be invariant to the specific context of the sources. Note that while the context was varied,  
473 the sources themselves were the same throughout the task, unless stated otherwise. Furthermore,  
474 in the majority of experiments the number of source signals and sensory stimuli was  $n_s = 2$ . A list  
475 of default parameters for the dynamic blind source separation task can be found in **Table 1**.



476 **Source signals**

As default source signals we used two compositions of two sines each ("chords") with a sampling rate of  $f_s = 8000\text{Hz}$  that can be written as

$$s_1(t) = \sin(2\pi f_{11}t/f_s) + \sin(2\pi f_{12}t/f_s) \quad (2)$$

$$s_2(t) = \sin(2\pi f_{21}t/f_s) + \sin(2\pi f_{22}t/f_s) \quad (3)$$

477 with frequencies  $f_{11} = 100\text{ Hz}$ ,  $f_{12} = 125\text{ Hz}$ ,  $f_{21} = 150\text{ Hz}$  and  $f_{22} = 210\text{ Hz}$ . Note that in our model  
 478 we measure time as the number of samples from the source signals, meaning that timescales are  
 479 relative and could be arbitrarily rescaled.

480 In **Figure 5**, we used pure sine signals with frequency  $f$  for visualisation purposes:  $s_i = \sin(2\pi ft/f_s)$ .  
 481 We also validated the model on signals that are not made of sine waves, as a sawtooth and a square  
 482 wave signal (**Figure 1–Figure Supplement 4**). Unless stated otherwise, the same signals were used  
 483 for training and testing the model.

484 **Time-varying contexts**

485 We generated the mixing matrix  $A$  for each context by drawing random weights from a uniform  
 486 distribution between 0 and 1, allowing only positive mixtures of the sources. Unless specified  
 487 otherwise, we sampled new contexts for each training batch and for the test data, such that the  
 488 training and test data followed the same distribution without necessarily being the same. The  
 489 dimension of the mixing matrices was determined by number of signals  $n_s$  such that  $A$  was of shape  
 490  $n_s \times n_s$ . To keep the overall amplitude of the sensory stimuli in a similar range across different  
 491 mixtures, we normalised the row sums of each mixing matrix to one. In the case of  $n_s = 2$ , this  
 492 implies that the contexts (i.e. the mixing matrices) are drawn from a 2-dimensional manifold (see  
 493 **Figure 8**, bottom left). In addition, we only used the randomly generated mixing matrices whose  
 494 determinant was larger than some threshold value. We did this to ensure that each signal mixture  
 495 was invertible and that the weights needed to invert the mixing matrix were not too extreme. A  
 496 threshold value of 0.2 was chosen based on visual inspection of the weights from the inverted  
 497 mixing matrix.

498 **Modulated feedforward network models**

499 Throughout this work, we modelled feedforward networks of increasing complexity. Common to all  
 500 networks was that they received the sensory stimuli  $\vec{x}$  and should provide an output  $\vec{y}$  that matches  
 501 the source signals  $\vec{s}$ . In the following, we first introduce the simplest model variant and how it is  
 502 affected by feedback from the modulatory system, and subsequently describe the different model  
 503 extensions.

504 **Modulation of feedforward weights by a recurrent network**

In the simplest feedforward network the network output  $\vec{y}(t)$  is simply a linear readout the sensory  
 stimuli  $\vec{x}(t)$ , with readout weights that are dynamically changed by the modulatory system:

$$\vec{y}(t) = (M(t) \odot W_0) \vec{x}(t) \quad (4)$$

**Table 1.** Default parameters of the dynamic blind source separation task.

parameter	symbol	value
number of signals	$n_s$	2
number of samples in context	$n_t$	1000
additive noise	$\sigma_n$	0.001
sampling frequency	$f_s$	8 kHz

505 where  $W_0$  are the baseline weights and  $M(t)$  the modulation provided by the modulatory system.  
 506  $M(t)$  is of the same shape as  $W_0$  and determines the element-wise multiplicative modulation of the  
 507 baseline weights. Because the task requires the modulatory system to dynamically infer the con-  
 508 text, we modelled it as a recurrent network – more specifically a long-short term memory network  
 509 (LSTMs; *Hochreiter and Schmidhuber, 1997*) – with  $N_h = 100$  hidden units. In particular, we used  
 510 LSTMs with forget gates (*Gers et al., 2000*) but no peephole connections (for an overview of LSTM  
 511 variants see *Greff et al. (2016)*).

In this work we treated the LSTM as a black-box modulatory system that receives the sensory stimuli and the feedforward network’s output and provides the feedback signal in return (*Figure 1a*). A linear readout of the LSTM’s output determines the modulation  $M(t)$  in *Equation 4*. In brief, this means that

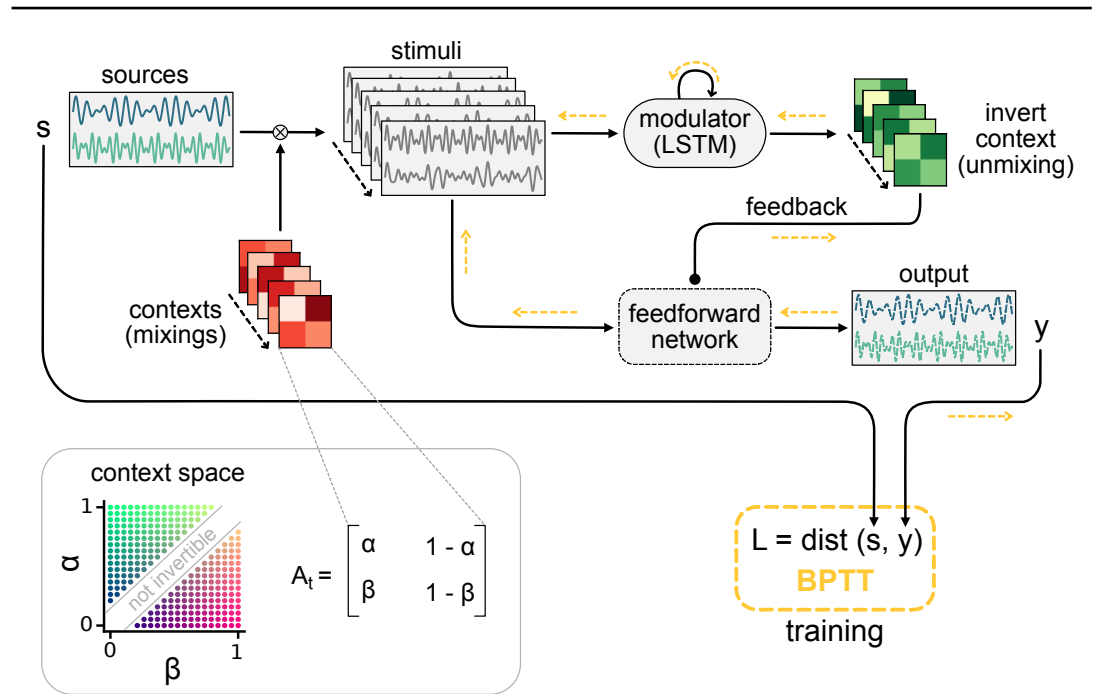
$$M(t) = \text{LSTM}(\vec{x}(t), \vec{y}(t)), \quad (5)$$

512 where  $\text{LSTM}(\cdot)$  is a function that returns the LSTM readout. For two-dimensional sources and sen-  
 513 sory stimuli, for instance,  $\text{LSTM}(\cdot)$  receives a concatenation of the two-dimensional vectors  $\vec{x}(t)$   
 514 and  $\vec{y}(t)$  as input and returns a two-by-two feedback modulation matrix – one multiplicative factor  
 515 for each weight in  $W_0$ . The baseline weights  $W_0$  were randomly drawn from the Gaussian distribu-  
 516 tion  $\mathcal{N}(1, 0.001)$  and fixed throughout the task. The LSTM parameters and readout were learned  
 517 during training of the model.

#### 518 Extension 1: Reducing the temporal specificity of feedback modulation

To probe our model’s sensitivity to the timescale of the modulatory feedback (*Figure 2*), we added a temporal filter to *Equation 5*. In that case the modulation  $M(t)$  followed the dynamics

$$\tau \frac{dM(t)}{dt} = -M(t) + \text{LSTM}(\vec{x}(t), \vec{y}(t)), \quad (6)$$



**Figure 8.** Schematic of the dynamic blind source separation task, the context space and the modulated feedforward network. Information flow is indicated by black arrows and the flow of the error during training with backpropagation through time (BPTT) is shown in yellow.

519 with  $\tau$  being the time constant of modulation. For small  $\tau$ , the feedback rapidly affects the feed-  
 520 forward network, whereas larger  $\tau$  imply a slowly changing modulatory feedback signal. The unit  
 521 of this timescale is the number of samples from the source signals. Note that the timescale of the  
 522 modulation should be considered relative to the timescale of the context changes  $n_i$ . As a default  
 523 time constant we used  $\tau = 100 < n_i$  (see **Table 2**).

#### 524 Extension 2: Reducing the spatial specificity of feedback modulation

To allow for spatially diffuse feedback modulation (**Figure 3**), we added an intermediate layer be-  
 tween the sensory stimuli and the network output. This intermediate layer consisted of a pop-  
 ulation of  $N_z = 100$  units that were modulated by the feedback, where neighbouring units were  
 modulated similarly. More specifically, the units were arranged on a ring to allow for a spatially con-  
 strained modulation without boundary effects. The population's activity vector  $\vec{z}(t)$  is described by

$$\vec{z}(t) = \vec{m}(t) \odot (W^x \vec{x}(t)), \quad (7)$$

with the sensory stimuli  $\vec{x}(t)$ , a weight matrix  $W^x$  of size  $N_z \times n_s$  and the vector of unit-specific  
 multiplicative modulations  $\vec{m}(t)$ . Note that the activity of the units was not constrained to be positive  
 here. The output of the network was then determined by a linear readout of the population activity  
 vector according to

$$\vec{y}(t) = W^{ro} \vec{z}(t) \quad (8)$$

525 with a fixed readout matrix  $W^{ro}$ .

The modulation to a single unit  $i$  was given by

$$\tau \frac{dm_i(t)}{dt} = -m_i(t) + \sum_{j=1}^{N_{FB}} K_{ij} l_j, \quad (9a)$$

$$\text{with } l_j = \text{LSTM}(x(t), y(t))_j. \quad (9b)$$

526 Here,  $\tau$  is the modulation time constant,  $K$  a kernel that determines the spatial specificity of mod-  
 527 ulation,  $\text{LSTM}(\cdot)_j$  the  $j$ -th feedback signal from the LSTM and  $N_{FB}$  the total number of feedback  
 528 signals. As in the simple model, the  $N_{FB}$  feedback signals were determined by a linear readout  
 529 from LSTM.

