

1 **The choice-wide behavioral association study: data-driven**
2 **identification of interpretable behavioral components**

3 David B. Kastner^{1,4}, Greer Williams¹, Cristofer Holobetz¹, Joseph P. Romano², Peter Dayan³

4 ¹*Department of Psychiatry and Behavioral Sciences, University of California, San Francisco, CA*
5 *94143, USA*

6 ²*Department of Statistics, Stanford University, Stanford, CA 94305, USA*

7 ³*Max Planck Institute for Biological Cybernetics, Tübingen 72076, Germany*

8 ⁴Lead Contact

9 Correspondence: David B. Kastner: david.kastner2@ucsf.edu

10

11 Abstract

12 **Behavior contains rich structure across many timescales, but there is a dearth of methods to**
13 **identify relevant components, especially over the longer periods required for learning and**
14 **decision-making. Inspired by the goals and techniques of genome-wide association studies,**
15 **we present a data-driven method—the choice-wide behavioral association study: CBAS—that**
16 **systematically identifies such behavioral features. CBAS uses powerful, resampling-based,**
17 **methods of multiple comparisons correction^{1–3} to identify sequences of actions or choices**
18 **that either differ significantly between groups or significantly correlate with a covariate of**
19 **interest. We apply CBAS to different tasks and species (flies⁴, rats⁵, and humans⁶) and find, in**
20 **all instances, that it provides interpretable information about each behavioral task.**

21 Understanding how behavior differs between different groups of humans or other
22 animals is critical for generating and testing hypotheses about the functional role of genes,
23 regions of the brain, and neural circuits, and is central to characterizing neurological and
24 psychiatric dysfunction⁷. However, behavior is highly complex, evolving over multiple
25 timescales and exhibiting substantial path dependencies due to individual experience^{8–12}. It is
26 increasingly possible to automate behavioral paradigms, and for computational methods to
27 revolutionize behavioral analyses^{13–22}. The latter come in two main flavors^{23,24}: model-based, or
28 top-down approaches, and data-driven, or bottom-up approaches. The former are substantially
29 more prevalent than the latter; we offer a partial corrective.

30 In model-based analyses, behavioral data, such as choices in a decision-making task, are
31 processed under the specific assumptions of a hypothesis or model. If the model or hypothesis
32 is correct, this is highly efficient, since large volumes of data can be reduced to a handful of
33 parameters that index semantically meaningful phenotypes, such as learning rates or
34 differential sensitivity to rewards or punishments. These parameters can then be compared
35 between the groups. However, even when substantial effort is put into building multiple
36 alternative models, and comparing them in a statistically rigorous manner, it remains possible
37 that the best fitting model nevertheless fails to characterize the behavior properly, rendering
38 nugatory any interpretation of group comparisons. Additionally, with such approaches,
39 confirmation bias²⁵ presents a significant challenge for accurate interpretation. Furthermore, in
40 the worst of cases, model and hypothesis-based analyses provides a post-hoc framework for
41 explaining any difference that can be found in a behavioral dataset²⁶. Modern machine learning
42 methods^{27,28} can provide useful lower bounds for how well hypothesis-driven models should fit,
43 but these approaches lack interpretability and data efficiency.

44 Data-driven analyses start from the other end, taking, and then characterizing, behavior
45 without relying on parametric assumptions, and making very few assumptions about how the
46 data is generated. Some of these approaches are unsupervised, for instance finding clusters in
47 behavioral space and defining them as ‘syllables’ which can subsequently be compared
48 between groups^{18,29,30}; others are more supervised, directly looking for discriminative
49 differences between populations or individuals^{13–16}. By not trying to force data into a limited set
50 of parameters, these methods should be more sensitive; however, comparisons between
51 groups pose severe statistical challenges, because of the complexity and dimensionality of
52 behavioral datasets.

53 Here, we present the choice-wide behavioral association study (CBAS), a data-driven
54 analysis method designed to identify relevant sequences of choices (or other discrete
55 behavioral features) made by subjects. CBAS has two components: 1) breaking down behavior
56 into a comprehensive language for comparison between two groups or correlation with a
57 covariate of interest; and 2) using rigorous, resampling-based, statistical corrections to account
58 for the resulting large number of comparisons and maintain statistical power despite
59 correlations in the data.

60 **Choice as a common discretization for behavior**

61 In developing our analysis, we were motivated by the data-driven approach of genome-
62 wide association studies (GWAS), whole exome sequencing (WES) and whole genome
63 sequencing (WGS)³¹. In some ways, the state of behavioral analysis in systems neuroscience
64 resembles the state of genetic analysis prior to GWAS/WES/WGS, where studies attempted to
65 associate candidate genes with phenotypes. Candidate gene studies were often underpowered
66 and failed to replicate^{32–34}, reminiscent of aspects of behavioral analyses³⁵.

67 GWAS/WES/WGS look for differences in base-pairs of the genome between groups of
68 subjects. The large numbers of base pairs being compared in these studies necessitates
69 statistical correction to enable reliable decisions about the significance of any differences
70 found. The discrete nature of base-pairs and the ability to compare that set across subjects
71 make GWAS/WES/WGS possible. To be able to develop a comparable method for behavioral
72 analyses it is critical to identify an appropriate discretization of behavioral tasks. In many cases,
73 choice provides just such a discretization (Fig 1). Indeed, there is evidence that, at a rather
74 fundamental level, behavior occurs through discrete choices^{18,36}. We use three tasks, in three
75 diverse species—flies, humans, and rats—to show the breadth of applicability of CBAS.

76 Using choice as the basis of the comparison for the behavioral analysis requires an
77 additional consideration beyond what is done for the genome with base-pairs. For
78 GWAS/WES/WGS, an individual base pair can be a meaningful unit of information (although this
79 is only a partial story^{37,38}) that can be compared between subjects. This need not be the case
80 for individual choices in behavior, whereby the choices that precede and follow a specific
81 choice can change the meaning of that choice. In this case, the relevant behavioral feature is a
82 whole sequence of choices. Therefore, CBAS does not just evaluate the occurrence of individual
83 choices, but, rather, the occurrence of all sequences up to a certain, user-defined, length. In
84 general, the longer the sequence length, the more data and computing time will be necessary.

85 Evaluating the occurrence of all sequences of choices up to a certain length, causes a
86 complexity, when it comes to correcting for the many comparisons, since, the sequences can be
87 highly correlated. Standard methods to correct statistically for multiple comparisons (e.g.
88 Bonferroni, Holm, Benjamini-Hochberg, Benjamini-Yekutieli) are either incorrect or
89 underpowered for correlated data. Therefore, we instead use an approach based on resampling
90 to correct for the multiple comparisons, which retains power in the face of correlations¹.

91 **CBAS identifies interpretable differences for fly y-maze**

92 The first task we considered involves drosophila walking on a y-maze (Fig 1a). The
93 movement of flies on the maze is tracked as they make the choice of going to either the left or

94 right arm after leaving the previous arm⁴. This left/right discretization of the task has enabled
95 many conclusions about the genetic nature of individual variability^{4,39–43}. Here, we compare the
96 choices of two outbred strains of fly. Analyses in the original paper⁴ identified some clear
97 indications about the difference; however, the data provide a useful testing ground for CBAS
98 since there are many subjects and the results are relatively low-dimensional.

99 When deciding to perform a data-driven genomics method, there are still decisions that
100 need to be made. For instance, deciding to focus on only exome sequences in WES, or SNPs in
101 GWAS. Similarly, to run CBAS, a few decisions need to be made about how to structure the data
102 (Fig 2a). These decisions are important, and are not normally pre-determined directly by the
103 data. Any conclusions drawn need to be interpreted in their light.

104 The first decision is the possible choices that will make up the sequences used in CBAS.
105 We refer to this as the basic language for the application of CBAS. For the fly task, the language
106 is left or right turns. Next, as described above, a decision needs to be made about the maximum
107 number of choices in a row that will make up all the sequences used in CBAS. For the fly task,
108 we chose a maximum sequence length of 10 choices in a row. That means that CBAS evaluates
109 all sequences from length 1 – 10 that exist in the dataset.

110 CBAS then works through evaluating the average sequence count for each sequence,
111 which we refer to as the rate of that sequence. The rate is calculated by counting the total
112 number of times that sequence occurs in the population divided by the number of individuals in
113 the population. To be precise about the occurrence of the sequence, a decision needs to be
114 made about the number of trials over which this is counted, i.e. a criterion. The same criterion
115 is applied to all subjects. For the fly task, we use the first 250 turns as that criterion. The last
116 decision that needs to be made is what to do with subjects that do not reach the criterion. For
117 the fly task, we exclude the subjects that do not reach criterion.

118 Given these decisions, we performed a CBAS comparing two outbred strains of flies,
119 Cambridge-A (CA) and w1118, from a publicly available dataset⁴. The w1118 strain is the
120 background strain for many transgenic flies⁴. For each sequence of left/right turns up to 10
121 turns long, the rate was calculated for the two strains (Fig 2b). Using the Romano-Wolf
122 resampling-based multiple comparisons correction^{1,2} (see methods), we determined which
123 sequences the two strains utilize significantly differently (Fig 2c). As a significance threshold, we
124 use median control of the false discovery proportion at 5%, which is comparable to 5% false
125 discovery rate (FDR) for the resampling-based method² (see methods).

126 CBAS identifies many sequences that differ significantly between the CA and w1118
127 lines, some of which are shown in Fig. 2d. Upon inspection of the sequences, a clear,
128 interpretable, difference is apparent between those sequences that occur more in the CA line
129 and those that occur more in the w1118 line (Fig 2d&e). The CA line utilizes sequences with
130 extended numbers of the same turn in a row, whereas the w1118 line utilizes sequences with
131 more frequent changes in turn direction. Therefore, CBAS not only identifies that there is a
132 difference between the two fly lines, but it also provides information to support an
133 interpretation as to the nature of the difference.

134 Data-driven methods invariably need more data than hypothesis-driven methods. To
135 estimate the sample size needed for appropriately powered experiments to detect any
136 difference between these two fly lines, we took advantage of the large dataset to resample
137 groups of smaller sizes from the data and recalculate the CBAS for each set of resampled
138 groups. We generated an estimate of the power for different number of flies per group (Fig 2f)
139 (see methods), by comparing the number of significant sequences identified by CBAS when
140 comparing CA to w1118 to the number of significant sequences identified by CBAS when
141 comparing the lines to themselves (Fig S1b). With 40 flies per group, CBAS has an estimated
142 power >80% to distinguish the CA and w1118 lines.

143 **Graceful decay of CBAS output with decreasing group size**

144 The power calculation concerns the ability of CBAS to distinguish between the two
145 strains. A more refined question is what the nature and comprehensiveness of the collection of
146 significantly different sequences would be as smaller numbers of flies per group are analyzed.
147 This would determine our ability to derive interpretations from fewer subjects.

