1    **Using top-down modulation to optimally balance shared versus separated task representations**

2

3    *Authors: Pieter Verbeke[1] & Tom Verguts[1]*

4    *Affiliations:* [1] Department of experimental psychology; Ghent University

5    ORCID IDs:        Pieter Verbeke: 0000-0003-2919-1528

6                          Tom Verguts: 0000-0002-7783-4754

7    Corresponding author email: pjverbek.verbeke@ugent.be

8

9    Keywords: Cognitive control, modulation, neural representations, generalization

10

11                              **Declarations**

12

19

20    **Data availability statement:** Code to simulate the model and analyze the resulting data is provided

21    in our GitHub repository: https://github.com/CogComNeuroSci/PieterV_public/tree/master/Gating

**Abstract**

22  Human adaptive behavior requires continually learning and performing a wide variety of tasks, often

23  with very little practice. To accomplish this, it is crucial to separate neural representations of different

24  tasks in order to avoid interference. At the same time, sharing neural representations supports

25  generalization and allows faster learning. Therefore, a crucial challenge is to find an optimal balance

26  between shared versus separated representations. Typically, models of human cognition employ top-

27  down modulatory signals to separate task representations, but there exist surprisingly little systematic

28  computational investigations of how such modulation is best implemented. We identify and

29  systematically evaluate two crucial features of modulatory signals. First, top-down input can be

30  processed in an additive or multiplicative manner. Second, the modulatory signals can be adaptive

31  (learned) or non-adaptive (random). We cross these two features, resulting in four modulation networks

32  which are tested on a variety of input datasets and tasks with different degrees of stimulus-action

33  mapping overlap. The multiplicative adaptive modulation network outperforms all other networks in

34  terms of accuracy. Moreover, this network develops hidden units that optimally share representations

35  between tasks. Specifically, different than the binary approach of currently popular latent state models,

36  it exploits partial overlap between tasks.

# 1. Introduction

Humans and other rational agents need to continually learn and perform an enormous number of complex tasks. Sometimes very similar contexts require totally different actions. For instance, while soccer and handball both require to put a ball in a goal which is guarded by a keeper and some defenders, soccer requires to manipulate the ball with the feet while handball requires to manipulate the ball with the hands. In such contexts, it is important to separate stimulus-action representations between the two tasks as much as possible in order to avoid interference. However, at other times, two different contexts nevertheless require partially similar actions. Despite the fact that tennis requires to play a ball over a low-hanging net and badminton requires to play a shuttle over a higher placed net, one can partially generalize the action of swinging the racket between the two sports. Thus, in these cases, an agent can significantly benefit from partially sharing knowledge between the two tasks.

Previous research (Baxter, 2019; Franklin & Frank, 2018; Musslick et al., 2017; Vaidya, Jones, Castillo, & Badre, 2021; Zambaldi et al., 2018) indeed illustrated that sharing task representations significantly improves learning and generalization across tasks, two hallmarks of human flexibility. However, sharing task representations in a neural network severely impacts the network's ability to perform more than one task at the same time (i.e., to multi-task; Alon et al., 2017; Musslick et al., 2017; Musslick, Saxe, Novick, Reichman, & Cohen, 2020). Moreover, shared task representations leave a network very vulnerable to overwriting previously learned information. This problem is known as catastrophic interference (French, 1999). In contrast, a network that develops separated task representations experiences less problems in multi-tasking (Musslick et al., 2020; Tsai, Saxe, & Cox, 2016) and can continually learn without forgetting (Kirkpatrick et al., 2017; Masse, Grant, & Freedman, 2018; McClelland, McNaughton, & O'Reilly, 1995; Verbeke & Verguts, 2019). However, such networks are less able to generalize and as a consequence must learn even very similar tasks (like tennis and badminton) from scratch. In sum, there exists a trade-off between sharing and separating task representations in neural networks (Musslick et al., 2017; Musslick & Cohen, 2020; Sagiv, Musslick, Niv, & Cohen, 2020).

One popular solution to deal with this sharing-separating trade-off are compositional task representations (Fidler, Boben, & Leonardis, 2009; Franklin & Frank, 2018; Lake et al., 2014; Sugita,

65    Tani, & Butz, 2011; Tubiana & Monasson, 2017; Yang, Joglekar, Song, Newsome, & Wang, 2019).

66    For instance, to a first approximation, knowledge of soccer can be decomposed in two basic building

67    blocks: the goal of the task (getting the ball past the goalkeeper) and the actions (kicking the ball). This

68    allows the agent to generalize the goal when learning to play handball but also to avoid interference by

69    separating the actions between both sports. Hence, a novel task can be learned quickly by recombining

70    building blocks from previously learned tasks. Indeed, generalizing information through compositional

71    task representations received considerable attention in several cognitive domains such as language

72    (Irsoy & Cardie, 2014; İrsoy & Cardie, 2015; Lake et al., 2014) and sensorimotor learning (Butz,

73    Achimova, Bilkey, & Knott, 2021; Butz, Bilkey, Humaidan, Knott, & Otte, 2019; Sugita et al., 2011).

74    Nevertheless, it is not clear which neural network configurations could learn such compositional

75    representations (Hupkes, Dankers, Mul, & Bruni, 2020; Lake & Baroni, 2018; Lake, Ullman,

76    Tenenbaum, & Gershman, 2017), and what the resulting compositional representations would look like.

77    The current work aims to build upon previous cognitive and computational work to investigate which

78    cognitive architectures can balance shared and separated task representations in typical cognitive tasks,

79    and what type of representations successful architectures would develop.

80          In cognitive science, the ability to perform one task while eliminating interference from other

81    tasks, is known as cognitive control. Influential theoretical work (Miller & Cohen, 2001), suggests that

82    cognitive control is implemented as a top-down modulatory signal that prioritizes relevant information

83    processing in other processing areas. Specifically, it has been suggested that the human prefrontal cortex

84    sends modulatory signals to more posterior processing areas; such signals excite task-relevant

85    processing pathways and inhibit task-irrelevant processing pathways (Aben, Calderon, Van den

86    Bussche, & Verguts, 2020). Hence, the prefrontal cortex can separate information by inhibiting all

87    processing that might interfere with the current task. This approach has proven fruitful to explain human

88    behavior in cognitively demanding tasks (Abrahamse, Braem, Notebaert, & Verguts, 2016; Botvinick,

89    Braver, Barch, Carter, & Cohen, 2001; Cohen, Dunbar, & McClelland, 1990; Verbeke & Verguts,

90    2019). However, as noted above, a complete separation between task representations would be

91    inefficient. Indeed, in some cases the network might benefit from information transfer between similar

92    tasks.

4

To study the balance between sharing and separation, we consider the nature of top-down signals. Interestingly, there exist some crucial differences in the literature with respect to how top-down modulation is implemented. For instance, while some research treats the top-down signal as any other input signal and add all inputs together (Cohen et al., 1990), other research has treated the top-down signal as multiplicative (Masse et al., 2018; O'Reilly & Frank, 2006), which allows to effectively shut down (multiply by zero) activity for irrelevant neurons. Additionally, while in most research the modulatory signal is adapted to the needs of the current task (Botvinick et al., 2001; Cohen et al., 1990; Verguts & Notebaert, 2008), other work (Bouchacourt & Buschman, 2019; Masse et al., 2018) illustrated that also random, non-adaptive modulatory signals can be sufficient to allow optimal performance on complex tasks. Hence, random signals can often meet performance of learned signals, while requiring far less computational constraints. Moreover, because less parameters need to be learned, these random modulatory signals are often faster in learning to process novel inputs. More generally, random signals have proven to be useful in constructing powerful neural networks (Lillicrap, Cowden, Tweed, & Akerman, 2016; Maass, Natschläger, & Markram, 2002). Thus, top-down modulation signals differ in whether they are additive or multiplicative and whether they are adaptive or non-adaptive.

The current work provides a systematic investigation of different types of modulation signals in balancing the trade-off between shared and separated representations. Specifically, we propose four approaches for modulation. In a first approach, non-adaptive additive modulation (N+ network) is applied. Here, for each task, a different random top-down signal contributes to the activity patterns in an additive manner. Second, in adaptive additive modulation (A+ network), top-down input is also added to the network. However, in this approach, the top-down input is treated like any other task-processing input in the sense that top-down weights are susceptible to the same (backpropagation) learning rules as the regular task-processing weights. Third, in non-adaptive multiplicative modulation (Nx network), the network inhibits and/or excites a random proportion of pathways in every task context by multiplying activation with zero (inhibition) or a random positive value (excitation). Fourth, in adaptive multiplicative modulation (Ax network), the network learns which processing pathways to excite or inhibit.

121     Since previous work illustrated that the impact of shared representations depends on the nature

122     of the task environment (Musslick et al., 2017), networks are tested on three different types of input

123     (discrete low-dimensional, continuous low-dimensional, and continuous high-dimensional; see also

124     Figure 1). For each input type, we consider a number of tasks that differ in the amount of overlap of

125     their stimulus-action mappings. Interestingly, for artificial agents, there is more catastrophic

126     interference when trained in a blocked fashion. In contrast, blockwise training appears beneficial for

127     human agents (Flesch, Balaguer, Dekker, Nili, & Summerfield, 2018). To evaluate each network's

128     ability to overcome interference, we thus trained our artificial networks in a blocked fashion. In sum,

129     we test the four modulation signals on a task that requires them to optimally balance the transfer

130     (sharing) and avoidance of interference (separating) between tasks. Network performance is evaluated

131     in terms of accuracy and the ability to find the optimal amount of sharing between task representations.

132

133                                 **2.   Methods**

134     **2.1 The network**

135         Our network (Figure 1a) consists of an Input, Hidden, and Output layer. Information flows in

136     a feedforward manner from Input to Hidden to Output layer. All neurons in each layer are fully

137     connected to all neurons in the next layer. Neurons in the Input layer are divided in a Stimulus group

138     and a Task group. Activation at the Hidden layer is a combination of input from the Stimulus group and

139     a modulatory signal from the Task group. In analogy to previous work, we use

$$H = f(SW_{s,h} + TW_{t,h}) \tag{1}$$

140

141     for additive modulation (e.g., Cohen et al., 1990) and

$$H = f(SW_{s,h}) \otimes g(TW_{t,h}) \tag{2}$$

142

143     for multiplicative modulation (e.g., Masse et al., 2018). However, see the Supplementary materials for

144     additional simulations with other implementations of modulation. In these equations, *H, S* and *T* are

145     vectors representing activation in the Hidden, Stimulus and Task group respectively. A weight matrix

146 between layers is represented by $\boldsymbol{W}$. The symbol $\otimes$ represents elementwise multiplication. The

147 functions $f()$ and $g()$ represent (elementwise) nonlinear activation functions. In the main text we

148 consider simulations in which $f()$ represents a sigmoid activation function:

149

$$sig(\boldsymbol{XW_{x,j}}) = \frac{1}{1 + e^{-(XW_{x,j})}} \tag{3}$$

150

151 and $g()$ represents a RELU activation function:

152

$$RELU(\boldsymbol{XW_{x,j}}) = \max(0, (\boldsymbol{XW_{x,j}})) \tag{4}$$

153

154 In these equations, $\boldsymbol{X}$ is a (row) vector representing activity in the sending layer and $\boldsymbol{W_{x,j}}$ represents a

155 weight matrix between the sending ($x$) and receiving ($j$) layer. Thus, Hidden neurons with a negative

156 weight from the active Task neuron are gated out (activation multiplied by 0), while activation of

157 Hidden neurons with a positive weight, are multiplied by a positive value. The Supplementary materials

158 present additional simulations in which different combinations of sigmoid and RELU functions are

159 explored. Activation at the Output layer ($\boldsymbol{O}$) simply follows $\boldsymbol{O} = sig(\boldsymbol{HW_{h,o}})$. After each trial, weights

160 are adapted by the backpropagation learning rule (Rumelhart, Hinton, & Williams, 1986):

161

$$\Delta \boldsymbol{W_{x,j}} = -\alpha \times \frac{\partial E}{\partial \boldsymbol{W_{x,j}}} \tag{5}$$

162

163 in which again $x$ represents the sending layer and $j$ represents the receiving layer. The parameter $\alpha > 0$

164 represents a learning rate, and $\frac{\partial E}{\partial W_{x,j}}$ is the (partial) derivative of the error ($E$) with respect to the weights
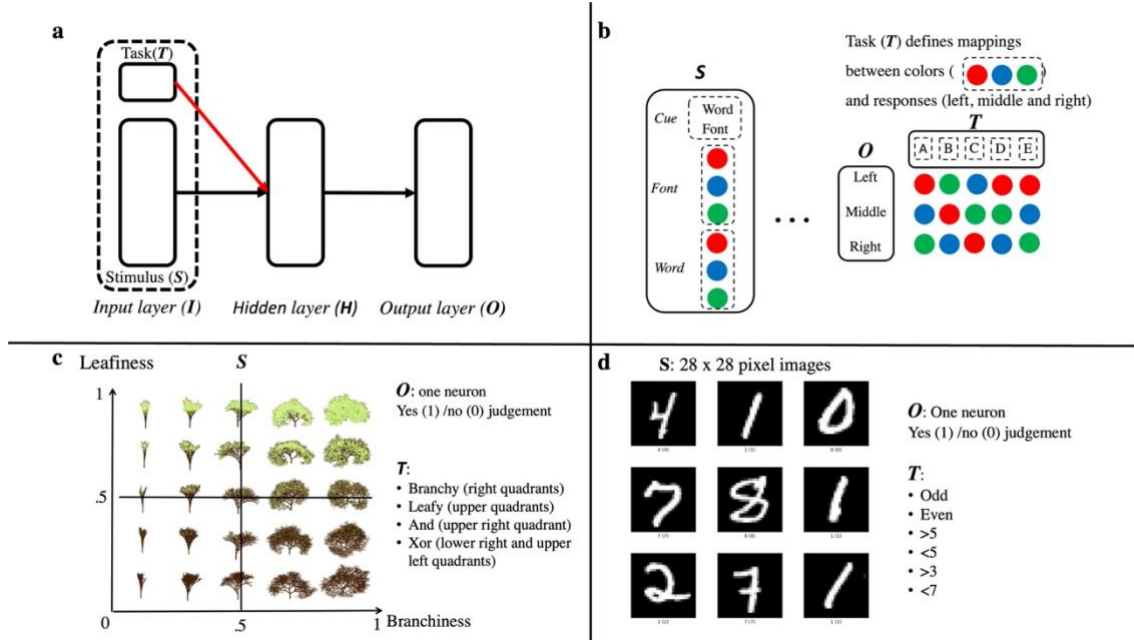
165 ($\boldsymbol{W}$).