The modulation kernel  $K$  was defined as a set of von Mises functions:

$$K_{ij} = \exp\left(\frac{1}{\sigma_m^2} \cos\left(z_i^{\text{loc}} - l_j^{\text{loc}}\right)\right), \quad (10)$$

530 where  $z_i^{\text{loc}} = \frac{2\pi i}{N_z} \in [0, 2\pi[$  represents the location of the modulated unit  $i$  on the ring and  $l_j^{\text{loc}}$  the  
 531 "preferred location" of modulatory unit  $j$ , i.e., the location on the ring that it modulates most ef-  
 532 fectively. These "preferred locations"  $l_j^{\text{loc}}$  of the feedback units were evenly distributed on the ring.  
 533 The variance parameter  $\sigma_m^2$  determines the spatial spread of the modulatory effect of the feedback  
 534 units, i.e., the spatial specificity of the modulation. Overall, the spatial distribution of the modu-  
 535 lation was therefore determined by the number of distinct feedback signals  $N_{FB}$  and their spatial  
 536 spread  $\sigma_m^2$  (see **Table 2** for a list of network parameters).

#### 537 Extension 3: Hierarchical rate-based network

We further extended the model with spatial modulation (**Equation 7–Equation 10**) to include a  
 two-stage hierarchy, positive rates and synaptic weights that obey Dale's law. Furthermore, we  
 implemented the feedback modulation as a gain modulation that scales neural rates but keeps  
 them positive. To this end, we modelled the feedforward network as a hierarchy of a lower-level  
 and a higher-level population. Only the higher-level population received feedback modulation.  
 Splitting the neural populations in this way allowed us to model the connections between them with  
 weights that follow Dale's law. Furthermore, the unmodulated lower-level population could serve

as a control for the emergence of context-invariant representations. The lower-level population consisted of  $N_L = 40$  rate-based neurons and the population activity vector was given by

$$\vec{z}^L(t) = [W^{Lx}\vec{x}(t)]_+ \quad , \quad (11)$$

where  $W^{Lx}$  is a fixed weight matrix,  $\vec{x}(t)$  the sensory stimuli and the rectification  $[\cdot]_+ = \max(0, \cdot)$  ensures that rates are positive. The lower-level population thus provides a neural representation of the sensory stimuli. The higher-level population consisted of  $N_H = 100$  rate-based neurons that received feedforward input from the lower-level population. The feedforward input consisted of direct excitatory projections as well as feedforward inhibition through a population of  $N_I = 20$  local inhibitory neurons. The activity vector of the higher-level population  $\vec{z}^H(t)$  was thus given by

$$\vec{z}^H(t) = [\vec{p}(t) \odot (W^{HL}\vec{z}^L(t) - W^{HI}\vec{z}^I(t))]_+ \quad (12)$$

$$\vec{z}^I(t) = [W^{IL}\vec{z}^L(t)]_+ \quad . \quad (13)$$

538 Here  $W^{HL}$ ,  $W^{HI}$  and  $W^{IL}$  are positive weight matrices,  $\vec{z}^I(t)$  the inhibitory neuron activities and  $\vec{p}(t)$   
 539 the neuron-specific gain modulation factors. As for the spatially modulated network of Extension  
 540 2, the network output  $\vec{y}(t)$  was determined by a fixed linear readout  $W^{ro}$  (see [Equation 8](#)). The  
 541 distributions used to randomly initialise the weight matrices are provided in [Table 3](#).

Again, the modulation was driven by feedback from the LSTM, but in this model variant we assumed inhibitory feedback, i.e., stronger feedback signals monotonically decreased the gain. More specifically, we assumed that the feedback signal targets a population of modulation units  $\vec{m}$ , which in turn modulate the gain in the higher-level population. The gain modulation of neuron  $i$  was constrained between 0 and 1 and determined by

$$p_i(t) = \frac{1}{1 + \exp(m_i(t))} \quad (14)$$

542 with  $m_i(t)$  being the activity of a modulation unit  $i$ , which follows the same dynamics as in [Equa-](#)  
 543 [tion 9a](#) (see [Figure 6a](#)).

#### 544 Training the model

We used gradient descent to find the model parameters that minimise the difference between the source signal  $\vec{x}(t)$  and the feedforward network's output  $\vec{y}(t)$ :

$$\mathcal{L} = \sum_{t=1}^{n_t} \text{dist}(\vec{x}(t), \vec{y}(t)) \quad (15)$$

545 with a distance measure  $\text{dist}(\cdot)$ . We used the machine learning framework PyTorch ([Paszke et al.,](#)  
 546 [2019](#)) to simulate the network model, obtain the gradients of the objective  $\mathcal{L}$  by automatic differen-  
 547 tiation and update the parameters of the LSTM using the Adam optimiser ([Kingma and Ba, 2014](#))  
 548 with a learning rate of  $\eta = 10^{-3}$ . As distance measure in the objective we used a smooth variant

**Table 2.** Default parameters of the network models.

parameter	symbol	value
number of hidden units in LSTM	$N_h$	100
number of units in middle layer $z$	$N_z$	100
number of distinct feedback signals	$N_{FB}$	4
number of neurons in lower-level population	$N_L$	40
number of neurons in higher-level population	$N_H$	100
number of inhibitory neurons	$N_I$	20
timescale of modulation	$\tau$	100
spatial spread of modulation	$\sigma_m^2$	0.2

**Table 3.** Distributions used for randomly initialised weight parameters

weights	distribution
$W_0$	$\mathcal{N}(1, 0.001)$
$W^x$	$\mathcal{N}(0, 0.5)$
$W^{Lx}$	$\mathcal{N}(0, 0.5)$
$W^{ro}$	$\mathcal{N}(0, 0.5)$
$W^{HL}$	$\mathcal{N}(1, 0.5) \cdot 20/N_H$
$W^{IL}$	$\mathcal{N}(1, 0.5)/N_I$
$W^{HI}$	$\mathcal{N}(1, 1) \cdot 20/N_H$
LSTM parameters	$\mathcal{U}(-\sqrt{1/N_H}, \sqrt{1/N_H})$
LSTM readout	$\mathcal{U}(-\sqrt{1/N_{FB}}, \sqrt{1/N_{FB}})$

549 of the L1 norm (PyTorch’s smooth L1 loss variant), because it is less sensitive to outliers than the  
550 mean squared error (*Huber, 1964*).

551 During training, we simulated the network dynamics over batches of 32 trials using forward Euler  
552 with a timestep of  $\Delta t = 1$ . Each trial consisted of  $n_t$  time steps (i.e. samples) and the context (i.e.  
553 mixing matrix) differed between trials. Since the model contains feedback and recurrent connections,  
554 we trained it using backpropagation through time (*Werbos, 1990*). This means that for each  
555 trial, we simulated the model and computed the loss for every time step. At the end of the trial  
556 we propagated the error through the  $n_t$  steps of the model to obtain the gradients and updated  
557 the parameters accordingly (*Figure 8*). Although the source signals were the same in every trial,  
558 we varied their phase independently across trials to prevent the LSTM from learning the exact signal  
559 sequence. To this end, we generated 16,000 samples of the source signals and in every batch  
560 randomly selected chunks of  $n_t$  samples independently from each source. Model parameters were  
561 initialised according to the distributions listed in *Table 3*.

562 In all model variants we optimised the parameters of the modulator (input, recurrent and read-  
563 out weights as well as the biases of the LSTM; see *Equation 5 & Equation 9b*). The parameters  
564 were initialised with the defaults from the corresponding PyTorch modules, as listed in *Table 3*.  
565 To facilitate the training in the hierarchical rate-based network despite additional constraints, we  
566 also optimised the feedforward weights  $W^{HL}$ ,  $W^{HI}$ ,  $W^{IL}$ ,  $W^{Lx}$  and  $W^{ro}$ . In principle, this allows to  
567 adapt the representation in the two intermediate layers such that the modulation is most effective.  
568 However, although we did not quantify it, we observed that optimising the network readout  $W^{ro}$   
569 facilitated the training the most, suggesting that a specific format of the sensory representations  
570 was not required for an effective modulation.

To prevent the gain modulation factor from saturating at 0 or 1, we added a regularisation term  $\mathcal{R}$  to the loss function *Equation 15* that keeps the LSTM’s output small:

$$\mathcal{R} = \lambda_{\text{out}} \sum_{t=1}^{n_t} \sum_{j=1}^{N_{FB}} \left| \text{LSTM}(x(t), y(t))_j \right| \quad (16)$$

571 with  $\lambda_{\text{out}} = 10^{-5}$ .

572 Gradient values were clipped between -1 and 1 before each update to avoid large updates. For  
573 weights that were constrained to be positive, we used their absolute value in the model. Each  
574 network was trained for 10,000 to 12,000 batches and for 5 random initialisations (*Figure 1-Figure  
575 Supplement 2*).

576 Testing and manipulating the model

577 We tested the network model performance on an independent random set of contexts (i.e. mixing  
578 matrices), but with the same source signals as during training. During testing, we also changed

579 the context every  $n_i$  steps, but the length of this interval was not crucial for performance (*Figure 1–*  
580 *Figure Supplement 1d*).

581 To manipulate the feedback modulation in the hierarchical rate-based network (*Figure 4*), we  
582 provided an additional input to the modulation units  $m$  in *Equation 9a*. We used an input of 3 or  $-3$   
583 depending on whether the modulation units were activated or inactivated, respectively. To freeze  
584 the feedback modulation (*Figure 6*), we discarded the feedback signal and held the local modula-  
585 tion  $p$  in *Equation 14* at a constant value determined by the feedback before the manipulation. The  
586 dynamics of the LSTM were continued, but remained hidden to the feedforward network until the  
587 freezing was stopped.

## 588 **Unmodulated feedforward network models**

### 589 **Linear regression.**

As a control, we trained feedforward networks with weights that were not changed by a modulatory system. First, we used the simplest possible network architecture, in which the sensory stimuli are linearly mapped to the outputs (*Figure 1–Figure Supplement 1a*):

$$y(t) = Wx(t). \quad (17)$$

590 It is intuitive that a fixed set of weights  $W$  cannot invert two different contexts (i.e. different mixing  
591 matrices  $A_1$  and  $A_2$ ). As an illustration we trained this simple feedforward network on one context  
592 and tested it on different contexts. To find the weights  $W$ , we used linear regression to minimise  
593 the mean squared error between the source signal  $s(t)$  and the network’s output  $y(t)$ . The training  
594 data consisted of 1024 consecutive time steps of the sensory stimuli for a fixed context, and the  
595 test data consisted of different 1024 time steps generated under a potentially different mixing.  
596 We repeated this procedure by training and testing a network for all combinations of 20 random  
597 contexts.