148 We first evaluated the number of significant sequences identified from smaller group
149 sizes as a fraction of the total number of significant sequences identified with the full dataset
150 (Fig 3a). As expected, with smaller group sizes, we find fewer overall sequences; however,
151 across the range of group sizes evaluated, the median fraction of sequences was larger than the
152 fraction of the population being used in the smaller group CBAS (Fig 3a). This indicates that, for
153 this dataset at smaller group sizes, the proportion of sequences identified by CBAS grows faster
154 than the proportion of subjects in the CBAS. This provides a rapid increase in the amount of
155 information provided by CBAS as sample size increases.

156 We next evaluated the fraction of the sequences identified by CBAS for the smaller sized
157 groups that are not identified in the CBAS on the full dataset. Since we control the false
158 discovery proportion and not the family-wise error rate, we do not expect this value to be zero;
159 however, it was consistently small across all the group sizes evaluated (Fig 3b). The medians
160 across the different repeats of the same sample sizes across the different sample sizes are all
161 less than 2%.

162 Then, we evaluated the similarity of the sequences identified by CBAS with
163 nonoverlapping sets of subjects with the same group size (Fig 3c). At smaller groups sizes, this
164 overlap could be quite low (median 24.4%), even though there was sufficient evidence to
165 discriminate the strains (Fig. 2f). This means that the specific sequences identified by CBAS can
166 be quite variable from one experiment to the next, especially at smaller group sizes. Given that,
167 we sought to understand if the structure of the sequences comported with the conclusion from
168 the full dataset—that the CA line uses more of the same turns in a row and the w1118 line
169 more frequent changes in direction of turn. Even though the overlap between CBAS on
170 different groups might not be that high (Fig. 3c), the sequences identified still follow the same
171 structure as the full dataset in regard to number of the same turns in a row, (Fig 3d). As the
172 group size increases the output of CBAS consistently becomes more similar to the results from
173 the full dataset (Fig 3e), licensing more generalized conclusions.

174 **CBAS provides evidence for testing model-based analysis**

175 The second task we considered, is the two-step task developed to test the interplay
176 between model-based and model-free (reinforcement) learning (Fig 1b)⁴⁴. In this task, subjects
177 choose between pairs of images at two different stages. The image choice at the first stage
178 governs the pair of images from which the subject can chose at the second stage. Following
179 image choice at stage two, reward is delivered based on dynamic probabilities associated with
180 each of the images. Model-based analyses of variants of this task have led to conclusions about
181 the algorithms used in different brains regions^{44,45} and differences that underlie psychiatric
182 symptoms⁶; however, the interpretations are not without controversy⁴⁶, and the complexity of
183 behavior makes it difficult to evaluate ways in which the model might be missing features in the
184 data.

185 To test the ability for CBAS to provide useful information in such a situation, we
186 evaluated the open-source dataset from Gillan et al⁶. In that work, a large population of human
187 subjects performed the two-step task, and also answered a series of psychiatric symptom
188 questionnaires. The authors performed factor analysis on the different questionnaires and
189 found that factor 2, which they associated with intrusive thoughts and compulsive behavior,
190 had significant association with the way subjects performed that task. They found that subjects'
191 factor 2 score was negatively associated with model-based decision making, leading to their
192 conclusion that the greater the loading on intrusive thoughts and compulsive behavior for a
193 person, the less likely they were to utilize model-based decision making on the two-step task.

194 There were two different groups to compare in the fly dataset. In this dataset the
195 relevant comparison is correlation with a covariate (i.e. factor 2 score from the psychiatric
196 symptom questionnaires)⁶. Therefore, we extended CBAS to identify the sequences that
197 significantly correlate with this particular factor score (see methods). As with the CBAS applied
198 to the fly data, there are decisions that need to be made to apply CBAS to this human dataset
199 (Fig 4a). For this CBAS, the language is comprised of 8 different units: choosing image 1 or 2;
200 making a choice within set A or set B, distinguishing A and B depending on whether a reward
201 was delivered (shown in the figure as bold and underlined); and making 'no choice' (either at
202 the first or second image), which occurs in the data, albeit rarely. Following Akam et al.⁴⁵ we
203 collapsed image 3 and 4 into set A and image 5 and 6 into set B because the specific images
204 within the set are not relevant for the critical decision at the first stage. We ran CBAS on all
205 sequences up to 4 choices long and evaluated the rate of the sequences over the 400 choices of
206 the dataset (200 trials of the two-step task). In the open-source dataset, there was no subject
207 who did not reach criterion.

208 CBAS evaluated the correlation between the usage of each sequence and the factor 2
209 score (Fig 4b&S2a). Multiple sequences were significantly correlated (Fig 4c&d). To understand
210 the output of CBAS, we review the expectation associated with the hypothesis-driven analysis
211 of this experiment. This task was designed to evaluate the interplay between model-based and
212 model-free decision making⁴⁴. Model-based decision making develops an understanding of the
213 structure of the world (i.e. the model) and makes choices based on that understanding. Model-
214 free decision making makes choices based on reinforced past successful actions, without
215 developing an understanding of the structure of the world. A way to see the difference
216 between these two decision-making schemes is when a reward is delivered at set A ('**A**' in our
217 language) after having come from image 2. A regular model-free learner will reinforce the visits

218 to both 2 and A and will therefore be more likely to choose 2 on the next trial. The model-based
219 learner will have an understanding that choosing image 1 on the next opportunity makes it
220 more likely to return to set A (because of the common/rare transition structure in the
221 environment) and will therefore chose object 1⁴⁴.

222 Instead of identifying correlations with sequences that relate to either the model-free or
223 model-based learning, CBAS instead identifies a positive correlation with many sequences
224 involving being rewarded at A and then selecting image 2 or being rewarded at B and then
225 selecting image 1 (Fig 4d). Just as it was helpful with the fly data to find complete sets of
226 sequences that were significant (e.g. all of the sequences that occurred in the data with 10 left
227 or right turns in a row were identified as happening more in CA than w1118, etc.) we were also
228 able to find practically complete sets of sequences that were identified as being positively
229 correlated with the factor 2 score. Practically all of the sequences in the dataset that contained
230 **1A2** or **2B1**, were significantly correlated with this score; whereas none of the sequences with
231 **1B1** was identified as significant (and there were no instances of **2A2** in the dataset)
232 (Fig4e&S2b).

233 Sequences **1A2** or **2B1** can be classified as anti-model-based decisions. The subjects get
234 rewarded after choosing from the common side (A from 1 or B from 2), but then selects the
235 image that will rarely bring them back to the previously rewarded side (2 from A or 1 from B).
236 CBAS therefore identifies that anti-model-based learning is a prominent feature that correlates
237 with intrusive thought and compulsive behavior loading. Further experimentation will be
238 needed to understand the interplay between this mode of decision making, the task, and these
239 symptoms scales.

240 As with the fly data, we could evaluate the sample size needed to reach a given
241 statistical power for identifying any significant correlation between the factor 2 score and this
242 task. We resampled smaller group sizes and compared the CBAS for the true relationship
243 between the sequence counts and factor 2 score to a randomly generated set of factor 2 score
244 values that was drawn from the same distribution as the data (Fig S2c). A sample size of 900
245 individual provides a power >80% to detect significant correlations between the factor 2 score
246 and the subject (Fig 4f), which compares favorably to the sample size of ~1,200 – 1,600 subjects
247 that Gillan et al. calculated to generate their dataset⁶.

248 **CBAS identifies a phenotype consistent with ASD in *Scn2a* haploinsufficient rats**

249 The third task we considered, involves spatial alternation behavior in rats. For this task,
250 to get reward, the rats must alternate between pairs of arms of a track whilst visiting a different
251 arm of the track in between (Fig 1c). Spatial alternation is a common behavioral paradigm for
252 phenotyping and neurophysiology, and the discretization of the behavior into the arms chosen
253 by the animals forms the basis of many of the conclusions from these studies^{47–51}. However, our
254 recent work calls into question the assumptions and hypothesis that motivate the standard
255 analysis for spatial alternation behavior^{5,52}. Even though we developed reinforcement learning
256 agents to fit individual behavior, those agents showed clear differences from the way the
257 animals learned⁵, limiting their use for phenotyping. Therefore, we applied CBAS.

258 We sought to discriminate wild-type (WT) rats from those haploinsufficient for *Scn2a*
259 (*Scn2a*^{+/-}), a high confidence, large effect, autism spectrum disorder (ASD) risk gene^{53,54}. We
260 collected a dataset of over 200 rats performing six different spatial alternation contingencies
261 using our previously described automated behavioral system⁵ (see methods). Each spatial
262 alternation contingency was defined by the three arms of the track where alternation needed
263 to occur to get reward. For example, if the contingency was at arms 2, 3, and 4, reward would
264 be provided for every arm visit in the sequence 3-4-3-2-3-4.

265 The language for the spatial alternation CBAS was visiting each arm of the track and not
266 getting reward and visiting the goal arms and getting reward. That means that within each
267 contingency there were a total of 9 possibilities—6 unrewarded arms and 3 rewarded arms. For
268 example, if the contingency was at arms 2, 3, and 4, the language is composed of visiting arms
269 1, 2, 3, 4, 5, or 6 and not getting rewarded or visiting arms 2, 3, or 4, and getting rewarded.
270 Each contingency was considered separately, as visiting a sequence of arms during one
271 contingency likely means something different than visiting those same arms in a different
272 contingency. We chose to evaluate all sequences up to six choices long. For the criterion, we
273 used the trial at which each animal reached 100 perfect performance sequences four choices
274 long for each contingency, and we included animals that did not reach the criterion (Fig 5a) (see
275 methods).

276 In running the CBAS on this dataset, we evaluated the rate difference of >86,000
277 sequences across the six different contingencies, with a total of 1,476 sequences being
278 identified as significantly different between the *Scn2a*^{+/-} and WT littermates (Fig 5b&c). To
279 interpret common features of those sequences, we sought categories of sequences for which a
280 substantial fraction was identified as being significantly different between the groups: the
281 strategy that was informative for the fly and human datasets (Fig 2e&4e). One common feature
282 of spatial alternation behavior is the 3-arm structure of each contingency (e.g. arms 2, 3, and 4
283 are the only arms with the potential to be rewarded during contingency A and E). Therefore, in
284 all contingencies, we evaluated all sequences that exclusively contained all sets of three arms.
285 For example, within a contingency we identified all sequences exclusively containing arms 2, 3,
286 and 4, which means that every sequence within that category only contains visits to arms 2, 3,
287 and 4 (but does not have to visit all of the arms). We did this with sets of three arms,
288 independent of whether the choice was rewarded or not, as well as the set of three arms in a
289 contingency that were all rewarded. There are 21 sets of three arms for each contingency, and,
290 consistent with the relevance of these sets for the behavior, across all contingencies they
291 contained only ~25% of all the sequences that the rats performed during the experiment, but
292 ~75% of all of the significant sequences (Fig S3a).

293 Some of these sets of 3 arms are clearly interpretable. One such is the set of three arms
294 from the current contingency, which contains all sequences where the animal exclusively visits
295 the arms of the contingency but could make errors in the order of those arm visits (Fig 5d; left).
296 A second interpretable set consists of three arms from the current contingency that are
297 rewarded, which contains all sequences where the animal consistently performs the task
298 correctly. A third and fourth are the sets of three arms exclusively containing the previous (Fig
299 5d; right), or the one before previous, contingency arms, which contain sequences where the
300 animals repeat prior actions that are no longer optimally rewarded.

