

# Extended Bayesian inference incorporating symmetry bias

Shuji Shinohara<sup>1\*</sup>, Nobuhito Manome<sup>1,2</sup>, Kouta Suzuki<sup>1,2</sup>, Ung-il Chung<sup>1</sup>, Tatsuji Takahashi<sup>1,3</sup>, Pegio-Yukio Gunji<sup>4</sup>, Yoshihiro Nakajima<sup>5</sup>, Shunji Mitsuyoshi<sup>1</sup>

<sup>1</sup>Department of Bioengineering, Graduate School of Engineering, The University of Tokyo, Tokyo, Japan

<sup>2</sup> Department of Research and Development, SoftBank Robotics Group Corp, Tokyo, Japan

<sup>3</sup>School of Science and Engineering, Tokyo Denki University, Saitama, Japan

<sup>4</sup>Department of Intermedia Art and Science, School of Fundamental Science and Technology, Waseda University, Tokyo, Japan

<sup>5</sup>Graduate School of Economics, Osaka City University, Osaka, Japan

\* Corresponding author

E-mail: [shinohara@bioeng.t.u-tokyo.ac.jp](mailto:shinohara@bioeng.t.u-tokyo.ac.jp) (SS)

## 1 **Abstract**

2           In this study, we start by proposing a causal induction model that incorporates symmetry bias. This  
3 model is important in two aspects. First, it can reproduce causal induction of human judgment with higher  
4 accuracy than conventional models. Second, it allows us to estimate the level of symmetry bias of subjects  
5 from experimental data. We further propose an inference method that incorporates the aforementioned causal  
6 induction model into Bayesian inference. In this method, the component of Bayesian inference, which  
7 updates the degree of confidence for each hypothesis, and the component of inverse Bayesian inference that  
8 modifies the model of the hypothesis coexist. Our study demonstrates that inverse Bayesian inference enables  
9 us to deal flexibly with unstable situations where the object of inference changes from time to time.

10

## 11 **Author summary**

12           We acquire knowledge through learning and make various inferences based on such knowledge and  
13 observational data (evidence). If the evidence is insufficient, then the certainty of the conclusion will decline.  
14 Moreover, even if the evidence is sufficient, the conclusion may be wrong if the knowledge is incomplete in  
15 the first place. In order to model such inference based on incomplete knowledge, we proposed an inference  
16 system that performs learning and inference simultaneously and seamlessly. Prepare two coins A and B with  
17 different probabilities of landing heads, and repeat the coin toss using either of them. However, the coin that  
18 is being tossed is also replaced repeatedly. The system observes only the result of coin toss each time, and  
19 estimates the probability of landing heads of coin tossed at the moment. In this task, it is necessary not only  
20 to estimate the probabilities of the landing heads of coin A and B, but also to estimate which coin is being  
21 used at the moment. In this paper, we show that the proposed system handles such tasks very efficiently by  
22 simultaneously performing inference and learning.

23

## 24 **Introduction**

25 As a cognitive bias observed in humans, the disposition to infer from ‘if P then Q’ to ‘if Q then P’  
26 or to ‘if not P, then not Q’ is well documented [1, 2, 3, 4, 5, 6, 7, 8]. The former is termed symmetry bias [9]  
27 and the latter is termed mutual exclusivity bias [4].

28 Consider a simple example. We tend to infer from ‘if you clean the room, then I will take you out’  
29 to ‘I will take you out if and only if you clean the room’ or ‘if you don’t clean the room, then I will not take  
30 you out’. Although these inferences are invalid according to classical logic, various people are inclined to  
31 make them regardless of age.

32 In contrast, among non-human animals, the symmetry bias has been reported in behaviour of only  
33 some California sea lions [10] and chimpanzees [11]. Although the symmetry bias produces wrong inferences  
34 from the classical logic point of view, humans do show some positive features apparently stemming from the  
35 symmetry bias. For instance, once you are able to respond to the question ‘What is this?’ with ‘This is an  
36 apple’ through learning, you will also be able to identify the correct object when asked ‘Which one is an  
37 apple?’. In other words, we automatically infer from ‘This is an apple’ to ‘An apple is this’ without any  
38 instruction. The symmetry bias has been studied in relation to stimulus equivalence in the field of comparative  
39 psychology [1, 2]. On the other hand, the mutual exclusivity bias has been studied primarily in the field of  
40 developmental psychology in the context of young children’s language acquisition [4]. Thus, although the  
41 symmetry and mutual exclusivity biases have been studied in different fields of psychology, since the  
42 contrapositive of ‘if Q then P’ is ‘if not P then not Q’, and since they are equivalent according to classical  
43 logic, the same implications can be associated with both biases. An example of such a shared implication is  
44 that both biases may be caused by the same neuroscientific factor.

45 Concurrently, in the field of cognitive psychology, experiments on causal induction were carried  
46 out, seeking to identify how humans evaluate the strength of causal relations between two events. In a regular  
47 conditional statement of the form ‘if  $p$  then  $q$ ’, the degree of confidence is considered to be proportional to  
48 the conditional probability  $P(q | p)$  which is the probability of occurrence of  $q$  following the occurrence of  
49  $p$  [12]. Further, in the case of causal relation, it has been experimentally demonstrated that humans have a

50 strong sense of causal relation when  $P(p|q)$  is high, as well as when  $P(q|p)$  is high, where  $P(p|q)$  is  
51 a conditional probability of the antecedent occurrence of  $p$ , given the occurrence of  $q$  [13].

52 Consider a simple causal induction model that infers the strength of a causal relation from the cause  
53 candidate of event  $C$  to the effect event  $E$  from four pieces of co-occurring information concerning  $C$  and  $E$ :  
54 the joint presence of  $C$  and  $E$ , absence of  $E$  given  $C$ , presence of the  $E$  given no  $C$ , and the joint absence of  
55  $C$  and  $E$ . The most representative model of causal induction is the  $\Delta P$  model [14]. It takes the difference  
56 between the conditional probability  $P(E|C)$  of occurrence of  $E$  given the occurrence of  $C$  and the  
57 conditional probability  $P(E|\neg C)$  of occurrence of  $E$  given non-occurrence of  $C$  (denoted by  $\neg C$ ) as an  
58 index for causal strength, that is,  $\Delta P = P(E|C) - P(E|\neg C)$ .