### 598 **Multi-layer nonlinear network.**

599 Since solving the task was not possible with a single set of readout weights, we extended the feed-  
600 forward model to include 3 hidden layers consisting of 32, 16 and 8 rectified linear units (*Figure 1–*  
601 *Figure Supplement 1d*). The input to this network was one time point from the sensory stimuli and  
602 the target output the corresponding time point of the sources. We trained the multi-layer network  
603 on 5000 batches of 32 contexts using Adam (learning rate 0.001) to minimise the mean squared  
604 error between the network output and the sources.

### 605 **Multi-layer network with sequences as input.**

606 Solving the task requires the network to map the same sensory stimulus to different outputs de-  
607 pending on the context. However, inferring the context takes more than one time point. To test  
608 if a feedforward network with access to multiple time points at once could in principle solve the  
609 task, we changed the architecture of the multi-layer network, such that it receives a sequence of  
610 the sensory stimuli (*Figure 1–Figure Supplement 1g*). The output of the network was a sequence  
611 of equal length. We again trained this network on 5000 batches of 32 contexts to minimise the  
612 error between its output and the target sources, where both the network input and output were  
613 sequences. The length of these sequences was varied between 1 and 150.

## 614 **Data analysis**

### 615 **Signal clarity**

To determine task performance, we measured how clear the representation of the source signals is in the network output. We first computed the correlation coefficient of each signal  $s_i$  with each output  $y_j$

$$r_{ij} = \frac{\sum_t (s_i(t) - \bar{s}_i)(y_j(t) - \bar{y}_j)}{\sigma_{s_i} \sigma_{y_j}}, \quad (18)$$

where  $\bar{x}_i$  and  $\bar{y}_j$  are the respective temporal mean and  $\sigma_{s,i}$  and  $\sigma_{y,j}$  the respective temporal standard deviations. The signal clarity in output  $y_j$  is then given by the absolute difference between the absolute correlation with one compared to the other signal:

$$c_j = | |r_{1j}| - |r_{2j}| | \quad . \quad (19)$$

616 By averaging over outputs we determined the overall signal clarity within the output. Note that  
 617 the same measure can be computed on other processing stages of the feedforward network. For  
 618 instance, we used the signal clarity of sources in the sensory stimuli as a baseline control.

#### 619 Signal-to-noise ratio

The signal-to-noise ratio in the sensory stimuli was determined as the variability in the signal compared to the noise. Since the mean of both the stimuli and the noise were zero, the signal-to-noise ratio could be computed by

$$\text{SNR} = \frac{\sigma_s^2}{\sigma_n^2} \quad ,$$

620 where  $\sigma_n$  was the standard deviation of the additive white noise and  $\sigma_s$  the measured standard  
 621 deviation in the noise-free sensory stimuli, which was around 0.32. As a scale of the signal-to-noise  
 622 ratio we used decibels (dB), i.e., we used  $\text{dB} = 10 \log_{10}(\text{SNR})$ .

#### 623 Linear decoding analysis

##### 624 Signal decoding.

625 We investigated the population-level invariance by using a linear decoding approach. If there was  
 626 an invariant population subspace, the source signals could be decoded by the same decoder across  
 627 different contexts. We therefore performed linear regression between the activity in a particular  
 628 population and the source signals. This linear decoder was trained on  $n_c = 10$  different contexts  
 629 with  $n_t = 1,000$  time points each, such that the total number of samples was 10,000. The linear  
 630 decoding was then tested on 10 new contexts and the performance determined using the  $R^2$  mea-  
 631 sure.

##### 632 Context decoding.

633 We took a similar approach to determine from which populations the context could be decoded.  
 634 For the dynamic blind source separation task the context is given by the source mixture, as de-  
 635 termined by the mixing matrix. Since we normalised the rows of each mixing matrix, the context  
 636 was determined by two context variables. We calculated the temporal average of the neuronal ac-  
 637 tivities within each context and performed a linear regression of the context variables onto these  
 638 averages. To exclude onset transients, we only considered the second half (500 samples) of every  
 639 context. Contexts were sampled from the two-dimensional grid of potential contexts. More specif-  
 640 ically, we sampled 20 points along each dimension and excluded contexts, in which the sensory  
 641 stimuli were too similar (analogously to the generation of mixing matrices), leaving 272 different  
 642 contexts (see **Figure 7c**, right). The linear decoding performance was determined with a 5-fold  
 643 cross-validation and measured using R-squared. Since the modulatory feedback signals depend  
 644 non-linearly on the context (**Figure 7c**), we tested two non-linear versions of the decoding approach.  
 645 First, we performed a quadratic expansion of the averaged population activity before a linear de-  
 646 coding. Second, we tested a linear decoding of the inverse mixing matrix (four weights) instead of  
 647 the two variables determining the context.

##### 648 Population subspace analysis

649 We visualised the invariant population subspaces by projecting the activity vector onto the two  
 650 readout dimensions and the first principal component. To measure how the orientation of the  
 651 subspaces changes when the context or feedback changes, we computed the angle between the  
 652 planes spanned by the respective subspaces. These planes were fitted on the three-dimensional

653 data described above using the least squares method. Since we were only interested in the relative  
654 orientation of the subspaces, we used a circular measure of the angles, such that a rotation of  
655 180 degrees corresponded to 0 degrees. This means that angles could range between 0 and 90  
656 degrees.

### 657 **Code availability**

658 The code for models and data analysis is publicly available under [https://github.com/sprekelerlab/  
659 feedback\\_modulation\\_Naumann22](https://github.com/sprekelerlab/feedback_modulation_Naumann22).

### 660 **Acknowledgments**

661 We thank Owen Mackwood for providing a code framework that manages simulations on a com-  
662 pute cluster, Loreen Hertäg and Johannes Letzkus for feedback on the manuscript, and the mem-  
663 bers of the Sprekeler lab for valuable discussions. No external funding was received for this work.

### 664 **References**

- 665 **Abs E**, Poorthuis RB, Apelblat D, Muhammad K, Pardi MB, Enke L, Kushinsky D, Pu DL, Eizinger MF, Conzelmann  
666 KK, et al. Learning-related plasticity in dendrite-targeting layer 1 interneurons. *Neuron*. 2018; 100(3):684-  
667 699.
- 668 **Alamia A**, Mozafari M, Choksi B, VanRullen R. On the role of feedback in visual processing: a predictive coding  
669 perspective. *arXiv preprint arXiv:210604225*. 2021; .
- 670 **Azimi Z**, Barzan R, Spoida K, Surdin T, Wollenweber P, Mark MD, Herlitze S, Jancke D. Separable gain control of  
671 ongoing and evoked activity in the visual cortex by serotonergic input. *Elife*. 2020; 9:e53552.
- 672 **Bang D**, Kishida KT, Lohrenz T, White JP, Laxton AW, Tatter SB, Fleming SM, Montague PR. Sub-second dopamine  
673 and serotonin signaling in human striatum during perceptual decision-making. *Neuron*. 2020; 108(5):999-  
674 1010.
- 675 **Bayer HM**, Glimcher PW. Midbrain dopamine neurons encode a quantitative reward prediction error signal.  
676 *Neuron*. 2005; 47(1):129-141.
- 677 **Bell AJ**, Sejnowski TJ. An information-maximization approach to blind separation and blind deconvolution.  
678 *Neural computation*. 1995; 7(6):1129-1159.
- 679 **Bengio Y**, Simard P, Frasconi P. Learning long-term dependencies with gradient descent is difficult. *IEEE trans-  
680 actions on neural networks*. 1994; 5(2):157-166.
- 681 **Branco T**, Staras K. The probability of neurotransmitter release: variability and feedback control at single  
682 synapses. *Nature Reviews Neuroscience*. 2009; 10(5):373-383.
- 683 **van den Brink RL**, Pfeffer T, Donner TH. Brainstem modulation of large-scale intrinsic cortical activity correla-  
684 tions. *Frontiers in human neuroscience*. 2019; 13:340.
- 685 **Brodbeck C**, Jiao A, Hong LE, Simon JZ. Neural speech restoration at the cocktail party: Auditory cortex recovers  
686 masked speech of both attended and ignored speakers. *PLoS biology*. 2020; 18(10):e3000883.
- 687 **Bronkhorst AW**. The cocktail-party problem revisited: early processing and selection of multi-talker speech.  
688 *Attention, Perception, & Psychophysics*. 2015; 77(5):1465-1487.
- 689 **Carandini M**, Heeger DJ. Normalization as a canonical neural computation. *Nature Reviews Neuroscience*.  
690 2012; 13(1):51-62.
- 691 **Cherry EC**. Some experiments on the recognition of speech, with one and with two ears. *The Journal of the  
692 acoustical society of America*. 1953; 25(5):975-979.
- 693 **Chittajallu R**, Pelkey KA, McBain CJ. Neurogliaform cells dynamically regulate somatosensory integration via  
694 synapse-specific modulation. *Nature neuroscience*. 2013; 16(1):13-15.
- 695 **Cichy RM**, Khosla A, Pantazis D, Torralba A, Oliva A. Comparison of deep neural networks to spatio-temporal  
696 cortical dynamics of human visual object recognition reveals hierarchical correspondence. *Scientific reports*.  
697 2016; 6(1):1-13.



698 **DiCarlo JJ**, Cox DD. Untangling invariant object recognition. *Trends in cognitive sciences*. 2007; 11(8):333–341.

699 **DiCarlo JJ**, Zoccolan D, Rust NC. How does the brain solve visual object recognition? *Neuron*. 2012; 73(3):415–  
700 434.

701 **Dipoppa M**, Ranson A, Krumin M, Pachitariu M, Carandini M, Harris KD. Vision and locomotion shape the  
702 interactions between neuron types in mouse visual cortex. *Neuron*. 2018; 98(3):602–615.

703 **Disney AA**, Aoki C, Hawken MJ. Gain modulation by nicotine in macaque v1. *Neuron*. 2007; 56(4):701–713.

704 **Douglas RJ**, Martin KA. Neuronal circuits of the neocortex. *Annu Rev Neurosci*. 2004; 27:419–451.

705 **Dubreuil A**, Valente A, Beiran M, Mastrogiuseppe F, Ostojic S. Complementary roles of dimensionality and  
706 population structure in neural computations. *bioRxiv*. 2020; .

707 **Failor SW**, Carandini M, Harris KD. Learning orthogonalizes visual cortical population codes. *bioRxiv*. 2021; .

708 **Felleman DJ**, Van Essen DC. Distributed hierarchical processing in the primate cerebral cortex. *Cerebral cortex*  
709 (New York, NY: 1991). 1991; 1(1):1–47.

710 **Ferguson KA**, Cardin JA. Mechanisms underlying gain modulation in the cortex. *Nature Reviews Neuroscience*.  
711 2020; 21(2):80–92.

712 **Finn C**, Abbeel P, Levine S. Model-agnostic meta-learning for fast adaptation of deep networks. In: *International*  
713 *Conference on Machine Learning* PMLR; 2017. p. 1126–1135.

