**SI1: Previous biological models of sameness detection**

Two broad classes of mechanisms have been proposed for determining whether two stimuli are identical (Grill-Spector, Henson, & Martin, 2006; Kumaran & Maguire, 2007). On the one hand, sequentially presented identical stimuli elicit less activation, due to neuronal "fatigue" or sharpening of the representations. As a result, novel, non-repeated representations have relatively higher levels of activation. However, such models cannot explain why sameness-relations can be generalized: after detecting the repetition in *babagu*, an item with new syllables and the same repetition-pattern (e.g., *wowofe*) will be just as unfamiliar as an item with new syllables and another repetition-pattern (e.g., *wofefe*). As a result, an explicit representation of sameness vs. difference (or match vs. non-match) is required (but see Cope et al., 2018, where generalization is observed under some circumstances).

The second class of mechanisms involves some kind of comparator between memory representations and sensory input, though there are few explicit and biologically realistic models of sameness matches. For example, it has been proposed that the hippocampal CA1 region (and maybe the CA3 region, depending on the studies) are crucial for detecting matches between memory representations and sensory input (while the CA3 regions might have an additional role in retrieving associations; Hasselmo, 2005; Lisman, 1999; Lisman & Otmakhova, 2001).

We will now discuss a number of representative models to illustrate these points.

SI1.1. Hasselmo and Wyble (1997), Carpenter and Grossberg (1987), Wen, Ulloa,

26        In their simulation of memory retrieval in the hippocampus, Hasselmo and

27        Wyble (1997) provide an explicit model of comparator-based sameness detection,

28        inspired by the anatomy of the trisynaptic circuit. Specifically, the hippocampus

29        receives sensory input from the entorhinal cortex, which in turn projects to region

30        CA3 (via the dentate gyrus). In contrast, CA1 receives input both from CA3 (via

31        the Schaffer collaterals) as well as directly from the entorhinal cortex. If

32        memories are encoded in CA3, the simultaneous input from the entorhinal cortex

33        and CA3 might allow CA1 to detect matches between sensory input (from the

34        entorhinal cortex) and memory representations (from CA3).

35        Specifically, during encoding of novel items, combined sensory and

36        memory (from CA1) input leads to novel self-organized representations in CA1.

37        To activate these representations during recognition (i.e., to enter the

38        corresponding attractor state), input from both the entorhinal cortex and CA3 is

39        required; sensory input alone does not activate the attractor state. In other words,

40        CA1 enters an attractor state only when the current sensory input matches

41        currently active memory representations in CA3 (see also Ludueña & Gros, 2013,

42        for a model that uses anti-Hebbian learning to configure a *mismatch* detector).

43        Relatedly, some working memory models detect matches between the

44        contents of working memory and current sensory input by *adding* input from

45        sensory input and WM (Carpenter & Grossberg, 1987; Wen et al., 2008). If an

46        item is in WM, it will provide an additional input. As a result, matches between

47        sensory input and WM can be detected using some threshold (though such a

48        mechanism might not be robust as it depends on the *absolute* fire firing rates;

49    Engel & Wang, 2011).

50                    <u>SI1.2. Engel and Wang (2011)</u>

51            While these models detect matches because the combined output of

52    memory representations and matching sensory input trigger the retrieval of other

53    representations, it is also possible to detect matches by *subtracting* sensory input

54    and memory representation. Such a model has been proposed in the context of

55    delayed-match-to-sample tasks. Specifically, Engel and Wang (2011) proposed a

56    biologically realistic model that detects matches through (i) a working memory

57    (WM) sub-network, (ii) a comparator sub-network, and (iii) a decision network

58    receiving input from the comparator network. Neurons in the WM network

59    receive sensory input (but only when attention is directed to the input) and can

60    maintain memory traces through self-excitation. Critically, the comparator

61    network is composed of two distinct populations. One receives both sensory input

62    and input from the WM network (hereafter called sensory+WM neurons). The

63    other receives *only* sensory input but no WM input (hereafter called sensory-only

64    neurons). Engel and Wang (2011) make two other critical assumptions. First, they

65    assume that the total level of excitation should be similar for matches in the

66    sensory+WM population, and for mismatches in the sensory-only population; as

67    the sensory+WM population has an extra excitatory input, they achieved this by

68    scaling down the synaptic excitation targeting the sensory+WM population.

69    Second, the comparator units show center-surround inhibition: there is a (limited

70    level of) excitation from similar stimuli, and much stronger inhibition from

71    dissimilar stimuli.