**Figure 1. Network architecture and simulations.** *a: General network architecture.* The network consists of three layers. Information flows in a feedforward manner from Input to Hidden to Output layer. The Input layer is divided in a Stimulus and Task group. The Task group sends a modulation signal (red arrow). We evaluate four different types of modulation and test the network on three types of stimulus datasets (b-d). *b: Stroop (discrete) dataset.* Stimuli are a combination of a cue, a font (color) and a (color) word. The cue indicates which other dimension (word or font) is relevant for responding. Five tasks are defined in which the mapping between three response options and three values for color (in font and word) are changed. Mappings are represented in the panel. *c: Trees (continuous low-dimensional) dataset.* The stimulus figure is adopted from Flesch et al. (2018). Current work has coded this dataset with two neurons that can take on any value between 0 and 1. Four tasks are defined in which the network needs to give a yes or no judgement (one response neuron). *d: MNIST (continuous high-dimensional) dataset.* Here, stimuli are 28 x 28 pixel images of handwritten digits from 0 to 9. Examples are taken from https://www.tensorflow.org/datasets/catalog/mnist. We defined six possible tasks in which again the network needs to give a yes/no judgement.

## 2.2 Datasets

The network was tested on three datasets (Figure 1b-d), allowing us to evaluate several combinations of input and task type. Note that the size of the Input, Output and Task layers were adapted depending on the input dataset. The first dataset had discrete (binary) low-dimensional input patterns. Specifically, we consider the classic cognitive control Stroop task (Stroop, 1935). See Figure 1b for an illustration of the network specifics for this dataset. Here, on each trial a (color) word is presented ("red", "green" or "blue") in a particular font (color) (red, green or blue). Additionally, a cue is provided

188   telling the agent to respond either to the word or to the font dimension. The task consists in learning

189   mappings between colors (red, green and blue) and response buttons (left, middle, right). Crucially,

190   both stimulus dimensions can provide congruent evidence (e.g. "red" presented in red) or incongruent

191   evidence ("red" presented in blue). In the latter case, the correct response depends on the cue dimension

192   (respond to red when cue is word and to blue when cue is font). In terms of the network, we consider 8

193   input neurons (2 cues, 3 words and 3 colors; see also Figure 1b). Here, each stimulus consists of the

194   activation (input value = 1) of 3 (a cue, a word and a font) out of 8 Stimulus neurons, resulting in 18 (2

195   instructions × 3 fonts × 3 words) possible stimuli. Additionally, we activate one Task neuron on every

196   trial, which determines the appropriate mappings between stimuli and responses. In this task, the Output

197   layer consists of 3 neurons. On each trial, the Output neuron with the highest activation (argmax($O$)) is

198   considered to be the network response. Depending on the task, each color value (red, green or blue) was

199   mapped to one of the neurons in the Output layer. Specifically, we define five tasks (see also Figure

200   1b). Here, tasks A, B and C share no stimulus-action mappings. Task D represents a mix of tasks A, B

201   and C. Specifically, D shares exactly 1/3 of stimulus-action mappings with all three other tasks. The

202   last task E shared all stimulus-action mappings with A but activated a different neuron in the Task group

203   (in a sense, A and E are synonyms). Note that we call this a Stroop task because the input consists of

204   font and word input where one dimension was relevant; we did not mimic the imbalance between color

205   naming and word reading that appears in typical Stroop tasks.

206         The second dataset is the Trees dataset (see also Flesch et al., 2018). Specifics for this dataset

207   are presented in Figure 1c. In the Trees dataset, there are two Stimulus neurons which can take on any

208   value in a range of 0 to 1 (continuous low-dimensional input). One Stimulus neuron represents the

209   'leafiness' of a tree and the other Stimulus neuron represents the 'branchiness' of the tree. The Output

210   layer contained only one neuron. For this dataset, the network has to make a yes (output = 1) or no

211   (output = 0) judgement. We defined four different tasks for this input type. One task was to respond yes

212   to leafy trees (leafiness >.5), a second task was to respond to branchy trees (branchiness >.5), a third

213   task required the network to respond to trees that were both leafy and branchy (AND task), and the

214   fourth task consisted of responding to trees that were either leafy or branchy (but not both; XOR task).

215   Note that for this dataset there were no completely (100%) dissimilar tasks. The Leafy and Branchy

216    task share 50% of stimulus-action mappings with each other but also with the AND and XOR taks. The

217    AND and XOR task share 25% of mappings with each other.

218            The third dataset consisted of images (continuous, high-dimensional input). More specifically,

219    we used the MNIST data set (LeCun, Cortes, & Burges, 2010) which contains grey-scaled images (28

220    x 28 pixels) of handwritten digits from 0 to 9. Again, the Output layer consisted of one neuron. For this

221    input type, 6 different tasks were provided. One task was to respond to odd digits (i.e., output = 1 for

222    odd digits; output = 0 for even digits); another task required a response to even digits. A third and fourth

223    task required the network to respond to digits that were respectively larger or smaller than 5. The fifth

224    task consisted of responding to digits larger than 3 and the sixth task was to respond to digits smaller

225    than 7. This resulted in a complex pattern of overlap between the different tasks. Tasks vary from 100%

226    dissimilar (odd and even), to only 20% dissimilar (>3 and >5; <5 and <7).

227    **2.3 Simulations**

228            As described before, four versions of the network were simulated. Activation at the Hidden

229    layer follows Equation (1) for additive modulation networks (N+ and A+), and Equation (2) for

230    multiplicative modulation networks (Nx and Ax). In adaptive modulation networks (A+ and Ax), the

231    weights between the Task group and Hidden layer are learned by the backpropagation rule (Equation

232    (5)), just like the other weights. In non-adaptive modulation networks (N+ and Nx), the weights between

233    the Task group and Hidden layer are fixed at their initial (random) values. All weights are initialized

234    with a random value drawn from the normal distribution N(0, 1). Only for the Ax network, modulating

235    weights (between Task and Hidden layer) had an initial random value drawn from the uniform

236    distribution U(0, 1), such that RELU($T$) > 0 (all gates open) at the first trial. This set up provides the

237    most optimal initialization for each network. We illustrate network performance with other weight

238    initialization distributions in the Supplementary materials.

239            All four versions of the network were tested on all three data sets. Additionally, we explored

240    different learning rates ($\alpha$) and shapes of Hidden layer. Also the shapes of the Input, Output and Task

241    layers were adapted depending on the input dataset. For the Stroop and Trees input datasets, $\alpha$ took on

242    6 values ranging from 0 to 1 in steps of .2. We explored the network with one Hidden layer of either 12

243 or 24 neurons. For the MNIST dataset we used lower learning rates. Here, α took on 6 values ranging

244 from 0 to .1 in steps of .02. For this data set, we explored performance with one Hidden layer of 400

245 neurons; and also with two Hidden layers (300 and 100 neurons respectively) and three Hidden layers

246 (200, 100 and 100 neurons respectively). Note that for this dataset, the total number of Hidden neurons

247 did not differ between architectures. Activation at the first Hidden layer ($H_1$) followed Equation (1) or

248 (2) for additive or multiplicative networks respectively. In standard simulations, activation at the second

249 and third Hidden layer followed: $H_i = sig(H_{i-1}W_{Hi-1,Hi})$, in which $i$ is the index of the Hidden layer.

250 Hence, the Task modulation signal was not sent directly to the deeper Hidden layer(s). However, for

251 completeness we also explored network performance (with two hidden layers) when the Task signal

252 was sent to only the second hidden layer, to both hidden layers or to none of the hidden layers. Results

253 of these simulations are presented in section 3.5.

254 For every combination of α and shape of the Hidden layer, 25 simulations ($N = 25$) were

255 performed for each dataset. For each simulation, 1200 or 12000 inputs were randomly sampled for the

256 Trees and MNIST datasets respectively. Since there were only 18 stimuli (input patterns) for the Stroop

257 dataset, we chose to repeat these 18 stimuli 75 times in each simulation, resulting in 1350 trials.

258 Additionally, we randomly shuffled the order of tasks before a simulation. In a next step, we divided

259 the sampled input patterns over 3 repetitions (450, 400 or 4000 trials per block for the Stroop, Trees

260 and MNIST dataset respectively). In every repetition, the network was trained (training phase)

261 blockwise on each task, using the predetermined input sample and order of tasks. Thus, each task was

262 repeated 3 times in a blocked fashion. At the end of the third block, weights were frozen, and the

263 network was tested (test phase). For this test phase a new order of tasks was generated. Each task was

264 tested for one block of trials. For this purpose, 100 or 500 new inputs were randomly sampled from the

265 Trees and MNIST datasets respectively. For the Stroop dataset, no new inputs could be generated so

266 we repeated the 18 possible inputs 5 times, resulting in 90 trials.

267 **2.4 Analyses**

268 *2.4.1 Accuracy*

269    To investigate whether networks suffered from catastrophic interference during learning, we

270    computed accuracy for each task repetition (averaged over all tasks). Networks that suffer from

271    catastrophic interference would need to relearn a task on every repetition because they would learn

272    other tasks in between. Hence, such a network would not improve over task repetitions.

273    Next, we investigated the network's ability to balance separating representations with sharing

274    representations. More specifically, we computed accuracy for each task during the test phase. For this

275    analysis we mainly focus on the Stroop dataset but we present results for the other datasets as well. The

276    Stroop dataset is optimally suited for this analysis since there is a larger variation in (dis)similarities

277    between tasks (see also the objective dissimilarity table in Figure 2) than is the case for the other

278    datasets. More specifically, five tasks were proposed for the Stroop dataset. As described in section 2.2,

279    three of them (A, B and C) did not share any stimulus-action mappings and thus can be totally separated.

280    Tasks A and E however, share all stimulus-action mappings and can be fully shared. Additionally, task

281    D shares 1/3 of its stimulus-action mappings with all other tasks. On the one hand, a full sharing of task

282    representations would allow the network to exploit the shared mappings between A, D and E but lead

283    to catastrophic interference in tasks B and C, illustrated by a strong decrease of accuracy in tasks B and

284    C compared to A, D and E. On the other hand, a full separation of tasks would improve accuracy of

285    tasks B and C (by less interference), but would also eliminate the advantage of the full overlap between

286    A and E and the partial overlap between D and the other tasks. Hence, accuracy would be the same for

287    all tasks. Importantly, when the network is only able to fully share or separate information it would

288    benefit from the full overlap between A and E but would not benefit from the partial overlap between

289    task D and the other tasks. In sum, an optimal network would find a balance between sharing and

290    separating, resulting in an improved accuracy for tasks A, D and E while minimizing the dip in accuracy

291    for tasks B and C.

292    To evaluate overall performance of the networks we also computed accuracy during the test

293    phase for all learning rates and all tasks.

294    *2.4.2 Representational dissimilarity*

295    In order to analyze to what extent the networks shared or separated stimulus-action mappings

296    across tasks, we computed how dissimilarity between tasks in terms of stimulus-action mappings was

297    represented in the network. This analysis considers several steps. An overview of these steps in the

298    context of the Stroop dataset is provided in Figure 2.

299          In a first step, we computed for each simulation the objective dissimilarity between stimulus-

300    action mappings across tasks. Specifically, we computed a matrix where rows and columns represent

301    the tasks, and each cell contains the proportion of stimuli that were matched with a different action

302    across the two respective tasks (row and column).

303          A second step was to compute the representational dissimilarity within the network. For this

304    purpose, we first computed the mean activation at Hidden layer for each stimulus (18 Stroop stimuli, 4

305    quadrants of branchy-leafy space and 10 digits in MNIST dataset) in each task across trials. Then the

306    difference between task representations was extracted by computing the mean Euclidean distance for

307    two tasks $T1$ and $T2$. Here,

308

$$Dissimilarity_{T1,T2} = \sum_{S}^{nStim} \left\| \boldsymbol{H}_{T1}^{S} - \boldsymbol{H}_{T2}^{S} \right\| / nStim \qquad (6)$$

309

310    in which $\boldsymbol{H}_{T1}^{S}$ and $\boldsymbol{H}_{T2}^{S}$ are vectors of length $nHidden$, representing the average activity for all Hidden

311    neurons when stimulus $S$ was presented to the network. Hence, we compute the Euclidean distance

312    (indicated by $\|\cdot\|$) for each stimulus ($S$) and each task pair ($T1$, $T2$). This distance is then averaged

313    over all possible stimuli ($nStim$) to obtain one dissimilarity matrix of Hidden representations between

314    tasks.

315          In a third and last step we compared the objective dissimilarity to the representational

316    dissimilarity. Specifically, we reshaped both matrices to vectors and computed the Spearman rank

317    correlation coefficient between these vectors. This resulted in one value of the dissimilarity correlation

318    between objective task dissimilarity and a network's representational task dissimilarity.
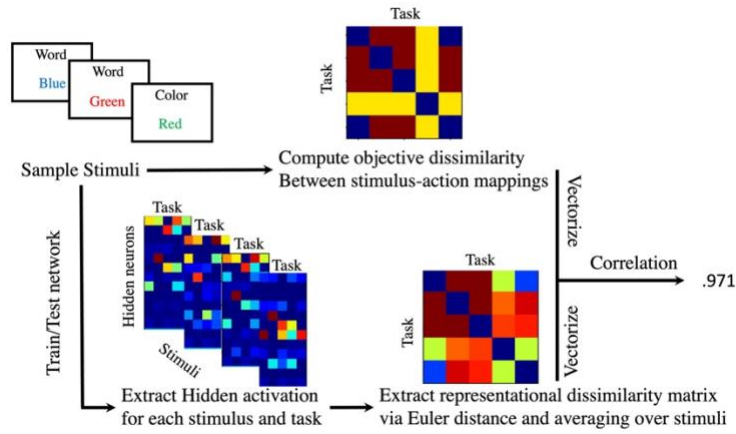
**Figure 2. Methods.** Illustration of the different steps in the representational dissimilarity analyses. Examples are shown for one simulation of the Stroop dataset with a learning rate of .6 and 12 Hidden neurons.

*2.4.3 Neural activation analyses*

We performed two additional analyses to gain more insight into how the different modulatory signals organize Hidden layer activity. For this purpose, we again computed for each task the mean activation at Hidden layer for each stimulus, resulting in a matrix with size (*nStim*, *nTask*, *nHidden*). First, we investigated the distribution of activation for all stimuli and tasks across the Hidden neurons. Second, in order to visualize the network representations for each task, we reduced Hidden layer dimensionality via principal component analysis. For this purpose, we entered the activation matrix with size (*nStim*, *nTask*, *nHidden*) into the principal component analysis and approximated it by a matrix of size (*nStim*, *nTask*, 2). As a result, we could plot the representation of each stimulus in each task in a two-dimensional space.