301 For all sets of three arms, in all the contingencies we tabulated the fraction of the total
302 sequences in the set that were significant either for WT > *Scn2a*^{+/-} or WT < *Scn2a*^{+/-}. We then
303 asked if any of these fractions were larger than would be expected from randomly distributing
304 the significant sequences across all possible sequence types within the entire category of sets
305 of three arms (see methods) and found consistent patterns in how the WT and *Scn2a*^{+/-} rats
306 differed (Fig 5e). In 4 out of the 6 contingencies, *Scn2a*^{+/-} rats showed increased usage of
307 sequences containing the current contingency, and in 3 out of the 6 contingencies, *Scn2a*^{+/-} rats
308 also showed increased usage of sequences related to prior contingencies (either 1 or 2
309 contingencies back). By contrast, in 4 out of the 6 contingencies the WT rats showed increased
310 usage of sequences containing the rewarded arms of the current contingency. This indicates
311 that WT rats are ultimately better at performing spatial alternation behavior than *Scn2a*^{+/-} rats,
312 and that *Scn2a*^{+/-} rats show difficulty transitioning from prior actions, repeating sequences
313 from prior contingencies, possibly consistent with restrictive and repetitive actions, which
314 forms one of the diagnostic criteria for ASD⁵⁵.

315 As with the fly and human data, we could evaluate the sample size needed to reach a
316 given statistical power to detect any difference between these genotypes. We resampled
317 smaller group sizes and compared the CBAS for the comparison between the WT and *Scn2a*^{+/-}
318 rats to the comparison between WT and itself and *Scn2a*^{+/-} and itself (Fig S3b). A sample size of
319 30 rats per group provides a power >80% to detect a difference between the two genotypes
320 (Fig 5f).

321 We have presented a general behavioral analysis method, CBAS, for identifying
322 interpretable behavioral components that is grounded in sequences of choices made by
323 subjects. There has been significant progress, in recent years, in tracking and analyzing, short
324 timescale actions of subjects¹³⁻²²; however, we lack methods for generally analyzing and
325 interpreting sequences of these actions or long-run choices and decisions of subjects during
326 behavior. CBAS provides just such a method and is applicable across a wide array of species and
327 different behavioral paradigms (Fig 1). It can be used to test models and hypotheses (Fig 2&4),
328 and, for instance, using generalizable simple principles about the relationship between
329 sequences and task contingencies, it can generate hypotheses in complex behaviors where
330 reliable computational understanding has yet to emerge (Fig 5). Through taking advantage of
331 large-scale data collection and rigorous statistical methods, CBAS has the potential to transform
332 our use of behavior in a comparable way to the ways that GWAS/WES/WGS changed the
333 paradigm for genomic studies.

334 **Methods**

335 **Animals:** All experiments on rats were conducted in accordance with University of California
336 San Francisco Institutional Animal Care and Use Committee and US National Institutes of Health
337 guidelines. Rats were fed standard rat chow (LabDiet 5001). To motivate the rats to perform
338 the task, reward was sweetened evaporated milk: 25 g of sugar per can (354 ml) of evaporated
339 milk (Carnation). The rats were food restricted to ~85% of their basal body weight.

340 The *Scn2a* mutant rats were generated due to funding from the Simons Foundation
341 Autism Research Initiative. Long Evans *Scn2a* mutant animals (LE-*Scn2a*^{em1Mcwi},

342 RRID:RGD_25394530) were generated at the Medical College of Wisconsin and shipped to the
343 University of California San Francisco for this study. Briefly, a single guide RNA targeting the
344 sequence GTGAAATCCAACCAATTCCA sequence within exon 5 of *Scn2a* was mixed with Cas9 (*S.*
345 *pyogenes*) protein (QB3 MacroLab, UC Berkeley) and injected into the pronucleus of fertilized
346 Long Evans (Crl:LE, Charles River Laboratories) embryos. Among the resulting offspring, a
347 mutant founder was identified harboring a net 4-bp deletion allele consisting of a 10-bp
348 deletion (rn7: chr3:50,364,411-50,364,420) along with a 6-bp insertion of TTCACT, inducing a
349 frameshift in the coding sequence predicted to truncate the normal protein after 193 amino
350 acids. The founder was backcrossed to the parental Crl:LE strain to establish a breeding colony.

351 **Spatial alternation behavior:** The automated behavior system for spatial alternation behavior
352 was previously described⁵. There are different symbols on each arm of the track serving as
353 proximal cues, and there are distal cues distinguishing the different walls of the room.
354 Pneumatic pistons (Clippard) open and close the doors. Python scripts, run through Trodes
355 (Spike Gadgets), control the logic of the automated system. The reward wells contain an
356 infrared beam adjacent to the reward spigot. The automated system uses the breakage of that
357 infrared beam to progress through the logic of the behavior. In addition to the infrared beam
358 and the spigot to deliver the reward, each reward well has an associated white light LED.

359 Each cohort of rats is divided into groups of four (or three) animals. The same groups
360 were maintained throughout the duration of the experiment. Within a group, a given rat is
361 always placed in the same rest box, and the four rats of a group serially perform the behavior.
362 The rats have multiple sessions on the track each day. Prior to beginning the first spatial
363 alternation contingency, the rats experience multiple days and sessions where they get
364 rewarded at any arm that they visit (provided it is not an immediate repeat). During this period
365 of the behavior, the duration of a session is defined by a fixed number of rewards, or a fixed
366 amount of time on the track (15 minutes), whichever came first. During the alternation task the
367 duration of a session was defined either by a fixed number of center arm visits and at least one
368 subsequent visit to any other arm, or a fixed amount of time on the track (15 minutes),
369 whichever came first.

370 The algorithm underlying the spatial alternation task is such that three arms on the track
371 have the potential for reward within a given contingency, for example during a contingency at
372 arms 2-3-4, arms 2, 3, and 4 have the potential to be rewarded, and arms 1, 5, and 6 do not. Of
373 those three arms we refer to the middle of the three arms as the center arm (arm 3 in the
374 above example) and the other two arms as the outer arms (arms 2 and 4 in the above example).
375 Reward is delivered at the center arms if and only if: 1) the immediately preceding arm whose
376 reward well infrared beam was broken was not the center arm. Reward was delivered at the
377 outer two arms if and only if: 1) the immediately preceding arm whose reward well infrared
378 beam was broken was the center arm, and 2) prior to breaking the infrared beam at the center
379 arm, the most recently broken outer arm infrared beam was not the currently broken outer
380 arm infrared beam. The one exception to the outer arm rules was at the beginning of a session,
381 if no outer arm infrared beam was broken prior to the first infrared beam break at the center
382 arm, then only the first condition had to be met.

383 For the running of the behavior, the infrared beam break determined an arm visit;
384 however, the rats sometimes go down an arm, get very close to the reward wells, but do not
385 break the infrared beam. Therefore, for all the analyses described for the rats, an arm choice is
386 defined as when a rat gets close to a reward well. These times were extracted from a video
387 recording of the behavior. These, effective, missed pokes were more frequent at the beginning
388 of a contingency. This proximity-based definition of an arm visit added additional arm visits to
389 those defined by the infrared beam breaks, and none of them could ever be rewarded, nor alter
390 the logic of the underlying algorithm. However, because of the non-Markovian nature of the
391 reward contingency, they could affect the rewards provided for subsequent choices.

392 A total of 121 WT (66 males, 55 females) and 120 *Scn2a*^{+/-} (66 males, 54 females) rats
393 were run on the spatial alternation task. WT and *Scn2a*^{+/-} rats were littermates and were
394 housed together prior to their being food restricted before the behavior. During the behavior
395 and food restriction the rats were single housed. 1 WT rat died after finishing the first spatial
396 alternation contingency, 1 WT rat died after finishing the fourth spatial alternation contingency,
397 and 1 *Scn2a*^{+/-} rat died after finishing the first spatial alternation contingency. The data from the
398 animals that died was included up until their expiration. For the first contingency: 5/121 WT
399 rats and 15/120 *Scn2a*^{+/-} rats did not reach the CBAS criterion. For the second contingency:
400 4/120 WT rats and 19/119 *Scn2a*^{+/-} rats did not reach the CBAS criterion; for the third
401 contingency: 4/120 WT rats and 12/119 *Scn2a*^{+/-} rats did not reach the CBAS criterion; for the
402 fourth contingency: 2/120 WT rats and 12/119 *Scn2a*^{+/-} rats did not reach the CBAS criterion;
403 for the fifth contingency: 4/119 WT rats and 14/119 *Scn2a*^{+/-} rats did not reach the CBAS
404 criterion; and for the sixth contingency: 3/119 WT rats and 12/119 *Scn2a*^{+/-} rats did not reach
405 the CBAS criterion.

406 **Romano-Wolf resampling based multiple comparisons correction:** We follow the terminology
407 and description laid out in Clarke et. al⁵⁶ to describe the Romano-Wolf multiple comparison
408 correction. First, we describe the way the method corrects for the family-wise error rate
409 (FWER) and then explain how the procedure is extended to provide median control of the false
410 discovery proportion. FWER control at a level of α means that across all comparisons there is a
411 α percent chance of having at least one false positive rejection of a null hypothesis. The
412 Romano-Wolf procedure provides FWER control through resampling the data¹. It tests a total of
413 S hypotheses.

414 It is not generally known if the Romano-Wolf procedure controls for type III, or
415 directional, errors. Type III errors are errors in the sign, or direction, of the conclusion. For
416 example, if a statistical test provided information to reject the null hypothesis $\theta_1 = \theta_2$, and you
417 then concluded that $\theta_1 > \theta_2$, when in fact $\theta_1 < \theta_2$. Therefore, instead of running a single two-
418 tailed test for each sequence, we run two one-tailed tests for each sequence. Therefore, the
419 total number of hypotheses tested, S , is twice the total number of sequences being compared.
420 For CBAS those hypotheses take one of two forms: 1) the rate of each sequence (r_s) is the same
421 between two groups $\Delta r_s = 0$, or 2) that there is no correlation (ρ_s) between each sequence and
422 a covariate of interest, $\rho_s = 0$. For the one-tailed versions of each hypothesis, we ask if $r_{s_1} -$
423 $r_{s_2} > 0$ and $r_{s_2} - r_{s_1} > 0$ for the comparison CBAS, where r_{s_N} is the rate of sequence s for
424 group N , or if $\rho_s > 0$ and $\rho_s < 0$ for the correlational CBAS.

425 The first step in the procedure is to create a studentized test statistic for each
 426 hypothesis. The studentization is different based on whether the CBAS is comparing two groups
 427 or calculating a correlation. In the case where two groups are being compared:

$$t_s = \frac{\Delta r_s}{\sigma_s} \quad (7)$$

428 where σ_s is the standard error of Δr_s , which we calculate by combining the standard error of
 429 the mean of the rate for each group using error propagation, i.e. $\sigma_s = \sqrt{\sigma_{s_1}^2 + \sigma_{s_2}^2}$, where σ_{s_N}
 430 is the standard error of the occurrence rate of sequence s for group N .