714 **Gers FA**, Schmidhuber J, Cummins F. Learning to forget: Continual prediction with LSTM. *Neural computation*.  
715 2000; 12(10):2451–2471.

716 **Gilbert CD**, Li W. Top-down influences on visual processing. *Nature Reviews Neuroscience*. 2013; 14(5):350–  
717 363.

718 **Glaser JI**, Benjamin AS, Chowdhury RH, Perich MG, Miller LE, Kording KP. Machine learning for neural decoding.  
719 *Eneuro*. 2020; 7(4).

720 **Greff K**, Srivastava RK, Koutník J, Steunebrink BR, Schmidhuber J. LSTM: A search space odyssey. *IEEE transac-*  
721 *tions on neural networks and learning systems*. 2016; 28(10):2222–2232.

722 **Halassa MM**, Sherman SM. Thalamocortical circuit motifs: a general framework. *Neuron*. 2019; 103(5):762–  
723 770.

724 **Hasselmo ME**, McGaughy J. High acetylcholine levels set circuit dynamics for attention and encoding and low  
725 acetylcholine levels set dynamics for consolidation. *Progress in brain research*. 2004; 145:207–231.

726 **Hochreiter S**, Schmidhuber J. Long short-term memory. *Neural computation*. 1997; 9(8):1735–1780.

727 **Hochreiter S**, Younger AS, Conwell PR. Learning to learn using gradient descent. In: *International Conference*  
728 *on Artificial Neural Networks* Springer; 2001. p. 87–94.

729 **Hong H**, Yamins DL, Majaj NJ, DiCarlo JJ. Explicit information for category-orthogonal object properties increases  
730 along the ventral stream. *Nature neuroscience*. 2016; 19(4):613–622.

731 **Huber PJ**. Robust Estimation of a Location Parameter. *The Annals of Mathematical Statistics*. 1964; 35(1):73–  
732 101.

733 **Hyvärinen A**, Oja E. Independent component analysis: algorithms and applications. *Neural networks*. 2000;  
734 13(4-5):411–430.

735 **Jarvis S**, Nikolic K, Schultz SR. Neuronal gain modulability is determined by dendritic morphology: a computa-  
736 tional optogenetic study. *PLoS computational biology*. 2018; 14(3):e1006027.

737 **Kar K**, Kubilius J, Schmidt K, Issa EB, DiCarlo JJ. Evidence that recurrent circuits are critical to the ventral stream’s  
738 execution of core object recognition behavior. *Nature neuroscience*. 2019; 22(6):974–983.

739 **Kawai H**, Lazar R, Metherate R. Nicotinic control of axon excitability regulates thalamocortical transmission.  
740 *Nature neuroscience*. 2007; 10(9):1168–1175.

741 **Keller GB**, Mrsic-Flogel TD. Predictive processing: a canonical cortical computation. *Neuron*. 2018; 100(2):424–  
742 435.

743 **Kietzmann TC**, Spoerer CJ, Sørensen LK, Cichy RM, Hauk O, Kriegeskorte N. Recurrence is required to capture  
744 the representational dynamics of the human visual system. *Proceedings of the National Academy of Sciences*.  
745 2019; 116(43):21854–21863.

746 **Kingma DP**, Ba J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*. 2014; .

747 **Kriegeskorte N**. Deep neural networks: a new framework for modeling biological vision and brain information  
748 processing. *Annual review of vision science*. 2015; 1:417–446.

749 **Kriegeskorte N**, Mur M, Bandettini PA. Representational similarity analysis-connecting the branches of sys-  
750 tems neuroscience. *Frontiers in systems neuroscience*. 2008; 2:4.

751 **Kuchibhotla KV**, Gill JV, Lindsay GW, Papadoyannis ES, Field RE, Sten TAH, Miller KD, Froemke RC. Parallel  
752 processing by cortical inhibition enables context-dependent behavior. *Nature neuroscience*. 2017; 20(1):62–  
753 71.

754 **Larkum ME**, Senn W, Lüscher HR. Top-down dendritic input increases the gain of layer 5 pyramidal neurons.  
755 *Cerebral cortex*. 2004; 14(10):1059–1070.

756 **Laviv T**, Riven I, Dolev I, Vertkin I, Balana B, Slesinger PA, Slutsky I. Basal GABA regulates GABABR conformation  
757 and release probability at single hippocampal synapses. *Neuron*. 2010; 67(2):253–267.

758 **LeCun Y**, Bengio Y, Hinton G. Deep learning. *nature*. 2015; 521(7553):436–444.

759 **Lillicrap TP**, Santoro A, Marris L, Akerman CJ, Hinton G. Backpropagation and the brain. *Nature Reviews*  
760 *Neuroscience*. 2020; 21(6):335–346.

761 **Lohani S**, Moberly AH, Benisty H, Landa B, Jing M, Li Y, Higley MJ, Cardin JA. Dual color mesoscopic imaging  
762 reveals spatiotemporally heterogeneous coordination of cholinergic and neocortical activity. *bioRxiv*. 2020;  
763 .

764 **Malina KCK**, Tsivourakis E, Kushinsky D, Apelblat D, Shtiglitz S, Zohar E, Sokoletsky M, Tasaka Gi, Mizrahi A,  
765 Lampl I, et al. NDNF interneurons in layer 1 gain-modulate whole cortical columns according to an animal's  
766 behavioral state. *Neuron*. 2021; .

767 **Mante V**, Sussillo D, Shenoy KV, Newsome WT. Context-dependent computation by recurrent dynamics in  
768 prefrontal cortex. *nature*. 2013; 503(7474):78–84.

769 **Markov NT**, Vezoli J, Chameau P, Falchier A, Quilodran R, Huissoud C, Lamy C, Misery P, Giroud P, Ullman  
770 S, et al. Anatomy of hierarchy: feedforward and feedback pathways in macaque visual cortex. *Journal of*  
771 *Comparative Neurology*. 2014; 522(1):225–259.

772 **Marques T**, Nguyen J, Fioreze G, Petreanu L. The functional organization of cortical feedback inputs to primary  
773 visual cortex. *Nature neuroscience*. 2018; 21(5):757–764.

774 **McAdams CJ**, Maunsell JH. Effects of attention on orientation-tuning functions of single neurons in macaque  
775 cortical area V4. *Journal of Neuroscience*. 1999; 19(1):431–441.

776 **McDermott JH**. The cocktail party problem. *Current Biology*. 2009; 19(22):R1024–R1027.

777 **Miller RJ**. Presynaptic receptors. *Annual review of pharmacology and toxicology*. 1998; 38(1):201–227.

778 **Molyneaux BJ**, Hasselmo ME. GABAB presynaptic inhibition has an in vivo time constant sufficiently rapid to  
779 allow modulation at theta frequency. *Journal of Neurophysiology*. 2002; 87(3):1196–1205.

780 **Naumann LB**, Sprekeler H. Presynaptic inhibition rapidly stabilises recurrent excitation in the face of plasticity.  
781 *PLoS Computational Biology*. 2020; 16(8):e1008118.

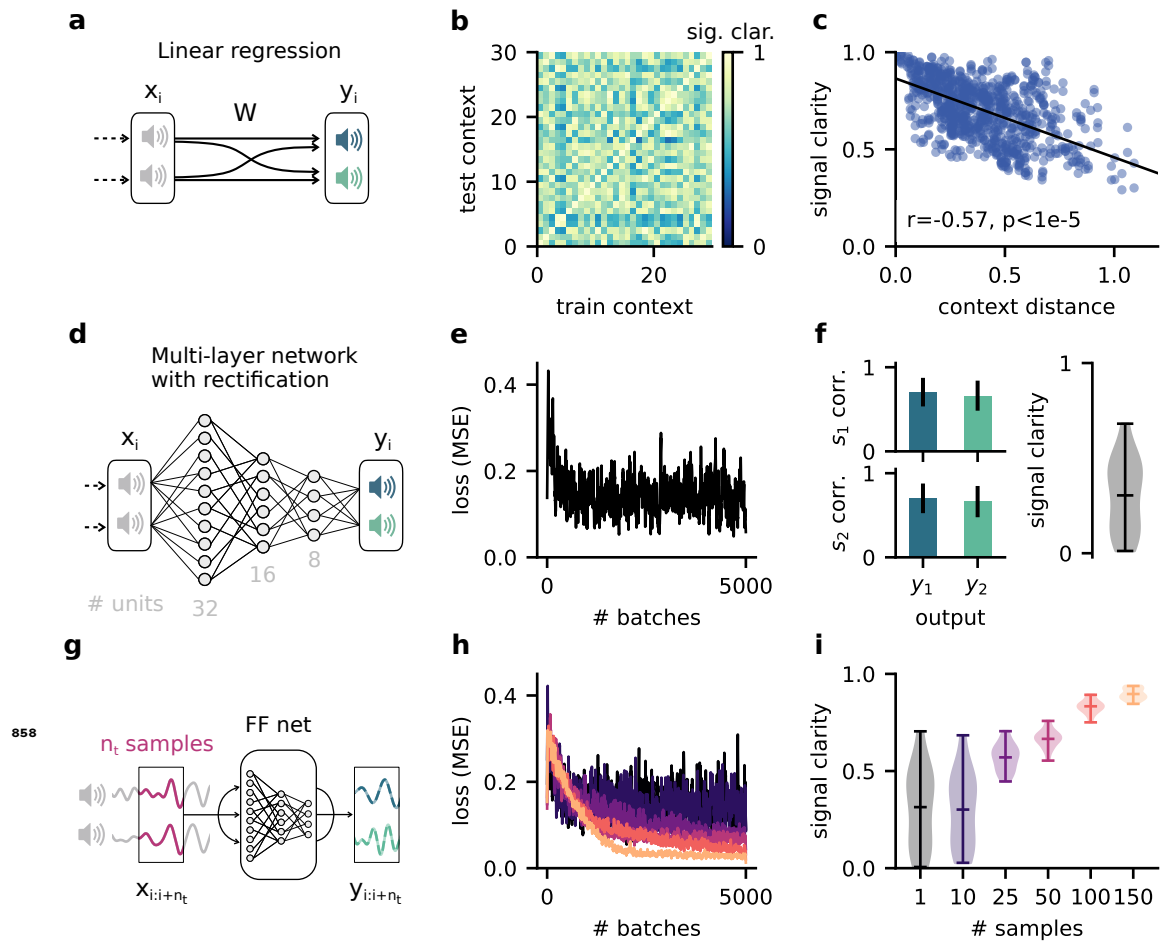
782 **Nayebi A**, Sagastuy-Brena J, Bear DM, Kar K, Kubilius J, Ganguli S, Sussillo D, DiCarlo JJ, Yamins DL. Goal-driven  
783 recurrent neural network models of the ventral visual stream. *bioRxiv*. 2021; .