59 Hattori and Oaksford proposed the dual-factor heuristic (*DFH*) model [13]. This model is based on  
60 the geometric mean of  $P(E|C)$ , which stands for the predictability of the effect from the cause, and its  
61 inverse  $P(C|E)$ , that is,  $DFH = \sqrt{P(E|C)P(C|E)}$ .

62 Both  $\Delta P$  and *DFH* models contain  $P(E|C)$ . In other words, given the occurrence of  $C$ , if the  
63 probability of occurrence of  $E$  following  $C$  is high, the chance of  $C$  to be the cause of  $E$  increases. Intuitively  
64 speaking, the strength of the causal relation does not seem to be solely determined by  $P(E|C)$ . The second  
65 item in the  $\Delta P$  model,  $-P(E|\neg C)$ , shows that even if the probability of occurrence of  $E$  is high given the  
66 occurrence of  $C$ , if the probability of occurrence of  $E$  is still high in the absence of the occurrence of  $C$ , that  
67 is, if the probability of occurrence of  $E$  is high irrespectively of the occurrence of  $C$ , then the chance of  $C$   
68 being the cause of  $E$  decreases.

69 Whereas for the *DFH* model, if the probability  $P(C|E)$ , which is the probability of the antecedent  
70 occurrence of  $C$  given the occurrence of  $E$ , is high, the chance of  $C$  being the cause of  $E$  increases. This can  
71 be understood as a probabilistic expression of the belief that where there is no cause, there is no effect.

72 We can also consider the  $\Delta P$  model and the *DFH* model in terms of biases. For the sake of  
73 simplicity,  $\Delta P$  and *DFH* are expressed as  $\Delta P(E|C)$  and  $DFH(E|C)$  respectively. Here, if we assign  $\neg C$

74 and  $\neg E$  to  $C$  and  $E$  respectively in  $\Delta P(E|C) = P(E|C) - P(E|\neg C)$ , we can obtain mutual exclusivity as  
75  $\Delta P(\neg E|\neg C) = P(\neg E|\neg C) - P(\neg E|\neg(\neg C)) = (1 - P(E|\neg C)) - (1 - P(E|C)) = P(E|C) - P(E|\neg C)$   
76  $= \Delta P(E|C)$ . If, on the other hand,  $C$  in  $DFH(E|C) = \sqrt{P(E|C)P(C|E)}$  is replaced by  $E$ , we can get  $DFH$   
77  $(C|E) = \sqrt{P(C|E)P(E|C)} = \sqrt{P(E|C)P(C|E)} = DFH(E|C)$  and the symmetry obtains.

78 Aside from the  $\Delta P$  and  $DFH$  models, Takahashi and colleagues [15] proposed the *pARIs*  
79 (proportion of assumed-to-be rare instances) as yet another model that has an unusually high affinity with the  
80 human causal induction judgment,  $pARIs(E|C) = P(C,E)/(P(C,E) + P(C,\neg E) + P(\neg C,E))$ , where  $P(x,y)$   
81 represents joint probability of  $x$  and  $y$ . If  $C$  in  $pARIs(E|C)$  is replaced by  $E$ , we get  $pARIs(C|E) =$   
82  $P(E,C)/(P(E,C) + P(E,\neg C) + P(\neg E,C)) = pARIs(E|C)$  and the symmetry obtains.

83 To view the relation between two events as a causal relation can therefore be understood as having  
84 both symmetry and mutual exclusivity biases.

85 Bayesian inference is based on the notion of conditional probability. Bayesian inference speculates  
86 the hidden cause behind an observation results from retrospectively applying statistical inferences. The  
87 relation between Bayesian inference and brain function has been attracting attention in recent years in the  
88 field of neuroscience [16, 17].

89 In Bayesian inference, the degree of confidence in a hypothesis is updated based on a model of  
90 predefined hypotheses and current observational data. In other words, Bayesian inference is a process of  
91 narrowing down hypotheses to one which best explains observational data. Changing the model of each  
92 hypothesis or adding new ones in the course of performing Bayesian inference is not allowed. In addition,  
93 Bayesian inference itself does not deal with alterations in the inference target during the inference or with its  
94 replacement. Therefore, such inference substantially needs to assume the identity of the target.

95 Note, however, that requirements of the invariability of the hypothetical model and the identity of  
96 the inference target stem from the theoretical framework, and they are not always met in actuality. For  
97 instance, if the object is unknown, it would be impossible to infer what it is without adding a new hypothetical

98 model. Moreover, it is likely that, under unsteady circumstances, the inference target undergoes alteration  
99 from time to time or is replaced by some other object.

100 In order to predetermine whether the object is replaced by another, one must first infer its identity.  
101 A correct inference depends on as much observational data as possible. However, in order to properly use  
102 accumulated observational data, it must be ensured that these data derive from the same object. In other  
103 words, to determine whether the object has been replaced or not, the object must be hypothesized not to have  
104 been replaced in the first place. In this sort of situation, it is necessary to infer what the object is while at the  
105 same time evaluating the legitimacy of the inference itself. How, then, could we model the inference under  
106 the situation described above?

107 Arecchi [18] proposed the concept of the inverse Bayesian inference where the hypothetical model,  
108 which is fixed in the traditional Bayesian inference, is modified according to circumstances. Gunji et al. [19,  
109 20] and Horry et al. [21] formulated the inverse Bayesian inference and demonstrated that animal herding  
110 and human decision-making can be satisfactorily modelled by combining Bayesian inference and inverse  
111 Bayesian inference. This framework can be said to seamlessly perform Bayesian inference, by picking up the  
112 optimal hypothesis from the predefined set of hypotheses, and simultaneously apply inverse Bayesian  
113 inference (learning), which creates a new hypothesis according to observational data. Although the inverse  
114 Bayesian inference was formulated by Gunji and others [19, 20], it is not necessarily linked with causal  
115 inference and symmetry bias.