431 In the case where the correlation is being calculated the studentized test statistic³ is:

$$t_s = \frac{\sqrt{n}\hat{\rho}_n}{\hat{t}_n} \quad (8)$$

432
 433 where:

$$\hat{\rho}_s = \frac{\sum X_i Y_i - n\bar{X}_s\bar{Y}}{\sqrt{\sum(X_i - \bar{X}_n)^2 \sum(Y_i - \bar{Y})^2}} \quad (9)$$

$$\hat{t}_s = \frac{\sqrt{\frac{1}{n}\sum_{i=1}^n(X_i - \bar{X}_s)^2(Y_i - \bar{Y})^2}}{\sqrt{\frac{1}{n}\sum_{i=1}^n(X_i - \bar{X}_s)^2} \sqrt{\frac{1}{n}\sum_{i=1}^n(Y_i - \bar{Y})^2}} \quad (10)$$

434 For eq. 8, 9, and 10, n is the number of subjects for which the correlation is being calculated, X_i
 435 and Y_i and the values of the metrics being correlated (in our case, the sequence count and
 436 factor 2 score for each individual respectively), \bar{X}_s is the mean sequence count for the specific
 437 sequence being considered, and \bar{Y} is the mean of the covariate of interest (factor 2 score, which
 438 is the same for any sequence).

439 When two groups are being compared, we resample from the entire population with
 440 replacement (separately for each group) and build up a null distribution by bootstrapping M
 441 times. The test statistic from the m^{th} bootstrap sample for $m = 1, \dots, M$ is:

$$t_s^{*,m} = \frac{\Delta r_s^{*,m}}{\sigma_s^{*,m}} \quad (11)$$

442 where, $\Delta r_s^{*,m}$ is the difference in the rate of each sequence whilst resampling, with
 443 replacement, from the entire population, ignoring the group labels. $\sigma_s^{*,m}$ is the accompanying
 444 standard error of the resampled groups. The resampled group sizes are the same as the two
 445 groups of interest.

446 In the case where the correlation is being calculated, the test statistic based on the m^{th}
 447 bootstrap sample is:

$$t_s^{*,m} = \frac{\sqrt{n}\hat{\rho}_s^{*,m}}{\hat{\tau}_s^{*,m}} \quad (12)$$

448 where, $\hat{\rho}_s^{*,m}$ is the correlations of each sequence whilst resampling, with replacement, from the
 449 entire population, ignoring the group labels. $\hat{\tau}_s^{*,m}$ is the accompanying normalization factor of
 450 the resampled groups. The resampled group size is the same as the original.

451 We used a value of $M = 10,000$. Importantly, for each individual resampling, m , the
 452 same resampled set is used for all sequences.

453 The test statistics, and their accompanying estimators are ordered from largest to
 454 smallest values. This creates a $M \times S$ matrix where each column contains all the estimators of
 455 the test statistics. The first column contains the estimators from the largest test statistic, the
 456 second column contains the estimators from the second largest, etc.

457 To define the distribution for which each test statistic is compared, which then
 458 determines the adjusted p-value, the following algorithm is used. The first sequence considered
 459 is the one with the maximum test statistic, t_s . Its comparison distribution is defined as the
 460 maximum value within each row of the matrix of estimators of the test statistic:

$$\max(t^{*,m}) = \max\{t_1^{*,m}, \dots, t_S^{*,m}\} \quad (13)$$

461 which provides a total of M values, $t^{*,m}$ (there is no longer an association with s , because these
 462 values can come from a resampling of any of the sequences). Using those M values, the
 463 adjusted p-value is calculated as follows:

$$p_s^{adj} = \frac{\#\{\max(t^{*,m}) \geq t_s\} + 1}{M + 1} \quad (14)$$

465 After calculating p_s^{adj} for the first sequence, the column with the test statistic
 466 estimators generated from the first sequence is removed from the matrix. This now leaves a
 467 matrix that is $M \times (S - 1)$. The above procedure is then used to calculate p_s^{adj} for the
 468 sequence with the second largest test statistic, and then its column of test statistic estimators is
 469 removed, etc until all p_s^{adj} have been calculated. Following the algorithm described in Clarke et
 470 al.⁵⁶, we enforce monotonicity of the p-values by resetting the p-value for each sequence:

$$p_s^{adj} = \max\{p_s^{adj}, p_{s-1}^{adj}\} \quad (15)$$

471 This is done prior to calculating the p-value for the next sequence.

472 To control the false discovery proportion (i.e., control the number of false positives
 473 divided by the total number of hypotheses rejected), the idea of k-FWER is introduced². For

474 control of the FWER, k is equal to 1, and that leads to an α percent chance that there is at least
475 1 false positive among all hypotheses rejected (most commonly $\alpha = 0.05$). If k equals 2, then
476 there is an α percent chance that there are at least 2 false positives among all hypotheses
477 rejected. Therefore, to get control of the false discovery proportion we need to find the k that
478 provides the proportion of interest given the number of hypotheses rejected. So, if we want a
479 false discovery proportion, γ , of 0.05, we need $k \sim 0.05 \times$ number of hypotheses rejected.

480 Romano and Wolf also derived an algorithm to do just that. The algorithm is as follows.
481 Start with $k = 1$. Apply the k -FWER procedure, and note the total number of hypotheses
482 rejected, N . If $N < \frac{k}{\gamma} - 1$, stop and you have identified k . Otherwise, increase k by 1, and
483 repeat². The way you determine k -FWER is in eq. 13, instead of taking the maximum value in
484 each row, you take the k^{th} largest value. Finally, to get median control of the false discovery
485 proportion, $\alpha = 0.5$. This means that 50% of the time you will get a value greater than γ , and
486 50% of the time you will get a value less than γ , leading to control of the median². This is a
487 similar decision to what is done when calculating the false discovery rate with Benjamini-
488 Hochberg or Benjamini-Yekutieli, except the false discovery rate controls the mean of the false
489 discovery proportion instead of the median.

490 **Power estimation for CBAS.** To estimate the statistical power of CBAS for a given sample size
491 we resampled the dataset without replacement and ran a CBAS to determine the number of
492 significant sequences. For each sample size we performed 20 repeats (Fig S1b, S2c, and S3b).
493 We also ran a CBAS comparing each group to itself (for the fly and rat datasets) (Fig S1b and
494 S3b), with 20 repeats for each group; or correlating the sequence counts with a randomly
495 generated set of factor 2 scores drawn from the same distribution as the actual factor 2 scores
496 (for the human dataset) (Fig S2c), with 40 repeats. Then for each sample size, the power is
497 estimated by identifying the largest 20 values of significant sequences, and determining the
498 fraction of those values that were generated by the comparison of the two groups or the
499 correlation with the data with the actual factor 2 scores.

500 **Spatial alternation category fraction determination of significance.** There are a total of 20
501 different sets of three arms, and 1 set of three arms with all choices being rewarded. For each
502 of the six contingencies the fraction of sequences that exclusively contains the set of three arms
503 that are significant is calculated separately for significant sequences greater in the WT and
504 greater in the *Scn2a*^{+/-}. That means that there a total of $21 \times 2 \times 6 = 252$ fractions across all
505 contingencies. To determine whether the fraction is significantly larger than expected by
506 chance we proceeded as follows: Within a contingency, we determined the sequences that
507 belonged to any of the 21 categories. Then, separately for sequences greater in WT and
508 sequences greater in *Scn2a*^{+/-}, we permuted the association between those sequences and
509 whether or not they were significant. With that permuted association we recalculated the
510 fraction of sequences in each category that were significant. We repeated that process 25,000
511 times and determined significance by calculating the number of permuted fractions that were
512 greater than or equal to the actual fraction value for each category and divided that by 25,001
513 (as with eq. 15). That number was then corrected for multiple comparisons using the
514 Bonferroni method and multiplied by 252, the number of tests being performed. Any category
515 whose corrected p-value was < 0.05 , was determined to be significantly larger than expected.

516 **Data availability.** Data will be made available upon reasonable request to the lead author.

517 **Code availability.** Code used to calculate CBAS will be posted to Github upon publication.

518 **Acknowledgements.** We thank Claire Gillan and Benjamin de Bivort for making their data
519 publicly available. We thank Thomas Akam and Adam Frank for helpful comments on the
520 manuscript. We thank the Simons Foundation Autism Research Initiative for supporting the
521 generation of the *Scn2a*^{+/-} rat line. This work was supported by grants from the Jane Coffin
522 Childs Memorial Fund for Medical Research (D.B.K.), the UCSF Physician Scientist Scholars
523 Program (D.B.K.), an NIH R25 (R25MH060482) (D.B.K.), a Simons Foundation Autism Research
524 Initiative grant (899599) (D.B.K.), the Max Planck Society (P.D) and the Humboldt Foundation
525 (P.D.). Some of the data used in this study were collected in the laboratory of Loren Frank at the
526 University of California, San Francisco.

527 **Author Contributions.** D.B.K. and P.D. designed the study, D.B.K., J.P.R., and P.D. developed the
528 analysis method, D.B.K., G.W., and C.H. collected the data, D.B.K. analyzed the data, D.B.K and
529 P.D. wrote the manuscript, and D.B.K., J.P.R. and P.D. edited the manuscript.

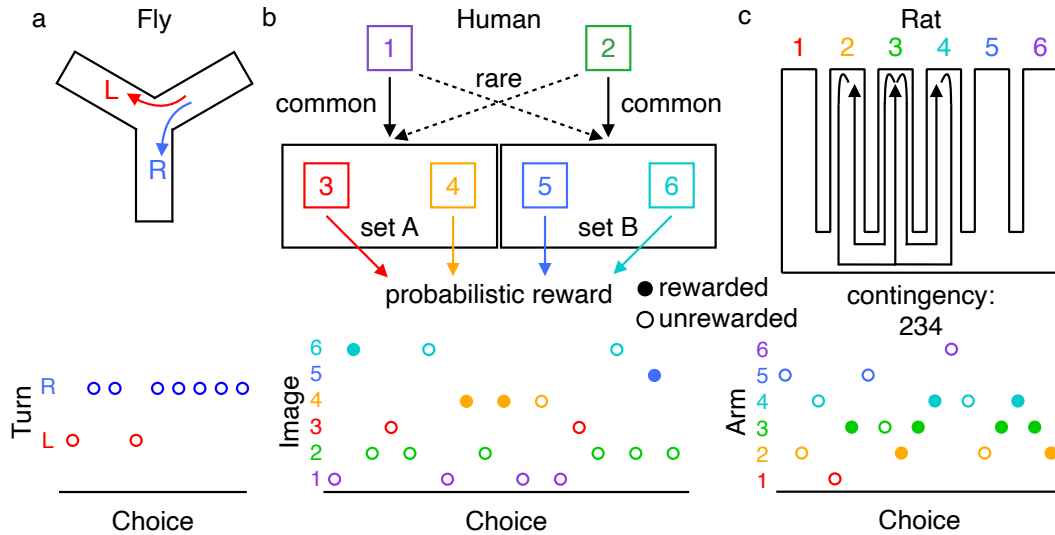
530 **Declaration of interests.** The authors declare no competing interests.