784 **Niell CM**, Stryker MP. Modulation of visual responses by behavioral state in mouse visual cortex. *Neuron*. 2010;  
785 65(4):472–479.

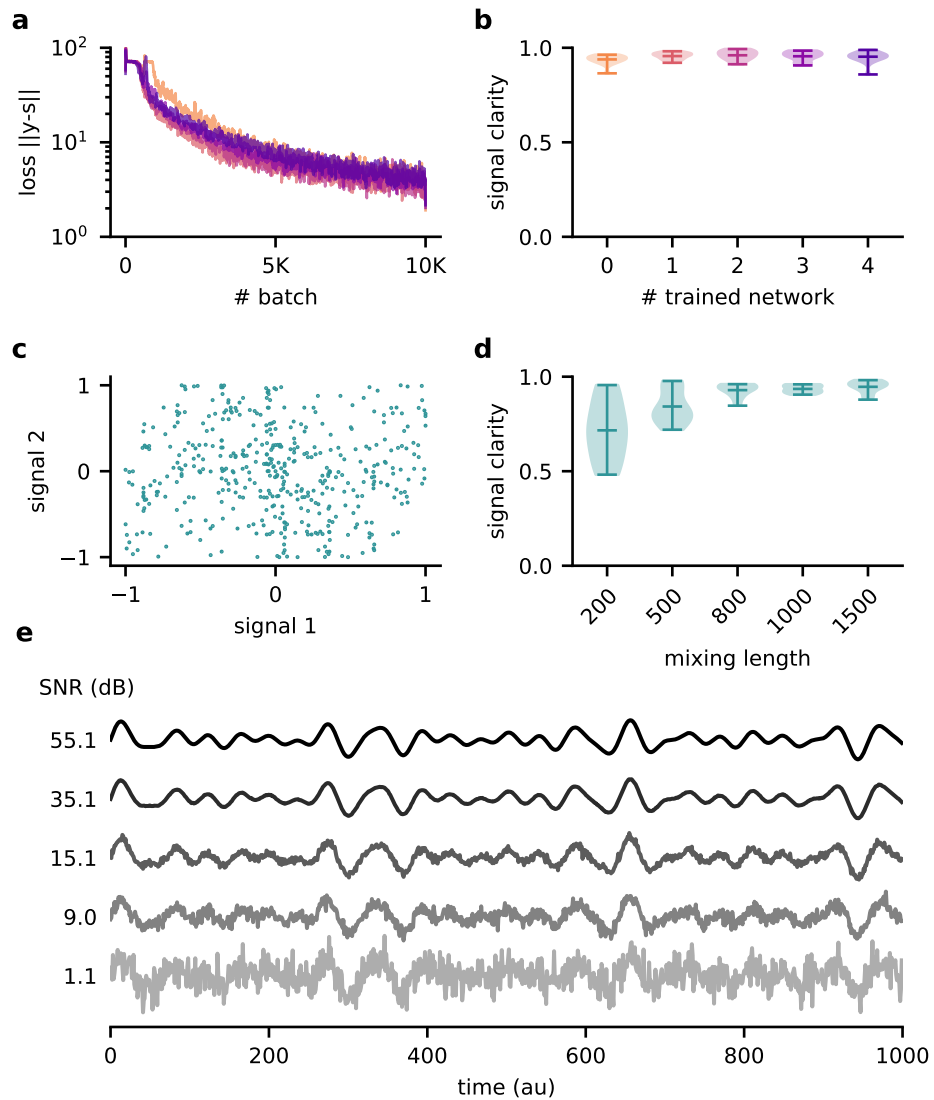
786 **Oberfeld D**, KloECKner-Nowotny F. Individual differences in selective attention predict speech identification at  
787 a cocktail party. *Elife*. 2016; 5:e16747.

- 788 **Olshausen BA**, Anderson CH, Van Essen DC. A neurobiological model of visual attention and invariant pattern  
789 recognition based on dynamic routing of information. *Journal of Neuroscience*. 1993; 13(11):4700–4719.
- 790 **Pardi MB**, Vogenstahl J, Dalmay T, Spanò T, Pu DL, Naumann LB, Kretschmer F, Sprekeler H, Letzkus JJ. A  
791 thalamocortical top-down circuit for associative memory. *Science*. 2020; 370(6518):844–848.
- 792 **Parthasarathy A**, Hancock KE, Bennett K, DeGruttola V, Polley DB. Bottom-up and top-down neural signatures  
793 of disordered multi-talker speech perception in adults with normal hearing. *Elife*. 2020; 9:e51419.
- 794 **Pascanu R**, Mikolov T, Bengio Y. On the difficulty of training recurrent neural networks. In: *International con-*  
795 *ference on machine learning* PMLR; 2013. p. 1310–1318.
- 796 **Paszke A**, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, Killeen T, Lin Z, Gimelshein N, Antiga L, et al. Pytorch:  
797 An imperative style, high-performance deep learning library. *Advances in neural information processing*  
798 *systems*. 2019; 32:8026–8037.
- 799 **Pinto L**, Goard MJ, Estandian D, Xu M, Kwan AC, Lee SH, Harrison TC, Feng G, Dan Y. Fast modulation of visual  
800 perception by basal forebrain cholinergic neurons. *Nature neuroscience*. 2013; 16(12):1857–1863.
- 801 **Polack PO**, Friedman J, Golshani P. Cellular mechanisms of brain state-dependent gain modulation in visual  
802 cortex. *Nature neuroscience*. 2013; 16(9):1331–1339.
- 803 **Poorthuis RB**, Bloem B, Schak B, Wester J, de Kock CP, Mansvelder HD. Layer-specific modulation of the pre-  
804 frontal cortex by nicotinic acetylcholine receptors. *Cerebral cortex*. 2013; 23(1):148–161.
- 805 **Purushothaman G**, Marion R, Li K, Casagrande VA. Gating and control of primary visual cortex by pulvinar.  
806 *Nature neuroscience*. 2012; 15(6):905–912.
- 807 **Quiroga RQ**, Reddy L, Kreiman G, Koch C, Fried I. Invariant visual representation by single neurons in the  
808 human brain. *Nature*. 2005; 435(7045):1102–1107.
- 809 **Reynolds JH**, Heeger DJ. The normalization model of attention. *Neuron*. 2009; 61(2):168–185.
- 810 **Riesenhuber M**, Poggio T. Hierarchical models of object recognition in cortex. *Nature neuroscience*. 1999;  
811 2(11):1019–1025.
- 812 **Roth MM**, Dahmen JC, Muir DR, Imhof F, Martini FJ, Hofer SB. Thalamic nuclei convey diverse contextual infor-  
813 mation to layer 1 of visual cortex. *Nature neuroscience*. 2016; 19(2):299–307.
- 814 **Sabatini BL**, Tian L. Imaging neurotransmitter and neuromodulator dynamics in vivo with genetically encoded  
815 indicators. *Neuron*. 2020; 108(1):17–32.
- 816 **Salinas E**, Abbott L. Invariant visual responses from attentional gain fields. *Journal of Neurophysiology*. 1997;  
817 77(6):3267–3272.
- 818 **Salinas E**, Sejnowski TJ. Book review: gain modulation in the central nervous system: where behavior, neuro-  
819 physiology, and computation meet. *The Neuroscientist*. 2001; 7(5):430–440.
- 820 **Salinas E**, Thier P. Gain modulation: a major computational principle of the central nervous system. *Neuron*.  
821 2000; 27(1):15–21.
- 822 **Sampathkumar V**, Miller-Hansen A, Sherman SM, Kasthuri N. Integration of signals from different cortical  
823 areas in higher order thalamic neurons. *Proceedings of the National Academy of Sciences*. 2021; 118(30).
- 824 **Sherman SM**. Thalamus plays a central role in ongoing cortical functioning. *Nature neuroscience*. 2016;  
825 19(4):533–541.
- 826 **Sherman SM**, Guillery R. On the actions that one nerve cell can have on another: distinguishing “drivers” from  
827 “modulators”. *Proceedings of the National Academy of Sciences*. 1998; 95(12):7121–7126.
- 828 **Shine JM**, Müller EJ, Munn B, Cabral J, Moran RJ, Breakspear M. Computational models link cellular mechanisms  
829 of neuromodulation to large-scale neural dynamics. *Nature neuroscience*. 2021; 24(6):765–776.
- 830 **Spoerer CJ**, McClure P, Kriegeskorte N. Recurrent convolutional neural networks: a better model of biological  
831 object recognition. *Frontiers in psychology*. 2017; 8:1551.
- 832 **Stroud JP**, Porter MA, Hennequin G, Vogels TP. Motor primitives in space and time via targeted gain modulation  
833 in cortical networks. *Nature neuroscience*. 2018; 21(12):1774–1783.

- 834 **Thorat S**, Aldegheri G, Kietzmann TC. Category-orthogonal object features guide information processing in  
835 recurrent neural networks trained for object categorization. *arXiv preprint arXiv:211107898*. 2021; .
- 836 **Thurley K**, Senn W, Luscher HR. Dopamine increases the gain of the input-output response of rat prefrontal  
837 pyramidal neurons. *Journal of neurophysiology*. 2008; 99(6):2985–2997.
- 838 **Urban-Ciecko J**, Fanselow EE, Barth AL. Neocortical somatostatin neurons reversibly silence excitatory trans-  
839 mission via GABA<sub>B</sub> receptors. *Current Biology*. 2015; 25(6):722–731.
- 840 **Vinck M**, Batista-Brito R, Knoblich U, Cardin JA. Arousal and locomotion make distinct contributions to cortical  
841 activity patterns and visual encoding. *Neuron*. 2015; 86(3):740–754.
- 842 **Wang JX**, Kurth-Nelson Z, Kumaran D, Tirumala D, Soyer H, Leibo JZ, Hassabis D, Botvinick M. Prefrontal cortex  
843 as a meta-reinforcement learning system. *Nature neuroscience*. 2018; 21(6):860–868.
- 844 **Wang J**, Narain D, Hosseini EA, Jazayeri M. Flexible timing by temporal scaling of cortical responses. *Nature*  
845 *neuroscience*. 2018; 21(1):102–110.
- 846 **Werbos PJ**. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*. 1990;  
847 78(10):1550–1560.
- 848 **Wiskott L**. How does our visual system achieve shift and size invariance. JL van Hemmen and TJ Sejnowski,  
849 editors. 2006; 23:322–340.
- 850 **Wiskott L**, Sejnowski T. Slow Feature Analysis: Unsupervised Learning of Invariances. *Neural Computation*.  
851 2002; 14:715–770.
- 852 **Yahav PHs**, Golumbic EZ. Linguistic processing of task-irrelevant speech at a Cocktail Party. *Elife*. 2021;  
853 10:e65096.
- 854 **Yamins DL**, DiCarlo JJ. Using goal-driven deep learning models to understand sensory cortex. *Nature neuro-*  
855 *science*. 2016; 19(3):356–365.
- 856 **Zhuang C**, Yan S, Nayebi A, Schrimpf M, Frank MC, DiCarlo JJ, Yamins DL. Unsupervised neural network models  
857 of the ventral visual stream. *Proceedings of the National Academy of Sciences*. 2021; 118(3).



