# eNeuro

# A Multilevel Computational Characterization of Endophenotypes in Addiction

**Vincenzo G. Fiore[1], Dimitri Ognibene[2,3], Bryon Adinoff[4,5] and Xiaosi Gu[1,6]**

[1]*School of Behavioral and Brain Sciences, University of Texas at Dallas, 800 W. Campbell Road, RichardsonTX 75080-3021, USA*

[2]*Department of Computer Science and Electronic Engineering, University of Essex, ColchesterCO4 3SQ, UK*

[3]*Department of Information and Communication Technologies, Universitat Pompeu Fabra, Barcelona08018, ES*

[4]*University of Texas Southwestern Medical Center, 5323 Harry Hines Blvd, DallasTX 75390, USA*

[5]*VA North Texas Health Care System, 4500 S. Lancaster Rd, DallasTX 75216, USA*

[6]*Department of Psychiatry, Icahn School of Medicine at Mount Sinai, 1 Gustave L. Levy Place, New Yorkny 10029-5674, USA*

**Correspondence should be addressed to** either Vincenzo G. Fiore, School of Behavioral and Brain Sciences, UT Dallas, 800 W. Campbell Road, Richardson, TX 75080-3021, USA. E-mail: vincenzo.g.fiore@gmail.com[MAIL] or Xiaosi Gu, School of Behavioral and Brain Sciences, UT Dallas, 800 W. Campbell Road, Richardson, TX 75080-3021, USA. E-mail: xiaosi.gu@utdallas.edu

1    **Manuscript Title:** A Multilevel Computational Characterization of Endophenotypes in

2    Addiction

3

4    **Abbreviated Title:** Computational Endophenotypes in Addiction

5

6    **List all Author Names and Affiliations in order as they would appear in the published**

7    **article**

8    Vincenzo G. Fiore[1], Dimitri Ognibene[2,3], Bryon Adinoff[4,5], Xiaosi Gu[1,6]

9    1 - School of Behavioral and Brain Sciences, University of Texas at Dallas, 800 W. Campbell

10   Road, Richardson, TX 75080-3021, USA

11   2 - Department of Computer Science and Electronic Engineering, University of Essex,

12   Colchester, CO4 3SQ, UK

13   3 - Department of Information and Communication Technologies, Universitat Pompeu Fabra,

14   08018, Barcelona, ES

15   4 - University of Texas Southwestern Medical Center, 5323 Harry Hines Blvd., Dallas, TX

16   75390, USA

17   5 - VA North Texas Health Care System, 4500 S. Lancaster Rd., Dallas, TX 75216, USA

18   6 - Department of Psychiatry, Icahn School of Medicine at Mount Sinai, 1 Gustave L. Levy

19   Place, New York, NY 10029-5674, USA

20

21   **Author Contributions:**

22   **VGF** and **DO**: Designed research, Performed research, Analyzed data, Wrote the paper. **BA** and

23   **XG** Analyzed data, Wrote the paper

24    **Correspondence should be addressed to (include email address)**

25    Vincenzo G. Fiore, vincenzo.g.fiore@gmail.com

26    Xiaosi Gu, xiaosi.gu@utdallas.edu

27    School of Behavioral and Brain Sciences, UT Dallas

28    800 W. Campbell Road, Richardson, TX 75080-3021, USA

29

30    **Number of Figures:** 5

31    **Number of Tables:** 5

32    **Number of Multimedia:** 1

33    **Number of words for Abstract:** 221

34    **Number of words for Significance Statement:** 107

35    **Number of words for Introduction:** 832

36    **Number of words for Discussion:** 1695

37

47  **Abstract**

48  Addiction is characterized by a profound intersubject (phenotypic) variability in the expression

49  of addictive symptomatology and propensity to relapse following treatment. However, laboratory

50  investigations have primarily focused on common neural substrates in addiction and have not yet

51  been able to identify mechanisms that can account for the multifaceted phenotypic behaviors

52  reported in the literature. To fill this knowledge gap theoretically, here we simulated phenotypic

53  variations in addiction symptomology and responses to putative treatments, using both a neural

54  model, based on cortico-striatal circuit dynamics, and an algorithmic model of reinforcement

55  learning. These simulations rely on the widely accepted assumption that both the ventral, model-

56  based, goal-directed system and the dorsal, model-free, habitual system are vulnerable to extra-

57  physiologic dopamine reinforcements triggered by addictive rewards. We found that

58  endophenotypic differences in the balance between the two circuit or control systems resulted in

59  an inverted U-shape in optimal choice behavior. Specifically, greater unbalance led to a higher

60  likelihood of developing addiction and more severe drug-taking behaviors. Furthermore,

61  endophenotypes with opposite asymmetrical biases among cortico-striatal circuits expressed

62  similar addiction behaviors, but responded differently to simulated treatments, suggesting

63  personalized treatment development could rely on endophenotypic rather than phenotypic

64  differentiations. We propose our simulated results, confirmed across neural and algorithmic

65  levels of analysis, inform on a fundamental and, to date, neglected quantitative method to

66  characterize clinical heterogeneity in addiction.

67

68

69

70 **Significance statement**

71 Addiction is known to encompass heterogeneity in its development, maintenance, and treatment

72 response**.** While previous work has mostly focused on the common mechanisms underlying

73 vulnerabilities in addiction at a group level, the neurocomputational causes for such intersubject

74 variability in addition are not well-understood. To fill this knowledge gap, we combine a neural

75 and a reinforcement learning model to reveal that the balance between neural circuits or

76 computational control modalities characterizes the presence of behavioral phenotypes in

77 addiction. The presence of converging effects, validated across neural and algorithmic levels of

78 analysis, informs on a quantitative method to characterize clinical heterogeneity, and potentially

79 helps future development of precision treatments.

80 **Introduction.**

81 Addiction is known to encompass a wide range of individual behavioral differences (i.e.

82 phenotypes) in development, maintenance and severity of symptoms, and treatment response

83 (Everitt and Robbins, 2016). Previous investigations into the mechanisms underlying this

84 heterogeneity of behaviors have identified two fundamental neurocomputational alterations

85 correlated with vulnerability in the development and severity of addictive behaviors (Garrison

86 and Potenza, 2014; Jupp and Dalley, 2014; Belin et al., 2016). These neural and computational

87 intersubject differentiations (i.e. endophenotypes) include 1) a dysregulation of D2 receptors in

88 the striatum (Morgan et al., 2002; Nader and Czoty, 2005; Dalley et al., 2007; Flagel et al., 2014)

89 and 2) an alteration of learning rates within a reinforcement-learning framework (Gutkin et al.,

90 2006; Piray et al., 2010). However, these endophenotypic differences are found across a wide

91 spectrum of dissociable phenotypes, so that the same neural or computational mechanism is used

92 to account for separable behavioral traits. For instance, different forms of striatal D2

93 dysregulation are found in individuals differing in terms of their impulsivity (Dalley et al., 2007;

94 Volkow et al., 2007), social dominance (Morgan et al., 2002; Gould et al., 2014), motor

95 reactivity or preference for novelty (Flagel et al., 2010; Flagel et al., 2014), or sensitivity to

96 rewards (Belcher et al., 2014). Each of these behavioral traits is separately correlated with

97 development of addiction, but they do not necessarily coexist in the same individuals (cf. novelty

98 seeking and impulsivity: Ersche et al., 2010; Molander et al., 2011; Belin and Deroche-Gamonet,

99 2012). This mismatch between few known endophenotypic differences and a wide variety of

100 multifaceted, dissociable, behavioral phenotypes suggests there are yet unknown neural and

101 computational mechanisms that are responsible, alone or in interaction, for the reported

102 behavioral differentiations. Finally, investigations into intersubject variability often emphasize

103 the initial stage of addiction development (but see e.g.: Belin et al., 2008; Economidou et al.,

104 2009; Pelloux et al., 2015). Yet, individual differences also exist in treatment response, resulting

105 in diverse relapse patterns among individuals showing similar severity of symptoms. These

106 differences have not been so far addressed in previous neural or computational models.

107

108 Here we propose a theoretical investigation into the interaction between ventral and dorsal

109 cortico-striatal circuits and the associated behavioral control modalities. Several studies have

110 emphasized that addiction is associated with alterations of ventral and dorsal cortico-striatal

111 circuits, and of motivations and habits (Volkow and Morales, 2015; Everitt and Robbins, 2016;

112 Koob and Volkow, 2016). However, the role played by the interaction between the two neural

113 circuits or between the two behavioral control modalities in generating intersubject variability in

114 addiction, has been so far neglected. To investigate this interaction, we use two models to

115 simulate neural dynamics and algorithmic (or normative) choice selections in a multiple-choice

116 task involving drug and non-drug rewards. Then we test these models under different conditions

117 of circuit or control modality dominance (i.e. simulated endophenotypes). Consistently with

118 previous models, we assume addictive substances hijack the healthy reward prediction error

119 signal (Schultz et al., 1997) by triggering extra-physiologic dopamine bursts (Nestler and

120 Aghajanian, 1997; Koob and Volkow, 2016). These dopamine activities signal the presence of an

121 aberrant unexpected reward, leading to the repetition of drug-related actions and escalation of

122 consumption (Redish et al., 2008; Dayan, 2009). In our neural model, this process of

123 reinforcement learning (RL, Sutton and Barto, 1998) is mediated by extra-physiologic changes in

124 cortico-striatal connectivity weights (Hyman et al., 2006; Haber, 2008; Koob and Volkow,

125 2016). These changes in turn aberrantly affect circuit gain and the stability of both ventral and

6

126    dorsal cortico-striatal circuits, disrupting their respective roles in encoding and selecting goal-

127    directed behaviors (Balleine, 2005; Balleine and O'Doherty, 2010; Gruber and McDonald, 2012)

128    and habitual responses (Yin et al., 2004; Balleine and O'Doherty, 2010). A similar effect is

129    assumed for our algorithmic model, where over-evaluation of drugs and related RL affect the

130    two control modalities, termed *model-based* and *model-free*, that approximate ventral/goal-

131    oriented and dorsal/habitual implementations (Dolan and Dayan, 2013; Voon et al., 2017). As a

132    result, and consistently with previous formulations of RL models of addiction (Redish et al.,

133    2008; Piray et al., 2010; Gillan et al., 2016), both the planned evaluation of known action-

134    outcome contingencies, represented in an *internal model* of the world, and the reactive

135    immediate motor responses are biased towards drug-related selections.

136

137    Based on these assumptions, our models show that phenotypic differentiation in addiction

138    development and treatment response can emerge as a function of the interaction between ventral

139    and dorsal circuits or model-based and model-free control modalities. Our simulated results offer

140    a proof-of-concept that this interaction is a candidate independent neural and computational

141    mechanism underlying addiction vulnerability, putatively characterizing three different

142    endophenotypes differing in the likelihood to develop addiction, severity of symptoms and

143    treatment response. We suggest this neurocomputational mechanism could interact with both

144    previously described D2 receptors dysregulation in the striatum (Dalley et al., 2007; Flagel et al.,

145    2014) and altered learning rates (Gutkin et al., 2006; Piray et al., 2010) to generate the variety of

146    dissociable behavioral traits reported in literature as associated with addiction vulnerabilities.