72            These assumptions conspire to yield stronger activation in the

73 sensory+WM population for matches, and stronger activation in the sensory-only

74 population for mismatches. As a result, to decide whether a probe matches the

75 target, the decision network just needs to compare the activation of the

76 sensory+WM population and the sensory-only populations. More specifically, in

77 the matching case, the sensory+WM neurons receive input both from the probe

78 and from the matching content of WM; in contrast, the sensory-only neurons

79 receive input only from the sensory representations of the probe. As a result, there

80 is stronger activation in the sensory+WM population. In contrast, in the case of a

81 mismatch, both populations receive input from the sensory representations of the

82 probe.

83 As Engel and Wang (2011) assume that excitatory input is stronger for the

84 sensory-only population, this population is expected to receive somewhat stronger

85 input than the sensory+WM population. Further, the sensory+WM population also

86 receives input from the (mismatching) WM representation; due to the center-

87 surround inhibition in the network, the probe and the (mismatching) target inhibit

88 each other, further reducing the activation in the sensory+WM population. The

89 decision network just has to decide whether similar orientations have stronger

90 presentation in the sensory+WM population or the sensory-only population.[1]

91 SI1.3. Johnson, Spencer, Luck, & Schöner (2009)

92 Another WM model that explicitly incorporates a same/different

93 distinction has been proposed by Johnson et al. (2009). In their model, sensory

---

[1] This model assumes that WM is mediated by self-sustained activity in a population of neurons. However, it has been questioned whether such self-sustained activation really plays a crucial role in WM (Rose et al., 2016; Stokes, 2015).

94    input excites (self-sustained) WM representations, which, in turn, inhibits the

95    corresponding sensory activation (with center-surround inhibition in all areas). As

96    a result, upon presentation of the sample stimulus, there is a self-sustained

97    representation of the memory items in WM, but little activation in sensory areas

98    due to the inhibitory input from WM. Hence, if a later sensory input matches the

99    items in memory, the sensory areas will remain largely silent. In contrast, if the

100   sensory input differs from the memory items, sensory input will be uninhibited.

101   Hence, in this model, "decision" neurons that receive excitation from WM will

102   respond to matches, while decision neurons receiving excitation from sensory

103   input will respond to mismatches, at least with mutual inhibition between these

104   decision populations.

105        However, there are a number of problems with this model. First, it has

106   been questioned whether WM really relies on self-sustained activity (Rose et al.,

107   2016; Stokes, 2015). Second, and crucially, items in (working) memory seem to

108   *attract* attention (Awh & Jonides, 2001; Downing, 2000; Fan & Turk-Browne,

109   2016) which seems inconsistent with the proposal that memory items suppress

110   perceptual input.

111        SI1.4. Difficulties of these models with generalizable repetition patterns

112        In their current instantiations such models are unlikely to account for the

113   generalization of sameness relatoions (nor were these models intended to do so).

114   For example, after exposure to *pupu*, they are unlikely to recognize *baba* over,

115   say, *bapu* when the syllables are novel. In Hasselmo and Wyble's (1997) model,

116   items like *baba* have no memory representation, and thus cannot trigger CA1-like

117   activation any more than *bapu* sequences. That said, a version of Hasselmo and

118   Wyble's (1997) model might act as a repetition-detector if each item undergoes

119   element-by-element encoding-retrieval cycles. For example, when processing the

120   item *pupu*, the network might first encode the first instance of the syllable *pu*;

121   upon presentation of the second instance of *pu*, a CA1-like structure might enter

122   an attractor state as the current sensory input matches an existing memory

123   representation. In contrast, for items like *bapu*, the second element has no

124   corresponding existing memory representation, and thus does not activate an

125   attractor in a CA1-like structure. Hence, a readout mechanism to a CA1-like

126   structure could, in principle, act as a repetition-detector.

127          However, there are a number of problems with such an architecture as

128   well. For example, after exposure to both *pupu* and *bapu*, the model might

129   classify *bapu* as a repetition, because the syllable *pu* has an existing memory

130   representation from a previous item. In other words, the model would show

131   proactive interference in sameness-detection, and there is no evidence that this

132   might be the case in real learners.

133          Likewise, such a model will face difficulties discriminating *ABB* patterns

134   as in *pulili* from *ABA* patterns as in *pulipu*. It is also unclear whether a memory-

135   based repetition-detector can detect the sameness of simultaneously (rather than

136   sequentially) presented items, and whether such models would "recognize" an

137   item in the presence of distractors; after all, a CA1-like region would receive

138   input from the distractors as well, which will bring up the total level of activity.