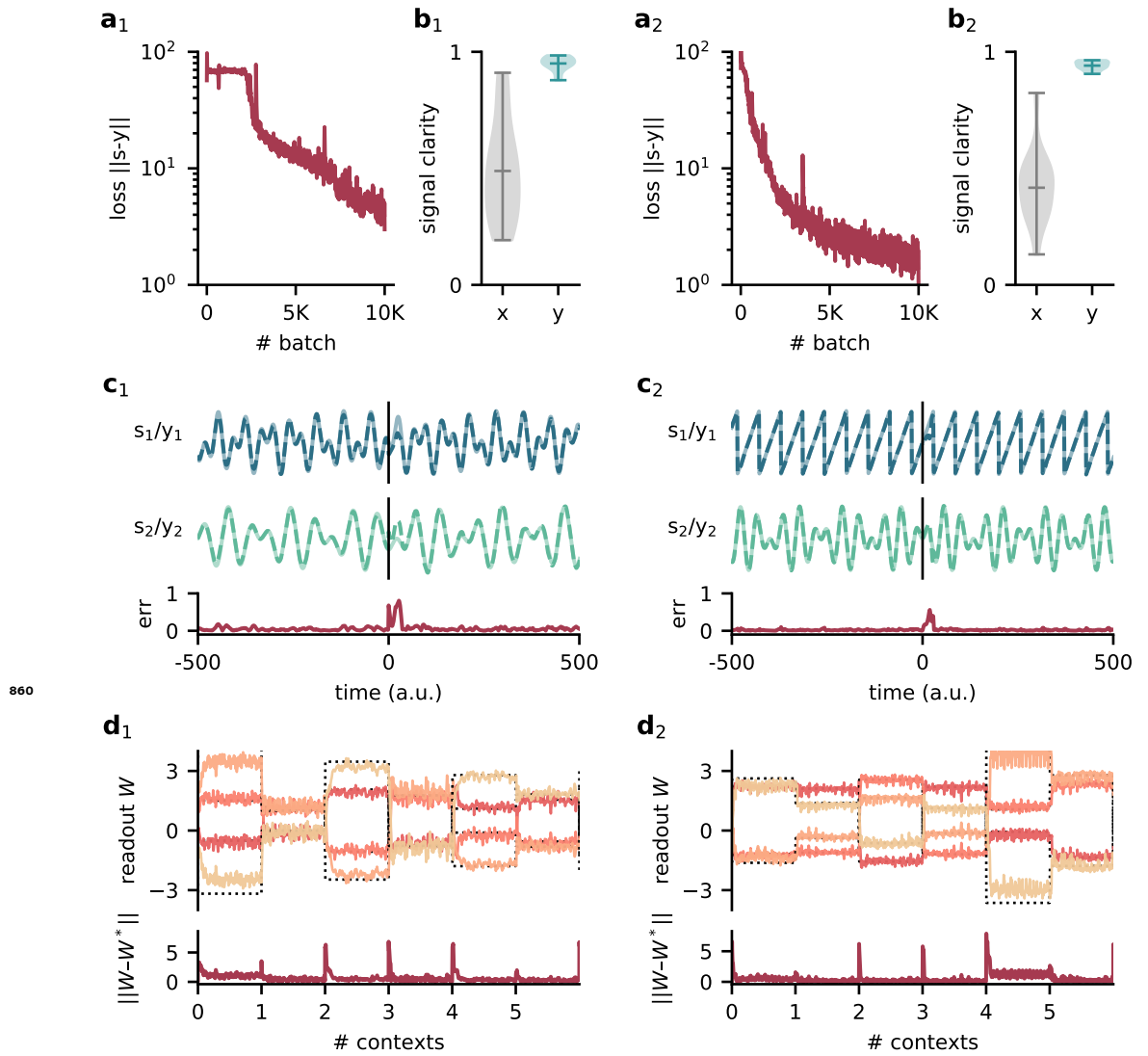
**Figure 1-Figure supplement 1. The dynamic blind source separation task cannot be solved with a feedforward network, unless the network receives a sequence of inputs at once. This would require an additional mechanism to retain information over time.** **a.** Schematic of a feedforward network consisting of a linear readout only. **b.** Pairwise signal clarity of one context when the network is trained on another context. **c.** Correlation between the distance between two contexts and their pairwise signal clarity (see (b)). **d.** Schematic of a multi-layer feedforward network with three hidden layers (32, 16 and 8 rectified linear units). **e.** Loss during training for the network in (d), measure by the mean squared error between the output and the sources. **f.** Network performance after training. Left: Correlation of the outputs with the sources over 20 contexts. Error bars indicate standard deviation. Right: Signal clarity across 20 contexts for the trained network. **g.** Schematic of network architecture and training setup when using a sequence of  $n_t$  samples as input to the multi-layer network. **h.** Same as (e) but for different number of samples. Color code corresponds to (i). **i.** Signal clarity for trained networks that receive different numbers of samples as input.



859

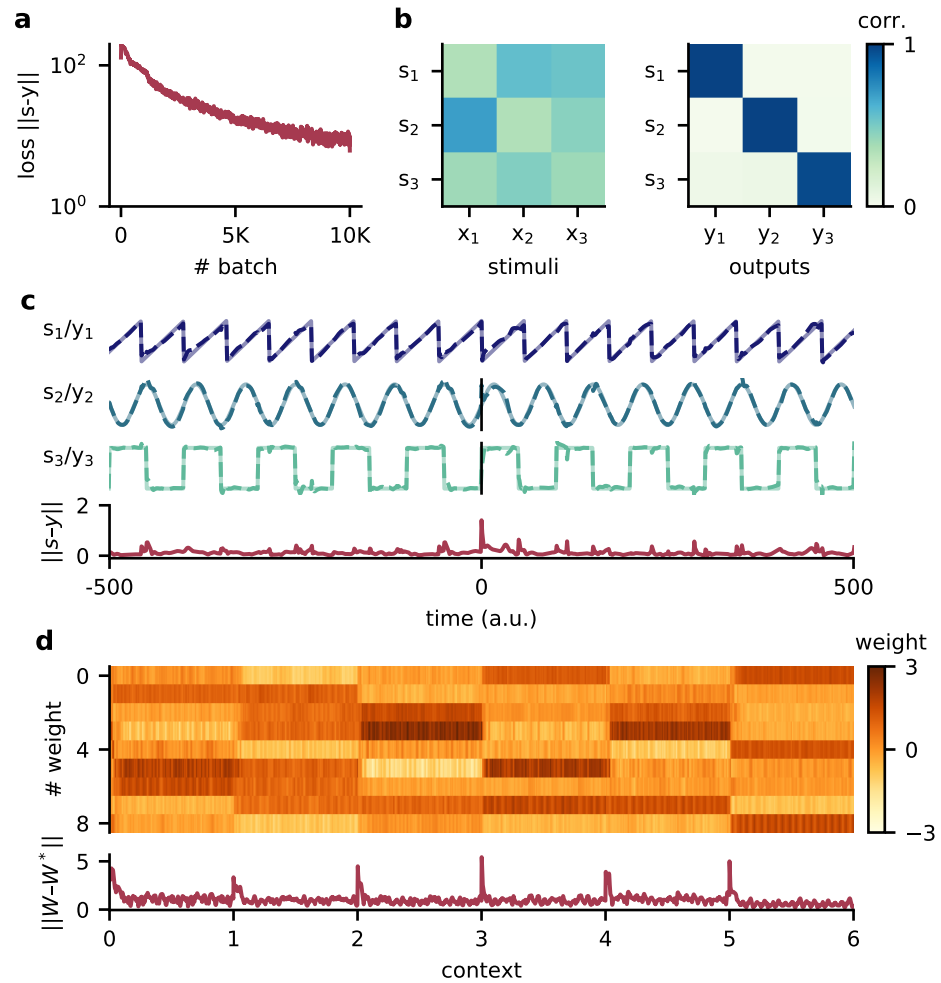
**Figure 1-Figure supplement 2. Robustness of the feedback-driven modulation mechanism.**

**a.** Loss over training for 5 different random initialisations of the model and **b.** signal clarity for 20 test contexts in the corresponding trained networks. The model performance is robust across model instantiations. **c.** Samples from the two default signals are uncorrelated. **d.** Signal clarity for different lengths of the context during testing. The length of the context interval is not crucial for performance, indicating that the network did not learn the interval by heart. **e.** Example traces of the sensory stimuli for different signal-to-noise ratios.



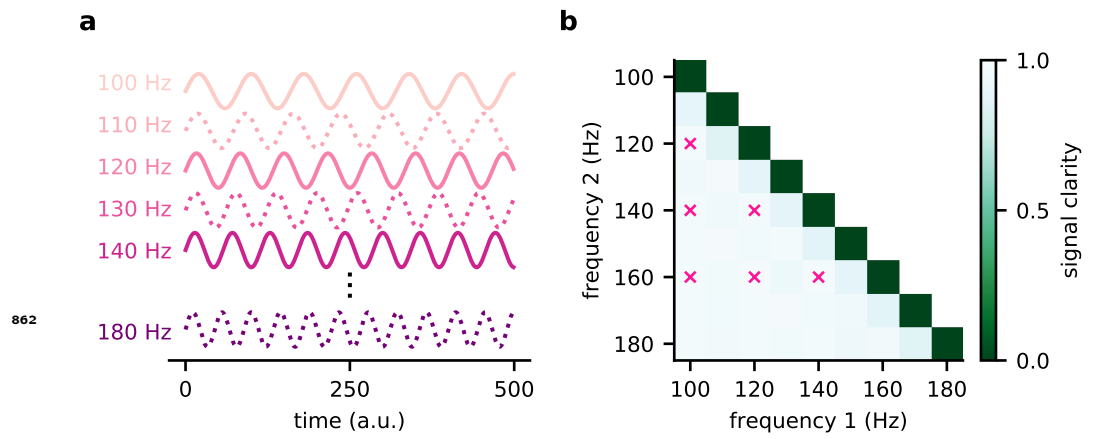
**Figure 1-Figure supplement 3. Model performance for two different sets of source signals.**

Left: Compositions of sines with  $f_{11} = 120$  Hz,  $f_{12} = 2.2$  Hz,  $f_{21} = 100$  Hz and  $f_{22} = 145$  Hz. Right: Sawtooth function with frequency 140 Hz and composed sine of 150 Hz and 210 Hz. **a**<sub>1/2</sub>. Loss over training. **b**<sub>1/2</sub>. Signal clarity for 20 test contexts measured in the sensory stimuli and the network output. **c**<sub>1/2</sub>. Example traces of the sources and the network output (top) and corresponding deviation between them (bottom). The context changes at time 0. **d**<sub>1/2</sub>. Top: Readout weights across 6 contexts; dotted lines indicate the optimal weights. Bottom: Deviation of readout from the optimal weights.