147

148    **Materials and Methods.**

149    In brief, we present two complementary models simulating endophenotypic differences and their

150    effects on addiction development and treatment response. In the models, intersubject differences

151    are expressed in terms of either neural circuit dominance (i.e. ventral or dorsal circuit) or control

152    modality dominance (i.e. model-based or model-free) in determining behavioral selections. The

153    resulting phenotypes are tested in environments granting free access to a simulated substance of

154    addiction, as usually implemented in laboratory studies. In particular, we compare our simulated

155    phenotypic variability with the results described in a recent study investigating individual

156    differences in rats self-administrating the stimulants cocaine or a designer drug, a dopamine- and

157    mixed dopamine-norepinephrine reuptake inhibitor, respectively (Gannon et al., 2017). We

158    selected this study because it highlights how different drugs, dosages, and tasks result in different

159    ranges of phenotypic differentiation. For instance, an initial acquisition phase, over a 10-day

160    period, shows compulsive behavior developed in up to 75% rats self-administering cocaine and

161    87.5% of those exposed to the designer drug. Furthermore, under a condition of fixed ratio (=5)

162    schedule, the study shows self-administration varied significantly among subjects. A subset of rat

163    population, termed high responders, self-administered cocaine up to 60% more times in

164    comparison with a different subset, termed low responders, depending on dosage (cf. figure 3 in:

165    Gannon et al., 2017). Importantly, the task setup chosen for both of our proposed models

166    involves the selection of a drug reward over explicit non-drug related alternatives; in contrast,

167    the chosen empirical study utilizes a time-out responding paradigm, where the only explicit non-

168    drug related behavior (a lever-press) is not rewarded. As for most studies simulating addiction

169    (e.g. see: Redish, 2004), we believe the choice to present our simulated agents with a richer set

170    of options (i.e. more than one) does not invalidate a parallel between simulated and real data. We

171    consider the simulated competing options as a proxy for the many conflicting stimuli and

172   associated behaviors that animals have access to, even in the limited environment of a standard

173   operant conditioning chamber. Thus, our focus is on perturbing the balance between the

174   dorsal/model-free and the ventral/model-based systems, to compare our simulated behavioral

175   differentiations in the escalation and compulsive selection of drug-related actions with the data

176   reported in the chosen laboratory study.

177

178   The two models comprise a neural mass model that has been validated and described in the

179   context of choice behavior and dopaminergic modulation (Fiore et al., 2016; Hauser et al., 2016;

180   Fiore et al., 2018) and a normative or algorithmic model based upon standard RL schemes

181   (Sutton and Barto, 1998). In the neural model, addiction and treatment response are modeled

182   through DA-dependent associative plasticity in both ventral and dorsal circuits. In the RL model,

183   aberrant learning is modeled using a duplex of model-based and model-free schemes that

184   competed for control over action selection. The model-based scheme entails learning a model of

185   the environment (in the form of probability transition matrices among states) that is used to

186   compute value functions under the Bellman optimality principle (Bellman, 1966). The equivalent

187   model-free scheme uses prediction error-based learning to directly acquire the value of state

188   action pairs. Both neural and RL models are tested under four successive stages or phases; 1)

189   before exposure to the simulated drug (termed *pre-drug*); 2) learning of addictive behavior

190   (termed *addiction*); 3) simulated ideal therapeutic interventions (termed *treatment*) that partially

191   revert the learning of the previous phase. Finally, 4) reinstated access to the simulated drug

192   following each treatment (termed *relapse*). The simulated treatments are conceived to emphasize

193   endophenotypic response and relapse differentiation and therefore they predominantly affect

194   only one control system, targeting either the goal-oriented/model-based or the habitual/model-

195    free. The former treatment is assumed to modify only the internal model of the environment and

196    related selection of action-outcome contingencies performed in the ventral circuit. The latter

197    treatment represents a condition in which the model of the world of the agent remains mainly

198    unaltered, but the acquired drug-related stimulus-response associations are disrupted, thus

199    preventing the agent from exhibiting habitual responses (cf. Doll et al., 2009).

200

201    The unique aspect of this complementary modeling approach is that converging results from

202    neural and algorithmic models can validate each other, as process and implementation theories

203    (i.e., synaptic and dynamical mechanisms) complement the normative principles formalized in

204    the RL model.

205

206    **Neural field model.**

207    *Basic model architecture and parameterization:* In cortico-striatal circuits, the signal processed

208    in the cortex is conveyed towards its respective area of the striatum, processed in basal ganglia

209    and finally relayed to the same cortical area where it originated, via thalamus (Haber, 2003;

210    Draganski et al., 2008; Jahanshahi et al., 2015). Thus, despite diverging in terms of the

211    information processed –e.g. sensorimotor or rewards and outcomes– these circuits are

212    characterized by similar computational dynamics (Obeso et al., 2014). Temporal responses in

213    recurrent neural networks co-occur with state transitions or input transformations that are often

214    described in terms of energy landscapes (**Figure 1A-C**). If multiple inputs or initial states

215    generate transitions towards the same final state, this is termed *attractor state* (Amit, 1989). In

216    recurrent networks such as cortico-striatal circuits, learning processes modulate the circuit gain,

217 thereby affecting the strength of the attractor states and the overall stability of the system (Fiore

218 et al., 2015; Fiore et al., 2016; Hauser et al., 2016).

219

220 We simulate the temporal responses in cortico-striatal circuits in a neural model (illustrative

221 representation of the neural architecture is represented in **Figure 1D**). This neural model

222 simulates mean-field activity (Deco et al., 2008) within multiple channels of both dorsal and

223 ventral cortico-striatal loops. A continuous-time differential equation simulates changes over

224 time ($\tau_g$) of the average action potential ($u_j$) of a pool of neurons (equation 1), and a positive

225 transfer function (equation 2) converts this action potential in the final activation of the pool ($y_j$).

226 Finally, the plasticity of the connections ($w_{ij}$) between cortex and striatum is characterized by

227 DA-dependent Hebbian learning, corrected with a constant threshold (*th*) as defined in equation

228 3. The resulting rule strengthens the connections among all active nodes in the cortex and those

229 active in the striatum, and weakens the connections among nodes showing opposite activation

230 status.

231

232 $\tau_g \dot{u}_j = -u_j + b_j + (\epsilon + \lambda d) \sum w_{ji} y_i$ (1)

233 $y_j = [\tanh(u_j - \theta)]^+$ (2)

234 $\Delta w_{ij} = \eta \left([y_i - th]^+ [y_j - th]^+ [d - th]^+\right) - \zeta \left([th - y_i]^+ [th - y_j]^+ [d - th]^+\right)$ (3)

235

236 The input ($\sum w_{ji} y_i$), reaching each node in the neural network is modulated by two coefficients $\lambda$

237 and $\epsilon$. These regulate the ratio between the signal affected by the presence of dopamine release $d$

238 and the amount of signal that is computed independent of dopamine release. For most units, the

239 values of the two coefficients are set to $\lambda = 0$ and $\epsilon = 1$, with the exception of the simulated

240     striatal units, where these parameters are set to $[\lambda = 1.4, \epsilon = 0.2]$ and $[\lambda = -0.5, \epsilon = 0.6]$, to

241     simulate the differential effect dopamine has, depending on the most prevalent receptor type

242     ($\lambda > 1$ and $\lambda < 0$ for D1 and D2 receptors, respectively). Due to the different effects the dopamine

243     receptors have on the activity of the simulated neurons, the drug-induced dopamine-dependent

244     Hebbian learning significantly affects D1-enriched units in the striatum, whilst having negligible

245     effects on D2-enriched units (Gerfen and Surmeier, 2011; Volkow and Morales, 2015).

246

247     *Simulating different addiction phenotypes and treatment effects*: Agents controlled by the neural

248     model are immersed in a simplified environment and can select among three arbitrary actions or

249     inactivity (cf. non-stationary three armed bandit environment). The selection of the actions is

250     carried out in the circuit simulating the dorsal cortico-striatal activity and it is considered

251     completed if the neural activity of any of the units in the external layer of the simulated cortex

252     (cf. **Figure 1D**) is maintained for at least 2 seconds. Ventral and dorsal circuits interact, both

253     ways, via cortico-cortical connectivity. Therefore, the activity in the simulated ventral circuit

254     biases action selection in the dorsal circuit and the selection of actions in the dorsal circuit biases

255     the activity in the ventral circuit. To test our hypothesis about the effect these reciprocal biases

256     have on choice behavior, we assumed cortico-cortical weights do not vary over time and we

257     tested eleven combinations for the parameters determining their weights, as $w_{ji}$=[0.02-0.2],

258     [0.03-0.17], [0.03-0.15], [0.05-0.15], [0.07-0.13], or [0.1-0.1] (and symmetrical). This spectrum

259     of weights describes the strength of the biases between the two major circuits, thereby

260     characterizing either a balanced condition or a dominance of one of the two circuits. We report

261     the effects in terms of behavioral responses for these putative endophenotypes and test each of

262     these with thirty noise seeds, random inputs and under four stages, to allow within phenotype

263    comparisons. The first stage, "pre-drug", represents an assessment of behavior before any drug

264    or reward is introduced, as the three available inputs randomly change their value to determine a

265    non-stationary order of preferences. Under the second stage, termed "addiction", one action is

266    associated with the administration of a simulated addictive substance, triggering DA phasic

267    responses and associated Hebbian learning in cortico-striatal connections of both ventral and

268    dorsal circuits. For the third stage, termed "treatment", we simulate the effects of deprivation

269    coupled with one of two hypothetical treatments targeting either the dorsal or the ventral cortico-

270    striatal circuits. The treatments are simulated by reverting the learning process in either the

271    dorsal or the ventral cortico-striatal circuit, respectively representing an intervention that would

272    block or extinguish either the habitual drug-related response (an ideal behavioral treatment) or

273    the drug-related emotional and value association (an ideal cognitive treatment). The dorsal

274    treatment brings back the pre-drug configuration in the dorsal circuit and keeps the configuration

275    reached under the addiction stage for the ventral circuit. The ventral treatment is achieved with

276    the opposite intervention. Finally, during the fourth stage, termed "relapse", we reintroduce

277    access to the simulated addictive substance, inducing relapse. For this stage, relapse time is

278    defined as the time required to reinstate the configuration of cortico-striatal weights found at the

279    end of the *addiction* stage.

280

281    **RL model.**

282    *Basic model architecture and parameterization*: In this model, we assume that the behavior of

283    the agent relies on a hybrid model (Daw et al., 2011) that learns and computes the value of

284    choices (actions, $a_t$) under each condition (state, $s_t$). Value is defined as a quantity that combines

285    short and long term expected rewards and negative outcomes when a specific strategy of action

286    is followed (policy, $\pi$). It is formally defined as:

287

288    $Q^{\pi}(s_t a_t) = r(s_t a_t) + E\left[\sum_i^{\infty} \gamma^i r(s_{t+i}, a_{t+i} = \pi(s_{t+i}))|s_t, a_t\right]$          (4)

289

290    In equation (4), $r(s,a)$ denotes the instantaneous reward received when action $a$ is performed in

291    state $s$. $\gamma$ is a discount factor, comprised between 0 and 1, which defines the trade-off between

292    immediate and long term rewards. The value of a state given the policy is defined as $V^{\pi}(s) =$

293    $\max_a Q^{\pi}(s,a)$. For each environment there is an optimal policy $\pi^*(s)$, which maximizes the

294    value $V^{\pi^*}(s)$ for every state (Sutton and Barto, 1998).