139   As a result, considerable computational and neuroscientific research is needed to

140   decide whether such an architecture might act as a repetition-detector.

141          Similar problems arise Engel and Wang's (Engel & Wang, 2011) model.

142    First, the model would show proactive interference, and falsely detect repetitions

143    in non-repetition sequences if the second item has been placed in working

144    memory on a considerable number of earlier occasions. While this is an empirical

145    prediction, it seems, at first sight, implausible.[2] Second, it is unclear whether a

146    memory-based repetition-detector would be able to detect the sameness of

147    simultaneously presented items. Third, Engel and Wang (2011) use supervised

148    training to teach units to subtract the activation of sensory+WM neurons and

149    sensory-only neurons, respectively. However, in all experiments on repetition

150    learning in infancy, learning is unsupervised; further, the reliance on supervised

151    training prevents the model from generalizing to items that are dissimilar from

152    those it has been trained on (Marcus, 1998a, 1998b).

153         There is another reason for which such memory-based repetition-detectors

154    are unlikely to support the kinds of generalizations reviewed here. Given how

155    widespread the ability to compute repetition-patterns is, one would expect it to

156    rely on fairly simple circuits. However, these memory-based models rely on the

157    interaction of different brain areas (the entorhinal cortex as well as CA1 and CA3

158    in the case of Hasselmo and Wyble's (1997) model, and a sensory as well as a

159    working memory system in the case of Engel and Wang's ( 2011) model).

---

[2] Engel and Wang's ( 2011) model can detect matches between the *A* items in *ABBA*
trials. However, they achieve this by assuming that the WM subnetwork receives sensory input
only when the input is attentionally encoded. As a result, only the first *A* from the *ABBA* items
ever reaches WM. However, this would predict that participants do not notice the repetition of the
*B* items. It thus seems that the WM component in Engel and Wang's (2011) has a similar function
as (pre) frontal regions in the recent models of inhibition (Egner & Hirsch, 2005; Erika-Florence,
Leech, & Hampshire, 2014; Hampshire & Sharp, 2015): it serves to highlight task-relevant
representations.

<u>SI1.5. Cope et al. (2018)</u>

161    Cope et al. (2018) proposed a model to explain the successful performance

162    of bees in delayed-match-to-sample tasks such as in Giurfa, Zhang, Jenett,

163    Menzel, and Srinivasan (2001). They used a model inspired by the architecture of

164    the bee mushroom body. At a conceptual level, the model comprises three

165    populations of neurons: (1) a population of input neurons encoding stimuli

166    (inspired by Kenyon cells); (2) a population of inhibitory neurons (inspired by the

167    protocerebellar tract); (3) and a population of output neurons (inspired by

168    extrinsic neurons), half of which code for a "go" response and half for a "no-go"

169    response.

170    The input population has excitatory connections (with fixed weights) to

171    both the output neurons and the inhibitory neurons; the inhibitory neurons project

172    to the output neurons as well, but, critically, with weights that are modifiable.

173    The critical assumption of the model is (an empirically observed) "fatigue

174    effect" in the input neurons: responses to repeated stimuli are weaker than to

175    novel stimuli. As these weaker activations are assumed to be insufficient to drive

176    the inhibitory population, novel and repeated stimuli play different roles in match-

177    to-sample tasks and non-match-to-sample tasks, respectively.

178    In *match-to-sample tasks*, repeated items fail to activate the inhibitory

179    neurons. As a result, the connection weight between the inhibitory neurons and

180    the output neurons is adjusted only when non-matching, novel items are

181    presented. Given that "go" responses to non-match items are not reinforced in

182    match-to-sample tasks, the strength of the connections between the inhibitory

183    neurons and *go* responses is increased relative to the strength of the connections

184   between the inhibitory neurons and *no go* responses. (The strength of the

185   connections between the inhibitory neurons and *no go* responses does not change

186   as no learning takes place if the bee refuses to "go" for a stimulus to begin with.)

187        In *non-match-to-sample tasks*, repeated items still fail to activate the

188   inhibitory neurons, so that learning occurs only with non-matching, novel items.

189   However, in *non-match-to-sample tasks*, connections between the inhibitory

190   neurons and *go* responses are weakened, relative to the connections between the

191   inhibitory neurons and *no go* responses.

192        In other words, the inhibitory population learns to select between go and

193   no-go responses, based on the frequency with which the responses are responses

194   are reinforced when it is activated by novel, non-matching stimuli. It thus detects

195   the correlation between the presence of rewards and input from non-repeated

196   stimuli.