- 531 1. Romano, J. P. & Wolf, M. Exact and Approximate Stepdown Methods for Multiple Hypothesis
532 Testing. *J Am Stat Assoc* **100**, 94–108 (2005).
- 533 2. Romano, J. P. & Wolf, M. Control of generalized error rates in multiple testing. *Ann Statistics*
534 **35**, 1378–1408 (2007).
- 535 3. DiCiccio, C. J. & Romano, J. P. Robust Permutation Tests For Correlation And Regression
536 Coefficients. *J. Am. Stat. Assoc.* **112**, 1211–1220 (2017).
- 537 4. Buchanan, S. M., Kain, J. S. & Bivort, B. L. de. Neuronal control of locomotor handedness in
538 *Drosophila*. *Proc National Acad Sci* **112**, 6700–6705 (2015).
- 539 5. Kastner, D. B. *et al.* Spatial preferences account for inter-animal variability during the
540 continual learning of a dynamic cognitive task. *Cell Reports* **39**, 110708–110708 (2022).
- 541 6. Gillan, C. M., Kosinski, M., Whelan, R., Phelps, E. A. & Daw, N. D. Characterizing a psychiatric
542 symptom dimension related to deficits in goal-directed control. *eLife* **5**, e11305 (2016).
- 543 7. Krakauer, J. W., Ghazanfar, A. A., Gomez-Marin, A., Maclver, M. A. & Poeppel, D.
544 Neuroscience Needs Behavior: Correcting a Reductionist Bias. *Neuron* **93**, 480–490 (2017).
- 545 8. Bialek, W. On the dimensionality of behavior. *P Natl Acad Sci Usa* **119**, e2021860119 (2022).
- 546 9. Meister, M. Learning, fast and slow. *Curr Opin Neurobiol* **75**, 102555 (2022).
- 547 10. Gomez-Marin, A. & Ghazanfar, A. A. The Life of Behavior. *Neuron* **104**, 25–36 (2019).
- 548 11. Trillmich, F., Günther, A., Müller, C., Reinhold, K. & Sachser, N. New perspectives in
549 behavioural development: adaptive shaping of behaviour over a lifetime? *Front Zool* **12**, S1
550 (2015).
- 551 12. Dayan, P., Roiser, J. P. & Viding, E. The first steps on long marches: the costs of active
552 observation. in *Rethinking Biopsychosocial Psychiatry*. (eds. Davies, W., Roache, R. & Savulescu,
553 J.) (Rethinking Biopsychosocial Psychiatry., 2018).
- 554 13. Hession, L. E., Sabnis, G. S., Churchill, G. A. & Kumar, V. A machine-vision-based frailty index
555 for mice. *Nat. Aging* **2**, 756–766 (2022).
- 556 14. Geuther, B. Q. *et al.* Action detection using a neural network elucidates the genetics of
557 mouse grooming behavior. *eLife* **10**, e63207 (2021).
- 558 15. Geuther, B. *et al.* High-throughput visual assessment of sleep stages in mice using machine
559 learning. *SLEEP* **45**, zsab260 (2021).

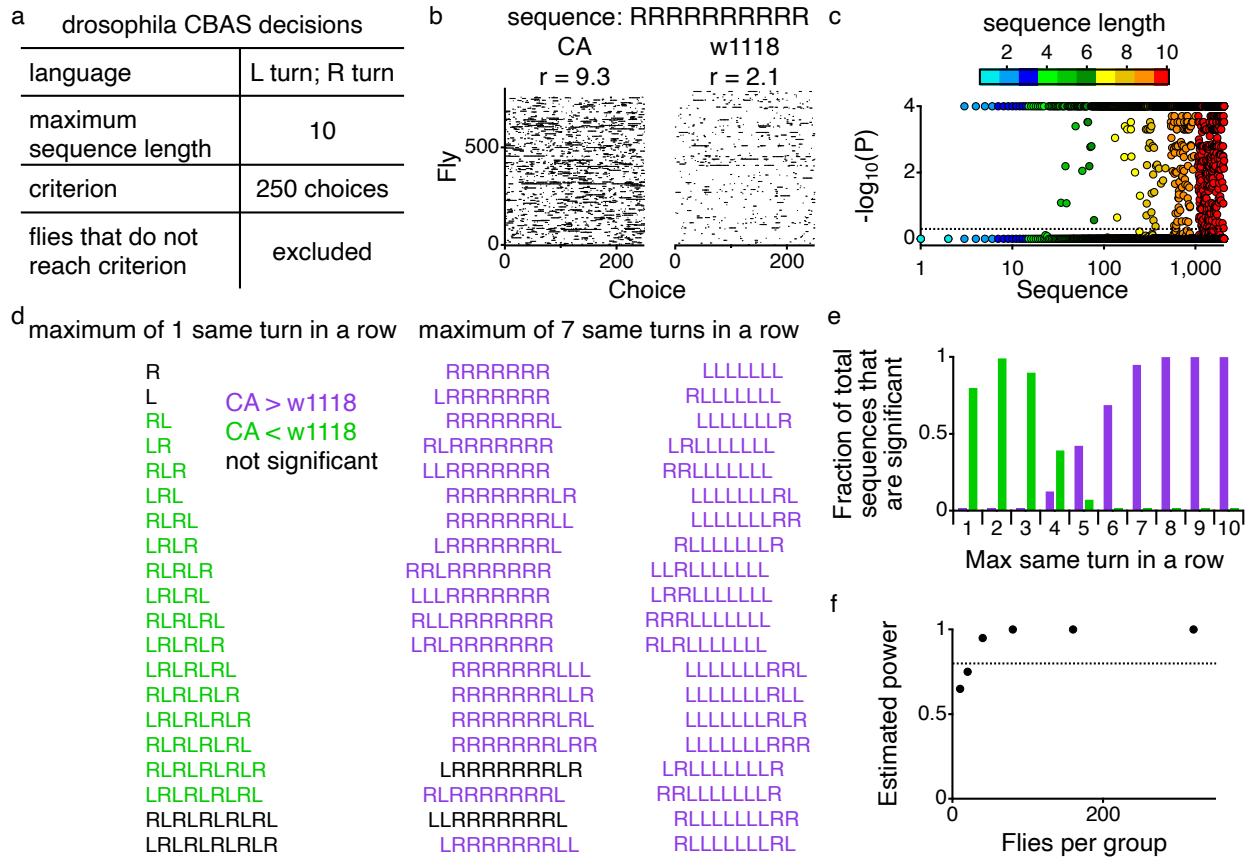
- 560 16. Sheppard, K. *et al.* Stride-level analysis of mouse open field behavior using deep-learning-
561 based pose estimation. *Cell Rep.* **38**, 110231 (2022).
- 562 17. Weinreb, C. *et al.* Keypoint-MoSeq: parsing behavior by linking point tracking to pose
563 dynamics. *bioRxiv* 2023.03.16.532307 (2023) doi:10.1101/2023.03.16.532307.
- 564 18. Wiltschko, A. B. *et al.* Mapping Sub-Second Structure in Mouse Behavior. *Neuron* **88**, 1121–
565 1135 (2015).
- 566 19. Pereira, T. D. *et al.* Fast animal pose estimation using deep neural networks. *Nat. Methods*
567 **16**, 117–125 (2019).
- 568 20. Pereira, T. D. *et al.* SLEAP: A deep learning system for multi-animal pose tracking. *Nat.*
569 *Methods* **19**, 486–495 (2022).
- 570 21. Marshall, J. D. *et al.* Continuous Whole-Body 3D Kinematic Recordings across the Rodent
571 Behavioral Repertoire. *Neuron* **109**, 420–437.e8 (2020).
- 572 22. Dunn, T. W. *et al.* Geometric deep learning enables 3D kinematic profiling across species
573 and environments. *Nat. Methods* **18**, 564–573 (2021).
- 574 23. Huys, Q. J. M., Maia, T. V. & Frank, M. J. Computational psychiatry as a bridge from
575 neuroscience to clinical applications. *Nat Neurosci* **19**, 404–413 (2016).
- 576 24. Kumar, P., Dayan, P. & Wolfers, T. From Complexity to Precision—Charting Decision-Making
577 Through Normative Modeling. *JAMA Psychiatry* **81**, (2024).
- 578 25. Nickerson, R. S. Confirmation bias: A ubiquitous phenomenon in many guises. *Rev Gen*
579 *Psychol* **2**, 175–220 (1998).
- 580 26. Kerr, N. L. HARKing: Hypothesizing After the Results are Known. *Pers Soc Psychol Rev* **2**,
581 196–217 (1998).
- 582 27. Dezfouli, A., Griffiths, K., Ramos, F., Dayan, P. & Balleine, B. W. Models that learn how
583 humans learn: The case of decision-making and its disorders. *PLoS Comput. Biol.* **15**, e1006903
584 (2019).
- 585 28. Eckstein, M. K., Summerfield, C., Daw, N. D. & Miller, K. J. Predictive and Interpretable:
586 Combining Artificial Neural Networks and Classic Cognitive Models to Understand Human
587 Learning and Decision Making. *bioRxiv* 2023.05.17.541226 (2023)
588 doi:10.1101/2023.05.17.541226.
- 589 29. Wiltschko, A. B. *et al.* Revealing the structure of pharmacobehavioral space through Motion
590 Sequencing. *Nat. Neurosci.* **23**, 1433–1443 (2020).

- 591 30. Gschwind, T. *et al.* Hidden behavioral fingerprints in epilepsy. *Neuron* **111**, 1440-1452.e5
592 (2023).
- 593 31. Uffelmann, E. *et al.* Genome-wide association studies. *Nat Rev Methods Primers* **1**, 59
594 (2021).
- 595 32. Farrell, M. S. *et al.* Evaluating historical candidate genes for schizophrenia. *Mol Psychiatr* **20**,
596 555–562 (2015).
- 597 33. Border, R. *et al.* No Support for Historical Candidate Gene or Candidate Gene-by-Interaction
598 Hypotheses for Major Depression Across Multiple Large Samples. *Am J Psychiat* **176**, 376–387
599 (2019).
- 600 34. Duncan, L. E., Ostacher, M. & Ballon, J. How genome-wide association studies (GWAS) made
601 traditional candidate gene studies obsolete. *Neuropsychopharmacol* **44**, 1518–1523 (2019).
- 602 35. Ioannidis, J. P. A. Why Most Published Research Findings Are False. *Plos Med* **2**, e124 (2005).
- 603 36. Markowitz, J. E. *et al.* Spontaneous behaviour is structured by reinforcement without
604 explicit reward. *Nature* **614**, 108–117 (2023).
- 605 37. Cui, T. *et al.* Gene–gene interaction detection with deep learning. *Commun. Biol.* **5**, 1238
606 (2022).
- 607 38. Fang, G. *et al.* Discovering genetic interactions bridging pathways in genome-wide
608 association studies. *Nat. Commun.* **10**, 4274 (2019).
- 609 39. Skutt-Kakaria, K., Reimers, P., Currier, T. A., Werkhoven, Z. & Bivort, B. L. de. A neural circuit
610 basis for context-modulation of individual locomotor behavior. *bioRxiv* 797126 (2019)
611 doi:10.1101/797126.
- 612 40. Ayroles, J. F. *et al.* Behavioral idiosyncrasy reveals genetic control of phenotypic variability.
613 *Proc National Acad Sci* **112**, 6706–6711 (2015).
- 614 41. Werkhoven, Z. *et al.* The structure of behavioral variation within a genotype. *eLife* **10**,
615 e64988 (2021).
- 616 42. Akhund-Zade, J., Ho, S., O’Leary, C. & Bivort, B. de. The effect of environmental enrichment
617 on behavioral variability depends on genotype, behavior, and type of enrichment. *J. Exp. Biol.*
618 **222**, jeb202234 (2019).
- 619 43. Bivort, B. de *et al.* Precise Quantification of Behavioral Individuality From 80 Million
620 Decisions Across 183,000 Flies. *Front. Behav. Neurosci.* **16**, 836626 (2022).