295

296    The environment can be completely characterized through the state transitions distributions

297    $p(s_{t+1} = s|s_t, a_t)$, and the expected rewards $E(r|s,a) = R(s,a)$. These two functions together

298    represent a model of the environment. Model-based behaviors compute $Q^{\pi}(s_t a_t)$ and the policy

299    relying on such functions, at each state, following the Bellman equation (Daw and Dayan, 2014):

300

301    $Q^{\pi^*}(s_t, a_t) = R(s_t, a_t) + \gamma \sum_s \left[p(s_{t+1}=s|s_t, a_t) \max_a Q^{\pi^*}(s,a)\right]$          (5)

302

303    The model-based component learns the transition distributions and the expected rewards during

304    the interaction with the environment. Thus, differently from other hybrid models (Daw et al.,

305    2005; Keramati et al., 2011; Pezzulo et al., 2013), the quality of Q value estimation at any given

306    moment depends on the experience the agent acquired up to that point in time. To compute value

307    estimation ($Q_{MB}$), this bounded (Gershman et al., 2015) component applies at each step the

308    Bellman equation (5) a limited number of times ($N_{PS} = 50$) to states sampled stochastically

309    following a heuristic for efficient state update selection. The algorithm is an early-interrupted

310    variation of the Prioritized Sweeping algorithm (Moore and Atkeson, 1993) with stochastic state

311    update selection. Crucially, our model-based component does not accumulate the variations of Q

312    values over time, and restarts the computation after each step (desJardins et al., 1999). This

313    choice is meant to instate a plausible bounded rationality for our model which can account for

314    the cognitive costs and ensuing limits of integrating old and new information about the

315    environment, whilst updating and extending a complex plan to navigate it. This implementation

316    is suitable for a bounded rational model-based component that shows controlled stochasticity of

317    deliberation performances in non-trivial environments. This choice allows to test the effects of

318    the hypothesized endophenotypic differentiation in an environment characterized by higher

319    degree of complexity in comparison with both the one chosen for the neural model and those

320    described in the literature of RL models of addiction. In particular, we consider drug

321    consumption to be associated with complex after-effects that make it difficult to predict the

322    overall result of pursuing the related action course.

323

324    In comparison with other hybrid models such as Dyna and Dyna2 (Sutton, 1990; Silver et al.,

325    2016), the proposed architecture does not share Q values between model-based and model-free

326    components, nor it requires that the two processes share the same state representations. The two

327    components separately represent their Q values and integrate them in a later phase. This

328    decoupling is assumed to result in a more biologically plausible agent (Daw & Dayan 2014), and

329    it is essential for the simulations of two separate treatments, essential requirement to establish a

330    comparison with the behavior simulated with the neural model. In contrast with previous work

331    using a hybrid Dyna-like architecture and Prioritized Sweeping algorithm, where the sharing of

332    the Q-values explained the appearance of model based drug oriented behavior (Simon and Daw,

333    2012), in our simulations this model based addiction emerges in independent model-free and

334    model based components. Thus, addiction behavior results from the joint effect of high reward

335    (i.e. the drug), a limited number of stochastically selected policy updates and limited knowledge

336    of the environment.

337

338    The model-free component has been implemented using the Q-Learning algorithm in tabular

339    form (Watkins and Dayan, 1992). Q-learning updates initial state value estimations as follows:

340

341    $Q_{MF}(s_t, a_t)_{new} = Q_{MF}(s_t, a_t)_{old} + \alpha \delta_t$         (6)

342    $\delta_t = R(s_t, a_t) + \gamma \max_{a'}[Q_{MF}(s_{t+1}, a')] - Q_{MF}(s_t, a_t)$     (7)

343

344    where α is a learning factor comprised between 0 and 1. Our hybrid model computes choice

345    values in a fashion that balances model-free (MF in the equations) and model-based (MB in the

346    equations) components depending on a parameter $\beta$. Six values (1, 0.8, 0.6, 0.4, 0.2, 0) are used

347    for this parameter to simulate different endophenotypes, on a spectrum between purely model-

348    based ($\beta$=1) and purely model-free ($\beta$=0) RL.

349

350    To allow exploration, the action to execute is selected randomly 10% of the times. This

351    exploration factor is kept constant to support adaptation to a changing environment (Singh et al.,

352    2000) and to simulate the continuous update of knowledge necessary to cope with ecological

353    environments. The remaining 90% of the times, actions are determined by maximizing $Q_{MX}(s,a)$

354 in a strategy defined as ε-greedy (ε=.1). These values are produced by combining the values

355 computed by the model-based and model-free components:

356

357 $$Q_{MX}(s,a) = \beta Q_{MB}(s,a) + (1-\beta)Q_{MF}(s,a) \tag{8}$$

358

359 The choice for a fixed balance between model-based and model-free requires minimal

360 assumptions on their interaction and has been used in recent reinforcement learning architectures

361 (Silver et al., 2016).

362

363 *Simulating different addiction phenotypes and treatment effects:* In comparison with the

364 simulations characterizing the neural model, a more complex environment is in use for the RL

365 model to highlight how our endophenotypic differentiations can also affect the likelihood to

366 develop addiction. This environment is characterized by a total of 20 states divided into four

367 different types (**Figure 2**): (i) healthy rewards (i.e. normal rewards that are not directly

368 associated with drugs); (ii) neutral states (no reward or negative outcome); (iii) drug-related

369 states, which give a high reward but are followed by multiple (iv) drug aftereffects, characterized

370 by small negative outcomes. Similar to the neural model investigations, the agent deals with

371 environment variations meant to simulate four phases of addiction: initial pre-drug phase (f1);

372 addiction (i.e. the drug becomes accessible for the first time, f2); treatment (f3); relapse (i.e.

373 second drug exposures, f4). Under the initial pre-drug phase ($d_{init}$=50 steps), the agent does not

374 receive any reward or negative outcome by entering the drug-related and aftereffects area, but a

375 moderate reward is assigned ($R_g$=1) by accessing the healthy reward state. Under the phases of

376 addiction and post-treatment addiction ($d_{tpy}$=1000 steps), the agent can also receive a high

377    reward, after accessing a drug-related state ($R_d$=10). The drug state always leads to a series of

378    randomized state transitions among the aftereffects states ($R_a$=-1.2) and simulates generic

379    negative consequences associated with addiction. The agent can occasionally leave this

380    aftereffect area of the environment (**Figure 2**) to reach a neutral state, at the price of a further

381    negative outcome ($R_a$=-4). Under the treatment phase ($d_{tpy}$=1000 steps), the drug-related state

382    results in a negative outcome (Rdt=-1, see **Tables 1** and **2**, column f3), thus increasing the

383    chances the agent stops pursuing this state. To allow for a comparison with the results in the

384    neural model, we simulate a model-based and model-free treatment by manipulating the learning

385    factor of the non-treated control modality, decreasing it: $\alpha_{Ctpy}$=0.01 * $\alpha$. Under the relapse phase,

386    we measure the simulated time required by the agents to reach at least 95% of drug-related action

387    preference as recorded under the addiction phase, after the drug is introduced again in the

388    environment. This threshold is used to measure the percentage of agents relapsing, as well as the

389    time required to complete the relapse, per endophenotype.

390

391    **Code Accessibility**

392    All models rely on custom code developed in Matlab (optimized for R2014b) that has been run

393    successfully on multiple OS (iOS, Linux and Windows) on different computers and local servers.

394    The code can be accessed at any time from the repository ModelDB (accession number: 239540;

395    https://senselab.med.yale.edu/modeldb/enterCode.cshtml?model=239540; title: 'Computational

396    Endophenotypes In Addiction: Source Code'). The downloadable archive file consists of two

397    folders (respectively for the neural model and the RL model), which include the entire source

398    code required to replicate the data reported in our Results section. Code available as Extended

399    Data 1.

400

**Results.**

**Simulations from the neural field model**.

403 During all stages, the three stimuli randomly change every few seconds, putatively representing a

404 dynamic fluctuation of values associated with perceived cues in a non-stationary environment.

405 This setup requires the agents to rapidly adapt to these changes, transiently triggering the motor

406 response associated with the most valuable cue, in order to achieve optimal behavior. During the

407 pre-drug stage, dorsal and ventral circuits perform unbiased selections, collaborating in the

408 generation of a near-optimal sequence of motor selections. All eleven endophenotypes show

409 uniform distributions of action selections, complying with the random distribution of the inputs

410 configurations (**Figure 3A**). This control stage allows the simulated network to generate

411 transient temporal responses that couple multiple initial states with multiple stable states, in a

412 transient *winner-take-all* or *winner-less* competition (Rabinovich et al., 2006; Afraimovich et al.,

413 2008).

414

415 During the simulated addiction stage, one of the actions is associated with drug administration

416 (**Figure 3B**, values represented in blue). Substance use triggers phasic dopamine bursts, leading

417 to Hebbian learning in cortico-striatal connections of both dorsal and ventral circuits (equation

418 3). In recurrent networks, circuit gain increases as a direct function of the weights of reentrant

419 synapses (Amit, 1989). A dopamine response triggered by healthy unexpected rewards would

420 create a bias towards the selection of the reinforced motor response to a perceived cue (Cohen

421 and Frank, 2009; Grahn et al., 2009; Baldassarre et al., 2013). However, drug consumption

422 triggers extra-physiologic dopamine-dependent learning, which in our model results in aberrantly

19

423  high circuit gain, compromising the ability of all affected circuits to discriminate among different

424  inputs and produce temporal transitions towards multiple stable states (cf. Fiore et al., 2014). The

425  cortico-striatal circuits become over-stable and resistant to perturbation caused by a change of

426  input or by noise as they are dominated by *parasitic attractors* (Hoffman and McGlashan, 2001)

427  (**Fig 1C**). In the ventral cortico-striatal circuit, a parasitic attractor sets and maintains the

428  selection of drug-related goals or outcomes, biasing the action-outcome assessments required for

429  planning. In the dorsal circuit, the same process determines over-stable selections of the

430  reinforced motor behavior, generating reactive responses and habits. Importantly, the learning

431  process simulated in our neural model leads to the generation of parasitic attractors in both

432  circuits across all endophenotypes, as all agents eventually reach a fixed threshold in cortico-

433  striatal neural plasticity. Despite the generation of a form of compulsive drug seeking behavior

434  across all endophenotypes, we observe significant differences in motor response patterns as a

435  function of the balance between ventral and dorsal circuits. Specifically, the endophenotypes

436  characterized by unbalanced dorsal or ventral control (i.e. **Figure 3B**, endophenotypes 1-3 and 9-

437  11) express distributions of motor selections that are significantly more compromised by drug-

438  related aberrant rewards, in comparison with balanced endophenotypes (i.e. **Figure 3B**,

439  endophenotypes 5-7). The presence of identical learning processes, and the associated attractor

440  formation in both ventral and dorsal circuits, ascribes all phenotypic differences univocally to the

441  only remaining independent variable, which controls cortico-cortical connectivity and therefore

442  the strength of the biases between circuits. Unbalanced agents are characterized by more frequent

443  drug-related selections as actions leading to drug consumption are selected more frequently than

444  in balanced endophenotypes, in a range between +3% and +45%. This result identifies all

20

445    phenotypes within the limits of individual differentiation described in the study chosen for

446    behavioral comparison (Gannon et al., 2017).