## 3. Results

### 3.1 Accuracy

First, we evaluated the networks' ability to separate task representations. For this purpose, we investigated whether average accuracy (during the training phase) increased over task repetitions. Networks that do not separate task sets suffer from catastrophic interference because they overwrite mappings of one task by the mappings of another task. As a result, such networks need to relearn the original task when it is presented again, and do not show any improvement over task repetitions. In Figure 3, it is observed that accuracy hardly improves for the additive networks (A+ and N+). Thus, these networks severely suffered from catastrophic interference. In contrast, for multiplicative

modulation networks, in particular for the Ax network, there was a significant improvement over task

repetitions. Thus, multiplicative modulation seems more efficient in separating task representations,

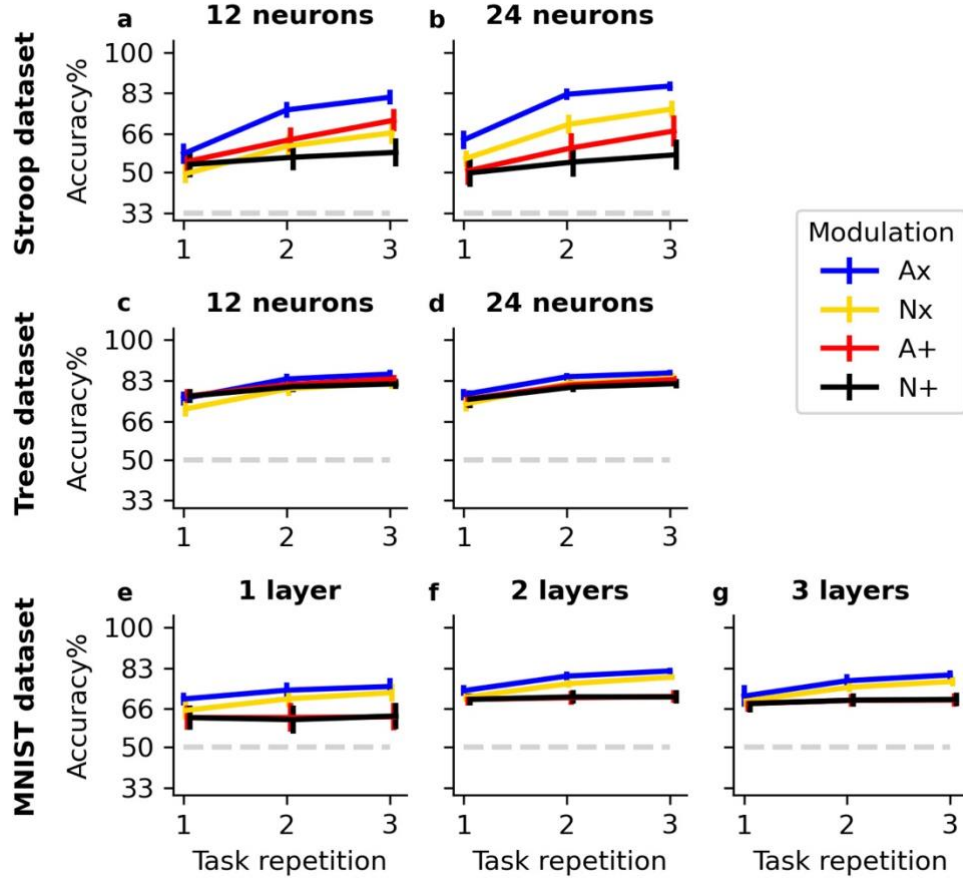rendering them less vulnerable to catastrophic interference during learning.



**Figure 3. Accuracy in training phase per task repetition.** Lines illustrate mean accuracy for each task
repetition during the training phase averaged across all learning rates ($\alpha$), all tasks and all simulations. Bars
indicate 95% confidence intervals over 25 simulations. The dashed lightgrey line indicates chance level
accuracy. Results are shown for different datasets (rows) and different shapes of Hidden layer (columns).

Second, we zoomed in on accuracy for each task during the test phase. Here, we focus mainly

on the Stroop dataset (see section 2.2) because the Stroop dataset has a broader range of dissimilarities

between tasks. Specifically, tasks A, B and C have completely dissimilar mappings, task D has a partial

overlap of 1/3 with all other tasks and task E shares all mappings with task A. An optimal network

would find a balance between sharing and separating, resulting in an improved accuracy for tasks A, D

and E while minimizing the dip in accuracy for tasks B and C (see also section 2.4.1). In Figure 4a,b

we observe a strong dip in accuracy for tasks B and C when the modulation signal was additive. This

suggests that additive modulating signals are well suited for sharing task representations, but less so for separating task representations. This dip in accuracy for tasks B and C is less strong for the multiplicative modulation networks. Importantly, the Nx network has an approximately equal accuracy for all tasks (Figure 4a,b). This suggests that the Nx network has strongly separated task representations, which did not allow that network to benefit from overlap between task mappings across the different tasks. The Ax network is clearly the network that was able to optimally balance the separation and sharing of task representations, showing an advantage in accuracy for all tasks compared to the other networks, and only a small dip for tasks A and B. Note that although task D only had a partial overlap with other tasks, accuracy is equally high as for tasks A and E which fully overlapped. Hence, the Ax network does not treat sharing or separation as an all-or-none process, but also captures partial overlap.
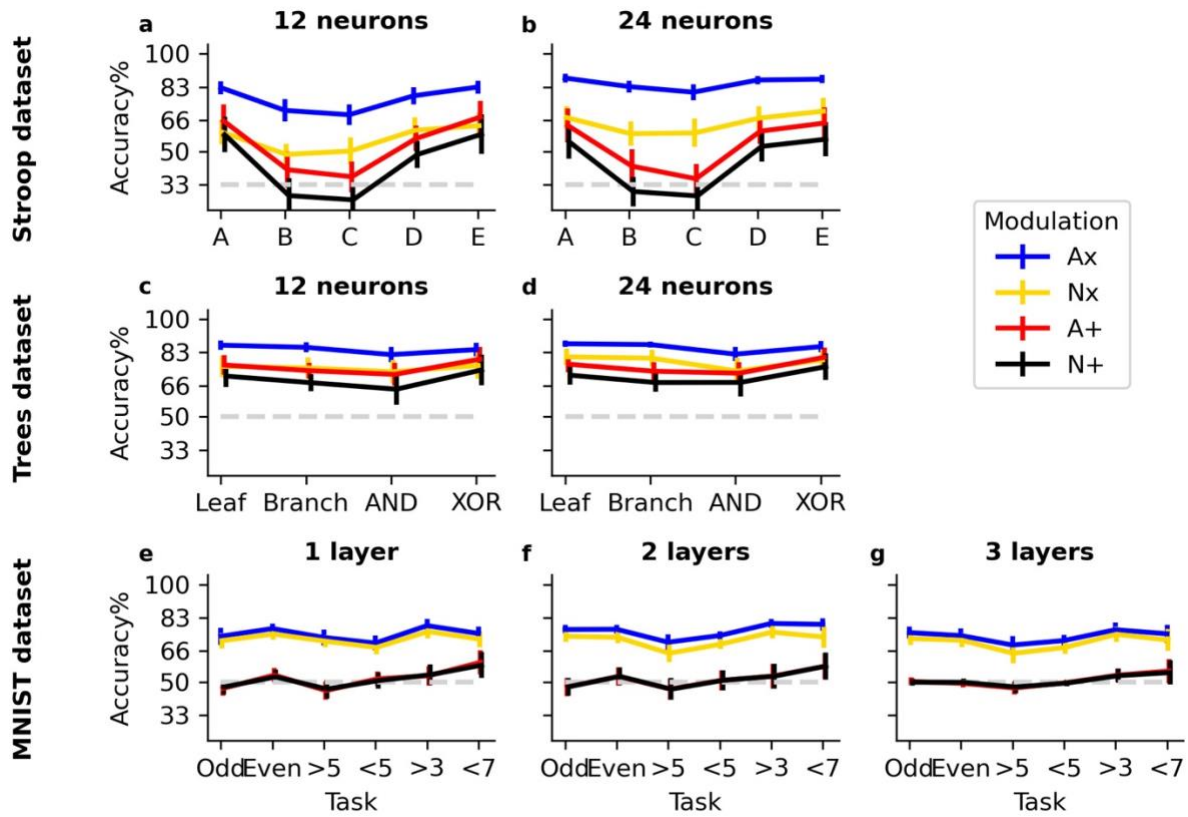


**Figure 4. Accuracy in test phase per task.** Lines illustrate mean accuracy during the test phase for each task, averaged across all learning rates ($\alpha$) and all simulations. Bars indicate 95% confidence intervals over 25 simulations. The dashed lightgrey line indicates chance level accuracy. Results are shown for different datasets (rows) and different shapes of Hidden layer (columns).

371       We also show results for the other datasets but these are less informative in this respect, because there is no strong variability in objective task dissimilarity (see section 2.2). For the Trees dataset (Figure 4 c,d) there was an average objective dissimilarity of 50% for the leafy and branchy task, and an average dissimilarity of 41.5% for the AND and XOR tasks. For the MNIST dataset (Figure 4e-g) there was a strong variability of dissimilarity between tasks themselves (20-100%), but when averaged over all tasks, the range of dissimilarity between one task and all other tasks was rather small (between 56 and 64%). Despite the absence of variability in overall dissimilarity between tasks, it is also clear for the Trees (Figure 4c, d) and MNIST (Figure 4e, f) dataset that the Ax network performs best.

380      This suggestion that the Ax network is best able to find a balance between sharing and separating task representations, is supported by the fact that this network reaches a higher accuracy overall when we analyze accuracy over all tasks during the test phase. In Figure 5, it is observed that the Ax network outperforms all other networks for all learning rates, all shapes of Hidden layer and all datasets. The non-adaptive additive (N+) network performs worst. The non-adaptive multiplicative modulation (Nx) network and the adaptive additive (A+) network perform in between these other networks. Here, the Nx modulation network seems to obtain an advantage over the A+ network when there are more Hidden neurons (Figure 5b,d) and for high-dimensional (MNIST; Figure 5e,f) input datasets. In sum, multiplicative modulation outperforms additive modulation, and adaptive modulation outperforms non-adaptive modulation.
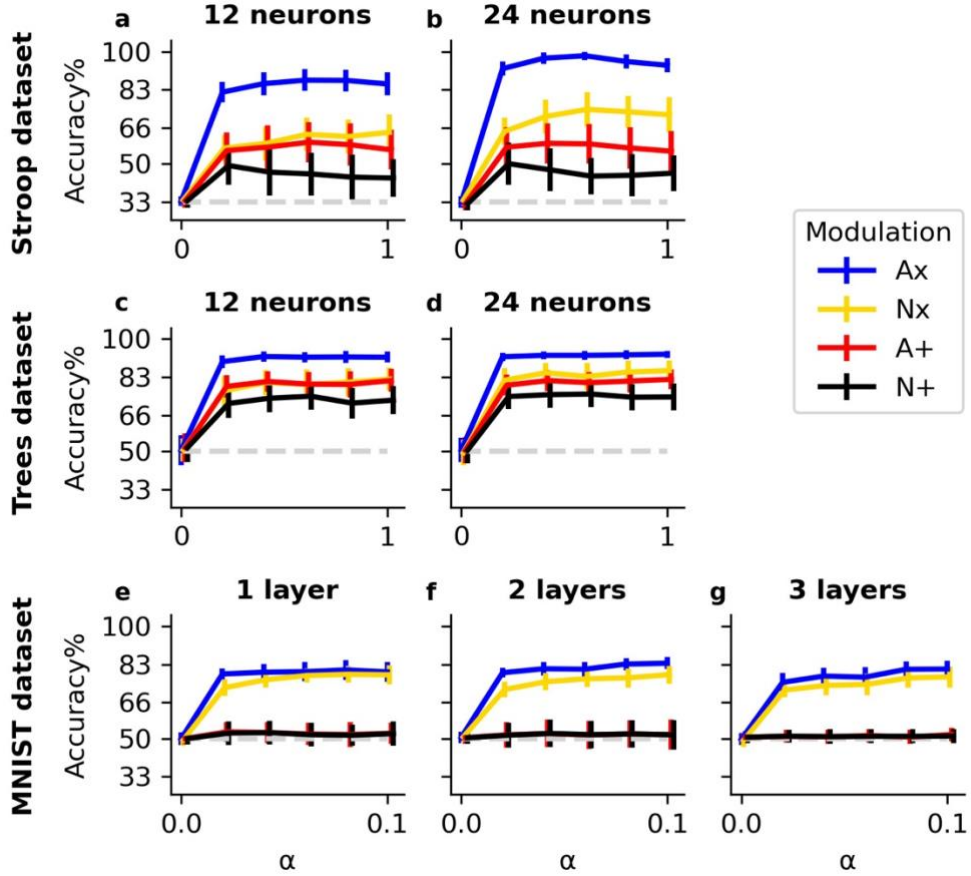
**Figure 5. Accuracy in test phase per learning rate (α).** Lines illustrate mean accuracy during the test phase for each value of α, averaged across all tasks and all simulations. Bars indicate 95% confidence intervals over 25 simulations. The dashed lightgrey line indicates chance level accuracy. Results are shown for different datasets (rows) and different shapes of Hidden layer (columns).

**3.2 Representational dissimilarity**

We next investigate whether objective dissimilarity of stimulus-action mappings between tasks was represented in the network. For this purpose, we computed for each network simulation a representational dissimilarity matrix and correlated this matrix elementwise with an objective dissimilarity matrix of stimulus-action mappings between tasks (see section 2.4.2 and Figure 2 for details). Results of this analysis are shown in Figure 6. As was already suggested by the accuracy results (section 3.1), the Ax network was clearly better at extracting the objective dissimilarity between tasks. Interestingly, the A+ network was also (although less strongly) able to capture the objective overlap between tasks for the Stroop and Trees input datasets (Figure 6a-d), but not for the MNIST dataset (Figure 6e-g).

405   For the MNIST dataset, we simulated a network with one Hidden layer (Figure 6e), one with

406   2 Hidden layers (Figure 6f) and one with three Hidden layers (Figure 6g). Overall, there does not seem

407   to be a strong benefit for dividing the (same number of) Hidden neurons over multiple layers in the

408   current setup. However, only the first Hidden layer received a modulation signal. In section 3.5 we

409   present results for simulations that modulated deeper and/or more hidden layers.
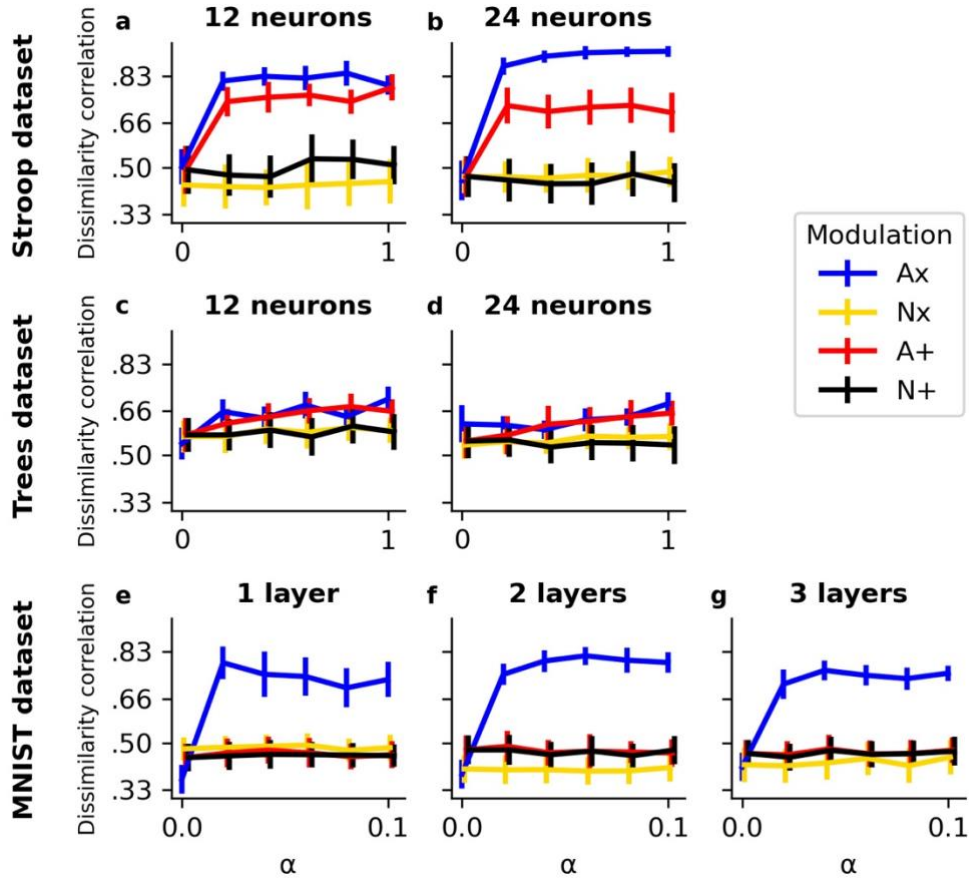


410

411   **Figure 6. Correlation of objective and representational dissimilarity between tasks.** Lines illustrate the

412   mean task dissimilarity correlation for each value of α across all simulations. Bars indicate 95% confidence

413   intervals over 25 simulations. Results are shown for different datasets (rows) and different shapes of Hidden

414   layer (columns).