197        Impressively, these simple computational principles are sufficient to allow

198   the model toe generalize the sameness-relations to untrained items; for example, if

199   the model is trained in delay (non-) match-to-sample task with, say, orientations,

200   it would transfer this learning to a task with, say, colors.

201        However, there are four situations that raise the question of whether this

202   model would appropriately account for sameness-detection in grammar-learning

203   situations. First, it is unclear to what extent this model can discriminate matching

204   from non-matching pairs when the elements of the pair are presented

205   simultaneously (Martinho & Kacelnik, 2016). This is because the model relies on

206   a decrease in representational strength of items presented repeatedly, and, if

207   identical items are presented simultaneously, no such decrease can occur (though

208    this issue might be solved if organisms attend to the items sequentially).

209         Second, and critically, humans and some other animals can learn

210    sameness-relations from positive evidence alone, in the absence of reinforcement

211    (Marcus, Vijayan, Rao, & Vishton, 1999).

212         Third, and relatedly, Cope et al.'s (2018) model learns in a fundamentally

213    different way from humans. Specifically, the model learns about *non*-matching

214    items. In match-to-sample tasks, it learns to increase the inhibition of "go"

215    responses to non-match stimuli; in non-match-to-sample tasks, it learns to

216    decrease the inhibition of "go" responses to non-match stimuli. In contrast,

217    humans learn predominantly about sameness rather than difference relations, and,

218    to the extent that they represent difference relations, they represent them as

219    negations of sameness relations (Hochmann, Carey, & Mehler, 2018; Hochmann,

220    Mody, & Carey, 2016).

221         Fourth, the model does not produce representations of sameness or

222    differences that can be used for further processing. For example, in Marcus et al.'s

223    (1999) discrimination between *AAB* and *ABB*, the critical distinction was not

224    whether the strings contained a repetition, but rather *where* in the strings the

225    repetition was located. As a result, learners had to bind the output of the sameness

226    detection computations to some kind of representation of sequential positions,

227    which seems beyond the representations produced by Cope et al.'s (2018) model.

228

229 **References**

230 Awh, E., & Jonides, J. (2001). Overlapping mechanisms of attention and spatial

231     working memory. *Trends in Cognitive Sciences*, *5*(3), 119–126.

232 Carpenter, G. A., & Grossberg, S. (1987). A massively parallel architecture for a

233     self-organizing neural pattern recognition machine. *Computer Vision,*

234     *Graphics, and Image Processing*, *37*(1), 54–115.

235     http://dx.doi.org/10.1016/S0734-189X(87)80014-2

236 Cope, A. J., Vasilaki, E., Minors, D., Sabo, C., Marshall, J. A. R., & Barron, A. B.

237     (2018). Abstract concept learning in a simple neural network inspired by

238     the insect brain. *PLoS Computational Biology*, *14*(9), e1006435.

239     https://doi.org/10.1371/journal.pcbi.1006435

240 Downing, P. E. (2000). Interactions between visual working memory and

241     selective attention. *Psychological Science*, *11*(6), 467–473.

242     https://doi.org/10.1111/1467-9280.00290

243 Egner, T., & Hirsch, J. (2005). Cognitive control mechanisms resolve conflict

244     through cortical amplification of task-relevant information. *Nature*

245     *Neuroscience*, *8*(12), 1784–1790. https://doi.org/10.1038/nn1594

246 Engel, T. A., & Wang, X.-J. (2011). Same or different? A neural circuit

247     mechanism of similarity-based pattern match decision making. *Journal of*

248     *Neuroscience*, *31*(19), 6982–6996.

249     https://doi.org/10.1523/JNEUROSCI.6150-10.2011

250 Erika-Florence, M., Leech, R., & Hampshire, A. (2014). A functional network

251     perspective on response inhibition and attentional control. *Nature*

252     *Communications*, *5*, 4073. https://doi.org/10.1038/ncomms5073

253  Fan, J. E., & Turk-Browne, N. B. (2016). Incidental biasing of attention from

254      visual long-term memory. *Journal of Experimental Psychology. Learning,*

255      *Memory, and Cognition*, *42*(6), 970–977.

256      https://doi.org/10.1037/xlm0000209

257  Giurfa, M., Zhang, S., Jenett, A., Menzel, R., & Srinivasan, M. V. (2001). The

258      concepts of "sameness" and "difference" in an insect. *Nature*, *410*(6831),

259      930–933. https://doi.org/10.1038/35073582

260  Grill-Spector, K., Henson, R. N., & Martin, A. (2006). Repetition and the brain:

261      neural models of stimulus-specific effects. *Trends Cogn Sci*, *10*(1), 14–23.