447

448    Next, we investigate how the simulated endophenotypes behave during the stages of treatment

449    and relapse. First, we measure the frequency of drug-related action selections during the stages

450    of addiction and treatment (**Figure 4A-B**). Both ventral (goal-oriented) and dorsal (habitual)

451    treatments effectively reduce the number of actions associated with drug consumption, in

452    comparison with baseline addiction. However, the dorsal treatment is more effective for dorsal-

453    dominated endophenotypes and the ventral treatment is more effective for ventral-dominated

454    endophenotypes. These endophenotype-specific treatment effects are further confirmed by our

455    analysis of individual differences under the relapse stage (**Figure 4C-D**): dorsal treatments are

456    more effective in elongating time to relapse for dorsal-dominated endophenotypes, whereas

457    ventral treatments are more successful in delaying relapse for ventral-dominated

458    endophenotypes. This analysis shows that simulated treatments focusing either on the dorsal

459    circuit (and therefore habitual responses) or the ventral circuit (and therefore motivational

460    responses) can have substantially different effects, depending on the balance between dorsal and

461    ventral circuits. Importantly, these differences emerge only after the treatment is applied, where a

462    pre-treatment comparison between compulsive behaviors expressed by the opposite unbalanced

463    endophenotypes (i.e. ventral-dominant or dorsal-dominant) does not show any significant

464    difference in choice selections (cf. **Figure 3B**, endophenotypes 1-3 and 9-11).

465

466    **Simulations from the RL model**.

21

467    By simulating explicit negative outcomes associated with drug consumption, the RL model

468    allows to measure the likelihood each agent has to develop addiction, as a function of its

469    endophenotype. In our analysis, addiction is defined as a behavior leading to drug selections

470    more frequently than the healthy alternative reward, under the addiction phase. The mean

471    percentage of these *addicted* agents (over 300 runs) was 43.05%, across endophenotypes, which

472    is consistent with the percentage of rats developing compulsive self-administration of cocaine, as

473    reported in the reference study (~40% over a period of 5 days, cf. Gannon et al., 2017).

474    Importantly, when considering endophenotype differentiation, the percentage varies

475    significantly: 60.3% for $\beta$ =0, 40.3% for $\beta$ =0.2, 30.1% for $\beta$ =0.4, 36.7% for $\beta$ =0.6, 39.3% for

476    $\beta$ =0.8, and 51.6% for β=1 (**Figure 5A-B**). This phenotypic differentiation is consistent with

477    well-established data from animal models. For instance, rat strains selectively bred for either

478    high or low voluntary running differ in the likelihood to develop addiction when given free

479    access to cocaine (respectively ~35% and ~60% of each strain develop addiction over a period of

480    5 days, cf. Smethells et al., 2016). Free access to substances of abuse does not necessarily lead to

481    compulsive behaviors (Piazza et al., 1989; Belin et al., 2011), as addiction varies as a function of

482    factors such as exposure extent, amount of drug delivered, and associated negative effects

483    (Pelloux et al., 2007; Jonkman et al., 2012). Our simulations suggest that endophenotypes with

484    lower chances of addiction are characterized by balanced control modalities. Note that an

485    optimal agent, knowing the environment structure and being able to compute the long-term

486    effects of drug, will never select drug-states (**Table 3**).

487

488    Finally, the simulations suggest that the hypothetical treatment targeting model-free control is

489    the most effective, reducing the likelihood to pursue drug-related behaviors for all

490    endophenotypes (**Figure 5A)**. In contrast, the model-based treatment appears to be less effective

491    for all endophenotypes, with the exception of the purely model-based one ($\beta$=1) (**Figure 5B)**.

492    Under the relapse phase, our data confirm that the simulated treatments significantly differ in

493    their effectiveness across the proposed endophenotypes, also suggesting the treatment targeting

494    model-free control is the most successful in prolonging relapse time (**Figure 5C-D**). Relapse

495    time after model-free treatment is mostly similar to the time required to develop addiction

496    behavior before any treatment (**Figure 5C**). At the opposite side of the control spectrum, the

497    model-based treatment shows a positive effect only for the purely model-based endophenotype.

498    All remaining endophenotypes show relapse times significantly shorter than those recorded for

499    the first development of addiction ($\beta$=1; **Figure 5D**).

500

501    **Discussion**

502    Individual differences in stress and anxiety responses (Dilleen et al., 2012; Jimenez and Grant,

503    2017), social dominance (Morgan et al., 2002; Covington and Miczek, 2005), aggressive

504    temperament (McClintick and Grant, 2016), preference for saccharine (Carroll et al., 2002),

505    sensation or novelty seeking (Suto et al., 2001; Nadal et al., 2002; Belin et al., 2011; Flagel et al.,

506    2014), impulsivity (Perry and Carroll, 2008; Verdejo-Garcia et al., 2008; Dalley et al., 2011),

507    and sensitivity to rewards (Belcher et al., 2014) have all been found in both animal models and

508    clinical studies in humans to be associated with addiction vulnerabilities, and in particular with

509    the likelihood to develop and maintain addiction, or to resist to treatment (Piazza et al., 1989;

510    Belin et al., 2016; Everitt and Robbins, 2016). However, investigations into the mechanisms

511    underlying this phenotypic differentiation in addiction has so far revealed few neural or

512    computational candidates, which are found to be associated with diverse and dissociable

23

513    behavioral traits. An important example is represented by the endophenotypic differentiation

514    reported in the expression and reactivity of striatal D2 dopaminergic receptors, which is found to

515    be negatively correlated with the traits of impulsivity (Dalley et al., 2007), social dominance

516    (Morgan et al., 2002), and sensitivity to rewards (Belcher et al., 2014) and non-linearly

517    correlated with novelty preference (Flagel et al., 2014). The overlap of this endophenotypic trait

518    across multiple, non-coexisting, phenotypes associated with addiction vulnerabilities suggests

519    other neural or computational mechanisms have yet to be identified to allow accounting for the

520    reported variety in behavioural traits.

521

522    Here we have presented a neural field model, augmented by an RL model, to expand on existing

523    neuropsychological and computational accounts of addiction. Our models propose a theoretical

524    investigation into the interaction among cortico-striatal circuits or behavioral control modalities,

525    and the effects this interaction has on addiction development and treatment response. As

526    described in classic models (Redish, 2004; Redish et al., 2008; Dayan, 2009), we have assumed

527    that over-evaluation of a drug leads to aberrant dopamine release and associated over-learning in

528    multiple DA targets (Volkow and Morales, 2015; Koob and Volkow, 2016). In the neural field

529    model, this mechanism results in the dysregulation of the circuit gain and associated dynamics of

530    both ventral and dorsal cortico-striatal circuits (Fiore et al., 2014; Hauser et al., 2016). In the

531    integrated model-based and model-free RL model, sequential choice behavior is confounded by

532    the presence of a high immediate reward (drug state). This leads to misrepresent the negative

533    outcomes following drug consumption, if their distribution across states and time is sufficiently

534    complex to escape the capabilities of the agent to correctly represent the environment (Doll and

535    Daw, 2016; Sadacca et al., 2016). We found that both models jointly indicate that the balance

24

536 between neural circuits or behavioral control modalities is a candidate neurocomputational

537 mechanism characterizing endophenotypes in addiction. The neural and RL models converge in

538 suggesting that individuals characterized by balanced behavioral control between reward-seeking

539 or planning (ventral circuit/model-based) and reactive or habitual responses (dorsal

540 circuit/model-free) would have a reduced chance to develop addiction and decreased severity of

541 symptoms if developing addiction. We propose that this neurocomputational mechanism may be

542 interacting with other known endophenotypic differentiations, such as alterations of D2 receptors

543 in the striatum (Morgan et al., 2002; Nader and Czoty, 2005; Dalley et al., 2007; Volkow et al.,

544 2007; Belcher et al., 2014; Flagel et al., 2014) or differences in learning rates (Gutkin et al.,

545 2006; Piray et al., 2010), to generate the multifaceted behavioral traits that have been reported in

546 literature to be associated with addiction vulnerabilities.

547

548 In our neural model, ventral and dorsal circuits are mostly in phase in their selections under the

549 pre-drug stage, exhibiting synchronous transient stability of neural activity and enhancing the

550 overall ability of the system to adapt to changing stimuli (i.e. the two circuits adapt to the input

551 changes with a similar pace and synchronize in their selection). Under the addiction stage, the

552 two circuits are mostly pulled towards the parasitic attractor state associated with drug

553 consumption, and they occasionally select the competing non-drug stimuli. If only one of the two

554 systems performs a selection outside of the attractor, the difference in selection generates a

555 dissonance or interference. In neural endophenotypes characterized by unbalanced control, this

556 dissonance is solved by one circuit taking the lead, so that both systems eventually converge on

557 the selection of the dominant circuit. These dynamics result in limited opportunities to generate

558 non-drug related responses to the external stimuli, as they can only be generated by the dominant

559 circuit. Conversely, in balanced control endophenotypes, if any of the two circuits ignores the

560 drug-stimulus and selects a competing option, the resulting dissonance can trigger a state

561 transition pulling out the parasitic attractor states associated with substance use. The

562 endophenotypes in our simulations vary only in the parameters regulating the balance between

563 circuits, as dopamine-driven learning processes established between cortex and striatum

564 (equation 3) do not vary across endophenotypes, resulting in identical habit formation and drug-

565 related biases in the outcome representations. Thus, our proposed phenotypic differentiation does

566 not interfere with the usual role ascribed to the ventral and dorsal circuits as respectively

567 implicated in the initial reward-seeking phase in addiction (Belin and Everitt, 2008; Willuhn et

568 al., 2012) and the subsequent consolidation of stimulus-response, habitual, association (Everitt

569 and Robbins, 2013, 2016). However, our simulated dynamics show that, after addiction is

570 developed, systemic over-stability can be reduced or further enhanced, depending on the cortico-

571 cortical biases between cortico-striatal circuits. In turn, this modulation of system stability can

572 foster or further impair input discrimination and motor response versatility, affecting addiction

573 symptomatology. As a result, our neural model shows phenotypic variability emerging after the

574 presentation of the reward simulating the drug and addiction is developed, in a gradient of over-

575 selection of drug-related actions.