116 We propose a causal induction model that primarily incorporates symmetry bias. First, we propose  
117 an extended model of degree of confidence that extends conditional probability by parametrising the mixed  
118 rate of  $P(q | p)$  and  $P(p | q)$ , i.e., the strength of symmetry bias. Second, we propose a realistic human  
119 inference model that incorporates the extended model into Bayesian inference, and we show that it  
120 necessarily involves inverse Bayesian inference. Specifically, we propose a framework of extended Bayesian  
121 inference which allows seamless and simultaneous learning and inference by replacing the conditional  
122 probability schema in Bayesian inference with the extended model of degree of confidence. Third, we explain

123 a conducted simulation, derived from the problem of inference, of the probability of getting heads in the  
 124 course of repetitive coin toss and how it helps verify the legitimacy of the Extended Bayesian Inference.

125

## 126 **Results**

### 127 **Proposal of extended confidence model**

128 We seek to establish an extended model of degree of confidence as the generalised weighted average  
 129 of  $P(q|p)$  and its inverse  $P(p|q)$  using parameters  $\alpha$  and  $m$ .

$$C(q|p) = [(1 - \alpha)P(q|p)^m + \alpha P(p|q)^m]^{1/m} \quad (1)$$

130

131 Hereinafter,  $C(q|p)$  will be termed the Extended Confidence Model. Here  $\alpha$  takes values in the range  
 132  $0.0 \leq \alpha \leq 1.0$  and denotes a weighted value of  $P(q|p)$  and  $P(p|q)$ , and  $m$  takes values in the range  
 133  $-\infty \leq m \leq \infty$  and denotes the manner of taking the mean. For example, suppose  $\alpha = 0.5$  and  $m = 1.0$ ,  
 134 then  $C(q|p) = 0.5P(q|p) + 0.5P(p|q)$  representing the arithmetic mean. Supposing  $m = 0.0$ , the  
 135 formula (1) is undefinable. If, however, we represent the mean value in the limit of  $m \rightarrow 0.0$ , we get  $C(q|p)$   
 136  $= P(q|p)^{1-\alpha}P(p|q)^\alpha$ , where if  $\alpha = 0.5$ , the geometric mean  $C(q|p) = \sqrt{P(q|p)P(p|q)}$  it coincides  
 137 with the *DFH* model. If  $\alpha = 0.5, m = -1.0$ , the formula represents the harmonic mean of  $P(q|p)$  and  $P$   
 138  $(p|q)$  and we get

$$C(q|p) = \frac{2P(q|p)P(p|q)}{P(q|p) + P(p|q)} = \frac{2P(p,q)}{2P(p,q) + P(p,-q) + P(-p,q)} = \frac{2pARIs}{1 + pARIs} \quad (2)$$

139 and  $C(q|p)$  can be expressed as a harmonic mean of 1 and  $pARIs$ . In other words,  $pARIs$  and  $C(q|p)$   
 140 are related by monotonically increasing functions that are in one-to-one correspondence. According to this,  
 141  $pARIs$  can be seen as a disguised form of  $C(q|p)$ . Here, the parameter  $\alpha$  can be regarded as a parameter

142 that controls that strength of the symmetry bias. When  $\alpha = 0$ ,  $C(q | p) = P(q | p)$  obtains irrespectively of  
 143 the value of  $m$ , and  $C(q | p)$  expresses a normal conditional probability without the symmetry bias.

144 Thus, the proposed model can be described as an extended model which accommodates the normal  
 145 conditional probability  $P$ ,  $DFH$  and  $pARIs$  as inner special cases.

146

### 147 **Evaluation of descriptive validity of extended confidence model**

148 In order to evaluate the descriptive validity of various models including the  $DFH$  model, Hattori and  
 149 Oaksford [13] performed meta-analysis using data from eight types of causal induction experiments. To test  
 150 the descriptive performance of the extended confidence model, we also performed the meta-analysis using  
 151 the same datasets as Hattori and Oaksford [13].

152 Generally, in a simple causal induction experiment, participants are given four types of co-  
 153 occurrence information concerning the cause  $C$  and the effect  $E$  (Table 1). Then, they are asked to assess  
 154 subjectively the strength of the causal relation between  $C$  and  $E$  using a number from 0 to 100. To measure  
 155 each model's fit to the data, we calculated the determination coefficient  $R^2$  from the pair of participants'  
 156 mean ratings of causal strength and the estimated value of each model computed from the same co-occurrence  
 157 information given to the participants in the experiments.

158 **Table 1. The  $2 \times 2$  contingency table for elemental causal induction.**

	effect ( $E$ )	no effect ( $\neg E$ )	marginal frequency
cause ( $C$ )	$N(C, E)$	$N(C, \neg E)$	$N(C)$
no cause ( $\neg C$ )	$N(\neg C, E)$	$N(\neg C, \neg E)$	$N(\neg C)$
marginal frequency	$N(E)$	$N(\neg E)$	

159  $N(x)$  denotes the frequency of occurrence of  $x$ , and  $N(x, y)$  denotes the frequency of co-occurrence  
 160 of  $x$  and  $y$ .

161

162 The experiment data used in the meta-analysis consists of the experiment I from [22], experiments  
 163 I and III from [23], experiments I and II from [13], experiments I and III from [24], and experiment II and  
 164 VI from [25]. The experiment data above will be abbreviated as AS95, BCC03.1, BCC03.3, HO07.1,



165 HO07.2, LS00, W03.2, and W03.6 respectively. See the methods section below for the content of each  
 166 experiment.

167 Values in parameters  $\alpha$  and  $m$  in formula (1) were shifted each with 0.05 increment in the interval  
 168  $[0.0,1.0]$  and  $[-2.0,2.0]$  and the determination coefficient  $R^2$  between assessment by participants and the  
 169 estimated value by the proposed model was calculated for each pair of parameters. Table 2 shows the pair of  
 170 parameters  $\alpha$  and  $m$  at which  $R^2$  becomes maximum in each experiment. As seen in Table 2, the  
 171 determination coefficients were greater than 0.9 for all experiments.  $\alpha$  was around 0.5 (0.25-0.6) and did  
 172 not reach 0.0, which stands for the normal conditional probability  $P$ . This suggests that symmetry bias was  
 173 deeply involved in causal induction. Moreover,  $m$  took a negative value in all experiments. This suggest  
 174 people show strong awareness of causal relations only if both  $P(E|C)$  and its inverse  $P(C|E)$  are large.