### 3.3 Neural activation analysis

416   To provide additional insight into how the different modulation signals organize Hidden layer

417   activity, Figure 7 shows the distribution of activation at the Hidden layer for all networks and datasets.

418   This is shown for the networks with 12 Hidden neurons (Stroop and Trees dataset) or 1 Hidden layer

419   (MNIST dataset). Here, it is observed that activation distributions are strongly bimodal with peaks

420　around 0 and 1. Note that, because of the RELU modulatory signal (see Equations (2) and (4)),

421　activation in the multiplicative modulation networks were theoretically not bound to 1. Nevertheless,

422　also these multiplicative modulation networks show a clear activation bound of 1 after learning.

423　Interestingly, the multiplicative modulation networks illustrate a strong asymmetrical distribution with

424　a higher peak of activity around zero. Especially the Nx network has a high zero-centered peak. This

425　suggests that the Nx network is learning more sparse representations and potentially creates different

426　groups or modules of neurons where each module learns (part of) one task. Hence, in line with what

427　was described before, the Nx network is well suited for separating task representations. The activity

428　distribution of the additive networks is clearly more symmetrical with a higher number of neurons that

429　exhibit strong activity for each stimulus and task. As a result, the additive networks will probably share

430　more neurons for representing stimuli and/or tasks. In line with what we described before, the Ax

431　network illustrates a mixture between the properties of the Nx and additive networks and is therefore

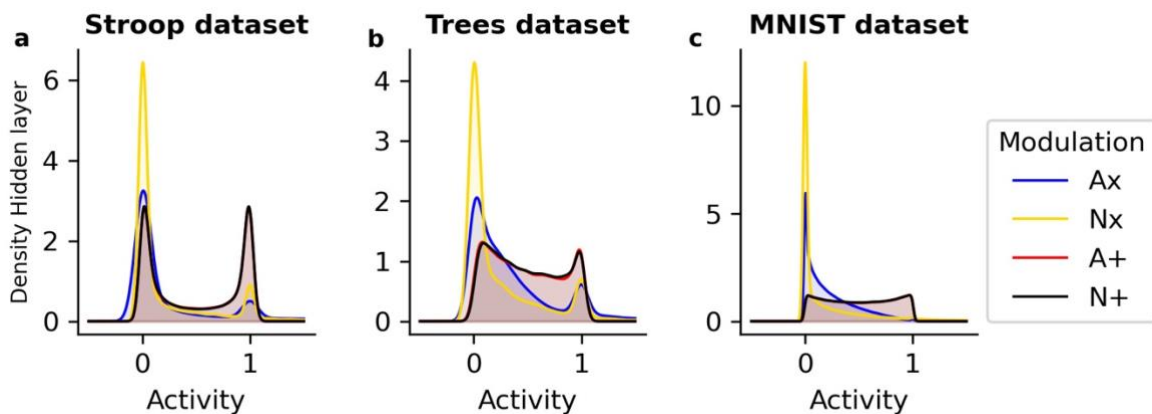432　optimally suited to balance shared and separated representations.

433



434　**Figure 7. Distribution of activity at Hidden layer.** The distribution of activation at the Hidden layer of the
435　different modulation networks is shown for each dataset.

436　　　For the next analyses, we reduced Hidden layer dimensionality into two principal components

437　with the highest eigenvalues. For the Stroop dataset these components explained on average over all

438　simulations, networks, learning rates and shapes of Hidden layer, 54.34% of the variance (SD = 6.41%).

439　For the Trees dataset, two components explained 71.29% of the variance (SD = 8.16%) and for the

440　MNIST dataset, two components explained 44.95% of the variance (SD = 4.68%). In Figure 8, we show

441　neural representations in the Hidden layer for each stimulus (dots) and each task (colors) on 2

442   dimensions. Notice that this analysis does not allow us to average over simulations. Hence, results are

443   shown for one representative simulation of the network with an intermediate learning rate of $\alpha = .6$ for

444   the Stroop and Trees dataset and $\alpha = .06$ for the MNIST dataset.

445         Generally, results are in line with our previous findings in accuracy and representational

446   dissimilarity analyses. For the Stroop dataset, we observe that the additive modulation networks show

447   a tendency to share neural representations across tasks (Figure 8g,j). In contrast, the Nx network (Figure

448   8d) effectively separates the different tasks but fails to share tasks A and E (blue and green dots) which

449   have no dissimilarities in their stimulus-action mappings. The Ax network however (Figure 8a), shows

450   a remarkable ability in discovering the overall relational structure between tasks. The network finds

451   three orthogonal axes for the three orthogonal tasks A, B and C. Task D which shares stimulus-action

452   mappings with all previous tasks, is placed in between (at the origin of the three axes) the

453   representations of A, B and C. Additionally, the network discovered that task E fully overlaps with task

454   A.

455         For the Trees dataset (Figure 8b,e,h,k), all networks were able to separate task representations.

456   Thus, in contrast to the Stroop input datasets, the additive networks were able to separate task

457   representations for the Trees dataset. This explains why the difference in accuracy between networks

458   was much smaller for the Trees dataset in comparison to other datasets (Figure 5). Note that in this

459   dataset the inputs were significantly less complex (only 2 dimensions) than for the other datasets (18 or

460   784 ($28^2$) dimensions for the Stroop and MNIST dataset respectively).

461         For the MNIST dataset, the additive networks again fail to separate task representations

462   (Figure 8i,l). Notice that the networks extracted separate representations for the 10 digits (0-9) but

463   shared the digit representations across all tasks. The Nx network (Figure 8f) was able to separate task

464   representations but did not extract a clear relational structure. Again, only the Ax network (Figure 8c)

465   was able to extract the full relational structure of the tasks. Here, two dimensions were extracted by the

466   network. One dimension (Dimension 2) was used to separate the odd from the even tasks, the other

467   dimension (Dimension 1) was used to separate the larger than (>5 and >3) tasks from the smaller than
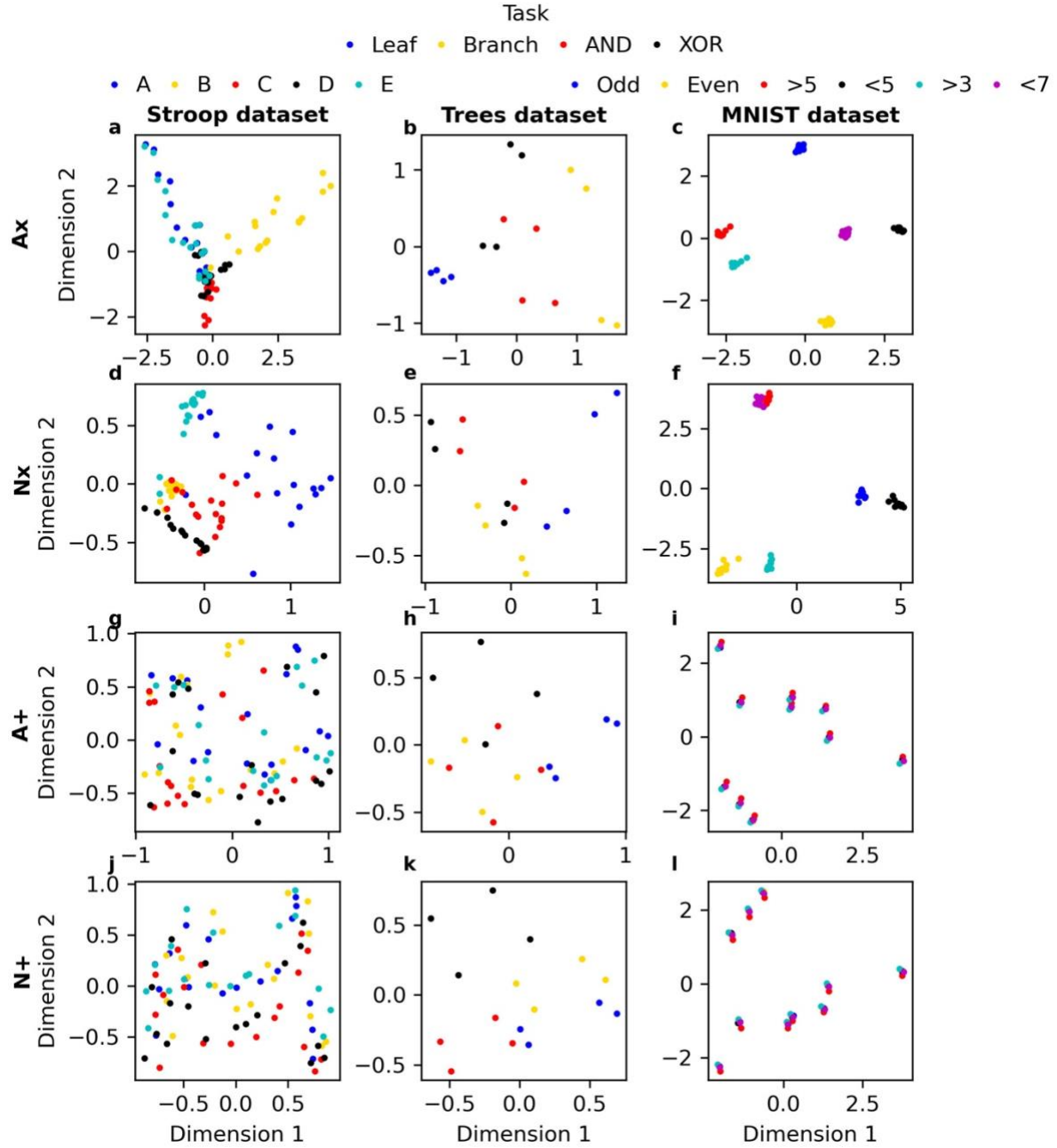
468   (<5 and <7) tasks.

**Figure 8. Task representations after principal component analysis.** The neural representation for each stimulus (dots) and each task (colors) are shown along the first two principal components. This is shown for a representative simulation and an intermediate learning rate, for all modulation networks (columns).

### 3.4. Multilayer networks

To provide more insight into the deeper (more than one Hidden layer) networks, we provide results of the representational dissimilarity analyses in each layer separately (Figure 9a-e). This is shown for the two and three Hidden layer networks that were tested on the MNIST dataset. Remarkably, while the modulation signal is only delivered to the first Hidden layer (Figure 9a,c), the other (second

478 and third but not the first) Hidden layer(s) represent the dissimilarity between tasks better (Figure

479 9b,d,e).

480        Intriguingly, we observe in Figure 9f-j that, in terms of activation, the differences between the

481 modulation networks are more pronounced at the first layer than at the second or third layer. This

482 contrasts with the previous result (Figure 9a-e) that task dissimilarity correlations are more pronounced

483 in the second and/or third Hidden layer(s). However, it is important to keep in mind that separating task

484 representations does not necessarily lead to higher dissimilarity correlations as the different tasks also

485 illustrate significant similarities. This is also emphasized by the fact that the Nx network clearly

486 separates all tasks but does not show a strong dissimilarity correlation. Although it deserves further

487 investigation, it might well be that task mappings are maximally separated in the first Hidden layer and

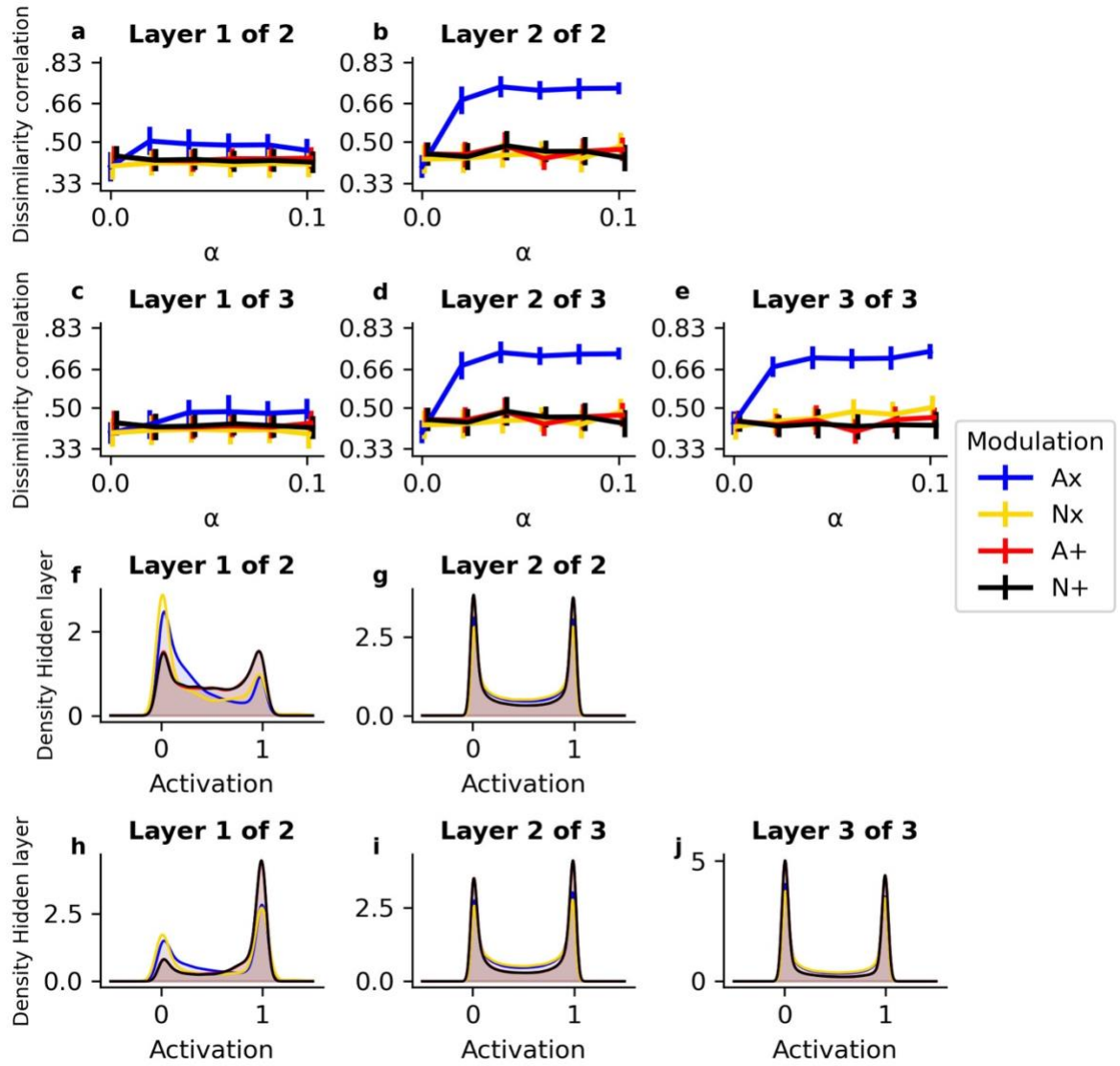488 then (compositionally) recombined in the deeper layers.