- 621 44. Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P. & Dolan, R. J. Model-Based Influences on
622 Humans' Choices and Striatal Prediction Errors. *Neuron* **69**, 1204–1215 (2011).
- 623 45. Akam, T. *et al.* The Anterior Cingulate Cortex Predicts Future States to Mediate Model-
624 Based Action Selection. *Neuron* **109**, 149-163.e7 (2020).
- 625 46. Silva, C. F. da & Hare, T. A. Humans primarily use model-based inference in the two-stage
626 task. *Nat. Hum. Behav.* **4**, 1053–1066 (2020).
- 627 47. Kim, S. M. & Frank, L. M. Hippocampal Lesions Impair Rapid Learning of a Continuous Spatial
628 Alternation Task. *Plos One* **4**, e5494 (2009).
- 629 48. Singer, A. C., Carr, M. F., Karlsson, M. P. & Frank, L. M. Hippocampal SWR Activity Predicts
630 Correct Decisions during the Initial Learning of an Alternation Task. *Neuron* **77**, 1163–1173
631 (2013).
- 632 49. Kay, K. *et al.* Constant Sub-second Cycling between Representations of Possible Futures in
633 the Hippocampus. *Cell* **180**, 552-567.e25 (2020).
- 634 50. Fernández-Ruiz, A. *et al.* Long-duration hippocampal sharp wave ripples improve memory.
635 *Science* **364**, 1082–1086 (2019).
- 636 51. Sigurdsson, T., Stark, K. L., Karayiorgou, M., Gogos, J. A. & Gordon, J. A. Impaired
637 hippocampal–prefrontal synchrony in a genetic mouse model of schizophrenia. *Nature* **464**,
638 763–767 (2010).
- 639 52. Kastner, D. B., Gillespie, A. K., Dayan, P. & Frank, L. M. Memory alone does not account for
640 the speed of learning of a simple spatial alternation task in rats. *J Neurosci* **40**, JN-RM-0972-20
641 (2020).
- 642 53. Sanders, S. J. *et al.* De novo mutations revealed by whole-exome sequencing are strongly
643 associated with autism. *Nature* **485**, 237–241 (2012).
- 644 54. Satterstrom, F. K. *et al.* Large-Scale Exome Sequencing Study Implicates Both
645 Developmental and Functional Changes in the Neurobiology of Autism. *Cell* **180**, 568-584.e23
646 (2020).
- 647 55. Lord, C. *et al.* Autism spectrum disorder. *Nat. Rev. Dis. Prim.* **6**, 5 (2020).
- 648 56. Clarke, D., Romano, J. P. & Wolf, M. The Romano-Wolf Multiple Hypothesis Correction in
649 Stata. *Ssrn Electron J* (2020) doi:10.2139/ssrn.3513687.

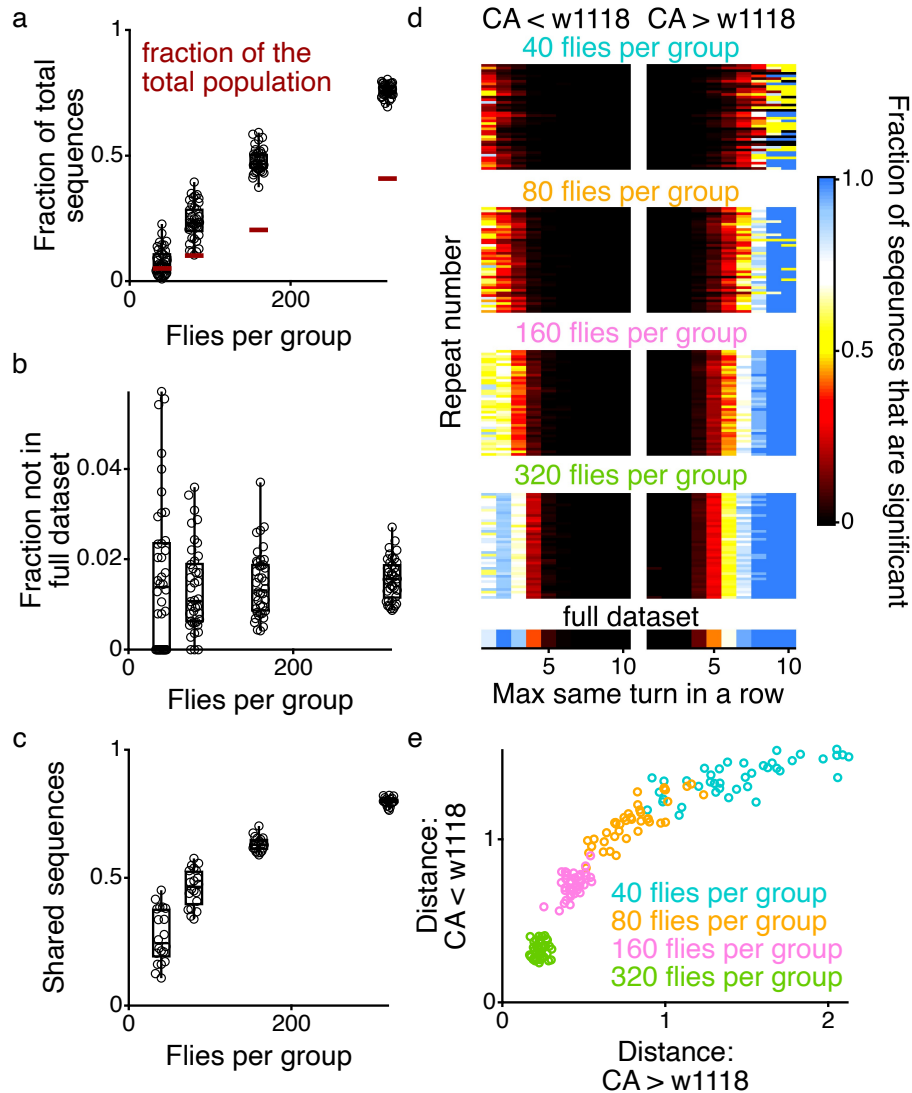


650 **Figure 1. Choice is a common discretization for behavior.** In the three behavioral tasks under
651 consideration in this work, the actions of the subjects are broken down into a series of choices.
652 (a) Fly y-maze (top), and example set of choices of left (red) and right (blue) turns (bottom) of
653 an individual fly. Data come from Buchanan et al.⁴. (b) Two-step task, performed by human
654 subjects (top), and an example set of choices of an individual subject (bottom). The colors
655 correspond to the different objects chosen (bottom). Data come from Gillan et al.⁶. (c) Spatial
656 alternation behavior, performed by rats (top), and an example set of choices of an individual rat
657 (bottom). The contingency is defined as the 3 arms that can be rewarded, and the 6-arm track
658 enables different sets of three arms to be rewarded. In b and c, filled in circles indicate that
659 reward was received.



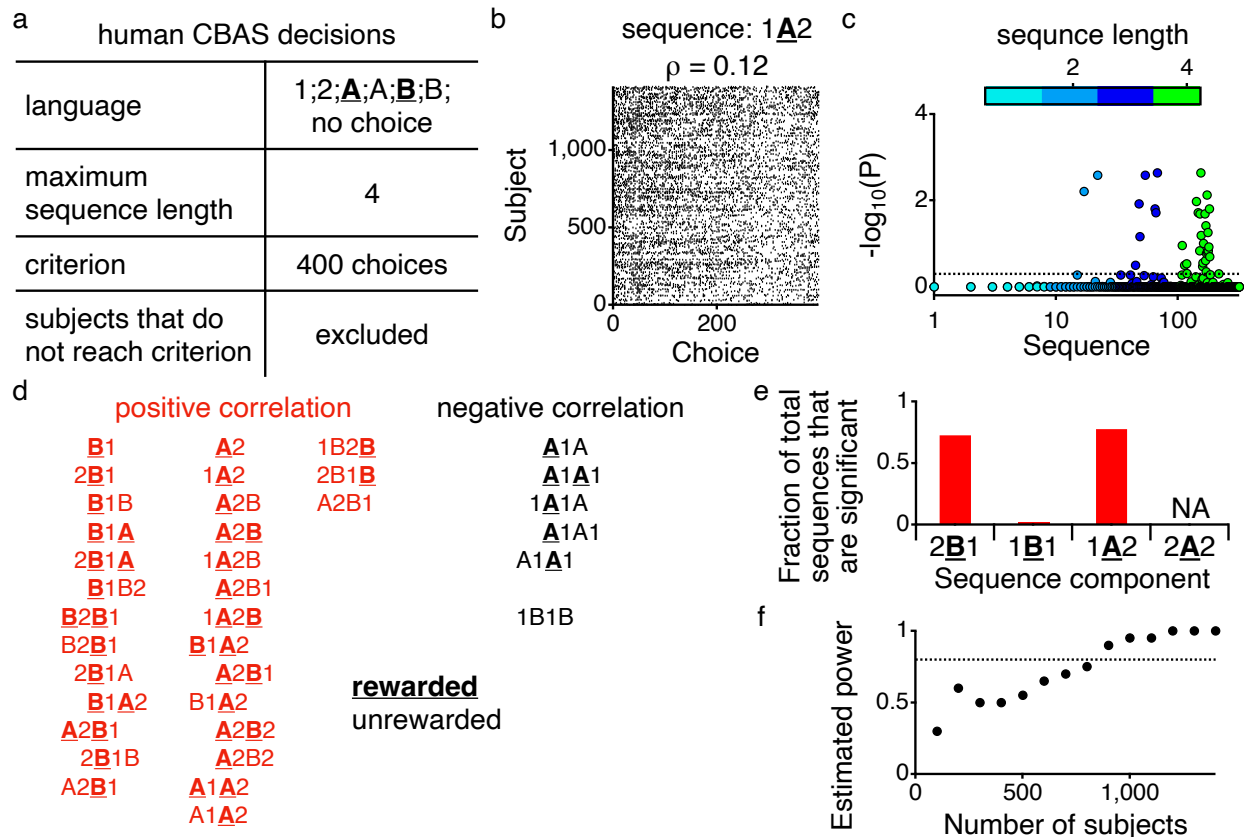
660 **Figure 2. CBAS applied to y-maze produces interpretable differences between fly strains.**
661 CBAS was applied to flies tracked on a y-maze (Fig 1a) (a) This table demarcates the decisions
662 made to perform the CBAS on this dataset. There was a total of 1,225 Cambridge-A (CA) and
663 1,372 w1118 flies in the dataset, 466 CA and 565 w1118 were excluded from analysis due to
664 not reaching a total of at least 250 turns. (b) The occurrence of a single sequence of turns (10
665 right turns in a row) in the two strains during the first 250 turns. The rows show the occurrence
666 of that sequence of turns for all individual flies from each strain. The CA strain has a rate of 9.3,
667 meaning that, on average, each fly utilizes this sequence 9.3 times. The W1118 strain has a rate
668 of 2.1, meaning that, on average, each fly utilizes this sequence 2.1 times. (c) CBAS Manhattan
669 plot displays the p-value for each sequence. Each sequence has two p-values on this plot, one
670 for CA > w1118, and the other for CA < w1118. The sequences are ordered based on the
671 number of choices in the sequence, and they are displayed on a log scale to make all sequence
672 lengths visible. Within a given sequence length, the sequences are ordered based on frequency
673 of occurrence in the entire dataset. The horizontal dotted line indicates the significance
674 threshold of 5% control of the median false discovery proportion. A total of 10,000 resamplings
675 were used to calculate the p-values making the maximum value on the plot 4.00004
676 $(-\log(\frac{1}{10,001}))$. (d) Every sequence that either has a maximum of 1 turn in a direction in a row
677 (left) or a maximum of 7 turns in the same direction in a row (right). Sequences are colored
678 based on whether CBAS identifies them as occurring significantly more in the CA strain (purple),
679 w1118 strain (green) or as not significantly different (black). The sequences on the right are
680 aligned to the 7 right or left turns. (e) The maximum number of the same turns in a row was