175 **Table 2. Performance evaluation of the extended confidence model based on the meta-analytic data**  
 176 **from Hattori and Oaksford (13).**

	AS95	BCC03.1	BCC03.3	HO07.1	HO07.2	LS00	W03.2	W03.6
$R^2$	0.97	0.98	0.99	0.99	0.99	0.91	0.97	1.00
$\alpha$	0.3	0.55	0.35	0.6	0.3	0.55	0.35	0.25
$m$	-0.65	-0.45	-0.25	-2.0	-2.0	-1.4	-1.6	-2.0
$N$	80	13	6	12	9	11	8	4

177  $R^2$  is a determination coefficient,  $\alpha$  and  $m$  are model parameters, and  $N$  is a number of combinations of  
 178 stimuli.

179  
 180 In these analyses, the optimal parameter value was calculated for each experiment. In what follows,  
 181 all experiments will be analysed comprehensively using common parameters. The determination coefficient  
 182 when parameters have fixed values will be calculated for each experiment along with their mean. The mean  
 183 value to be calculated is the weighted average using Fisher's Z conversion. This procedure was repeated by  
 184 changing the parameter values with 0.05 increment. Fig 1 shows the mean coefficient of determination for  
 185 each pair of parameter values.

186  
 187 **Fig 1. Mean coefficient of determination for each pair of parameter values.** The mean value for the  
 188 determination coefficient peaked when  $(\alpha,m)=(0.35,-1.15)$  and the value was  $\overline{R^2} = 0.93$ .

189

190 Hattori and Oaksford [13] demonstrated that their *DFH* model (without parameters) showed the best  
 191 fit compared to 33 models without parameters and seven models with parameters. We used eight types of  
 192 experimental data to compare the performance of our proposed model with that of other models including  
 193 *DFH*, *pARIs*,  $\Delta P$  model, and conditional probability *P*. The results are shown in Table 3.

194 **Table 3. A replication trial of the meta-analysis of Hattori and Oaksford (13).**

	AS95	BCC03.1	BCC03.3	HO07.1	HO07.2	LS00	W03.2	W03.6	$\overline{R^2}$
<i>C</i>	<b>0.92</b>	0.90	<b>0.97</b>	<b>0.97</b>	0.96	0.72	<b>0.91</b>	<b>0.93</b>	<b>0.93</b>
<i>P</i>	0.82	0.68	0.82	0.80	0.0	0.17	0.20	0.54	0.75
$\Delta P$	0.78	0.84	0.70	0.50	0.0	0.77	0	NA	0.73
<i>DFH</i>	0.91	0.95	0.91	0.93	0.96	<b>0.80</b>	0.69	0.80	0.91
<i>pARIs</i>	0.89	<b>0.96</b>	0.94	0.93	<b>0.99</b>	0.80	0.78	0.88	0.90

195 Values in the table represent the determination coefficient between the estimated value of each model and  
 196 the assessment value of the subject of each experiment. *C* represents the extended confidence model when  
 197  $\alpha=0.35$ ,  $m=-1.15$ . *P* represents the conditional probability. The rightmost column represents the weighted  
 198 average of the determination coefficient for each experiment. Bold-faced numbers represent the value of the  
 199 model that marked the greatest determination coefficient for each experiment.

200

201 While the mean determination coefficient exceeded 0.9 for our proposed model, *DFH* and *pARIs*,  
 202 it did not for the  $\Delta P$  model or the conditional probability *P*. Particularly, the proposed model recorded the  
 203 highest determination coefficient in five out of eight experiments, as well as in the mean of all experiments.  
 204 Thus, it was shown that introducing symmetry bias into the conditional probability could significantly  
 205 improve the determination coefficient with human assessment.

206

### 207 **Proposal of extended Bayesian inference**

208 In this section, we propose the extended Bayesian inference where the conditional probability in  
 209 Bayesian inference is replaced by the extended confidence model.

210 First, we describe Bayesian inference. This study deals with the problem of inferring a generative  
 211 model (probability distribution) from observational data. To this end, in what follows, the hypothesis *h* and  
 212 data *d* will be used on behalf of *p* and *q*. Moreover, discrete models will be considered.

213 Bayesian inference first defines several hypotheses  $h_i$  and provides a model for each hypothesis  
 214 (probability distribution of data) in the form of conditional probability  $P(d|h_i)$ . When data are fixed and  
 215 regarded as a function of a hypothesis, this conditional probability is termed likelihood. The confidence  $P$   
 216 ( $h_i$ ) for each hypothesis is given as a prior probability.

217 We can take  $P(d|h_i)$  and  $P(h_i)$  as initial values and calculate the posterior probability  $P(h_i|\mathbf{d})$   
 218 when observing data  $\mathbf{d}$  using Bayes' theorem as follows.

$$P(h_i|\mathbf{d}) = P(h_i) \frac{P(\mathbf{d}|h_i)}{P(\mathbf{d})} \quad (3)$$

219 Hereinafter, data observed at a point in time are represented by the bold  $\mathbf{d}$ . Afterwards, we can replace the  
 220 posterior probability with the prior probability using Bayesian updating.

$$P(h_i) \leftarrow P(h_i|\mathbf{d}) \quad (4)$$

221 By combining formulas (3) and (4), we get

$$P(h_i) \leftarrow P(h_i) \frac{P(\mathbf{d}|h_i)}{P(\mathbf{d})} \quad (5)$$

222 Whenever new data are observed,  $P(h_i)$  in the formula (5), i.e., confidence in each hypothesis, is  
 223 updated and the inference continues. The inference distribution during this procedure can be expressed as

$$P(\mathbf{d}) = \sum_i P(h_i) P(\mathbf{d}|h_i) \quad (6)$$

224 Note that in Bayesian inference, while the probability  $P(h_i)$  of each hypothesis changes over time,  
 225 the probability of the model of each hypothesis  $P(\mathbf{d}|h_i)$  does not. Fig 2 (a) shows an overview of the  
 226 processing flow of Bayesian inference.