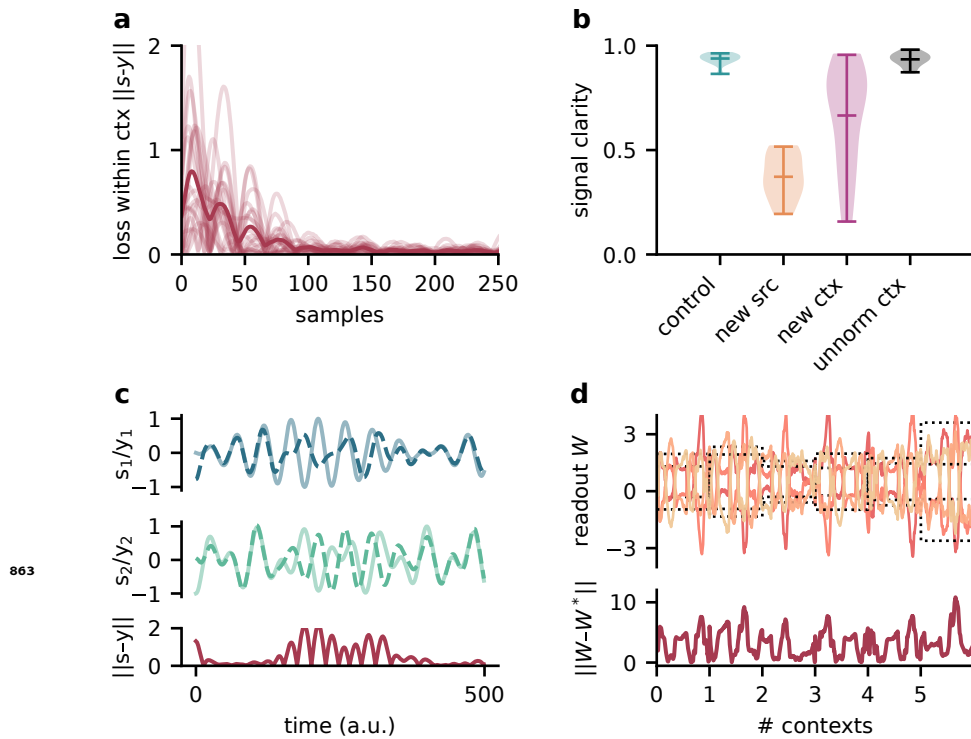


**Figure 1-Figure supplement 4. Model performance for three source signals.** **a.** Loss over training. **b.** Correlation of the sources with the mixed sensory stimuli (left) and with the network outputs (right). **c.** Example traces of the three source signals and network outputs (top) and corresponding deviation between them (bottom). The context changes at time 0. The source signals are a sawtooth of frequency 140 Hz, a sine wave of frequency 120 Hz and a square wave signal of 80 Hz. **d.** Top: Readout weights across 6 contexts. Bottom: Deviation of readout from the optimal weights.

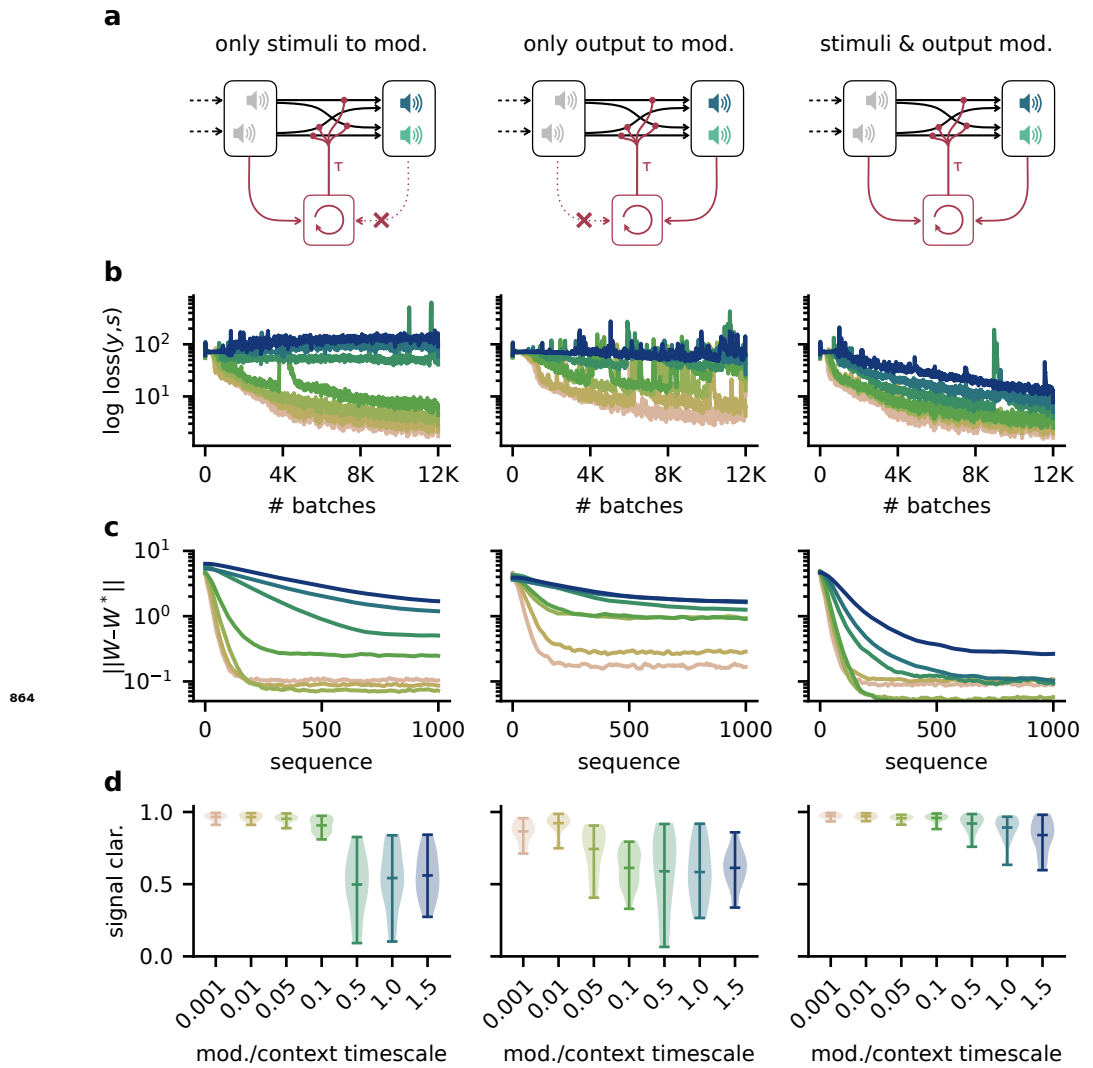




**Figure 1-Figure supplement 5. The modulated network model generalises across frequencies.** **a.** Illustration of the source signals used during training (solid lines) and only during testing (dotted lines). During the training, the model experiences only a subset of potential signals. **b.** Signal clarity for different combinations of test frequencies. Combinations used during training are marked with a pink cross.

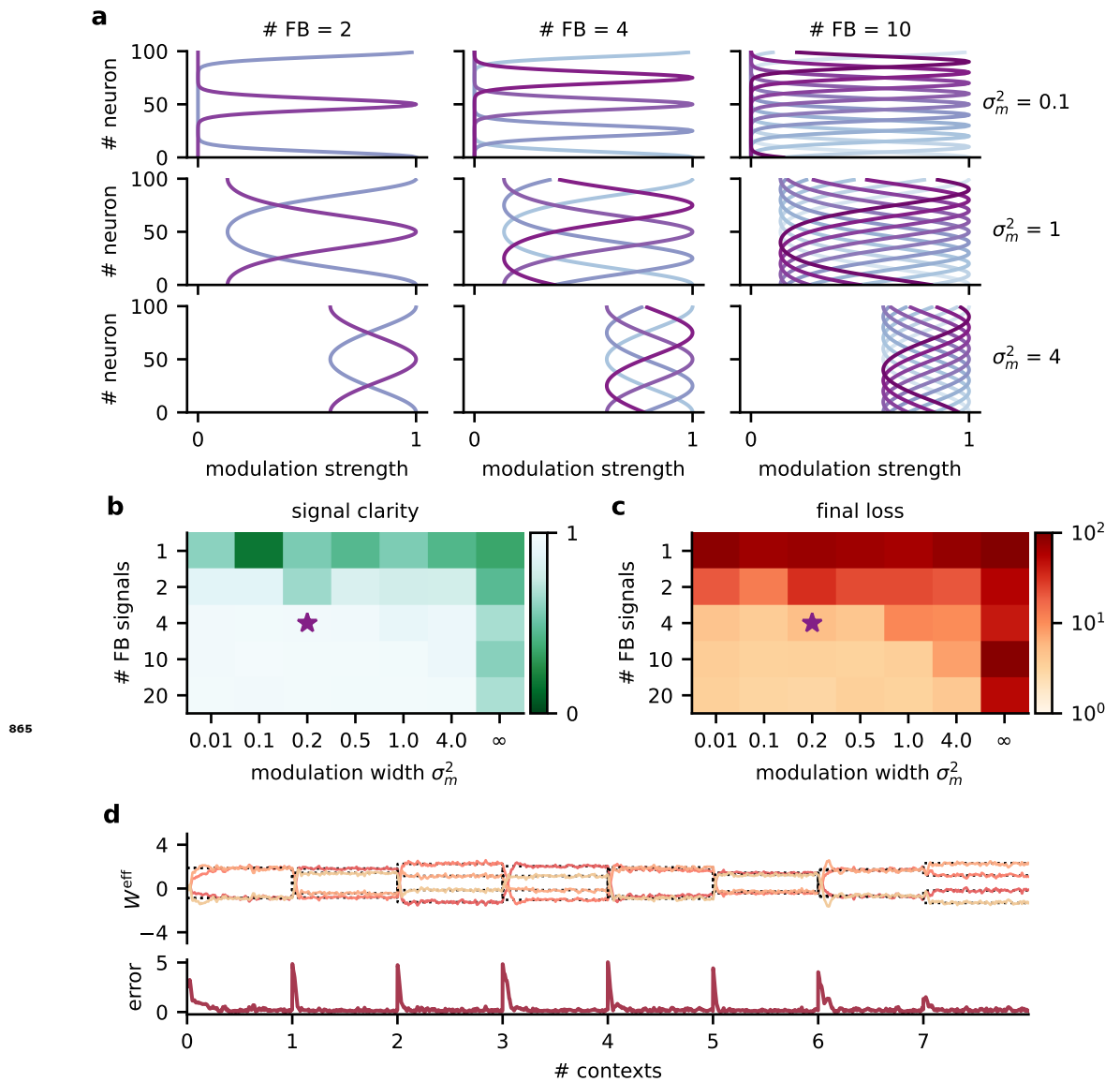


**Figure 1-Figure supplement 6. The modulator learns a model of the sources and contexts, and infers the current context from the stimuli. Testing the network on sources and contexts with different statistics than during training thus impairs its performance. a.** Deviation of network output from sources within contexts. Average across contexts shown in dark red. **b.** Signal clarity for different test cases: same sources and same context statistics as during training ("control"), new sources ("new src"), same sources but different context statistic (i.e. unnormalized mixing matrices, "new ctx"), and different context statistics but when training the network on them ("unnorm ctx"). **c.** Top: Sources ( $s_{1,2}$ ) and network output ( $y_{1,2}$ ) for a context when testing on new sources. Bottom: Deviation of outputs from the sources. **d.** Top: Modulated readout weights across 6 contexts when testing on new sources; dotted lines indicate the inverse of the current mixing. Bottom: Deviation of readout from target weights.