576

577 With the RL model, we investigate whether the balance between model-based and model-free

578 modalities would also increase the robustness of the system against the selection of drug-states in

579 a more complex environment and in presence of explicit negative outcomes. Similar to the neural

580 model, a system with balanced control modalities introduces more diversity in action selection

581 during exploration, reducing (yet not cancelling) the chances of developing maladaptive reactive

582    responses. This increased diversity and overall reliability are likely to be induced by a higher

583    redundancy and diversification of the system. While both components may fail, the causes of

584    failures are not necessarily correlated. The model-based system can fail due to its sensitivity to

585    cognitive resources but it is more efficient in encoding previous experience of the agent. On the

586    other hand, the model-free component is affected by limited exploration but it is reliable in its

587    selections, which are not affected by the availability of cognitive resources. Consistent with the

588    neural model, differentiations in behaviors among endophenotypes emerge in an inverted-U

589    shape, where unbalanced control system are the most vulnerable to developing addiction.

590

591    The phenomenon of relapse is more elusive and the two models do not fully converge on this

592    aspect. To investigate this phenomenon, we have adapted the complexity of real world

593    treatments to the capabilities of our simulated agents and environments, where we can easily

594    manipulate or extinguish consolidated memory, but we cannot engage all other aspects

595    commonly involved in addiction treatment, such as cognitive or emotional functions or

596    developing new behavioral strategies to compete with drug-related habits. Therefore, we

597    implemented two compartmentalized treatments that we consider as ideal reference models that

598    target only a single decision system or circuit. These putatively represent treatments capable of

599    affecting only drug-related emotional/value or habitual/motor associations. In the neural model,

600    balanced dorsal and ventral endophenotypes respond well to both types of simulated treatments.

601    For the unbalanced endophenotypes, however, only the appropriate treatment, targeting the

602    dominant neural circuit, is effective. The simulations in the RL model do not show the same

603    symmetric effects for the two treatments: the model-free treatment is effective for most

604    endophenotypes, whereas the model-based treatment is mostly unsuccessful, with short relapse

27

605  times across all endophenotypes, but the purely model-based one. The latter result is possibly

606  due to the learning process characterizing the model-based component, which is affected by

607  conflicting information as drug use is associated with both positive and negative outcomes,

608  experienced by the agent when entering the drug state under different phases.

609

610  It is worth noting that habitual and goal-oriented behaviors have neural representations in the

611  dorsal and ventral cortico-striatal circuits respectively, but they do not fully overlap with model-

612  based and model-free control modalities in RL (Dolan and Dayan, 2013). Nonetheless, the neural

613  and RL models independently simulate choices among competing options in addiction.  Thus, we

614  have been able to test our hypothesis of endophenotypic differentiation under two

615  complementary levels in Marr's tri-level of analysis: the neural implementation and the

616  algorithmic level (Marr and Poggio, 1976). This multilevel modeling approach has been often

617  used in computational psychiatry (Maia and Frank, 2011; Montague et al., 2012; Adams et al.,

618  2016; Hauser et al., 2016; Huys et al., 2016) to highlight model convergence and associate

619  specific neural structure and dynamics with mathematical formalizations of optimal and

620  suboptimal behavior in RL. The convergence of neural and RL models on important predictions

621  also provides more confidence in the reliability of the identified computational mechanisms

622  underlying addiction and the associated characterization of endophenotypes. Specifically, both

623  models indicate individuals with unbalanced cortico-striatal activity or control modality are at

624  higher risk of developing addiction and relapse after any treatment. Thus, independent of

625  phenotypic-specific treatments, our results suggest that individuals with these traits would

626  require a prolonged or more intense treatment, in comparison with balanced endophenotypes.

627  Finally, when considering phenomena that are divergent across both models (e.g. **response**

28

628  across endophenotypes to our simulated treatments), our findings still demonstrate that important

629  endophenotypic features might remain undetected in terms of pre-treatment observable behavior.

630  The models showed that opposite unbalanced agents resulted in similar addictive behaviors and

631  vulnerabilities, but diverged in treatment response, potentially informing the development of

632  precision interventions. Further studies will be required to provide empirical validation of our

633  models. For example, computational analysis of fMRI data can be used to test effective

634  connectivity among cortico-striatal circuits (e.g. Friston et al., 2003), in conjunction with

635  cognitive tasks targeting the model-based and model-free control systems.

636     **References**

637     Adams RA, Huys QJ, Roiser JP (2016) Computational Psychiatry: towards a mathematically
638             informed understanding of mental illness. J Neurol Neurosurg Psychiatry 87:53-63.
639     Afraimovich V, Tristan I, Huerta R, Rabinovich MI (2008) Winnerless competition principle
640             and prediction of the transient dynamics in a Lotka-Volterra model. Chaos
641             18:043103.
642     Amit DJ (1989) Modeling brain function: the world of attractor neural networks. Cambridge
643             ; New York: Cambridge University Press.
644     Baldassarre G, Mannella F, Fiore VG, Redgrave P, Gurney K, Mirolli M (2013) Intrinsically
645             motivated action-outcome learning and goal-based action recall: a system-level bio-
646             constrained computational model. Neural Netw 41:168-187.
647     Balleine BW (2005) Neural bases of food-seeking: affect, arousal and reward in
648             corticostriatolimbic circuits. Physiol Behav 86:717-730.
649     Balleine BW, O'Doherty JP (2010) Human and rodent homologies in action control:
650             corticostriatal determinants of goal-directed and habitual action.
651             Neuropsychopharmacology 35:48-69.
652     Belcher AM, Volkow ND, Moeller FG, Ferre S (2014) Personality traits and vulnerability or
653             resilience to substance use disorders. Trends Cogn Sci 18:211-217.
654     Belin D, Everitt BJ (2008) Cocaine seeking habits depend upon dopamine-dependent serial
655             connectivity linking the ventral with the dorsal striatum. Neuron 57:432-441.
656     Belin D, Deroche-Gamonet V (2012) Responses to novelty and vulnerability to cocaine
657             addiction: contribution of a multi-symptomatic animal model. Cold Spring Harbor
658             perspectives in medicine 2.
659     Belin D, Belin-Rauscent A, Everitt BJ, Dalley JW (2016) In search of predictive
660             endophenotypes in addiction: insights from preclinical research. Genes, brain, and
661             behavior 15:74-88.
662     Belin D, Mar AC, Dalley JW, Robbins TW, Everitt BJ (2008) High impulsivity predicts the
663             switch to compulsive cocaine-taking. Science 320:1352-1355.
664     Belin D, Berson N, Balado E, Piazza PV, Deroche-Gamonet V (2011) High-novelty-
665             preference rats are predisposed to compulsive cocaine self-administration.
666             Neuropsychopharmacology 36:569-579.
667     Bellman R (1966) Dynamic programming. Science 153:34-37.
668     Carroll ME, Morgan AD, Lynch WJ, Campbell UC, Dess NK (2002) Intravenous cocaine and
669             heroin self-administration in rats selectively bred for differential saccharin intake:
670             phenotype and sex differences. Psychopharmacology (Berl) 161:304-313.
671     Cohen MX, Frank MJ (2009) Neurocomputational models of basal ganglia function in
672             learning, memory and choice. Behavioural brain research 199:141-156.
673     Covington HE, 3rd, Miczek KA (2005) Intense cocaine self-administration after episodic
674             social defeat stress, but not after aggressive behavior: dissociation from
675             corticosterone activation. Psychopharmacology (Berl) 183:331-340.
676     Dalley JW, Everitt BJ, Robbins TW (2011) Impulsivity, compulsivity, and top-down
677             cognitive control. Neuron 69:680-694.
678     Dalley JW, Fryer TD, Brichard L, Robinson ES, Theobald DE, Laane K, Pena Y, Murphy ER,
679             Shah Y, Probst K, Abakumova I, Aigbirhio FI, Richards HK, Hong Y, Baron JC, Everitt

680   BJ, Robbins TW (2007) Nucleus accumbens D2/3 receptors predict trait impulsivity
681    and cocaine reinforcement. Science 315:1267-1270.
682 Daw ND, Dayan P (2014) The algorithmic anatomy of model-based evaluation. Philos Trans
683    R Soc Lond B Biol Sci 369.
684 Daw ND, Niv Y, Dayan P (2005) Uncertainty-based competition between prefrontal and
685    dorsolateral striatal systems for behavioral control. Nat Neurosci 8:1704-1711.
686 Daw ND, Gershman SJ, Seymour B, Dayan P, Dolan RJ (2011) Model-based influences on
687    humans' choices and striatal prediction errors. Neuron 69:1204-1215.
688 Dayan P (2009) Dopamine, reinforcement learning, and addiction. Pharmacopsychiatry 42
689    Suppl 1:S56-65.
690 Deco G, Jirsa VK, Robinson PA, Breakspear M, Friston K (2008) The dynamic brain: from
691    spiking neurons to neural masses and cortical fields. PLoS Comput Biol 4:e1000092.
692 desJardins ME, Durfee EH, Ortiz J, Charles L., Wolverton MJ (1999) A survey of research in
693    distributed, continual planning. AI Magazine 20:13-22.
694 Dilleen R, Pelloux Y, Mar AC, Molander A, Robbins TW, Everitt BJ, Dalley JW, Belin D (2012)
695    High anxiety is a predisposing endophenotype for loss of control over cocaine, but
696    not heroin, self-administration in rats. Psychopharmacology (Berl) 222:89-97.
697 Dolan RJ, Dayan P (2013) Goals and habits in the brain. Neuron 80:312-325.
698 Doll BB, Daw ND (2016) The expanding role of dopamine. Elife 5.
699 Doll BB, Jacobs WJ, Sanfey AG, Frank MJ (2009) Instructional control of reinforcement
700    learning: a behavioral and neurocomputational investigation. Brain Res 1299:74-94.
701 Draganski B, Kherif F, Kloppel S, Cook PA, Alexander DC, Parker GJ, Deichmann R,
702    Ashburner J, Frackowiak RS (2008) Evidence for segregated and integrative
703    connectivity patterns in the human Basal Ganglia. J Neurosci 28:7143-7152.
704 Economidou D, Pelloux Y, Robbins TW, Dalley JW, Everitt BJ (2009) High impulsivity
705    predicts relapse to cocaine-seeking after punishment-induced abstinence. Biol
706    Psychiatry 65:851-856.
707 Ersche KD, Turton AJ, Pradhan S, Bullmore ET, Robbins TW (2010) Drug addiction
708    endophenotypes: impulsive versus sensation-seeking personality traits. Biol
709    Psychiatry 68:770-773.
710 Everitt BJ, Robbins TW (2013) From the ventral to the dorsal striatum: devolving views of
711    their roles in drug addiction. Neurosci Biobehav Rev 37:1946-1954.
712 Everitt BJ, Robbins TW (2016) Drug Addiction: Updating Actions to Habits to Compulsions
713    Ten Years On. Annu Rev Psychol 67:23-50.
714 Fiore VG, Dolan RJ, Strausfeld NJ, Hirth F (2015) Evolutionarily conserved mechanisms for
715    the selection and maintenance of behavioural activity. Philos Trans R Soc Lond B
716    Biol Sci 370.
717 Fiore VG, Nolte T, Rigoli F, Smittenaar P, Gu X, Dolan RJ (2018) Value encoding in the globus
718    pallidus: fMRI reveals an interaction effect between reward and dopamine drive.
719    Neuroimage 173:249-257.
720 Fiore VG, Rigoli F, Stenner MP, Zaehle T, Hirth F, Heinze HJ, Dolan RJ (2016) Changing
721    pattern in the basal ganglia: motor switching under reduced dopaminergic drive. Sci
722    Rep 6:23327.
723 Fiore VG, Sperati V, Mannella F, Mirolli M, Gurney K, Friston K, Dolan RJ, Baldassarre G
724    (2014) Keep focussing: striatal dopamine multiple functions resolved in a single
725    mechanism tested in a simulated humanoid robot. Front Psychol 5:124.