227 The extended Bayesian inference is an inference that has  $\alpha$  and  $m$  as parameters and  
 228 accommodates normal Bayesian inference as its special case when  $\alpha = 0$ . Specifically, it is constructed by  
 229 the following two update formulas.

$$C(h_i) \leftarrow \frac{C(h_i)C(\mathbf{d}|h_i)}{[(1-\alpha)(C(\mathbf{d}))^{-m} + \alpha(C(h_i))^{-m}]^{1/m}} \quad (7)$$

$$C(\mathbf{d}|\mathbf{h}) \leftarrow \frac{C(\mathbf{h})C(\mathbf{d}|\mathbf{h})}{[(1-\alpha)(C(\mathbf{h}))^{-m} + \alpha(C(\mathbf{d}))^{-m}]^{-1/m}} \quad (8)$$

230

231 In formula (8), the bold-faced  $\mathbf{h}$  represents the hypothesis that has the highest confidence. See the methods  
 232 section below for a detailed derivation of the update formulas. Here, we can see that, supposing  $\alpha = 0$  in  
 233 formula (7), the right side shows the same form as that of Bayesian inference seen in formula (5).

$$C(h_i) \leftarrow \frac{C(h_i)C(\mathbf{d}|h_i)}{C(\mathbf{d})} \quad (9)$$

234

235 Now, in the case of Bayesian inference, the model  $P(\mathbf{d}|h_i)$  was invariable. We can ask if this is the  
 236 same for extended Bayesian inference. Here, if we look closely to the right side of formula (8), supposing  
 $\alpha = 0$ , formula (8) becomes a tautology as shown below.

$$C(\mathbf{d}|\mathbf{h}) \leftarrow \frac{C(\mathbf{h})C(\mathbf{d}|\mathbf{h})}{C(\mathbf{h})} = C(\mathbf{d}|\mathbf{h}) \quad (10)$$

237

238 In other words, if  $\alpha = 0$ , then formula (8) substantially disappears. Conversely, if  $\alpha > 0$ ,  $C(\mathbf{d}|\mathbf{h})$   
 239 is subject to the denominator  $C(\mathbf{d})$  in the right side, that is, the estimated value of the data. Following Gunji  
 240 et al. [19], the process shown in formula (8) is termed Inverse Bayesian Inference.

241 In what follows we show the processing flow of Extended Bayesian inference. First, we take  $P(\mathbf{d}|h_i$   
 242 ) and  $P(h_i)$  as initial values and substitute them with  $C$ .

$$\begin{aligned} C(\mathbf{d}|h_i) &= P(\mathbf{d}|h_i) \\ C(h_i) &= P(h_i) \end{aligned} \quad (11)$$

243

244 Second, we calculate the degree of confidence for each hypothesis  $C(h_i)$  and the model  $C(\mathbf{d}|\mathbf{h})$   
 245 using the formula (7) and (8) whenever  $\mathbf{d}$  is observed. Following the application of formulas (7) and (8),  
 246 we can normalise  $C(h_i)$  and  $C(\mathbf{d}|\mathbf{h})$ .

$$C(h_i) \leftarrow \frac{C(h_i)}{\sum_k C(h_k)} \quad (12)$$

$$C(d|h_i) \leftarrow \frac{C(d|h_i)}{\sum_j C(d_j|h_i)}$$

247 Finally, we can calculate the estimated distribution values as with Bayesian inference.

$$C(d) = \sum_i C(h_i)C(d|h_i) \quad (13)$$

248 Fig 2(b) shows an overview of the processing flow of extended Bayesian inference.

249

250 **Fig 2. A flow chart comparing Bayesian Inference to Extended Bayesian inference.** (a) An overview of  
 251 the processing flow of Bayesian inference; and (b) An overview of the processing flow of extended Bayesian  
 252 inference. In (b), for simplicity, the portion that belongs to the normalisation process is omitted. In (b), if we  
 253 suppose  $\alpha=0$ , the portion of inverse Bayesian inference disappears, corresponding to the Bayesian inference  
 254 in (a).

255

## 256 **Performance evaluation of extended Bayesian inference using simulation**

257 To observe the behaviour of extended Bayesian inference, a simulation was performed. Specifically,  
 258 a coin was tossed repeatedly, using a simulator, to observe the results and estimate the probability of getting  
 259 heads using extended Bayesian inference. The probability of landing heads at the  $t^{th}$  trial was designated  $p^t$   
 260 , and the probability of it landing tails was designated as  $1 - p^t$  to handle cases where the probability  
 261 changes over time. In each trial, a uniformly distributed random number was generated from interval  $[0.0,$   
 262  $1.0]$ . Numbers equal to or less than predefined  $p^t$  were regarded as heads; numbers larger than  $p^t$ , were  
 263 tails.

264 Whenever a coin toss result is observed, the correct probability of landing heads ( $p^t$ ) is estimated  
 265 by extended Bayesian inference. Additionally, for comparison with extended Bayesian inference, estimation

266 using only inverse Bayesian inference and estimation using Exponential Moving Average (EMA) were also  
267 carried out.

268 First, let heads be expressed as *HEAD* and tails as *TAIL*. Second, we prepare  $N$  hypotheses ( $h_1,$   
269  $h_2, \dots, h_N$ ) and define the probability of heads and the probability of tails in each hypothesis  $h_i$  as follows.

$$\begin{aligned} C(\text{HEAD}|h_i) &= 0.5 \\ C(\text{TAIL}|h_i) &= 1.0 - C(\text{HEAD}|h_i) = 0.5 \end{aligned} \quad (14)$$

270

271 That is, the models for all hypotheses are the same, and it means that this system has substantially no  
272 model of hypothesis at the initial stage.

273 Further, we must suppose that the prior probability for each hypothesis is equal.

$$C(h_i) = \frac{1}{N} \quad (15)$$

274

275 Whenever a coin toss result is observed, by performing extended Bayesian inference using formulas (7),  
276 (8), (12), and (13),  $C(\text{HEAD})$  is successively updated. In the simulations,  $N = 3$ . For the simplicity of  
277 subsequent analysis, in the following simulations, the parameter  $m$  was fixed to  $-1$  in the extended  
278 Bayesian inference.

279 When updating the degree of confidence  $C(h_i)$  for each hypothesis using formula (7), we set the  
280 minimum value  $\varepsilon$  to impose a restriction so that the degree of confidence will not be zero.

$$C(h_i) \leftarrow \max \left( \varepsilon, \frac{C(h_i)C(\mathbf{d}|h_i)}{[(1-\alpha)(C(\mathbf{d}))^{-m} + \alpha(C(h_i))^{-m}]^{-1/m}} \right) \quad (16)$$

281 Where  $\max(x, y)$  is a function whose output is a larger value of the two arguments  $x, y$ . In the simulation,  
282  $\varepsilon$  was set to 0.00001.