262      https://doi.org/10.1016/j.tics.2005.11.006

263  Hampshire, A., & Sharp, D. J. (2015). Contrasting network and modular

264      perspectives on inhibitory control. *Trends in Cognitive Sciences*, *19*(8),

265      445–452. https://doi.org/10.1016/j.tics.2015.06.006

266  Hasselmo, M. E. (2005). The Role of Hippocampal Regions CA3 and CA1 in

267      Matching Entorhinal Input With Retrieval of Associations Between

268      Objects and Context: Theoretical Comment on Lee et al. (2005).

269      *Behavioral Neuroscience*, *119*(1), 342–345.

270  Hasselmo, M. E., & Wyble, B. P. (1997). Free recall and recognition in a network

271      model of the hippocampus: simulating effects of scopolamine on human

272      memory function. *Behavioural Brain Research*, *89*(1–2), 1–34.

273  Hochmann, J.-R., Carey, S., & Mehler, J. (2018). Infants learn a rule predicated

274      on the relation same but fail to simultaneously learn a rule predicated on

275      the relation different. *Cognition*, *177*, 49–57.

276      https://doi.org/10.1016/j.cognition.2018.04.005

277    Hochmann, J.-R., Mody, S., & Carey, S. (2016). Infants' representations of same

278         and different in match- and non-match-to-sample. *Cognitive Psychology*,

279         *86*, 87–111. https://doi.org/10.1016/j.cogpsych.2016.01.005

280    Johnson, J. S., Spencer, J. P., Luck, S. J., & Schöner, G. (2009). A dynamic

281         neural field model of visual working memory and change detection.

282         *Psychological Science*, *20*(5), 568–577. https://doi.org/10.1111/j.1467-

283         9280.2009.02329.x

284    Kumaran, D., & Maguire, E. A. (2007). Which computational mechanisms

285         operate in the hippocampus during novelty detection? *Hippocampus*,

286         *17*(9), 735–748. https://doi.org/10.1002/hipo.20326

287    Lisman, J. E. (1999). Relating hippocampal circuitry to function: recall of

288         memory sequences by reciprocal dentate-CA3 interactions. *Neuron*, *22*(2),

289         233–242.

290    Lisman, J. E., & Otmakhova, N. A. (2001). Storage, recall, and novelty detection

291         of sequences by the hippocampus: elaborating on the SOCRATIC model

292         to account for normal and aberrant effects of dopamine. *Hippocampus*,

293         *11*(5), 551–568. https://doi.org/10.1002/hipo.1071

294    Ludueña, G. A., & Gros, C. (2013). A self-organized neural comparator. *Neural*

295         *Computation*, *25*(4), 1006–1028. https://doi.org/10.1162/NECO_a_00424

296    Marcus, G. F. (1998a). Can connectionism save constructivism? *Cognition*, *66*(2),

297         153–182.

298    Marcus, G. F. (1998b). Rethinking eliminative connectionism. *Cognit Psychol*,

299         *37*(3), 243–82.

300    Marcus, G. F., Vijayan, S., Rao, S. B., & Vishton, P. (1999). Rule learning by

301   seven-month-old infants. *Science*, *283*(5398), 77–80.

302 Martinho, A., & Kacelnik, A. (2016). Ducklings imprint on the relational concept

303   of "same or different." *Science*, *353*(6296), 286–288.

304   https://doi.org/10.1126/science.aaf4247

305 Rose, N. S., LaRocque, J. J., Riggall, A. C., Gosseries, O., Starrett, M. J.,

306   Meyering, E. E., & Postle, B. R. (2016). Reactivation of latent working

307   memories with transcranial magnetic stimulation. *Science*, *354*(6316),

308   1136–1139. https://doi.org/10.1126/science.aah7011

309 Stokes, M. G. (2015). "Activity-silent" working memory in prefrontal cortex: a

310   dynamic coding framework. *Trends in Cognitive Sciences*, *19*(7), 394–

311   405. https://doi.org/10.1016/j.tics.2015.05.004

312 Wen, S., Ulloa, A., Husain, F., Horwitz, B., & Contreras-Vidal, J. L. (2008).

313   Simulated neural dynamics of decision-making in an auditory delayed

314   match-to-sample task. *Biological Cybernetics*, *99*(1), 15–27.

315   https://doi.org/10.1007/s00422-008-0234-0

316

317