**Figure 9. Hidden layer activity for multilayer networks.** The upper two rows (panels a-e) illustrate the mean task dissimilarity correlation for each value of α across all simulations. Bars indicate 95% confidence intervals across 25 simulations. This is shown separately for each Hidden layer of the two Hidden layer network in (a-b) and for the three Hidden layer network (c-e). The lower two rows (f-j) illustrate the distribution of activation at each Hidden layer of the two-layer network (f-g) and the three-layer network (h-j).

## 3.5. The location of modulation

To gain insight in how network performance is influenced by the location of modulation we performed additional simulations of the two Hidden layer network on the MNIST dataset. Here, we explored performance when the network received no modulatory input from the Task layer, when modulation was applied at the first Hidden layer (as before), at the second Hidden layer, or at both

Hidden layers. As can be observed in Figure 10a, all networks perform at chance level when no modulation is applied. When modulation is applied at one or more Hidden layers, the multiplicative modulating networks, and in particular the Ax network, outperforms the additive networks. Although network performance seems more reliable (narrow confidence intervals) with modulation at deeper and/or more Hidden layers, the increase in mean accuracy is very small.
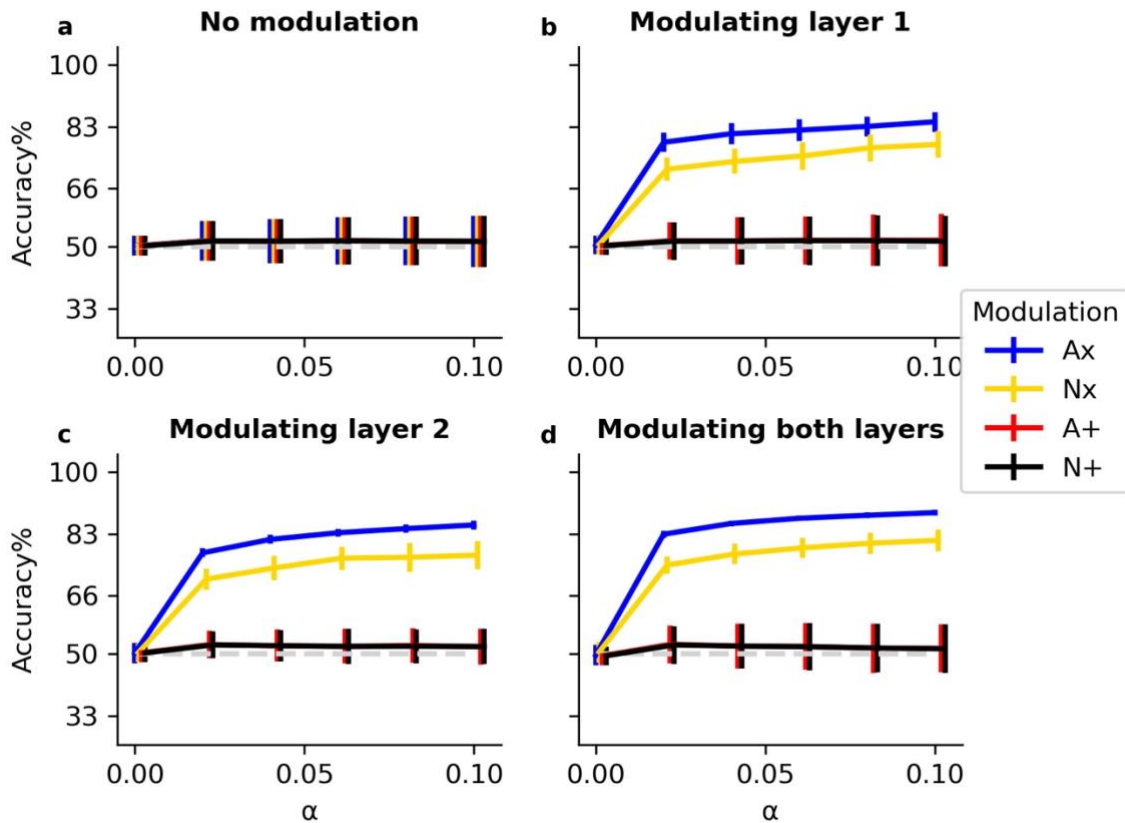


**Figure 10. Accuracy in test phase for different locations of modulation.** Lines illustrate mean accuracy during the test phase for each value of $\alpha$, averaged across all tasks and all simulations. Bars indicate 95% confidence intervals over 25 simulations. The dashed lightgrey line indicates chance level accuracy. Results are shown for different datasets (rows) and different shapes of Hidden layer (columns).

Also the results of task representational dissimilarity correlations seem highly similar for all locations of modulation. Note that this analysis could not be performed for the network that used no modulation. Since there was no task modulation, task representations were completely similar and the representational dissimilarity matrices for these networks were constant. Hence, no correlation could be computed with the objective dissimilarity matrix.

516        In sum, there does not seem to be a significant difference in network performance depending

517    on where modulation is applied. In this specific case, modulation at layer 2 could be considered as most

518    optimal since it requires to only learn 100 weights from each Task neuron to that second layer, compared

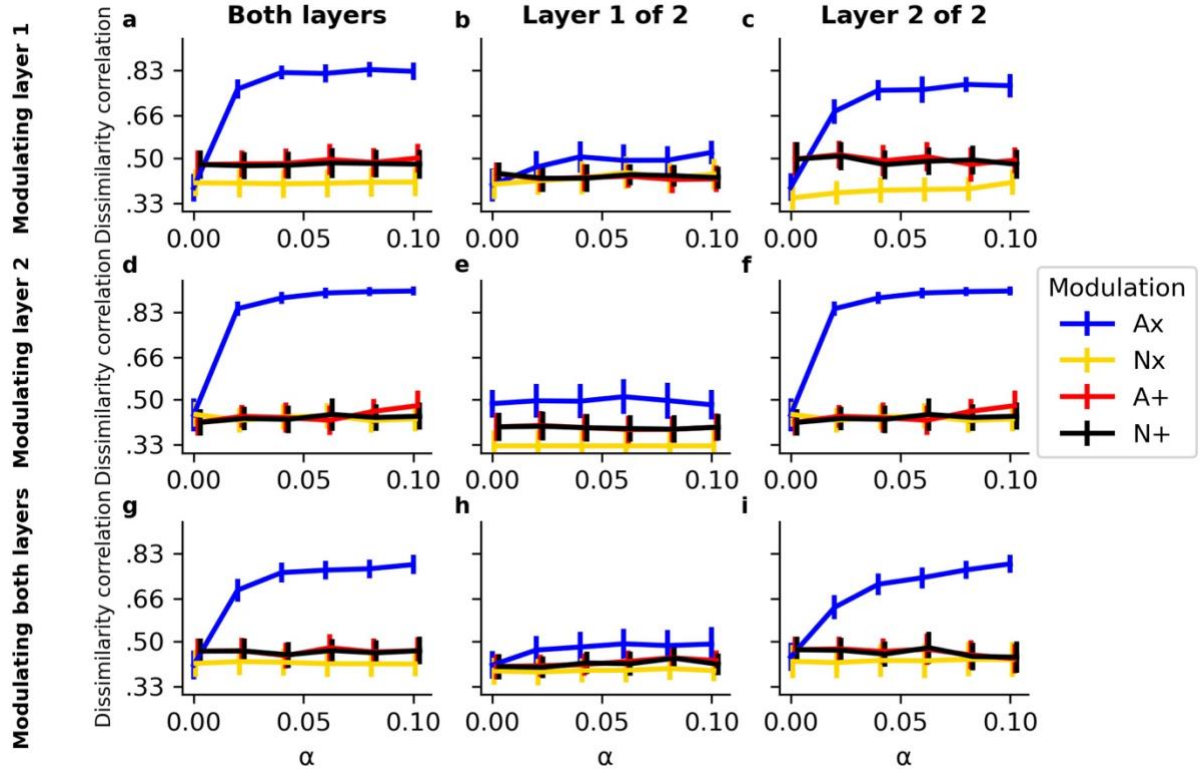519    to 300 for the first layer, and 400 for both layers.



520

521    **Figure 11. Correlation of objective and representational dissimilarity depending on location of**

522    **modulation.** Lines illustrate the mean task dissimilarity correlation for each value of α across all

523    simulations. Bars indicate 95% confidence intervals over 25 simulations. Results are shown for different

524    locations of modulation (rows) and different Hidden layer(s) (columns).

525                       **4.  Discussion**

526        Current work investigated how neural networks can optimally balance the trade-off between

527    avoiding interference via separating task representations and generalizing information via sharing

528    representations. For this purpose, we identified and systematically investigated two crucial features of

529    modulation signals. First, the modulation signals can be additive or multiplicative. Multiplicative

530    signals were better suited for separating task representations. The multiplicative networks were less

531    vulnerable to catastrophic interference than the additive networks. Second, the modulation signals could

532    be adaptive (learned) or non-adaptive (random). Adaptive modulation signals provided a clear

533 advantage over non-adaptive modulation signals in terms of both accuracy and balancing

534 representations. Hence, the adaptive multiplicative (Ax) network was able to optimally balance the

535 trade-off between sharing and separating task representations. This Ax network can avoid interference

536 but also generalize across tasks which resulted in an overall better accuracy compared to the other

537 networks.

538        Crucially, multiplicative signals modulated task-specific input more strongly. A Hidden

539 neuron that receives input from many bottom-up Stimulus neurons, needs a strong (negative) additive

540 modulation signal in order to be inhibited. In contrast, our multiplicative signal followed a RELU

541 activation function (Equation (4); see Supplementary materials for an investigation of activation

542 functions), which means that a small negative weight was sufficient to shut down (multiply by zero) a

543 Hidden neuron activation. As a result, multiplicative signals developed sparser representations (Figure

544 7) which is optimal for separating task representations. This advantage was especially present when

545 there were many Stimulus neurons (i.e., the MNIST task). When there were only 2 Stimulus neurons,

546 as in the Trees task (see Figure 1c), the additive network was also able to separate task representations

547 (see Figure 8h,k). Thus, multiplicative modulation is more efficient than additive signals, especially for

548 separating high-dimensional inputs. Specifically, the Ax network warped the representational space in

549 order to effectively organize tasks that obey similar mappings as well as tasks that obey dissimilar

550 mappings within one neural architecture. In this representational space, dissimilar tasks were placed at

551 the edges of a regular grid (see Figure 8c) and tasks that were similar were placed closer together. Such

552 a geometrical organization of task rules is optimally suited for generalizing task rules (Bernardi et al.,

553 2020; Kim, Pitt, & Myung, 2013). Consistent with the current analysis, Kim et al., (2013) demonstrated

554 how backpropagation shapes hidden space to accommodate quasi-regularities in language processing,

555 thus to accommodate both regular and irregular stimuli (e.g., orthography-phonology mappings). They

556 demonstrated that after training by backpropagation, both regular and irregular stimuli could be placed

557 in hidden neuron space at the edges of a slightly deformed grid; sufficiently grid-like to process the

558 regular stimuli (and profit from generalization), but sufficiently deformed to cope with irregular

559 mappings as well. Moreover, previous work has illustrated a similar systematicity of task

representations in neural activation of prefrontal areas and hippocampus of monkeys (Bernardi et al., 2020).

An extensive amount of work describes how humans share or separate representations by extracting latent states in the environment (Collins & Frank, 2016; Franklin & Frank, 2018; Gershman & Niv, 2012; Wilson, Takahashi, Schoenbaum, & Niv, 2014; Yu, Wilson, & Nassar, 2020). Here, the agent decides on every new experience whether to categorize it as belonging to a new state or as belonging to a state that it has experienced before. Each latent state would then develop its own representations. An important disadvantage of the latent state approach is that it uses a dichotomous decision on whether an object belongs to the state or not. Such an approach is less suited to capture partial overlap between tasks. In the example of the Stroop dataset, a latent state approach would correctly assign tasks A, B and C to three different states because the mappings are completely dissimilar. The latent state approach would also correctly assign task E to the same latent state as A because they fully share the stimulus-action mappings. However, task D shares 1/3 of the mappings with all four other tasks. In this case, D would be optimally handled as a combination of the other mappings that are already learned. This is problematic for a latent state approach since it can only decide to assign D to a new latent state or to one of the previous ones.

To accommodate this limitation of the latent state approach, previous work has proposed compositionality (Fidler et al., 2009; Franklin & Frank, 2018; Lake et al., 2014; Sugita et al., 2011; Tubiana & Monasson, 2017; Yang et al., 2019). Specifically, the latent state approach could allow mixed overlap between tasks by representing a task as multiple states, each representing a subset of mappings (Franklin & Frank, 2020; Griffiths & Ghahramani, 2011). Nevertheless, this approach could significantly increase the number of possible states, which can be problematic in very complex task environments. This raises the question in how many states/dimensions the agent should cluster its experiences in order to optimally balance generalization and interference (Badre, Bhandari, Keglovits, & Kikumoto, 2021). As we have shown in Figure 4a,b, the current Ax network could benefit equally from the partial overlap in D and the full overlap between A and E without the need of extracting latent states. This is consistent with previous work, showing that multiplicative network interactions lead to useful compositional task representations (İrsoy & Cardie, 2015; Sugita et al., 2011).

588    Multiplicative modulation is also sometimes called gating (Masse et al., 2018; O'Reilly &

589    Frank, 2006; Rougier, Noelle, Braver, Cohen, & O'Reilly, 2005). A crucial question that remains is

590    how multiplicative signals are mechanistically implemented in the human brain. In this respect, we

591    point to recent work that described an important role for neural oscillations in organizing functional

592    networks. For example, it has been proposed that neural oscillations at alpha frequency (8-12 Hz) reflect

593    gating by inhibition (Jensen & Mazaheri, 2010). Here, GABAergic inhibition provided by the task-

594    irrelevant areas would be reflected by stronger alpha activity in those areas. Another oscillatory

595    frequency that is known to organize functional networks in the brain is the theta frequency (4-8 Hz).