283 In case only inverse Bayesian inference is performed, the hypothesis is limited to only  $\mathbf{h}$ , the process of  
284 formula (7) is not performed, and  $C(\mathbf{h})$  is always set to 1.0.

285 In this paper, we deal with a task in which the probability of heads can take two values, and they are  
286 replaced by the probability  $\theta$ . If a uniformly distributed random number generated from interval  $[0.0, 1.0]$   
287 at the  $t^{th}$  trial is denoted as  $rnd^t$ , the probability of heads is expressed by the following formula.

288

$$p^{t+1} = \begin{cases} p^t & \text{if } rnd^t > \theta \\ 1 - p^t & \text{else} \end{cases} \quad (17)$$

289

290 In this simulation,  $\theta$  was set to 0.0001. The initial value of the probability of heads  $p^0$  was set to 0.85.  
291 That is, the probability of heads can take two values of 0.15 and 0.85.

292 EMA is calculated as a weighted average between the current estimated value and the observed data as  
293 follows.

$$s^{t+1} = \beta d^t + (1 - \beta)s^t = s^t + \beta(d^t - s^t) \quad (18)$$

294

295 Here,  $s^t$  and  $s^{t+1}$  represent estimated values of the probability of heads at  $t^{th}$  and  $t + 1^{th}$  trial,  
296 respectively.  $\alpha$  represents a learning rate, which takes a value of interval  $[0.0, 1.0]$ .  $d^t$  represents the  
297 result of coin toss at  $t^{th}$  trial; in the case of *HEAD*,  $d^t = 1$ , and in the case of *TAIL*,  $d^t = 0$ . The weight  
298 of each data decreases exponentially as it goes to the past, and it is expressed by  $\beta(1 - \beta)^x d^{t - (x + 1)}$ . In the  
299 simulation, the value of  $\beta$  was shifted from 0.0005 to 0.0063 with 0.0002 increment. That is, the  
300 estimations by EMA were performed using  $\beta$  of the 30 patterns.

301 Fig 3 (a), (b), and (c) show the results of extended Bayesian inference, inverse Bayesian inference  
302 and estimations by EMA. However, for EMA, only three results of  $\beta = 0.0005$ ,  $\beta = 0.0021$  and  $\beta =$   
303 0.0063 are shown for easy viewing.

304 As can be seen from Fig 3(c), in the estimations by EMA, although rapid change can be followed as the  
305 learning rate  $\alpha$  increases, the fluctuation in the stable period becomes larger. In other words, there is a  
306 trade-off between the ability to follow change and the accuracy in the stable period.

307           The result of only inverse Bayesian inference is very similar to the estimation result by EMA with  $\beta$   
308       = 0.0021. On the other hand, although extended Bayesian inference takes time to follow changes as in the  
309       case of inverse Bayesian inference initially, the ability to follow gradually improves, and it becomes possible  
310       to respond rapidly to sudden changes.

311

312       **Fig 3. Time progress of the estimated values for the probability of head landing.** The figure includes the  
313       correct probability. (a) Estimated values by extended Bayesian inference. (b) Estimated values by only  
314       inverse Bayesian inference. (c) Estimated values by EMA.

315           Fig 4 shows the internal state of the extended Bayesian inference. Fig 4 (a) shows the time progress  
316       of the probability of head landing for each hypothesis. Initially, the probabilities for all hypotheses were 0.5  
317       by definition, but learning by inverse Bayesian inference gradually formed hypothesis models. After the  
318       middle stage, the probabilities of head landing for three hypotheses  $h_1$ ,  $h_2$  and  $h_3$  became approximately  
319       0.15, 0.85, and 0.5, respectively. Here, 0.15 and 0.85 correspond to two correct values in this simulation, as  
320       shown in formula (17).

321           Figure 4(b) shows the time progress of the hypothesis with the greatest degree of confidence. As can  
322       be observed from the figure, in the second half of the simulation, extended Bayesian inference switches the  
323       hypotheses quickly when the probability of head landing is replaced. That is, extended Bayesian inference  
324       tries to respond to changes by learning using inverse Bayesian inference in the first half of the simulation,  
325       while in the second half of the simulation, abrupt changes are dealt with by switching the hypotheses formed  
326       by the learning.

327

328       **Fig 4. Internal state of extended Bayesian inference.** (a) Time progress of the probability of head landing  
329       for each hypothesis . (b) Time progress of the hypothesis with the greatest degree of confidence.

330



331 We show the relationship between the ability to follow the sudden change and the accuracy of the  
332 estimation in the stable period. Fig 5 shows the results of estimation by extended Bayesian inference,  
333 inverse Bayesian inference, and EMA in an enlarged manner between 40914<sup>th</sup> trial at which the  
334 probability of heads suddenly changed from 0.15 to 0.85 and the 60913<sup>th</sup> trial.

335

336 **Fig 5. Time progress of the estimated values for the probability of head landing.** The figure includes the  
337 correct probability.

338

339 This period was divided into the first half interval from 40914<sup>th</sup> trial to 50913<sup>th</sup> trial and the second  
340 half interval from 50914<sup>th</sup> trial to 60913<sup>th</sup> trial, and the differences between the correct values and the  
341 estimated values were calculated using root-mean-square error (RMSE) in each interval. RMSE is defined  
342 as follows.

$$RMSE = \sqrt{\frac{\sum_{t=k}^{T+k-1} (\hat{x}_t - x_t)^2}{T}} \quad (19)$$

343

344 Here,  $\hat{x}_t$  and  $x_t$  represent the correct value and the estimated value in  $t^{\text{th}}$  trial, respectively.  $T$   
345 represents the length of the interval.

346 We use the RMSE of the first half as a measure of the ability to follow rapid changes, and the RMSE of  
347 the second half as a measure of the accuracy of the estimation in the stable period.

348 Fig 6 shows the relationship between the followability and the estimation accuracy in the extended  
349 Bayesian inference, the inverse Bayesian inference, and EMA estimations.

350 As can be seen from the figure, there was a trade-off relationship of the accuracy being lost if the  
351 followability improved in EMA estimation. The regression curve for EMA data is also shown in this figure.

352 The data of the inverse Bayesian inference was located slightly lower left on this trade-off curve. On the  
 353 other hand, the data of the extended Bayesian inference was almost the same as the data of the inverse  
 354 Bayesian inference with regards to the accuracy, but the followability was greatly improved.