681 calculated for every sequence in the dataset (2,046 total sequences), and for every number of
682 maximum turns in a row, the fraction of sequences that whose prevalence was significantly
683 greater in the CA strain (purple) or the w1118 strain (green) is plotted. The CA and w1118
684 strains showed separation with other related metrics even for the max turns that show both
685 significance for both direction of comparison (e.g. max same turn in a row of 4) (Fig S1a). (f)
686 Power estimate for different sized groups of flies in each strain (see methods). Horizontal
687 dotted line shows a value of 80%.

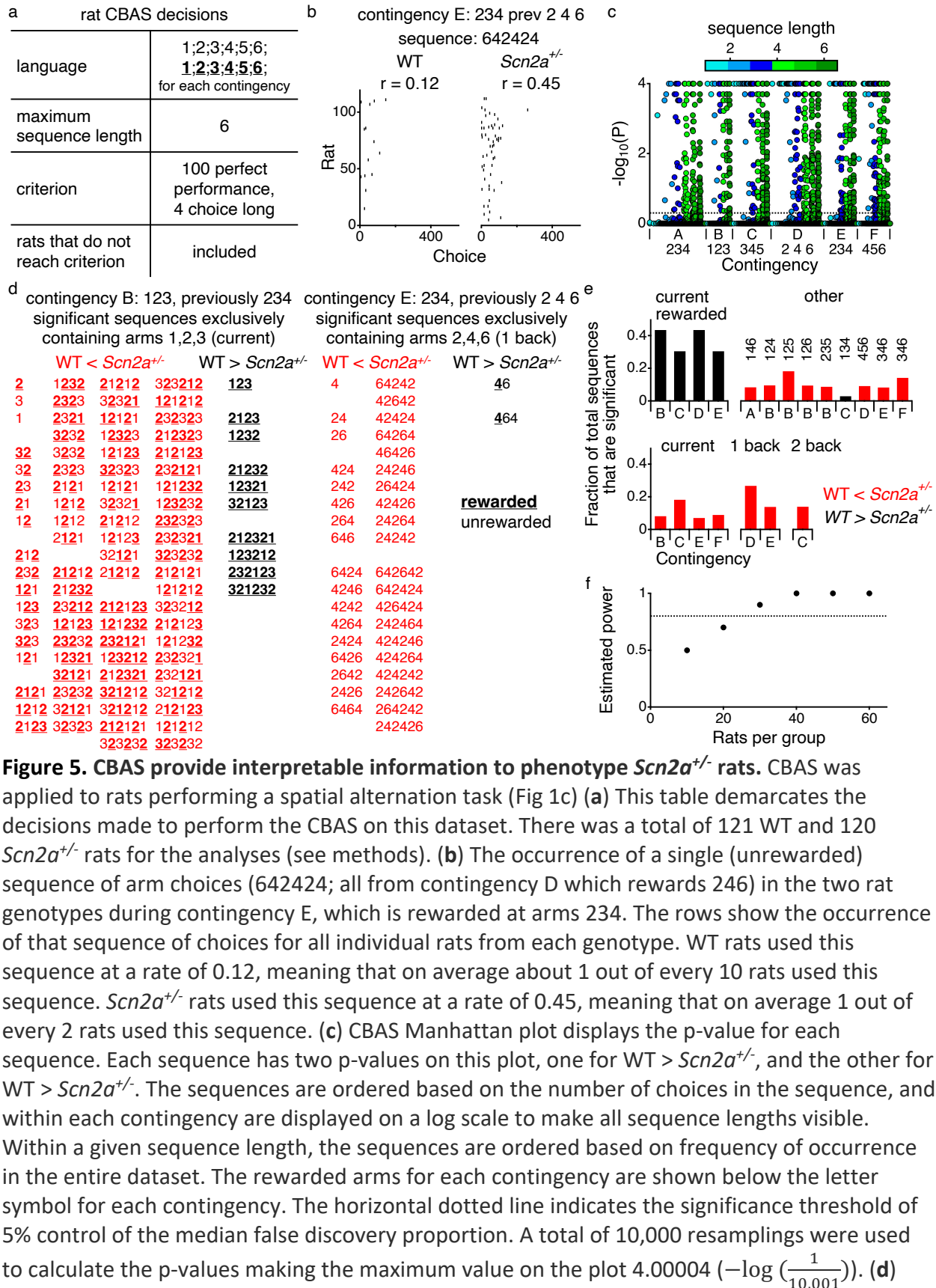


688 **Figure 3. Common, but degraded output with CBAS on smaller sample sizes.** For all panels in
689 this figure, many repeats of smaller sample sizes were generated from the fly data used in Fig 2
690 by resampling subjects (without replacement within a group) from the full dataset. (a) Each
691 CBAS run on the smaller sample size identified some number of significant sequences. Those
692 sequences were compared to the sequences identified in the CBAS on the full dataset, and the
693 graph shows the ratios of the number of sequences identified by the smaller sample size that
694 were also in the full dataset to the total number of sequences identified in the full dataset. As a
695 comparison, the ratio of the number of flies in the smaller CBAS to the total number of flies in
696 the dataset is plotted in the maroon horizontal lines. (b) The number of significant sequences
697 identified in the smaller sample size that were not also identified in the full dataset is plotted
698 over the total number of significant sequences identified in the smaller sample size. (c) In
699 creating the smaller samples sizes, 20 paired sets of animals were generated that had no
700 overlapping individuals. The number of the same significant sequences that were identified
701 with these nonoverlapping sets of flies is plotted over the average number of significant
702 sequences in the pair. (d) For each repeat of each sample size, all of the sequences were

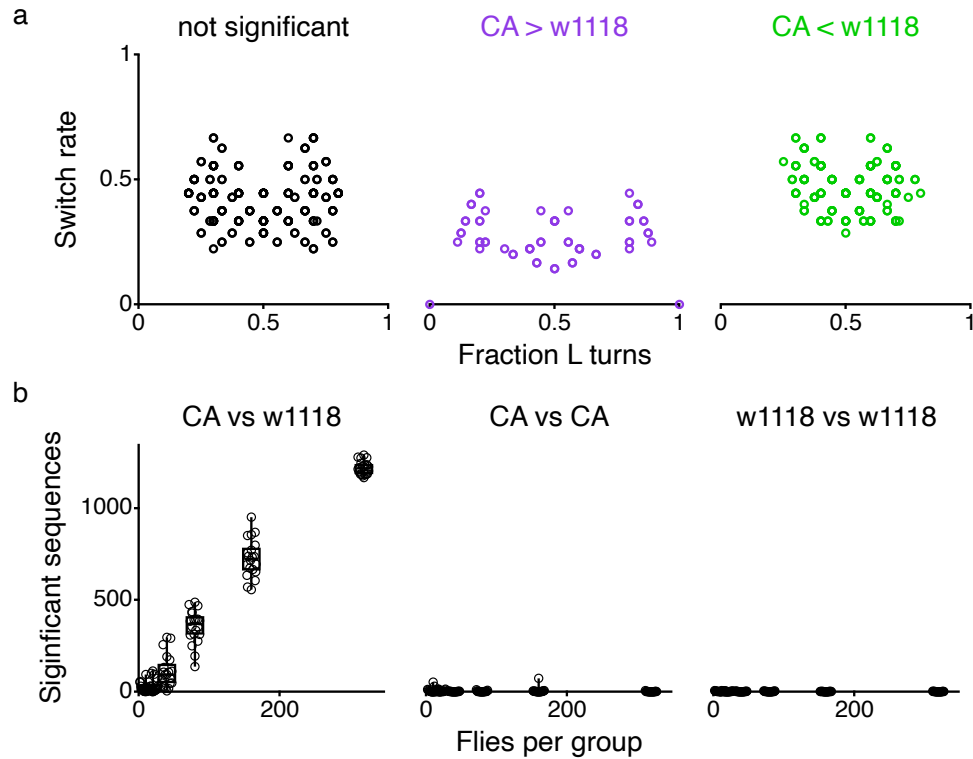
703 categorized based on the maximum number of turns in the same direction and the fraction of
704 significant sequences within those categories are plotted (as in Fig 2e). The bottom row of this
705 plot is from the full dataset and is identical to Fig 2e. (e) The Euclidean distance between each
706 row from panel **d** and the full dataset row is plotted. Colors correspond to the sample sizes as
707 shown in panel **d**. For panels **a** – **c** data points are overlaid by box and whisker plots. The
708 center line of the box displays the median of the data, the top and bottom lines of the box
709 show the 25th and 75th quartiles, respectively, and the end of the whiskers show the full range
710 of the data.