596    More specifically, recent theoretical and empirical work (Helfrich & Knight, 2016; Lisman & Jensen,

597    2013; Verbeke et al., 2021; Verbeke & Verguts, 2019; Verguts, 2017) has proposed that prefrontal theta

598    activity functions to (de)synchronize gamma (>40 Hz) activity in posterior processing areas. Here,

599    synchronization leads to effective communication (gates open) between processing areas while

600    desynchronization eliminates effective communication (gates closed) between processing areas (Fries,

601    2005, 2015). Thus, previous work has described how oscillatory interactions within and between

602    different frequency bands (theta, alpha and gamma) reflect gating processes that could biologically

603    implement the multiplicative modulation signals that were used in the current networks.

604    Previous work has also described how neurotransmitters such as dopamine and noradrenaline

605    can modulate neural activation in a way that mimics multiplicative modulation (O'Reilly & Frank,

606    2006; Servan-Schreiber, Printz, & Cohen, 1990). Moreover, current work observed that additive

607    modulation can reach similar performance as multiplicative modulation when stimulus dimensionality

608    was low (Trees dataset). Hence, combining additive modulation with weight regularization between

609    Stimulus and Hidden layer but not between Task and Hidden layer could allow the additive modulation

610    networks to overcome high-dimensional inputs and reach similar performance as the multiplicative

611    modulation networks. Future work should further explore biologically plausible implementations of

612    multiplicative modulation.

613    In analogy to previous work (e.g., Cohen et al., 1990), the current networks used a low-

614    dimensional Task layer which sent modulation signals to modulate a higher-dimensional Hidden layer.

615    Typically, information in the Task layer has been considered to correspond to dorsolateral prefrontal

616    cortex (DLPFC) and the Hidden layer to posterior task-related (e.g., visual and motor) processing

617    pathways (Miller & Cohen, 2001). Hence, we propose that modulation signals are implemented by

618    DLPFC.

619    A detailed investigation of activation at Hidden layer (Figure 7) illustrated that

620    (multiplicative) networks separate tasks by developing sparse representations (see also Bowers,

621    Vankov, Damian, & Davis, 2014). Here, for every task, only a small subset of neurons is active. In the

622    current context, this suggests that the network develops groups or modules of neurons that each become

623    specialized for a given task. Developing specialized modules might be beneficial for cognition in

624    various ways (Bullinaria, 2007; Clune, Mouret, & Lipson, 2013; Coltheart, 1999; Fodor, 1983;

625    Meunier, Lambiotte, Fornito, Ersche, & Bullmore, 2009). Moreover, recent work in reinforcement

626    learning has also described hierarchical forms of modularity (Botvinick, Niv, & Barto, 2009; Dietterich,

627    2000; Holroyd & Verguts, 2021; Krueger & Dayan, 2009). Here, it is proposed that the dorsal part of

628    the anterior cingulate cortex processes prediction errors related to specific events while the rostral part

629    processes prediction errors related to the context in which these events occur (Alexander & Brown,

630    2015). However, modularity of processes in a single task requires integration of information across

631    stages of processing. Hence, exploring the trade-off between sharing and separating representations at

632    different levels of processing is an important avenue for future research.

633    In the current work, modulation signals are employed to guide learning over trials. Here, (for

634    the adaptive networks) weights between the Task and Hidden layer are adapted to learn representations

635    of different tasks. In contrast, a lot of previous work considered how modulation is adapted to guide

636    online performance. Here, the intensity of the modulation signal is typically increased in response to

637    some evaluation of the cost-benefit structure of the task context (Shenhav, Botvinick, & Cohen, 2013).

638    These networks typically adapt the intensity of the modulation signal by changing activity in the Task

639    layer instead of changing the weights (Botvinick et al., 2001; Verbeke & Verguts, 2020; Verguts, 2017).

640    Consistent with this approach, research on visual attention has proposed that activity can be modulated

641    (in a multiplicative manner) at the level of the Stimulus layer (Martinez-Trujillo & Treue, 2004; Treue

642    & Martínez Trujillo, 1999). Alternatively, previous work (Cheadle et al., 2014) suggested that decisions

643    can be guided via adaptive (multiplicative) gain functions in the transfer from input to output. Hence,

644     there is a potentially important functional distinction between modulation in learning versus

645     performance (see also Lindsay & Miller, 2017).

646          The trade-off between shared and separated representation also impacts performance

647     (Musslick et al., 2020) Specifically, a wide range of empirical observations of interference during multi-

648     tasking (performing multiple tasks at the same time) can be explained by a tendency to share task

649     representations. Musslick et al. (2020) suggest that for generalization purposes, representations of novel

650     tasks should strongly overlap with other task representations; unfortunately, such overlap leads to strong

651     interference when these tasks need to be performed at the same time. However, with extensive training,

652     the networks will gradually separate task representations, which leads to less interference. Hence, future

653     work should consider a more extensive exploration of modulation signals in performance as well.

654          We evaluated the ability of different modulation signals to balance shared and separated

655     representations by investigating how the networks could overcome catastrophic interference.

656     Importantly, while modulation has proven to be efficient (Masse et al., 2018; Verbeke & Verguts, 2019),

657     previous work also described other methods to avoid catastrophic interference. For instance, sharing or

658     separation of task sets might be implemented in complementary learning systems (O'Reilly & Norman,

659     2002). Alternatively, machine learning literature has introduced methods such as synaptic intelligence

660     (Kirkpatrick et al., 2017) in which weights learn whether they should be specific for one task and hence

661     not change during further learning or whether they can be shared across all tasks. Future work should

662     further address the differences between these approaches.

663          Several additional tests and extensions can be made to the network. First, previous work

664     (Flesch et al., 2018) has pointed to an important distinction between interference in artificial and human

665     agents. While artificial agents show more interference when they learn in a blocked fashion, human

666     agents exhibit more interference when they learn in an interleaved fashion. The current work trained

667     artificial agents in a blocked fashion. However, the different types of modulation signals should also be

668     evaluated for interleaved learning. Potentially this approach could yield more insight in the existing

669     distinction of learning benefits for artificial compared to human agents. Second, the current networks

670     did not learn which task features were relevant for modulation. Here, the input layer was divided a-

671     priori in a Stimulus group and a Task group. Hence, an important next step for the network would be to

672     learn a hierarchical structure in input features in order to extract which inputs are relevant for

673     modulation and which are relevant for basic stimulus-action mappings (Rougier et al., 2005).

674     Alternatively, previous work has illustrated that compositional representations can even develop

675     without providing task input if they are considered useful as situational signals (Butz et al., 2021). Third,

676     although we illustrated that the Ax network was able to significantly benefit (in terms of accuracy) from

677     shared mappings between tasks (Figure 4a,b), we did not perform a direct test of generalization. That

678     is, we did not evaluate whether newly learned mappings in Task A of the Stroop task were transferred

679     to task E without further training in task E. That is, we did not test whether the learned task relations

680     were also suited for few-shot learning (Lake et al., 2014; Sylvain, Petrini, & Hjelm, 2020). Note that

681     for an exact transfer between A and E, the weight matrix between Task neuron A and the Hidden layer

682     should be exactly the same as the weight matrix between Task neuron E and the Hidden layer, which

683     seems a strong requirement. Yet, it is not clear either whether two contexts that require the same

684     stimulus-action mappings would also function in exactly the same way at behavioral and neural levels.

685     Nor is this computationally desirable: In the natural environment, two labels that lead to the same

686     stimulus-action contingencies, may still suggest (subtle) differences. Consider looking for lunch and

687     seeing a "restaurant" versus a "snack bar" in the distance: Both tell you that you will be able to eat

688     there, but expectations will differ at least slightly. As we discussed before, extracting a relational

689     structure between contexts that allows for partial generalization might be more optimal than using a

690     dichotomous same or different decision.

691         In sum, efficient human learning and performance requires to balance a trade-off between

692     sharing representations to allow generalization and separating representations to avoid interference. We

693     evaluated four different modulation signals and found that an adaptive multiplicative modulation signal

694     was best suited to balance the sharing/separation trade-off. This modulation signal allowed the Hidden

695     layer of the network to make a geometrical abstraction of the relational structure between tasks.

696     Importantly, our work opens several avenues for future work to increase the understanding of the

697     sharing/separation trade-off in both artificial and human agents.

698

699     **5. References**

700  Aben, B., Calderon, C. B., Van den Bussche, E., & Verguts, T. (2020). Cognitive Effort Modulates

701    Connectivity between Dorsal Anterior Cingulate Cortex and Task-Relevant Cortical Areas. *The*

702    *Journal of Neuroscience*, *40*(19). https://doi.org/10.1523/jneurosci.2948-19.2020

703  Abrahamse, E., Braem, S., Notebaert, W., & Verguts, T. (2016). Grounding Cognitive Control in

704    Associative Learning. *Psychological Bulletin*, *142*(7), 693–728.

705    https://doi.org/10.1037/bul0000047

706  Alexander, W. H., & Brown, J. W. (2015). Hierarchical error representation: A computational model

707    of anterior cingulate and dorsolateral prefrontal cortex. *Neural Computation*, *27*(11), 2354–

708    2410. https://doi.org/10.1162/NECO

709  Alon, N., Reichman, D., Shinkar, I., Wagner, T., Musslick, S., Cohen, J. D., … Ozcimder, K. (2017).

710    A graph-theoretic approach to multitasking. In *Advances in Neural Information Processing*

711    *Systems* (pp. 2097–2106).

712  Badre, D., Bhandari, A., Keglovits, H., & Kikumoto, A. (2021). The dimensionality of neural

713    representations for control. *Current Opinion in Behavioral Sciences*, *38*, 20–28.

714    https://doi.org/10.1016/j.cobeha.2020.07.002

715  Baxter, J. (2019). Learning Internal Representations. In *Proceedings of the eigth annual conference*

716    *on computational learning theory* (pp. 311–320).

717    https://doi.org/10.7551/mitpress/2906.003.0006

718  Bernardi, S., Benna, M. K., Rigotti, M., Munuera, J., Fusi, S., & Salzman, C. D. (2020). The

719    Geometry of Abstraction in the Hippocampus and Prefrontal Cortex. *Cell*, *183*(4), 954–967.

720    https://doi.org/10.1016/j.cell.2020.09.031

721  Botvinick, M. M., Braver, T. S., Barch, D. M., Carter, C. S., & Cohen, J. D. (2001). Conflict

722    monitoring and cognitive control. *Psychological Review*, *108*(3), 624–652.

723    https://doi.org/10.1037/0033-295X.108.3.624

724  Botvinick, M. M., Niv, Y., & Barto, A. C. (2009). Hierarchically organized behavior and its neural

725    foundations : A reinforcement learning perspective. *Cognition*, *113*(3), 262–280.

726    https://doi.org/10.1016/j.cognition.2008.08.011

727  Bouchacourt, F., & Buschman, T. J. (2019). A Flexible Model of Working Memory. *Neuron*, *103*(1),

147–160. https://doi.org/10.1016/j.neuron.2019.04.020

Bowers, J. S., Vankov, I. I., Damian, M. F., & Davis, C. J. (2014). Neural networks learn highly selective representations in order to overcome the superposition catastrophe. *Psychological Review*, *121*(2), 248–261. https://doi.org/10.1037/a0035943

Bullinaria, J. A. (2007). Understanding the advantage of modularity in neural systems. *Cognitive Science*, *31*(4), 673–695. https://doi.org/10.1080/15326900701399939

Butz, M. V., Achimova, A., Bilkey, D., & Knott, A. (2021). Event-Predictive Cognition: A Root for Conceptual Human Thought. *Topics in Cognitive Science*, *13*(1), 10–24. https://doi.org/10.1111/tops.12522

Butz, M. V., Bilkey, D., Humaidan, D., Knott, A., & Otte, S. (2019). Learning, planning, and control in a monolithic neural event inference architecture. *Neural Networks*, *117*, 135–144. https://doi.org/10.1016/j.neunet.2019.05.001

Cheadle, S., Wyart, V., Tsetsos, K., Myers, N., deGardelle, V., HerceCastañón, S., & Summerfield, C. (2014). Adaptive gain control during human perceptual choice. *Neuron*, *81*(6), 1429–1441. https://doi.org/10.1016/j.neuron.2014.01.020

Clune, J., Mouret, J. B., & Lipson, H. (2013). The evolutionary origins of modularity. *Proceedings of the Royal Society B: Biological Sciences*, *280*(1755). https://doi.org/10.1098/rspb.2012.2863

Cohen, J. D., Dunbar, K., & McClelland, J. L. (1990). On the control of automatic processes: a parallel distributed processing account of the Stroop effect. *Psychological Review*, *97*(3), 332–361. https://doi.org/10.1037/0033-295X.97.3.332

Collins, A. G. E., & Frank, M. J. (2016). Neural signature of hierarchically structured expectations predicts clustering and transfer of rule sets in reinforcement learning. *Cognition*, *152*, 160–169. https://doi.org/10.1016/j.cognition.2016.04.002

Coltheart, M. (1999). Modularity and cognition. *Trends in Cognitive Sciences*, *3*(3), 115–120. https://doi.org/10.1016/S1364-6613(99)01289-9

Dieterich, T. G. (2000). Hierarchical reinforcement learning with the MAXQ value function decomposition. *Journal of Artificial Intelligence Research*, *13*, 227–303. https://doi.org/10.1613/jair.639

756     Fidler, S., Boben, M., & Leonardis, A. (2009). Learning hierarchical compositional representations of

757         object structure. *Object Categorization: Computer and Human Vision Perspectives*,

758         *9780521887*, 196–215. https://doi.org/10.1017/CBO9780511635465.012

759     Flesch, T., Balaguer, J., Dekker, R., Nili, H., & Summerfield, C. (2018). Comparing continual task

760         learning in minds and machines. *Proceedings of the National Academy of Sciences of the United*

761         *States of America*, *115*(44). https://doi.org/10.1073/pnas.1800755115

762     Fodor, J. A. (1983). *The modularity of mind*. MIT Press/Bradford Books.

763     Franklin, N. T., & Frank, M. J. (2018). Compositional clustering in task structure learning. *PLoS*

764         *Computational Biology*, *14*(4), 1–25. https://doi.org/10.1371/journal.pcbi.1006116

765     Franklin, N. T., & Frank, M. J. (2020). Generalizing to generalize: Humans flexibly switch between

766         compositional and conjunctive structures during reinforcement learning. *PLoS Computational*

767         *Biology*, *16*(4), 1–33. https://doi.org/10.1371/journal.pcbi.1007720

768     French, R. M. (1999). Catastrophic forgetting in connectionist networks. *Trends in Cognitive*

769         *Sciences*, *3*(4), 128–135.

770     Fries, P. (2005). A mechanism for cognitive dynamics: neuronal communication through neuronal

771         coherence. *Trends in Cognitive Sciences*, *9*(10), 474–480.