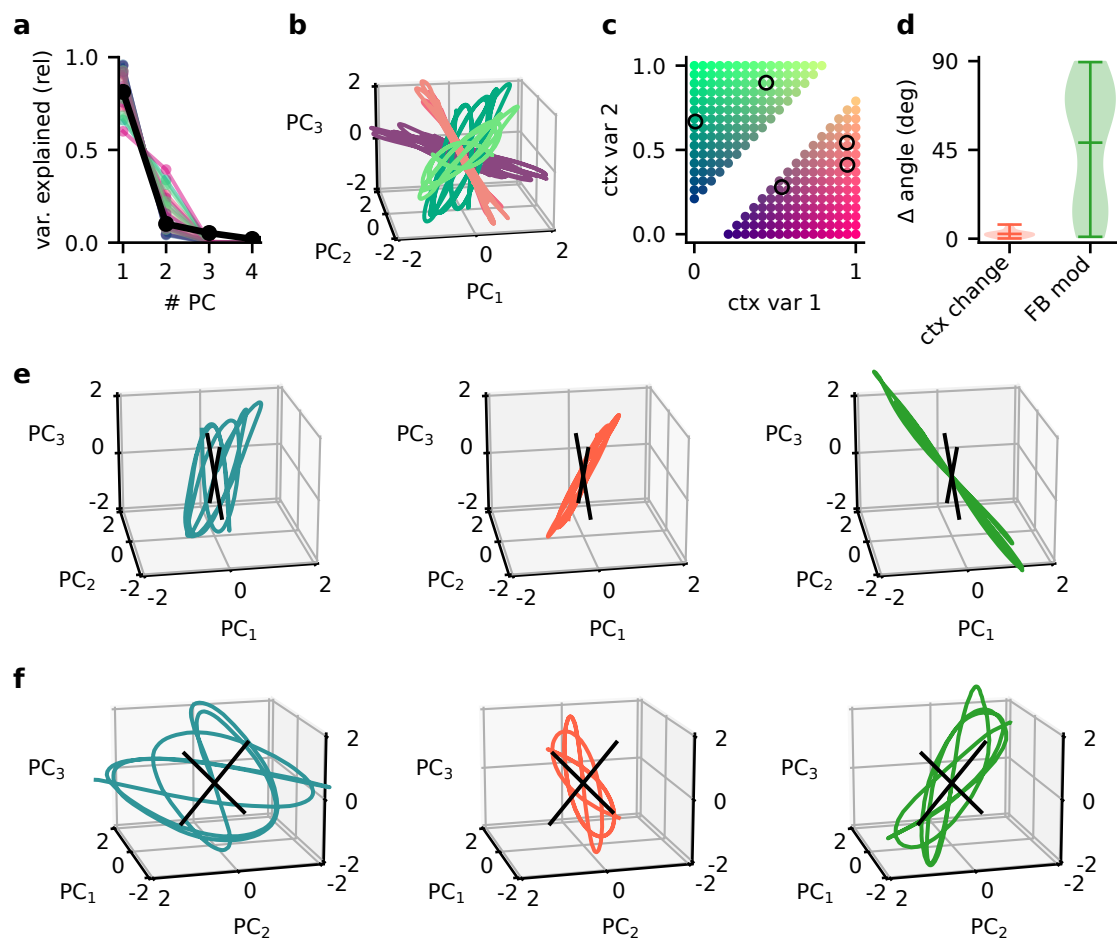
**Figure 2-Figure supplement 1. Robustness to slow feedback modulation depends on the inputs to the modulatory system.**

**a.** Illustration of different input configurations: the modulatory system receives only the sensory stimuli as feedforward input (left), only the network output as feedback input (right) or both (right). **b.** Loss over training for different timescales. Colours correspond to values shown in (d). **c.** Deviation of the readout weights from the optimal weights over the duration of a context for different modulation timescales, averaged across 20 contexts. Colours correspond to values shown in (d). **d.** Signal clarity for different timescales of the modulatory feedback signal.

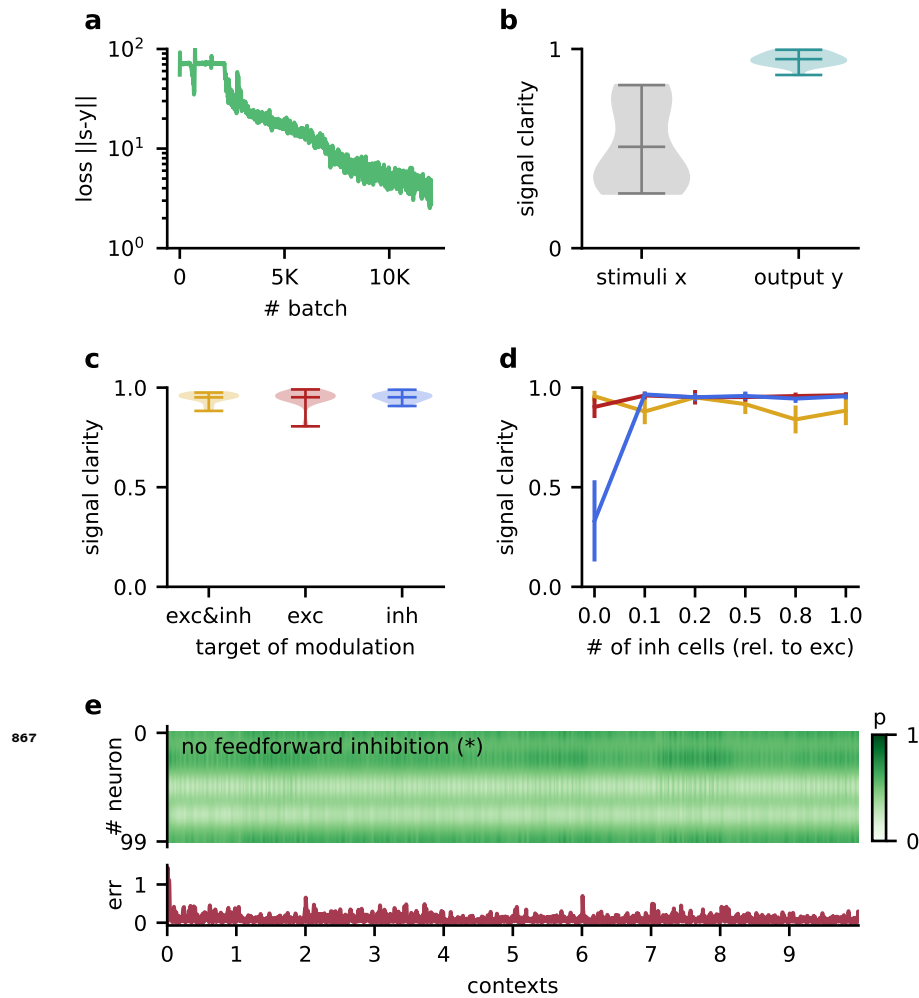


**Figure 3-Figure supplement 1. Robustness to the spatial scale of feedback modulation. a.**

Examples of the spatial extent of feedback modulation for different numbers of feedback signals (# FB) and spatial spread ( $\sigma_m^2$ ). **b.** Signal clarity and **c.** final log loss in network models with different parameters determining the spatial scale of feedback modulation. Signal clarity was averaged across 20 contexts. Final loss was averaged across the last 200 batches during training. The purple star indicates default values used in the main results. Modulation width of " $\infty$ " corresponds to a homogeneous modulation over the whole population. **d.** Top: Effective weights from stimuli to network output over 8 contexts. Effective weights are computed as the modulated weights from stimuli to neural population, multiplied with the readout weights. Dotted lines indicate inverse of mixing. Bottom: Deviation of effective weights from the inverse.



**Figure 5-Figure supplement 1. Principal component analysis captures the low-dimensional population subspaces and the subspace re-orientation with feedback.** **a.** Fraction of variance explained by principal component analyses on single contexts (coloured lines) and across all contexts (black line). **b.** Population activity in the space of the first 3 PCs for 5 contexts. Colour indicates the location of the contexts in context space as shown in **(c)**. **d.** Violin plot of the angle change between original subspace and the subspace for context changes (ctx change) and feedback modulation (FB mod). **e.** Population activity in the space of the first 3 PCs in different stages of the experiment. Left: context 1 with intact feedback, center: context 2 with frozen feedback, right: context 2 with intact feedback. Black lines indicate the readout vectors. **f.** Same as **(e)** but from a different viewpoint to show the readout space.



**Figure 6-Figure supplement 1. The Dalean network can learn to solve the dynamic blind source separation task, and the performance does not depend on specifics of the model architecture.** **a.** Loss over training. **b.** Violin plot of the signal clarity for 20 test contexts measured in the sensory stimuli and the network output. **c.** Violin plot of signal clarity for models in which excitatory, inhibitory or both types of synapses are modulated by feedback; measured over 20 contexts. **d.** Mean signal clarity across 20 contexts for different numbers of inhibitory neurons  $N_i$  (relative to the number of neurons in the higher-level population). Colours correspond to the targets of modulation from (c). Error bars indicate standard deviation. The yellow arrow indicates the default parameter used in the main results. The star indicates networks without feedforward inhibition (see (e)). **e.** Top: Modulation of neurons in the higher-level population across 10 contexts without feedforward inhibition. The modulation does not switch with the context but fluctuates on a faster timescale. Bottom: Corresponding deviation of the network output from the sources.