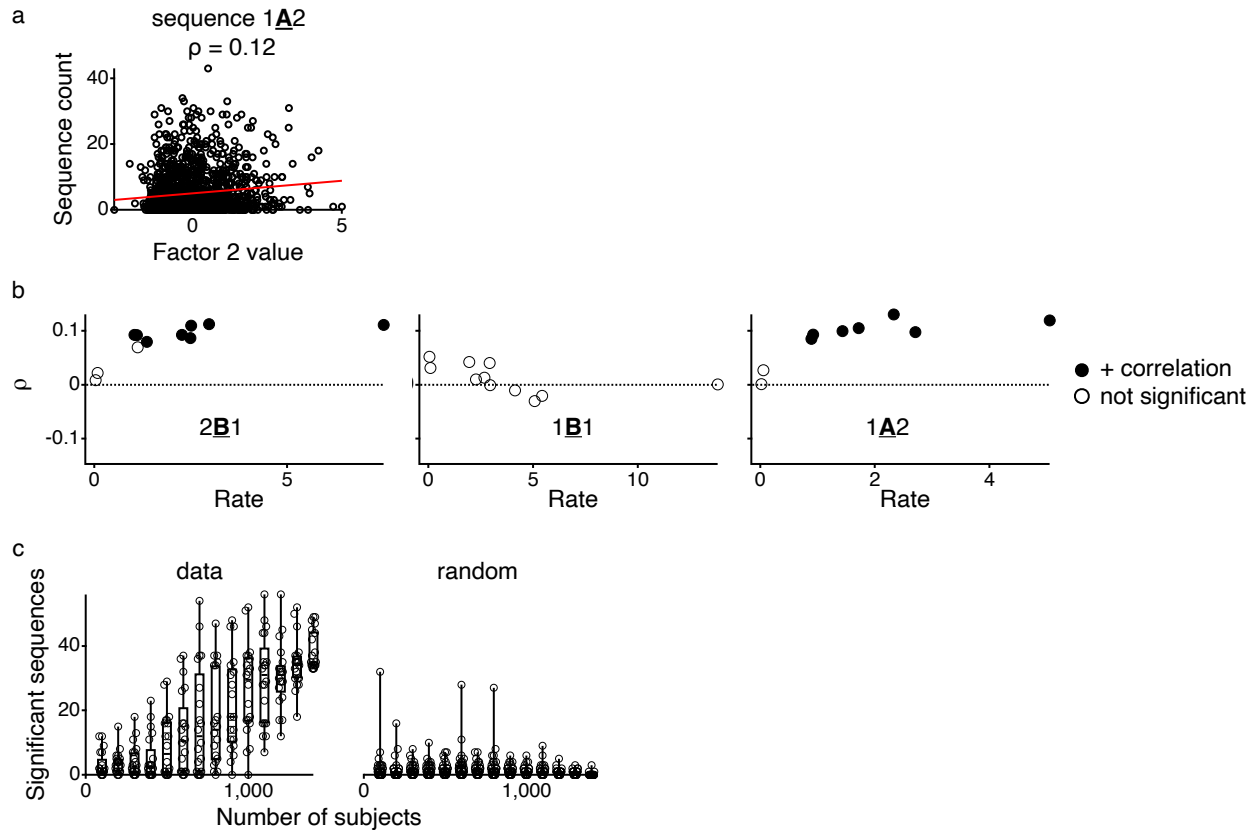
711 **Figure 4. CBAS identifies unexpected sequences in human dataset.** CBAS was applied to
712 humans performing the two-step task (Fig 1b) (a) This table demarcates the decisions made to
713 perform the CBAS on this dataset. There was a total of 1,413 human subjects in the dataset. (b)
714 The occurrence of a single sequence of choices ('1A2', meaning choosing object 1, then getting
715 rewarded for a choice in set A, then choosing object 2) across all subjects in the dataset.
716 Subjects are ordered based on their factor 2 ("intrusive thoughts and compulsive behaviors")
717 score. This sequence shows a correlation (ρ) of 0.12 with factor 2 score from the questionnaire
718 factor analysis (see S2a). (c) CBAS Manhattan plot displays the p-value for each sequence. Each
719 sequence has two p-values on this plot, one for positive correlation and one for negative
720 correlation. The sequences are ordered based on the number of choices in the sequence and
721 are displayed on a log scale to make all sequence lengths visible. Within a given sequence
722 length, the sequences are ordered based on frequency of occurrence in the entire dataset. The
723 horizontal dotted line indicates the significance threshold of 5% control of the median false
724 discovery proportion. A total of 10,000 resamplings were used to calculate the p-values making
725 the maximum value on the plot $4.00004 \left(-\log\left(\frac{1}{10,001}\right)\right)$. (d) Every sequence that was
726 significantly positively (left) or negatively (right) correlated with factor 2 score. The sequences
727 on the left are aligned to B1 or A2, and sequences on the right are aligned to A1. (e) For the
728 categories listed (2B1, 1B1, 1A2, 2A2), the number of significant sequences that contain the
729 category is plotted over total number of sequences that exist in the dataset that contain the
730 category. For 2A2 no subject in the dataset performed that sequence. (f) Power estimate for
731 different sized groups of subjects (see methods). Horizontal dotted line shows a value of 80%.



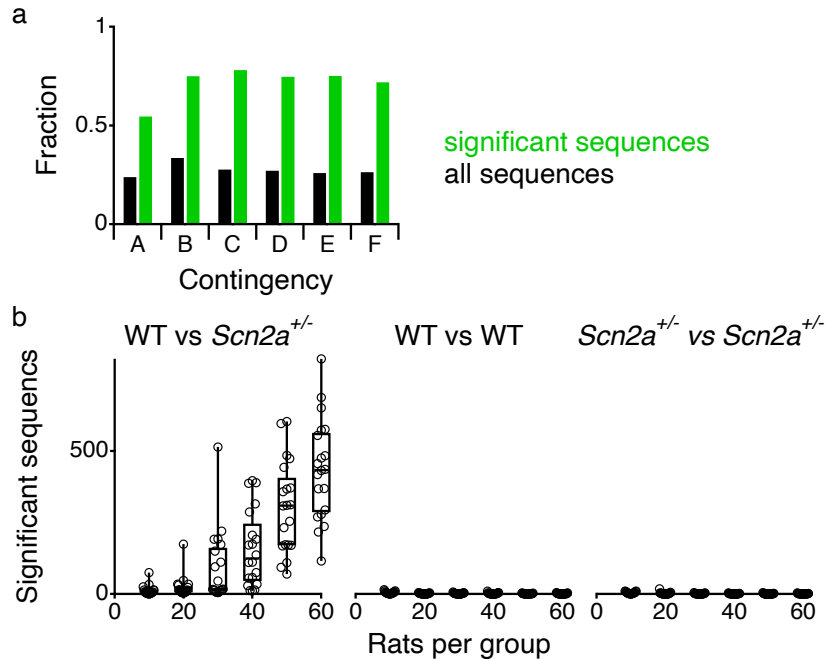
750 Every significant sequence that exclusively contains arms 1, 2, 3 (irrespective of reward) during
751 contingency B (left) or every significant sequence that exclusively contains arms 2, 4, 6
752 (irrespective of reward) during contingency E (right). Sequences that occur significantly more in
753 *Scn2a^{+/-}* are shown in red, and those that occur significantly more in WT are shown in black. (e)
754 The sets of 3 arms which show a significantly greater fraction of significant sequences than
755 would be expected by chance (see methods). The categories above each group of bar plots
756 indicate the structure of the set of 3 arms. “Current” indicates the 3 arms from the current
757 contingency irrespective of reward, “current rewarded” indicates the 3 arms from the current
758 contingency all of which are rewarded, “1 back” indicates the 3 arms from the prior
759 contingency, “2 back” indicate the 3 arms from 2 contingencies prior, and “other” are
760 categories that show significant fractions that do not fit the other categories. Of note, the
761 sequences in the “current rewarded” category are a subset of the sequences in the “current”
762 category, as can be seen in panel **d**, left. (f) Power estimate for different sized groups of rats in
763 each genotype (see methods). Horizontal dotted line shows a value of 80%.



764 **Supplementary Figure 1.** (a) All sequences with a maximum length of 4 turns in the same
765 direction in a row (Fig 2e) were evaluated for the fraction of L turns in the sequence
766 (encompassing $L = 4 - 8$) and the rate of switching (number of changes in turn direction divided
767 by the total number of turns in the sequence). These two metrics show a separation in the
768 sequences that were significantly greater in the CA strain compared to those that were
769 significantly greater in the w1118 strain: the CA strain had sequences with more extreme
770 fraction of left turns and lower switch rate, consistent with more turns in the same direction.
771 (b) Left: the number of significant sequences when randomly resampling the populations
772 without replacement and calculating CBAS on the smaller sample sizes. Middle: the number of
773 significant sequences when comparing smaller samples sizes of the CA strain to itself, with
774 nonoverlapping individuals in each group. Right: the number of significant sequences when
775 comparing smaller samples sizes of the w1118 strain to itself, with nonoverlapping individuals
776 in each group. Data points are overlaid by box and whisker plots. The center line of the box
777 displays the median of the data, the top and bottom lines of the box show the 25th and 75th
778 quartiles, respectively, and the end of the whiskers show the full range of the data.



779 **Supplementary Figure 2.** (a) Correlation between the sequences count for each human subject
780 and their factor 2 score for sequence: 1A2, which means choosing object 1, then making a
781 choice in set A and getting rewarded, and then choosing object 2. Factor 2 reflects the intrusive
782 thoughts and obsessive behavior loading on the psychiatric symptom questionnaires. Red line
783 shows the linear fit to the data. (b) The correlation plotted as a function of the average
784 sequence rate for all sequences containing the unit 2B1 (left), 1B1 (middle), or 1A2 (right).
785 Filled in circles indicate those sequences that show a significant positive correlation. (c) Left:
786 the number of significant sequences when randomly resampling the populations without
787 replacement and calculating CBAS on the smaller sample sizes. Right: the number of significant
788 sequence when randomly resampling the populations without replacement and comparing it to
789 randomly generated factor 2 scores drawn from a distribution imputed from the original factor
790 2 scores and calculating CBAS on the smaller sample sizes. Data points are overlaid by box and
791 whisker plots. The center line of the box displays the median of the data, the top and bottom
792 lines of the box show the 25th and 75th quartiles, respectively, and the end of the whiskers show
793 the full range of the data.



794 **Supplementary Figure 3.** (a) For the sequences that exclusively contain all sets of 3 arms (Fig
795 5d&e) the fraction of the total number of sequences is plotted in black for each contingency, and the fraction of significant sequences that are a part of the sequences that exclusively
796 contain all sets of 3 arms compared to the total number of significant sequences in each
797 contingency is plotted in green. (b) Left: the number of significant sequences when randomly
798 resampling the populations without replacement and calculating CBAS on the smaller sample
799 sizes. Middle: the number of significant sequence when comparing smaller samples sizes of the
800 WT genotype to itself, with nonoverlapping individuals in each group. Right: the number of
801 significant sequence when comparing smaller samples sizes of the *Scn2a*^{+/-} genotype to itself,
802 with nonoverlapping individuals in each group. Data points are overlaid by box and whisker
803 plots. The center line of the box displays the median of the data, the top and bottom lines of
804 the box show the 25th and 75th quartiles, respectively, and the end of the whiskers show the full
805 range of the data.
806