772         https://doi.org/10.1016/j.tics.2005.08.011

773     Fries, P. (2015). Rhythms for Cognition: Communication through Coherence. *Neuron*, *88*(1), 220–

774         235. https://doi.org/10.1016/j.neuron.2015.09.034

775     Gershman, S. J., & Niv, Y. (2012). Exploring a latent cause theory of classical conditioning. *Learning*

776         *and Behavior*, *40*(3), 255–268. https://doi.org/10.3758/s13420-012-0080-8

777     Griffiths, T. L., & Ghahramani, Z. (2011). The Indian buffet process: An introduction and review.

778         *Journal of Machine Learning Research*, *12*, 1185–1224.

779     Helfrich, R. F., & Knight, R. T. (2016). Oscillatory Dynamics of Prefrontal Cognitive Control. *Trends*

780         *in Cognitive Sciences*, *20*(12), 916–930. https://doi.org/10.1016/j.tics.2016.09.007

781     Holroyd, C. B., & Verguts, T. (2021). The best laid plans: Computational principles of ACC. *Trends*

782         *in Cognitive Sciences*, *25*(4), 316–329. https://doi.org/10.1016/j.tics.2021.01.008

783     Hupkes, D., Dankers, V., Mul, M., & Bruni, E. (2020). Compositionality decomposed: How do neural

784     networks generalise? *Journal of Artificial Intelligence Research*, *67*, 757–795.

785     Irsoy, O., & Cardie, C. (2014). Deep recursive neural networks for compositionality in language.

786         *Advances in Neural Information Processing Systems*, *3*(January), 2096–2104.

787     İrsoy, O., & Cardie, C. (2015). Modeling compositionality with multiplicative recurrent neural

788         networks. *3rd International Conference on Learning Representations, ICLR 2015 - Conference*

789         *Track Proceedings*, (2013), 1–10.

790     Jensen, O., & Mazaheri, A. (2010). Shaping functional architecture by oscillatory alpha activity:

791         Gating by inhibition. *Frontiers in Human Neuroscience*, *4*, 1–8.

792         https://doi.org/10.3389/fnhum.2010.00186

793     Kim, W., Pitt, M. A., & Myung, J. I. (2013). How Do PDP Models Learn Quasiregularity?

794         *Psychological Review*, *120*(4), 903–916. https://doi.org/10.1037/a0034195

795     Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., … Hadsell, R.

796         (2017). Overcoming Catastrophic Forgetting in Neural Networks. *Proceedings of the National*

797         *Academy of Sciences*, *114*(13), 3521–3526. https://doi.org/10.1073/pnas.1611835114

798     Krueger, K. A., & Dayan, P. (2009). Flexible shaping: How learning in small steps helps. *Cognition*,

799         *110*(3), 380–394. https://doi.org/10.1016/j.cognition.2008.11.014

800     Lake, B., & Baroni, M. (2018). Generalization without systematicity: On the compositional skills of

801         sequence-to-sequence recurrent networks. *35th International Conference on Machine Learning,*

802         *ICML 2018*, *7*, 4487–4499.

803     Lake, B., Lee, C.-Y., Glass, J., Lake, B. M., Glass, J. R., & Tenenbaum, J. B. (2014). One-shot

804         learning of generative speech concepts Publication Date One-shot learning of generative speech

805         concepts. *Proceedings of the Annual Meeting of the Cognitive Science Society*, (36), 803–808.

806         Retrieved from https://cloudfront.escholarship.org/dist/prd/content/qt3xf2n3vc/qt3xf2n3vc.pdf

807     Lake, B., Ullman, T., Tenenbaum, J., & Gershman, S. (2017). Building machines that learn and think

808         like people. *Behavioral and Brain Sciences*, *40*(2017).

809         https://doi.org/10.1017/S0140525X16001837

810     LeCun, Y., Cortes, C., & Burges, C. J. (2010). MNIST handwritten digit database.

811     Lillicrap, T. P., Cownden, D., Tweed, D. B., & Akerman, C. J. (2016). Random synaptic feedback

weights support error backpropagation for deep learning. *Nature Communications*, *7*, 1–10.

https://doi.org/10.1038/ncomms13276

Lindsay, G. W., & Miller, K. D. (2018). How biological attention mechanisms improve task

performance in a large-scale visual system model. *ELife*, *7*, 1–29.

https://doi.org/10.7554/eLife.38105.001

Lisman, J. E., & Jensen, O. (2013). The Theta-Gamma Neural Code. *Neuron*, *77*(6), 1002–1016.

https://doi.org/10.1016/j.neuron.2013.03.007

Maass, W., Natschläger, T., & Markram, H. (2002). Real-time computing without stable states: A

new framework for neural computation based on perturbations. *Neural Computation*, *14*(11),

2531–2560. https://doi.org/10.1162/089976602760407955

Martinez-Trujillo, J. C., & Treue, S. (2004). Feature-Based Attention Increases the Selectivity of

Population Responses in Primate Visual Cortex. *Current Biology*, *14*, 744–751.

https://doi.org/10.1016/j.cub.2004.04.028

Masse, N. Y., Grant, G. D., & Freedman, D. J. (2018). Alleviating catastrophic forgetting using

context-dependent gating and synaptic stabilization. *Proceedings of the National Academy of

Sciences*, *115*(44), 1–12. https://doi.org/10.1073/pnas.1803839115

McClelland, J. L., McNaughton, B. L., & O'Reilly, R. C. (1995). Why There Are Complementary

Learning Systems in the Hippocampus and Neo-cortex: Insights from the Successes and Failures

of Connectionists Models of Learning and Memory. *Psychological Review.*, *102*(3), 419–457.

https://doi.org/10.1037/0033-295X.102.3.419

Meunier, D., Lambiotte, R., Fornito, A., Ersche, K. D., & Bullmore, E. T. (2009). Hierarchical

modularity in human brain functional networks. *Frontiers in Neuroinformatics*, *3*, 1–12.

https://doi.org/10.3389/neuro.11.037.2009

Miller, E. K., & Cohen, J. D. (2001). An Integrative Theory of Prefrontal Cortex Function. *Annual

Review of Neuroscience*, *24*, 167–202. https://doi.org/10.1146/annurev.neuro.24.1.167

Musslick, S., & Cohen, J. D. (2020). Rationalizing constraints on the capacity for cognitive control.

*PsyArxiv*, 45. https://doi.org/10.31234/osf.io/vtknh

Musslick, S., Saxe, A. M., Ozcimder, K., Dey, B., Henselman, G., & Cohen, J. D. (2017).

840  Multitasking Capability Versus Learning Efficiency in Neural Network Architectures. In *Annual*

841  *meeting of the Cognitive Science Society* (pp. 829–834).

842  Musslick, S., Saxe, A., Novick, A., Reichman, D., & Cohen, J. D. (2020). On the rational

843  boundedness of cognitive control: Shared versus separated representations. *PsyArXiv*.

844  https://doi.org/10.31234/osf.io/jkhdf

845  O'Reilly, R. C., & Frank, M. J. (2006). Making Working Memory Work : A Computational Model of

846  Learning in the Prefrontal Cortex and Basal Ganglia. *Neural Computation*, *18*(2), 283–328.

847  https://doi.org/10.1162/089976606775093909

848  O'Reilly, R. C., & Norman, K. A. (2002). Hippocampal and neocortical contributions to memory:

849  Advances in the complementary learning systems framework. *Trends in Cognitive Sciences*,

850  *6*(12), 505–510. https://doi.org/10.1016/S1364-6613(02)02005-3

851  Rougier, N. P., Noelle, D. C., Braver, T. S., Cohen, J. D., & O'Reilly, R. C. (2005). Prefrontal cortex

852  and flexible cognitive control: Rules without symbols. *Proceedings of the National Academy of*

853  *Sciences of the United States of America*, *102*(20), 7338–7343.

854  https://doi.org/10.1073/pnas.0502455102

855  Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-

856  propagating errors. *Nature*, *323*, 533–536. https://doi.org/10.1038/323533a0

857  Sagiv, Y., Musslick, S., Niv, Y., & Cohen, J. D. (2020). Efficiency of learning vs. processing:

858  Towards a normative theory of multitasking. In *Proceedings of the 40th Annual Meeting of the*

859  *Cognitive Science Society,* (p. 1004—1009).

860  Servan-Schreiber, D., Printz, H., & Cohen, J. D. (1990). A network model of catecholamiine effects:

861  Gain, signal-to-noise ratio, and behavior. *Science*, *249*(4971), 892–895.

862  https://doi.org/10.1126/science.2392679

863  Shenhav, A., Botvinick, M. M., & Cohen, J. D. (2013). The expected value of control: An integrative

864  theory of anterior cingulate cortex function. *Neuron*, *79*(2), 217–240.

865  https://doi.org/10.1016/j.neuron.2013.07.007

866  Stroop, J. R. (1935). Studies of interference in serial verbal reactions. *Journal of Experimental*

867  *Psychology*, *18*(6), 643–662. https://doi.org/10.1037/h0054651

868  Sugita, Y., Tani, J., & Butz, M. V. (2011). Simultaneously emerging braitenberg codes and

869       compositionality. *Adaptive Behavior*, *19*(5), 295–316.

870       https://doi.org/10.1177/1059712311416871

871  Sylvain, T., Petrini, L., & Hjelm, R. D. (2020). Zero-Shot Learning from scratch (ZFS): leveraging

872       local compositional representations. Retrieved from http://arxiv.org/abs/2010.13320

873  Treue, S., & Martínez Trujillo, J. C. (1999). Feature-based attention influences motion processing

874       gain in macaque visual cortex. *Nature*, *399*, 575–579. https://doi.org/10.1038/21176

875  Tsai, C. Y., Saxe, A., & Cox, D. (2016). Tensor switching networks. *Advances in Neural Information*

876       *Processing Systems*, (Nips), 2046–2054.

877  Tubiana, J., & Monasson, R. (2017). Emergence of Compositional Representations in Restricted

878       Boltzmann Machines. *Physical Review Letters*, *118*(13), 1–5.

879       https://doi.org/10.1103/PhysRevLett.118.138301

880  Vaidya, A. R., Jones, H. M., Castillo, J., & Badre, D. (2021). Neural representation of abstract task

881       structure during generalization. *ELife*, (10:e63226.), 1–22. https://doi.org/10.7554/eLife.63226

882  Verbeke, P., Ergo, K., De Loof, E., Verguts, T., Loof, E. De, & Verguts, T. (2021). Learning to

883       synchronize: Midfrontal theta dynamics during rule switching. *Journal of Neuroscience*, *41*(7),

884       1–13. https://doi.org/10.1523/JNEUROSCI.1874-20.2020

885  Verbeke, P., & Verguts, T. (2019). Learning to synchronize: How biological agents can couple neural

886       task modules for dealing with the stability-plasticity dilemma. *PLoS Computational Biology*,

887       *15*(8). https://doi.org/10.1371/journal.pcbi.1006604

888  Verbeke, P., & Verguts, T. (2020). Neural synchrony for adaptive control. *Psyarxiv*, 1–37.

889       https://doi.org/10.31234/osf.io/523x9

890  Verguts, T. (2017). Binding by random bursts: A computational model of cognitive control. *Journal*

891       *of Cognitive Neuroscience*, *29*(6), 1103–1118. https://doi.org/10.1162/jocn

892  Verguts, T., & Notebaert, W. (2008). Hebbian Learning of Cognitive Control : Dealing With Specific

893       and Nonspecific Adaptation. *Psychological Review*, *115*(2), 518–525.

894       https://doi.org/10.1037/0033-295X.115.2.518

895  Wilson, R. C., Takahashi, Y. K., Schoenbaum, G., & Niv, Y. (2014). Orbitofrontal cortex as a

896     cognitive map of task space. *Neuron*, *81*(2), 267–279.

897     https://doi.org/10.1016/j.neuron.2013.11.005

898     Yang, G. R., Joglekar, M. R., Song, H. F., Newsome, W. T., & Wang, X. J. (2019). Task

899     representations in neural networks trained to perform many cognitive tasks. *Nature*

900     *Neuroscience*, *22*(2), 297–306. https://doi.org/10.1038/s41593-018-0310-2

901     Yu, L. Q., Wilson, R. C., & Nassar, M. R. (2020). Adaptive learning is structure learning in time.

902     *Psyarxiv*, 1–27. https://doi.org/10.31234/osf.io/r637c

903     Zambaldi, V., Raposo, D., Santoro, A., Bapst, V., Li, Y., Babuschkin, I., … Battaglia, P. (2018).

904     Relational Deep Reinforcement Learning, (2), 1–15. Retrieved from

905     http://arxiv.org/abs/1806.01830

906

**Supplementary materials**

**S.1. Exploration of activation functions**

In the main text (see Equation (2)), multiplicative modulation was established by combining two nonlinear transformations of the Stimulus and Task input via $f()$ and $g()$. Here, $f()$ represented the sigmoid activation function (see Equation (3)) and $g()$ represented a RELU activation function (see Equation (4)). Additive modulation (see Equation (1)) was implemented by transforming both the Task and Stimulus input with $f()$. Here, we present 4 novel simulations in which we explored all combinations of activation functions. Specifically, we tested networks in which both $f()$ and $g()$ corresponded to the sigmoid function (i.e. Sig ($\otimes$ Sig)), we tested when $f()$ corresponded to the sigmoid function and $g()$ to the RELU function (i.e. Sig ($\otimes$ RELU); as in the main text), we tested when $f()$ corresponded to the RELU function and $g()$ to the sigmoid function (i.e. RELU ($\otimes$ Sig)) and we tested the networks when both $f()$ and $g()$ corresponded to the RELU function (i.e. RELU ($\otimes$ RELU)). Note that for the additive modulation networks, only $f()$ is relevant. This is why our notation shows the second function (i.e., $g()$) between brackets.

For these additional simulations, the networks were tested on the Stroop and Trees dataset with 12 Hidden neurons. We again explored different values of $\alpha$ ranging from 0 to 1 in steps of .2. Again, 25 simulations were performed in which we shuffled task contexts and trained networks for three context repetitions after which weights were frozen and the networks were tested again on each context.