31

726    Flagel SB, Waselus M, Clinton SM, Watson SJ, Akil H (2014) Antecedents and consequences
727        of drug abuse in rats selectively bred for high and low response to novelty.
728        Neuropharmacology 76 Pt B:425-436.
729    Flagel SB, Robinson TE, Clark JJ, Clinton SM, Watson SJ, Seeman P, Phillips PE, Akil H (2010)
730        An animal model of genetic vulnerability to behavioral disinhibition and
731        responsiveness to reward-related cues: implications for addiction.
732        Neuropsychopharmacology 35:388-400.
733    Friston KJ, Harrison L, Penny W (2003) Dynamic causal modelling. Neuroimage 19:1273-
734        1302.
735    Gannon BM, Galindo KI, Rice KC, Collins GT (2017) Individual Differences in the Relative
736        Reinforcing Effects of 3,4-Methylenedioxypyrovalerone under Fixed and
737        Progressive Ratio Schedules of Reinforcement in Rats. The Journal of pharmacology
738        and experimental therapeutics 361:181-189.
739    Garrison KA, Potenza MN (2014) Neuroimaging and biomarkers in addiction treatment.
740        Current psychiatry reports 16:513.
741    Gerfen CR, Surmeier DJ (2011) Modulation of striatal projection systems by dopamine.
742        Annu Rev Neurosci 34:441-466.
743    Gershman SJ, Horvitz EJ, Tenenbaum JB (2015) Computational rationality: A converging
744        paradigm for intelligence in brains, minds, and machines. Science 349:273-278.
745    Gillan CM, Kosinski M, Whelan R, Phelps EA, Daw ND (2016) Characterizing a psychiatric
746        symptom dimension related to deficits in goal-directed control. Elife 5.
747    Gould RW, Duke AN, Nader MA (2014) PET studies in nonhuman primate models of cocaine
748        abuse: translational research related to vulnerability and neuroadaptations.
749        Neuropharmacology 84:138-151.
750    Grahn JA, Parkinson JA, Owen AM (2009) The role of the basal ganglia in learning and
751        memory: neuropsychological studies. Behavioural brain research 199:53-60.
752    Gruber AJ, McDonald RJ (2012) Context, emotion, and the strategic pursuit of goals:
753        interactions among multiple brain systems controlling motivated behavior. Front
754        Behav Neurosci 6:50.
755    Gutkin BS, Dehaene S, Changeux JP (2006) A neurocomputational hypothesis for nicotine
756        addiction. Proc Natl Acad Sci U S A 103:1106-1111.
757    Haber S (2008) Parallel and integrative processing through the Basal Ganglia reward
758        circuit: lessons from addiction. Biol Psychiatry 64:173-174.
759    Haber SN (2003) The primate basal ganglia: parallel and integrative networks. J Chem
760        Neuroanat 26:317-330.
761    Hauser TU, Fiore VG, Moutoussis M, Dolan RJ (2016) Computational Psychiatry of ADHD:
762        Neural Gain Impairments across Marrian Levels of Analysis. Trends Neurosci 39:63-
763        73.
764    Hoffman RE, McGlashan TH (2001) Neural network models of schizophrenia.
765        Neuroscientist 7:441-454.
766    Huys QJ, Maia TV, Frank MJ (2016) Computational psychiatry as a bridge from
767        neuroscience to clinical applications. Nat Neurosci 19:404-413.
768    Hyman SE, Malenka RC, Nestler EJ (2006) Neural mechanisms of addiction: the role of
769        reward-related learning and memory. Annu Rev Neurosci 29:565-598.
770    Jahanshahi M, Obeso I, Rothwell JC, Obeso JA (2015) A fronto-striato-subthalamic-pallidal
771        network for goal-directed and habitual inhibition. Nat Rev Neurosci 16:719-732.

772    Jimenez VA, Grant KA (2017) Studies using macaque monkeys to address excessive alcohol
773            drinking and stress interactions. Neuropharmacology 122:127-135.
774    Jonkman S, Pelloux Y, Everitt BJ (2012) Drug intake is sufficient, but conditioning is not
775            necessary for the emergence of compulsive cocaine seeking after extended self-
776            administration. Neuropsychopharmacology 37:1612-1619.
777    Jupp B, Dalley JW (2014) Behavioral endophenotypes of drug addiction: Etiological insights
778            from neuroimaging studies. Neuropharmacology 76 Pt B:487-497.
779    Keramati M, Dezfouli A, Piray P (2011) Speed/accuracy trade-off between the habitual and
780            the goal-directed processes. PLoS Comput Biol 7:e1002055.
781    Koob GF, Volkow ND (2016) Neurobiology of addiction: a neurocircuitry analysis. The
782            lancet Psychiatry 3:760-773.
783    Maia TV, Frank MJ (2011) From reinforcement learning models to psychiatric and
784            neurological disorders. Nat Neurosci 14:154-162.
785    Marr D, Poggio T (1976) From understanding computation to understanding neural
786            circuitry. Cambridge: Massachusetts Institute of Technology, Artificial Intelligence
787            Laboratory.
788    McClintick MN, Grant KA (2016) Aggressive temperament predicts ethanol self-
789            administration in late adolescent male and female rhesus macaques.
790            Psychopharmacology (Berl) 233:3965-3976.
791    Molander AC, Mar A, Norbury A, Steventon S, Moreno M, Caprioli D, Theobald DE, Belin D,
792            Everitt BJ, Robbins TW, Dalley JW (2011) High impulsivity predicting vulnerability
793            to cocaine addiction in rats: some relationship with novelty preference but not
794            novelty reactivity, anxiety or stress. Psychopharmacology (Berl) 215:721-731.
795    Montague PR, Dolan RJ, Friston KJ, Dayan P (2012) Computational psychiatry. Trends Cogn
796            Sci 16:72-80.
797    Moore A, Atkeson CG (1993) Prioritized sweeping: Reinforcement learning with less data
798            and less time. Machine Learning 13:103-130.
799    Morgan D, Grant KA, Gage HD, Mach RH, Kaplan JR, Prioleau O, Nader SH, Buchheimer N,
800            Ehrenkaufer RL, Nader MA (2002) Social dominance in monkeys: dopamine D2
801            receptors and cocaine self-administration. Nat Neurosci 5:169-174.
802    Nadal R, Armario A, Janak PH (2002) Positive relationship between activity in a novel
803            environment and operant ethanol self-administration in rats. Psychopharmacology
804            (Berl) 162:333-338.
805    Nader MA, Czoty PW (2005) PET imaging of dopamine D2 receptors in monkey models of
806            cocaine abuse: genetic predisposition versus environmental modulation. Am J
807            Psychiatry 162:1473-1482.
808    Nestler EJ, Aghajanian GK (1997) Molecular and cellular basis of addiction. Science 278:58-
809            63.
810    Obeso JA, Rodriguez-Oroz MC, Stamelou M, Bhatia KP, Burn DJ (2014) The expanding
811            universe of disorders of the basal ganglia. Lancet 384:523-531.
812    Pelloux Y, Everitt BJ, Dickinson A (2007) Compulsive drug seeking by rats under
813            punishment: effects of drug taking history. Psychopharmacology (Berl) 194:127-
814            137.
815    Pelloux Y, Murray JE, Everitt BJ (2015) Differential vulnerability to the punishment of
816            cocaine related behaviours: effects of locus of punishment, cocaine taking history
817            and alternative reinforcer availability. Psychopharmacology (Berl) 232:125-134.

818     Perry JL, Carroll ME (2008) The role of impulsive behavior in drug abuse.
819             Psychopharmacology (Berl) 200:1-26.
820     Pezzulo G, Rigoli F, Chersi F (2013) The mixed instrumental controller: using value of
821             information to combine habitual choice and mental simulation. Front Psychol 4:92.
822     Piazza PV, Deminiere JM, Le Moal M, Simon H (1989) Factors that predict individual
823             vulnerability to amphetamine self-administration. Science 245:1511-1513.
824     Piray P, Keramati MM, Dezfouli A, Lucas C, Mokri A (2010) Individual differences in nucleus
825             accumbens dopamine receptors predict development of addiction-like behavior: a
826             computational approach. Neural Comput 22:2334-2368.
827     Rabinovich MI, Varona P, Selverston AI, Abarbanel HDI (2006) Dynamical principles in
828             neuroscience. Reviews of Modern Physics 78.
829     Redish AD (2004) Addiction as a computational process gone awry. Science 306:1944-
830             1947.
831     Redish AD, Jensen S, Johnson A (2008) A unified framework for addiction: vulnerabilities in
832             the decision process. Behav Brain Sci 31:415-437; discussion 437-487.
833     Sadacca BF, Jones JL, Schoenbaum G (2016) Midbrain dopamine neurons compute inferred
834             and cached value prediction errors in a common framework. Elife 5.
835     Schultz W, Dayan P, Montague PR (1997) A neural substrate of prediction and reward.
836             Science 275:1593-1599.
837     Silver D, Huang A, Maddison CJ, Guez A, Sifre L, van den Driessche G, Schrittwieser J,
838             Antonoglou I, Panneershelvam V, Lanctot M, Dieleman S, Grewe D, Nham J,
839             Kalchbrenner N, Sutskever I, Lillicrap T, Leach M, Kavukcuoglu K, Graepel T,
840             Hassabis D (2016) Mastering the game of Go with deep neural networks and tree
841             search. Nature 529:484-489.
842     Simon DA, Daw ND (2012) Dual-system learning models and drugs of abuse. In:
843             Computational Neuroscience of Drug Addiction, 1 Edition (Gutkin B, Ahmed SH,
844             eds), pp 145–161. New York: Springer-Verlag.
845     Singh S, Jaakkola T, Littman ML, Szepesvári C (2000) Convergence Results for Single-Step
846             On-Policy Reinforcement-Learning Algorithms. Machine Learning 38:287-308.
847     Smethells JR, Zlebnik NE, Miller DK, Will MJ, Booth F, Carroll ME (2016) Cocaine self-
848             administration and reinstatement in female rats selectively bred for high and low
849             voluntary running. Drug and alcohol dependence 167:163-168.
850     Suto N, Austin JD, Vezina P (2001) Locomotor response to novelty predicts a rat's
851             propensity to self-administer nicotine. Psychopharmacology (Berl) 158:175-180.
852     Sutton RS (1990) Integrated architecture for learning, planning, and reacting based on
853             approximating dynamic programming. In: Proceedings of the seventh international
854             conference (1990) on Machine learning, pp 216-224. Austin, Texas, USA: Morgan
855             Kaufmann Publishers Inc.
856     Sutton RS, Barto AG (1998) Reinforcement Learning: An Introduction. Cambridge, MA: MIT
857             Press.
858     Verdejo-Garcia A, Lawrence AJ, Clark L (2008) Impulsivity as a vulnerability marker for
859             substance-use disorders: review of findings from high-risk research, problem
860             gamblers and genetic association studies. Neurosci Biobehav Rev 32:777-810.
861     Volkow ND, Morales M (2015) The Brain on Drugs: From Reward to Addiction. Cell
862             162:712-725.