355 That is, it can be seen that the extended Bayesian inference broke the trade-off found in EMA estimation.

356

357 **Fig 6. Relationship between the followability and the estimation accuracy in the extended Bayesian**  
 358 **inference, the inverse Bayesian inference, and EMA estimations.**

359

360 In the extended Bayesian inference, the inverse Bayesian inference could be applied only to the  
 361 hypothesis  $\mathbf{h}$  which has the greatest degree of confidence. Moreover, we set  $m = -1$ . Because of this, we  
 362 can rewrite formula (8) for inverse Bayesian inference as follows.

$$C(\mathbf{d}|\mathbf{h}) \leftarrow \frac{C(\mathbf{h})}{(1-\alpha)C(\mathbf{h}) + \alpha C(\mathbf{d})} \cdot C(\mathbf{d}|\mathbf{h}) \quad (20)$$

363 Here, we can see that the denominator on the right side is the weighted average of  $C(\mathbf{h})$  and  $C(\mathbf{d})$ ,  
 364 if  $C(\mathbf{h}) > C(\mathbf{d})$ ,  $C(\mathbf{h}) > (1-\alpha)C(\mathbf{h}) + \alpha C(\mathbf{d})$ , so  $C(\mathbf{d}|\mathbf{h})$  increases. At this point, the increment of  $C$   
 365  $(\mathbf{d}|\mathbf{h})$  is larger if the degree of confidence  $C(\mathbf{h})$  is higher.

366 Conversely, if  $C(\mathbf{h}) < C(\mathbf{d})$ ,  $C(\mathbf{h}) < (1-\alpha)C(\mathbf{h}) + \alpha C(\mathbf{d})$ , so  $C(\mathbf{d}|\mathbf{h})$  decreases greatly if the  
 367 degree of confidence  $C(\mathbf{h})$  is lower.

368 Let us turn to the analysis of the stable period where  $C(\mathbf{h}) = 1$ . When  $C(\mathbf{h}) = 1$ , the total sum of  
 369 confidence of all hypotheses is 1, and for any hypothesis  $h_i$  other than  $\mathbf{h}$ ,  $C(h_i) = 0$ . Hence, formula (13)  
 370 can be rewritten as follows.

$$C(\mathbf{d}) = \sum_i C(h_i)C(\mathbf{d}|h_i) = C(\mathbf{d}|\mathbf{h}) \quad (21)$$

371 At this step, formula (20) for inverse Bayesian inference can also be rewritten as follows.

$$C(\mathbf{d}|\mathbf{h}) \leftarrow \frac{1}{(1-\alpha) \cdot 1 + \alpha C(\mathbf{d}|\mathbf{h})} \cdot C(\mathbf{d}|\mathbf{h}) \quad (22)$$

372

373 The right side of this formula shows the weighted harmonic average of 1 and  $C(\mathbf{d}|\mathbf{h})$ . Since  $0 \leq C$   
 374  $(\mathbf{d}|\mathbf{h}) \leq 1$ , the denominator is necessarily less than 1, and the likelihood  $C(\mathbf{d}|\mathbf{h})$  increases whenever  
 375 updated. In other words, when certain data are observed, the connection between the data and the hypothesis  
 376 with the highest degree of confidence at that time is reinforced. Conversely, unobserved data, i.e., for  $d_j$   
 377 other than  $\mathbf{d}$ ,  $C(d_j|\mathbf{h})$  can be standardised using formula (12), hence decreasing by the increment amount  
 378 in  $C(\mathbf{d}|\mathbf{h})$ . Here, the rate of increase for  $C(\mathbf{d}|\mathbf{h})$  depends on  $\alpha$ . When  $\alpha = 1$ , regardless of the presence  
 379 value of  $C(\mathbf{d}|\mathbf{h})$ , the right side of formula (22) is 1. As  $\alpha$  gets smaller, the increase rate lowers, and when  
 380  $\alpha = 0$ , it coincides with Bayesian inference and  $C(\mathbf{d}|\mathbf{h})$  becomes invariable.

381 These considerations suggest that  $\alpha$  becomes larger according to increase of updates to the model.  
 382 In this sense, we can say that formula (20) for inverse Bayesian inference during the steady period represents  
 383 a process of learning, and  $\alpha$  corresponds to the rate of learning.

384 With respect to the portion that corresponds to Bayesian inference, suppose  $m = -1$  in the formula  
 385 (7), then we can rewrite it as:

$$C(h_i) \leftarrow \frac{C(h_i)C(\mathbf{d}|h_i)}{[(1-\alpha)(C(\mathbf{d}))^{-m} + \alpha(C(h_i))^{-m}]^{-1/m}} = \frac{C(h_i)C(\mathbf{d}|h_i)}{(1-\alpha)C(\mathbf{d}) + \alpha C(h_i)} \quad (23)$$

386 Through careful observation of this formula we can note that  $\alpha$  becomes larger, while  $C(\mathbf{d})$ , i.e.,  
 387 the effect of observation data, gets smaller. Where  $\alpha = 1$ , the denominator  $C(\mathbf{d})$  disappears so  $C(h_i)$  in  
 388 the numerator and denominator is cancelled, and the formula can be expressed as follows.

$$C(h_i) \leftarrow C(\mathbf{d}|h_i) \quad (24)$$

389  
 390 This means that when  $\alpha = 1$ , the degree of confidence in each hypothesis  $C(h_i)$  does not even  
 391 consider the past observation history and is seen to be identical to the likelihood at that point in time. This  
 392 coincides with the maximum likelihood estimation. In contrast, when  $\alpha = 0$ , the present formula coincides  
 393 with the Bayesian inference expressed in formula (5).

394 A comparison of formula (5) with formula (24) reveals that their difference lies in the presence or  
395 absence of  $C(h_i)$  on the right side, because  $P(\mathbf{d})$  can be regarded as a constant. In this sense, the difference  
396 between them can be said to lie in the extent to which they accept history in order to determine degree of  
397 confidence.