Accuracy and dissimilarity correlations were analysed during the test phase. As can be observed (Figure S1), the network that is presented in the main text Sig ($\otimes$ RELU) is most optimal for both the multiplicative and additive modulation networks. Although the multiplicative modulation seemed less efficient when implemented via a sigmoid function (Figure S1a,c,e,g,i,k,m,o), the Ax network still outperformed the other modulation methods. Presumably, the RELU function is more efficient for modulation because it operates at a much larger scale $[0, +\infty]$ than the sigmoid function which has bounds at 0 and 1.
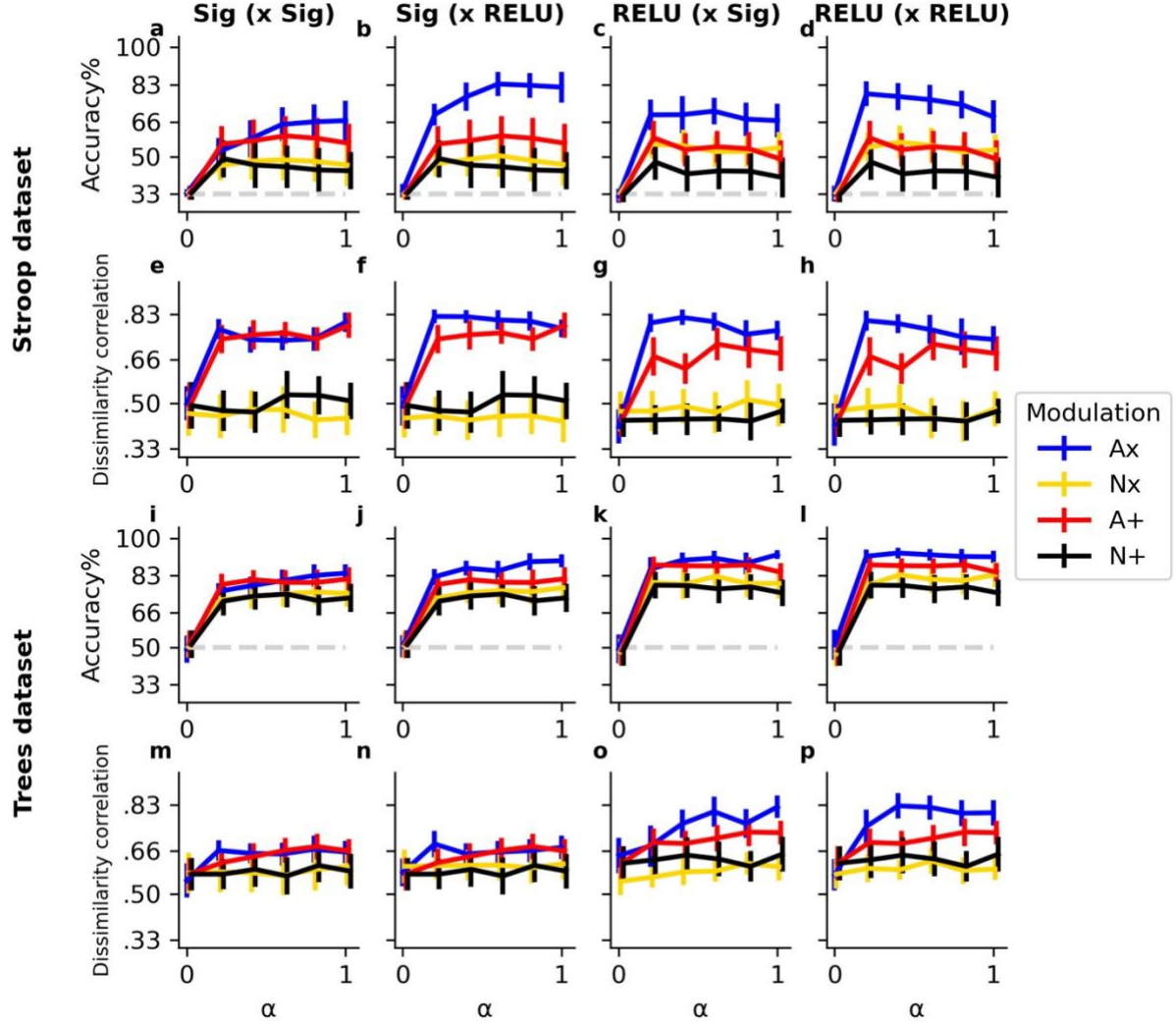
**Figure S1. Activation function exploration.** Lines illustrate mean accuracy/dissimilarity correlations for each value of α across all tasks and all simulations. Bars indicate 95% confidence intervals over 25 simulations. The dashed lightgrey line indicates chance level accuracy. Results are shown for different datasets (rows) and different combinations of activation functions (columns).

### S.2. Concatenated versus separated input transformations

In the main text (see Equation (2)), multiplicative modulation was established by combining two separated nonlinear transformations of the Stimulus and Task input via *f(SW)⊗ g(TW)* while Additive modulation (see Equation (1)) was implemented by transforming both the Task and Stimulus input with *f(SW+TW)*. Hence, for additive modulation, inputs were concatenated in one transformation, while for multiplicative modulation inputs were separated in two transformations. We choose these different implementations because they are most closely related to previous work (e.g., Cohen, Dunbar, & McClelland, 1990; Masse, Grant, & Freedman, 2018 for additive and multiplicative modulation

respectively). However, for completeness, we also explored other implementations of multiplicative and additive modulation. Specifically, we tested networks in which for both types of modulation (additive and multiplicative) the inputs were concatenated. This results in $f(SW) \otimes g(TW)$ for multiplicative modulation and $f(SW) + g(TW)$ for additive modulation. Additionally, we tested networks in which the inputs were separated. This resulted in $f(SW \otimes TW)$ for multiplicative modulation and $f(SW + TW)$ for additive modulation.

Again, the networks were tested on the Stroop and Trees dataset with 12 Hidden neurons. We explored different values of $\alpha$ going from 0 to 1 in steps of .2 and 25 simulations were performed in which we shuffled task contexts and trained networks for three context repetitions after which weights were frozen and the networks were tested again on each context.

Accuracy and dissimilarity correlations were analysed during the test phase. As illustrated in (Figure S2), for none of the modulation methods there is a clear difference when using separated versus concatenated input transformations.
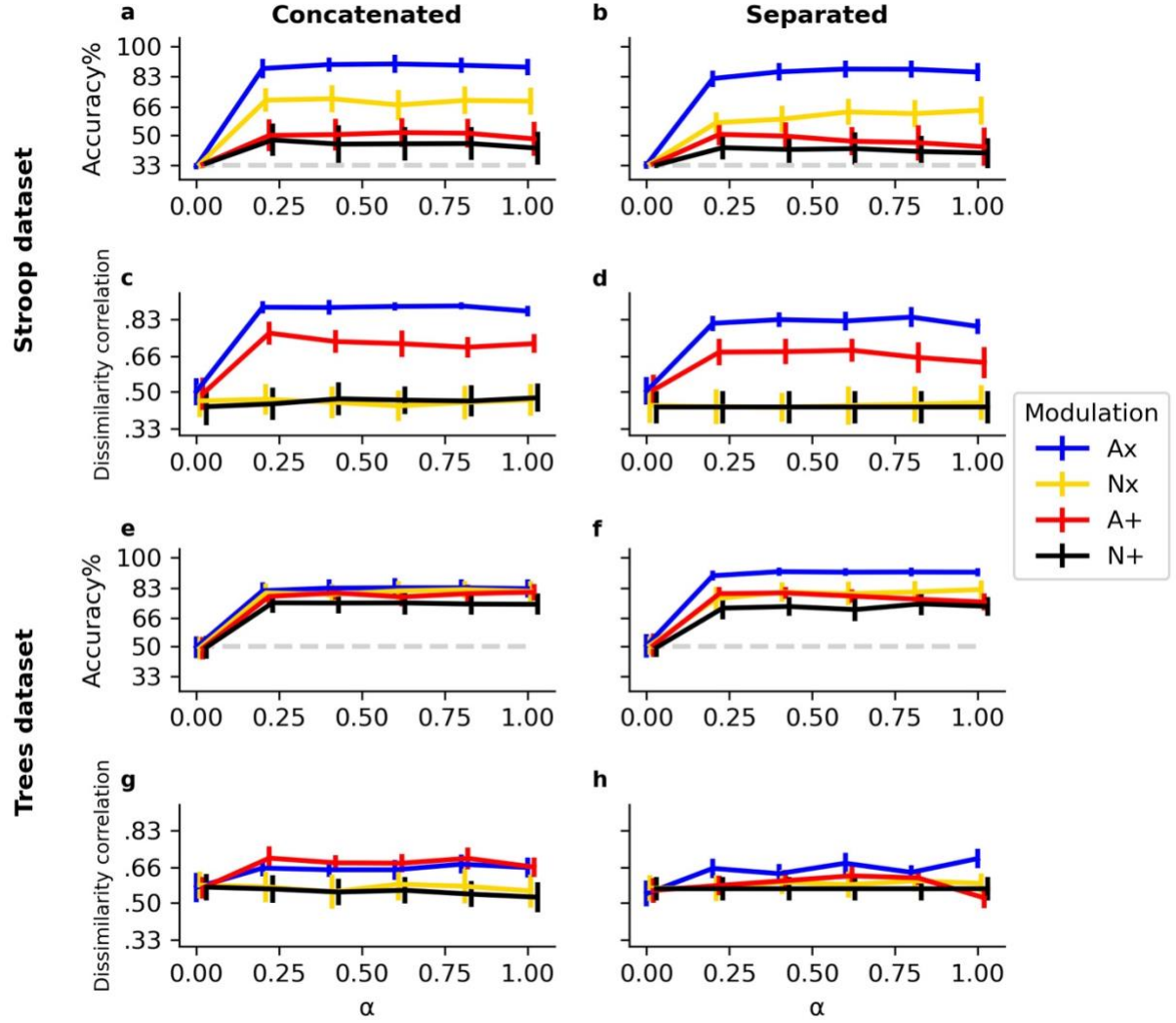
**Figure S2. Concatenated versus separated input transformations.** Lines illustrate mean accuracy/dissimilarity correlations for each value of α across all tasks and all simulations. Bars indicate 95% confidence intervals over 25 simulations. The dashed lightgrey line indicates chance level accuracy. Results are shown for different datasets (rows) and different concatenated versus separated input transformations (columns).

## S.3. Exploration of weight initialization

In the main text we described that all weights are initialized with a random value drawn from the normal distribution $N(0, 1)$. Only for the Ax network, modulating weights (between Task and Hidden layer) had an initial random value drawn from the uniform distribution $U(0, 1)$, such that $RELU(T) > 0$ (all gates open) at the first trial. However, previous work has demonstrated that the way in which weights are initialized is not trivial (e.g., Flesch, Juechems, Dumbalska, & Saxe, 2021). To illustrate that our results were not solely driven by this choice of weight initialization, we performed

additional simulations in which we tested a normal initialization for all modulating weights and compared this to a uniform initialization for all modulating weights.

The networks were again tested on the Stroop and Trees dataset with 12 Hidden neurons. We explored different values of $\alpha$ ranging from 0 to 1 in steps of .2. Again, 25 simulations were performed in which we shuffled task contexts and trained networks for three context repetitions after which weights were frozen and the networks were tested again on each context.

As shown in Figure S3, the Ax network indeed performs a bit worse when the weights are initialized from the normal distribution. However, it is also the case that for the other three networks the uniform initialization was a bit less optimal. Hence, the main text describes simulations in which the optimal weight initialization was chosen for each of the four modulation methods.
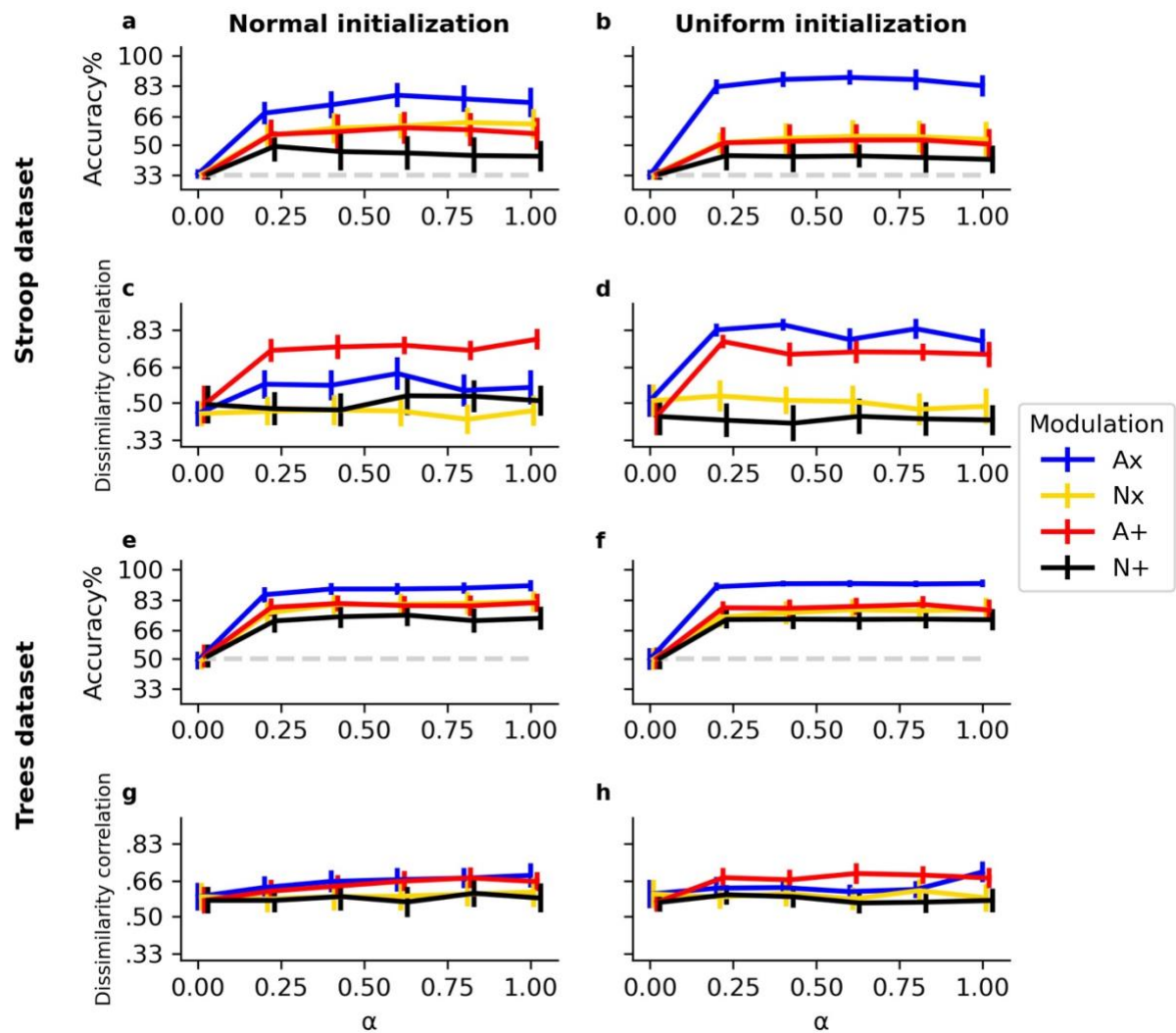
**Figure S2. Weight initialization exploration.** Lines illustrate mean accuracy/dissimilarity correlations for each value of α across all tasks and all simulations. Bars indicate 95% confidence intervals over 25 simulations. The dashed lightgrey line indicates chance level accuracy. Results are shown for different datasets (rows) and for different initializations of the modulation weights (columns).

## References

Cohen, J. D., Dunbar, K., & McClelland, J. L. (1990). On the control of automatic processes: a parallel distributed processing account of the Stroop effect. *Psychological Review*, *97*(3), 332–361. https://doi.org/10.1037/0033-295X.97.3.332

Flesch, T., Juechems, K., Dumbalska, T., & Saxe, A. (2021). Rich and lazy learning of task representations in brains and neural networks. *BioRxiv*.

Masse, N. Y., Grant, G. D., & Freedman, D. J. (2018). Alleviating catastrophic forgetting using context-dependent gating and synaptic stabilization. *Proceedings of the National Academy of Sciences*, *115*(44), 1–12. https://doi.org/10.1073/pnas.1803839115