863     Volkow ND, Fowler JS, Wang GJ, Swanson JM, Telang F (2007) Dopamine in drug abuse and
864             addiction: results of imaging studies and treatment implications. Archives of
865             neurology 64:1575-1579.
866     Voon V, Reiter A, Sebold M, Groman S (2017) Model-Based Control in Dimensional
867             Psychiatry. Biol Psychiatry 82:391-400.
868     Watkins CJ, Dayan P (1992) Q-Learning. Machine Learning 8:279-292.
869     Willuhn I, Burgeno LM, Everitt BJ, Phillips PE (2012) Hierarchical recruitment of phasic
870             dopamine signaling in the striatum during the progression of cocaine use. Proc Natl
871             Acad Sci U S A 109:20703-20708.
872     Yin HH, Knowlton BJ, Balleine BW (2004) Lesions of dorsolateral striatum preserve
873             outcome expectancy but disrupt habit formation in instrumental learning. Eur J
874             Neurosci 19:181-189.
875

876 **Extended Data Code File 1.** To access the source code of both models, visit the ModelDB

877 website (https://senselab.med.yale.edu/modeldb/enterCode.cshtml?model=239540), and

878 download the archive. The source code shows its structure in the commented main files

879 "separate_test.m" and "RunExperimentLearning96.m", respectively in the folder "neural_model"

880 and "RL_model".

881

882

883 **Figure 1. Illustrative representation of energy landscapes and neural architecture of the**

884 **model. A-C.** These representations of energy landscapes are meant to illustrate differences in the

885 temporal responses provided by neural systems. Depending on the energy landscape, three

886 arbitrary inputs (magenta dots) are transformed into different stable states (grey dots). Learning

887 processes increase or decrease the strength of the connections among nodes in a network, thereby

888 altering its energy landscape and reshaping temporal responses towards existing attractors.

889 Attractors are defined as low energy states (bottom of the basins) at the end point of the temporal

890 responses to multiple starting inputs. **(A)** The landscape is characterized by multiple shallow

891 attractors: these allow slow temporal responses, transforming multiple inputs into multiple

892 weakly stable states. Noise and changes in the incoming input easily determine new responses

893 towards different attractors. **(B)** In this second illustrative configuration, steep and vast attractors

894 characterize the energy landscape, allowing quick state transitions towards two equilibrium

895 points. This new configuration is able to resist noise and minor changes in the incoming input,

896 and, at the same time, allows a differentiation of inputs in two broad categories. **(C)** Finally, the

897 third energy landscape illustrates the presence of a parasitic attractor, exemplifying the condition

898 of addiction: all inputs fall now at the bottom of a single steep basin. Under this condition, noise

36

899   and changes in the incoming input determine temporal responses that keep falling in the same

900   attractor, therefore preventing the system from executing different behaviors. **(D)** Neural

901   architecture used to simulate neural dynamics and behavior for the mean field neural model. The

902   activity in the dorsal cortico-striatal circuit is responsible for the motor output of the system (left

903   circuit), whilst activity in the ventral cortico-striatal circuit is responsible for goal selections

904   (right circuit). The two systems bias each other via cortico-cortical connectivity and learning

905   processes affect the weights of the connections between the two cortical outputs and the striatum

906   in their corresponding circuits. The components in the architecture are labeled as follows: cortex

907   (Cx), thalamus (Th), globus pallidus pars externa and interna (GPe and GPi), substantia nigra

908   pars reticulata (SNr), sub-thalamic nucleus (STN) and striatum (Str), divided into two areas

909   enriched by either D1 or D2 receptors.

910

911   **Figure 2**. **Illustrative representation of the environment used for the RL model of**

912   **addiction.** The states are disposed in a linear arrangement: on one extreme is a healthy reward

913   state (1), on the opposite side a drug state (8) followed by twelve aftereffects states (9-22).

914   Healthy reward and drug states are separated by 6 neutral states (2-7). The agent can traverse

915   between nearby neutral states. From the two borders of the central segment of neutral states, an

916   agent can enter the healthy reward state (from state 2), securing a moderate reward (Rg=1), or

917   the drug state (from state 7), receiving an initial high reward (Rd=10, during the phase of

918   addiction) and a series of sparse but temporally extended negative outcomes, characterizing the

919   aftereffects states. The presence of negative outcomes makes entering the drug and aftereffects

920   area suboptimal during all experimental phases (see optimal policy in **table 3**). From both the

921   goal state and the drug/aftereffects segment the agent is then returned to the middle of the neutral

37

922   segment. In this representation, we explicitly portray the transitions related to states 1 (healthy

923   reward), 4 (neutral), and 15 and 20 (drug aftereffects) for illustrative purposes. Line width

924   represents related transition probability value. Line and text color represent the action class ($a_s$,

925   $a_g$, $a_w$, $a_d$). Neutral states are navigable with actions $a_{s2-7}$ which are deterministic for adjacent

926   state while have high chance of failing for distant states. From the neutral states the agent can

927   reach: (i) the healthy reward, if executing action $a_g$ when in state 2; and (ii) the drug state (8) and

928   aftereffects area (state 9 to 22), if executing action $a_d$, when in state 7. From the healthy reward

929   area the agent can issue again $a_g$, receiving a reward of 1 and going back to the center of the

930   neutral area, state 4. By entering the drug area, the agent receives a reward of 10. Action results

931   in the drug/aftereffect area are probabilistic: the agent can reach a nearby state in the area or

932   leave the area and reach the center of the neutral state. Leaving the drug/aftereffects area has a

933   cost of -4, whereas every other transition inside the area costs -1.2. For a full description of

934   transitions and their probability distribution in the environment (see **Tables 1-2,4-5**).

935

936   **Figure 3. Distribution of action selections across endophenotypes controlled by the neural**

937   **model.** Histograms show how the distribution of simulated action selections changes depending

938   on the endophenotype (11 variations in cortico-cortical connectivity weights). 30 random

939   seeds/inputs are used per endophenotype, tested under two stages: pre-drug **(A)** and addiction

940   **(B)**. The three colors represent the occurrence of selections of three arbitrary actions. Under the

941   pre-drug stage, no reward is provided and action selections are triggered by random fluctuation

942   in values of competing sensory inputs. The simulations show the agents adapt to the changes in

943   sensory stimuli and therefore exhibit a near-uniform distribution of action selections.

944   Conversely, under the addiction stage, the action represented in blue is associated with

945     administration of the simulated drug, triggering DA-dependent Hebbian learning in cortico-

946     striatal connectivity, and consequently over-selection. Under addiction, the differences among

947     endophenotypes clearly emerge in the selection frequency of the action leading to drug

948     consumption. Asymmetric control (endophenotypes 1-3 and 9-11) leads to a stronger over-

949     selection in comparison with balanced control (endophenotypes 4-7), despite identical learning

950     processes and reward encoding.

951

952     **Figure 4. Severity of addiction and relapse time across endophenotypes controlled by the**

953     **neural model.** Shaded error bars report mean and standard error for 30 simulated agents across

954     endophenotypes (11 variations in cortico-cortical connectivity weights). Panels A and B show

955     the selections of actions leading to substance consumption, as a percentage of the overall number

956     of action selections. In the first case **(A)** we compare the values recorded during the addiction

957     stage with those recorded during the stage of dorsal treatment jointly with abstinence (i.e. drug-

958     related actions do not trigger self-administration of a drug and the treatment targets the dorsal

959     circuit). In the second case **(B)** the comparison involves addiction and ventral treatment

960     (treatment targeting the ventral circuit, during abstinence). Panels C and D compare the

961     simulated time required by the 11 endophenotypes to reach an arbitrary threshold of cortico-

962     striatal connectivity during the stage of addiction and during the stage of relapse after either

963     dorsal **(C)** or ventral **(D)** treatment. Within the time of a simulation run, all simulated agents

964     reached the addiction threshold. The two treatments are simulated by restoring either the

965     dorsal/motor (A-C) or the ventral/outcome circuit (B-D) to the configuration characterizing the

966     pre-drug stage. The percentage of the action selections shows the dorsal treatment is more

967     effective in endophenotypes characterized by high dorsal dominance **(A)**, whereas the ventral

968    treatment only has an effect in endophenotypes characterized by high ventral dominance **(B)**.

969    Similarly, dorsal and ventral treatments result in long relapse times in endophenotypes

970    characterized by high dorsal and high ventral dominance, respectively. (*) indicates significant

971    difference: p<0.05.

972

973    **Figure 5. Likelihood to develop addiction and relapse time across endophenotypes**

974    **controlled by the RL model.** Shaded error bars report mean and standard error for ~100

975    simulated agents across 6 endophenotypes (differential balance between model-based and model-

976    free control modalities, $\beta$=[0, 0.2, 0.4, 0.6, 0.8, 1]). Panels A and B show the percentage of

977    agents developing addiction (i.e. drug-related choices are more frequent than healthy reward-

978    related choices), per endophenotype, under the addiction and treatment phases. In the first case

979    **(A)** the comparison involves data recorded during the phase of addiction and those recorded

980    during the phase of model-free treatment. In the second case **(B)** the comparison involves the

981    phases of addiction and model-based treatment. Panels C and D illustrate the simulated time

982    required by the 6 endophenotypes to reach 95% of action preference towards the drug state, in

983    comparison with action preference recorded during the phase of addiction (f2). In the first case

984    **(C)** the comparison involves the phases of addiction and relapse after model-free treatment,

985    whereas in the second case **(D)** the comparison involves the phases of addiction and relapse after

986    model-based treatment. In terms of action selection ratio, the simulated results show both

987    treatments have a significant effect only on those phenotypes characterized by strong unbalance

988    of control **(A-B)**. In terms of relapse, the results show the model-free treatment is on average

989    more successful than the model-based one, as 5 endophenotypes show no significant difference

990    between the phases of addiction and post-treatment addiction (i.e. the time required to relapse is

40

991     not significantly different than the time required to develop addiction the first time). Each

992     endophenotype, or parameter selection, was simulated 100 times across the four phases (3050

993     steps per simulation). Results depend on the statics of the environment, but over similar

994     environments the results were qualitatively similar. (*) indicates significant difference: $p<0.05$.

995

996

997     **Table 1** Environment transition probabilities across endophenotypes controlled by the RL model.

998     Changes during phases in italic.

| Transition Description | Probability for each phase | | | | |
|---|---|---|---|---|---|
| | P (f1) | P (f2) | P (f3) | P (f4) | |
| $P(s=i\|s=i,a=a^{s=i})$, i neutral state | 1 | 1 | 1 | 1 | |
| $P(s=i+j\|s=i,a=a^{s=i+j})$,j=+1/-1, i neutral state, i+j neutral state | 0.99 | 0.99 | 0.99 | 0.99 | |
| $P(s=i\|s=i,a=a^{s=i+j})$ ,j=+1/-1, i neutral state, i+j neutral state | 0.01 | 0.01 | 0.01 | 0.01 | |
| $P(s=i+k\|s=i,a=a^{s=i+k})$,k!=+1/-1, i neutral state, i+k neutral state | 0.0001 | 0.0001 | 0.0001 | 0.0001 | From Neutral States |
| $P(s=i\|s=i,a=a^{s=i+k})$,k!=+1/-1, i neutral state, i+k neutral state | 0.9999 | 0.9999 | 0.9999 | 0.9999 | |
| $P(s=i\|s=i,a=a^{w})$, i neutral state | 1 | 1 | 1 | 1 | |
| $P(s=1\|s=2,a=a^{g})$ | 1 | 1 | 1 | 1 | |
| $P(s=i\|s=i,a=a^{g})$, i!=2 neutral state | 1 | 1 | 1 | 1 | |
| $P(s=8\|s=7,a=a^{d})$ | 1 | 1 | 1 | 1 | |

41

| Transition Description | | | | | |
|---|---|---|---|---|---|
| P(s=i\|s=i,a=aᵈ), i!=7 neutral state | 1 | 1 | 1 | 1 | |
| P(s=i\|s=i,a=aᵍ), i drug/aft state | 0.999 | 0.999 | *0.8* | 0.999 | |
| P(s=4\|s=i,a=aᵍ), i drug/aft state | 0.001 | 0.001 | *0.2* | 0.001 | |
| P(s=i\|s=i,a=aˢ⁼*), i drug/aft state | 0.999 | 0.999 | *0.8* | 0.999 | |
| P(s=4\|s=i,a=as=*), i drug/aft state | 0.001 | 0.001 | *0.2* | 0.001 | |
| P(s=j\|s=i,a=aʷ), i!=15 drug/aft state, j next or previous drug/aft state | 0.4995 | 0.4995 | *0.4* | 4.995 | From Drug/aft States |
| P(s=4\|s=i,a=aʷ), i!=15 drug/aft state | 0.001 | 0.001 | *0.2* | 0.001 | |
| P(s=14/16\|s=15,a=aʷ) | 0.2 | 0.2 | *0.15* | 0.2 | |
| P(s=4\|s=15,a=aʷ) | 0.6 | 0.6 | *0.7* | 0.6 | |
| P(s=j\|s=i,a=aᵈ), i drug/aft state, j next drug/aft state | 0.745 | 0.745 | *0.6* | 0.745 | |
| P(s=j\|s=i,a=aᵈ), i drug/aft state, j previous drug/aft state | 0.245 | 0.245 | *0.2* | 0.245 | |
| P(s=4\|s=i,a=aᵈ), i drug/aft state | 0.01 | 0.01 | *0.2* | 0.01 | |
| P(s=4\|s=1,a=aᵍ) | 1 | 1 | 1 | 1 | |
| P(s=1\|s=1,a=aˢ⁼*) | 1 | 1 | 1 | 1 | Goal |
| P(s=1\|s=1,a=aʷ) | 1 | 1 | 1 | 1 | |
| P(s=1\|s=1,a=aᵈ) | 1 | 1 | 1 | 1 | |

999

1000 **Table 2** Environment rewards across endophenotypes controlled by the RL model. Changes

1001 during phases in italic.

| Transition Description | Probability for each phase | |
|---|---|---|

| | P (f1) | P (f2) | P (f3) | P (f4) | |
|---|---|---|---|---|---|
| $T(s=i|s=i,a=a^{s=i})$, i neutral state | 0 | 0 | 0 | 0 | **From States** |
| $T(s=i+j|s=i,a=a^{s=i+j})$, j=+1/-1, i neutral state, i+j neutral state | 0 | 0 | 0 | 0 | |
| $T(s=i|s=i,a=a^{s=i+j})$, j=+1/-1, i neutral state, i+j neutral state | 0 | 0 | 0 | 0 | |
| $T(s=i+k|s=i,a=a^{s=i+k})$,k!=+1/-1, i neutral state, i+k neutral state | -0.3 | -0.3 | -0.3 | -0.3 | |
| $T(s=i|s=i,a=a^{s=i+k})$,k!=+1/-1, i neutral state, i+k neutral state | 0 | 0 | 0 | 0 | |
| $T(s=i|s=i,a=a^{w})$, i neutral state | 0 | 0 | 0 | 0 | |
| $T(s=1|s=2,a=a^{g})$ | 0 | 0 | 0 | 0 | |
| $T(s=i|s=i,a=a^{g})$, i!=2 neutral state | 0 | 0 | 0 | 0 | |
| $T(s=8|s=7,a=a^{d})$ | 0 | *10* | *-1* | *10* | |
| $T(s=i|s=i,a=a^{d})$, i!=7 neutral state | 0 | 0 | 0 | 0 | |
| $T(s=i|s=i,a=a^{g})$, i drug/aft state | -0.3 | *-1.2* | *-1.2* | *-1.2* | **From Drug/aft States** |
| $T(s=4|s=i,a=a^{g})$, i drug/aft state | -4 | -4 | -4 | -4 | |
| $T(s=i|s=i,a=a^{s=*})$, i drug/aft state | -0.3 | *-1.2* | *-1.2* | *-1.2* | |
| $T(s=4|s=i,a=a^{s=*})$, i drug/aft state | -4 | -4 | -4 | -4 | |
| $T(s=j|s=i,a=a^{w})$, i!=15 drug/aft state, j next or previous drug/aft state | -0.3 | *-1.2* | *-1.2* | *-1.2* | |
| $T(s=4|s=i,a=a^{w})$, i!=15 drug/aft state | -4 | -4 | -4 | -4 | |
| $T(s=14/16|s=15,a=a^{w})$ | -0.3 | *-1.2* | *-1.2* | *-1.2* | |

| | | | | | |
|---|---|---|---|---|---|
| T(s=4|s=15,a=a$^w$) | -4 | -4 | -4 | -4 | |
| T(s=j|s=i,a=a$^d$), i drug/aft state, j next drug/aft state | -0.3 | *-1.2* | *-1.2* | *-1.2* | |
| T(s=j|s=i,a=a$^d$), i drug/aft state, j previous drug/aft state | -0.3 | *-1.2* | *-1.2* | *-1.2* | |
| T(s=4|s=i,a=a$^d$), i drug/aft state | -4 | -4 | -4 | -4 | |
| T(s=4|s=1,a=a$^g$) | 1 | 1 | 1 | 1 | |
| T(s=1|s=1,a=a$^{s=*}$) | 0 | 0 | 0 | 0 | **Goal** |
| T(s=1|s=1,a=a$^w$) | 0 | 0 | 0 | 0 | |
| T(s=1|s=1,a=a$^d$) | 0 | 0 | 0 | 0 | |

1002

1003 **Table 3** Optimal policy across endophenotypes controlled by the RL model (2nd Drug phase)

| State Id | State Type | Action | Q value |
|---|---|---|---|
| 1 | goal | a$^g$ | 2.8967 |
| 2 | neutral | a$^g$ | 2.607 |
| 3 | neutral | a$^{s=2}$ | 2.3439 |
| 4 | neutral | a$^{s=3}$ | 2.1074 |
| 5 | neutral | a$^{s=4}$ | 1.8948 |
| 6 | neutral | a$^{s=5}$ | 1.7036 |
| 7 | neutral | a$^{s=6}$ | 1.5317 |
| 8 | drug | a$^d$ | -10.1134 |
| 9 | drug-aft | a$^d$ | -10.3781 |
| 10 | drug-aft | a$^w$ | -10.4882 |

44

| 11 | drug-aft | $a^w$ | -10.2809 |
| 12 | drug-aft | $a^w$ | -9.7099 |
| 13 | drug-aft | $a^w$ | -8.6469 |
| 14 | drug-aft | $a^w$ | -6.8532 |
| 15 | drug-aft | $a^w$ | -3.9265 |
| 16 | drug-aft | $a^d$ | -5.2928 |
| 17 | drug-aft | $a^d$ | -6.4251 |
| 18 | drug-aft | $a^d$ | -7.3633 |
| 19 | drug-aft | $a^d$ | -8.1408 |
| 20 | drug-aft | $a^d$ | -8.7849 |
| 21 | drug-aft | $a^d$ | -9.318 |
| 22 | drug-aft | $a^d$ | -9.7575 |

1004

1005 **Table 4** Agent model parameters across endophenotypes controlled by the RL model

| Name | Description | Value |
|------|-------------|-------|
| $\alpha$ | MF learning factor | 0.05 |
| $\gamma$ | Discount factor | 0.9 |
| $d_{MB}$ | MB decay factor | 0.01 |
| $N_{PS}$ | MB number of updates | 50 |
| $T_{MB}$ | Temperature for stochastic state update selection | 1 |
| $\varepsilon$ | Exploration factor | 0.1 |
| $\alpha_{Ctpy}$ | Cognitive therapy MF learning factor | 0.0001, 0.0005, 0.001 |

1006

1007    **Table 5** Environment parameters across endophenotypes controlled by the RL model.

| Name | Description | Value |
|------|-------------|-------|
| $N_T$ | Number of states | 22 |
| $N_G$ | Number Goal States | 1 |
| $N_D$ | Number Drug/aft States | 15 |
| $N_n$ | Number Neutral States | 6 |
| $N_a$ | Number of actions | 9 |
| $S_0$ | Starting state | 4 |
| $R_p$ | Punishment end of drug/aft consumption | -4 |
| $R_c$ | Punishment in drug/aft area | -1.2 |
| $R_{dd}$ | Reward at init drug consumption (f2,f4) | 10 |
| $R_{dt}$ | Reward at init drug consumption in therapy | -1 |
| $R_g$ | Reward when entering goal state | 1 |
| $d_{init}$ | Duration initial (no drug) phase | 50 |
| $d_{drug1}$ | Duration first drug phase | 1000 |
| $d_{tpy}$ | Duration therapy phase | 1000 |
| $d_{drug2}$ | Duration second drug phase | 600 |

1008

**A** Pre-Drug

**B** Addiction

Distribution of action selections expressed as a percentage

Phenotype dominance:
high dorsal (left) - high ventral (right)

Action 1, never reinforced
Action 2, never reinforced
Action 3, reinforced with addictive substance, under addiction condition

Pre-drug (f1) → Addiction (f2) → Model-free treatment (f3) → Relapse after MF treat. (f4)

Addiction (f2) → Model-based treatment (f3) → Relapse after MB treat. (f4)

**A** **Model-free Treatment**

%

Percentage of agents developing addiction, 100 agents, 300 runs

60
50
40
30
20

* * *

1 2 3 4 5 6

Phenotype dominance:
model-free (left) – model-based (right)

**B** **Model-based Treatment**

%

Percentage of agents developing addiction, 100 agents, 300 runs

60
50
40
30
20

* *

1 2 3 4 5 6

Phenotype dominance:
model-free (left) – model-based (right)

**C** **Model-free Treatment**

Time to develop addiction or relapse per 100 simulated agents

*

1 2 3 4 5 6

Phenotype dominance:
model-free (left) – model-based (right)

**D** **Model-based Treatment**

Time to develop addiction or relapse per 100 simulated agents

* * * * *

1 2 3 4 5 6

Phenotype dominance:
model-free (left) – model-based (right)