398 As shown above, in extended Bayesian inference, symmetry bias plays two roles. First, the strength  
399 of symmetry bias indicates the rate of learning in portions of inverse Bayesian inference. In other words, the  
400 stronger the symmetry bias, the greater the degree of modification to the hypothesis model based on  
401 observational data. Second, the strength of the symmetry bias indicates how much the model takes into  
402 account past history in portions of Bayesian inference. In other words, as symmetry bias gets stronger,  
403 confidence in each hypothesis is updated based solely on more recent observational data.

404

## 405 **Discussion**

406 In this study, we first proposed a different causal induction model. This model can replicate human  
407 judgments concerning causal induction with higher accuracy than previous models. Then we formulated an  
408 inference model that incorporates the said causal induction model into Bayesian inference. We noticed this  
409 inference model necessarily involves inverse Bayesian inference, which allows for flexibility to handle  
410 unsteady situations where the inference target changes from time to time. Finally, we demonstrated how this  
411 model can work well with unknown situations by forming new hypotheses through inverse Bayesian  
412 inference.

413 The causal induction model that we are proposing has two parameters that control the strength of  
414 symmetry bias. Conditional probability and causal induction models like *DFH* and *pARIs* can be shown to  
415 be special cases where particular values are assigned to parameters in our model. In other words, using the  
416 proposed model allows us to seamlessly express degree of confidence in those statements of the forms ‘if P  
417 then Q’, which stand for prediction, as well as those in the form of ‘P therefore Q’, which stand for causal

418 relation in a single model. The results of the meta-analysis of causal induction experiments revealed that the  
419 proper incorporation of symmetry bias in the proposed model allows it to replicate human judgment with  
420 high accuracy. However, as shown in Table 2, the values of the parameters are different for each experiment.  
421 It is known that the interpretation of conditionals largely change depending on the type and the contents of  
422 the conditionals, as well as subject's age [26]. Further studies are necessary to determine how parameters  
423 change according to type and contents of conditionals, as well as age.

424 The performance of the present model and the catalogue performance of *DFH* and *pARIs* were  
425 compared using eight sorts of experiment data. It turned out that the proposed model recorded the best  
426 determination coefficient in five of these experiments and the mean of all experiments. Thus, the present  
427 model is important in two ways. First, regarding human causal induction judgments, it has a capability that  
428 outperforms *DFH* and *pARIs*, which had hitherto demonstrated the best catalogue performance. Second, by  
429 using the extended confidence model, it is possible to determine the parameters  $\alpha$  and  $m$  that can best explain  
430 the participant's judgement from the data of simple causal induction experiments. In other words, we can  
431 measure the strength of symmetry bias the participant has. It would be possible to compare the strength of  
432 symmetry bias in patients with a mental illness against that of healthy individuals. For example, 'the Von  
433 Domarus principle' applies to the speech of schizophrenic patients [27] and refers to an inference of the form  
434 'Men die. Weed die. Therefore, men are weed'. There is a widely observed tendency in schizophrenic patients  
435 to identify two things as the same when they share a common property – a mechanism said to underlie  
436 delusion [28]. Logically, it is wrong to conclude from 'A is C' and 'B is C' that 'A is B'. However, if the  
437 symmetry bias allows derivation of 'C is B' from 'B is C', then from 'A is C' and 'C is B', we can derive 'A  
438 is B'. In other words, one influence on the delusion of schizophrenic patients may be a strong susceptibility  
439 to symmetry bias. To test this hypothesis, it is possible to use the proposed model to estimate and compare  
440 the strength of symmetry bias in both patients with schizophrenia and healthy people. This is a research goal  
441 we can trace.

442 Parameters that denote the strength of symmetry bias indicate, in the case of extended Bayesian  
443 inference, the strength of inverse Bayesian inference, that is, the rate of learning. At the same time, they  
444 indicate how much the history is taken into account when updating the degree of confidence in each  
445 hypothesis. In this way, learning and inference become interlocked via parameters that denote the strength of  
446 symmetry bias in such way that it takes account only of more recent observational data in inference as the  
447 rate of learning becomes greater.

448 As a third form of inference following induction and deduction, various authors mention abduction.  
449 Abduction is an inference from knowledge or from known theories, for example, that ‘if it rains, people open  
450 an umbrella’ and to determine that from ‘people open an umbrella’ we can infer ‘it is raining’. Abduction  
451 can be seen as a procedure of selecting hypothesis that best explain the observational data. In this regard,  
452 abduction is akin to maximum likelihood estimation and Bayesian inference. They differ, however, in that  
453 while the latter two inferences proceed by extracting the optimal hypothesis from the existing ones based on  
454 observational data, abduction focuses in the formation of a new hypothesis. An example is Newton, who  
455 introduced the law of universal gravitation to explain free fall of physical bodies.

456 Whereas the models for maximum likelihood estimation and Bayesian inference remain constant,  
457 the model for extended Bayesian inference is modified by virtue of inverse Bayesian inference enabling the  
458 model to match observation data. In other words, a new hypothesis that better explains the fact is formed in  
459 each case. In this sense, extended Bayesian inference that accompanies inverse Bayesian inference can be  
460 said to be akin to abduction.

461 In interpersonal communication, it is important to mutually estimate emotions of others. Further, for  
462 future studies on human-machine interaction this sort of information is essential. When estimating the  
463 emotion of others, since we cannot directly observe private internal states, there is no way to estimate their  
464 emotion other than using external clues (observational data) such as facial expression and tone of voice. In  
465 general, to perform proper estimation under these circumstances, the more observational data the better.  
466 However, emotions are not always constant, it is a variable that changes from time to time. Under these

467 circumstances, we must infer emotions while considering whether observational data to be used for  
468 estimation derive from the same emotions. In the future it would be interesting to apply extended Bayesian  
469 inference to tasks like these.

470         Suppose that we have knowledge (based on models of others) of the kinds of emotions the other has  
471 and how they are expressed in him. Of course, one cannot attain perfect or complete knowledge because his  
472 mental states have some degree of privacy. Suppose that when estimating based on incomplete knowledge  
473 that the other is pleased, his expression suddenly changes. At this moment, one may think that his emotion  
474 has changed or that this is another way of expressing pleasure. The former is an inference based on knowledge  
475 and corresponds to Bayesian inference. The latter, on the other hand, is a modification of knowledge or an  
476 addition of new knowledge and it corresponds to inverse Bayesian inference.

477         Incorporating the function of inverse Bayesian inference may help to develop robots that  
478 autonomously learn and make various human-like inferences.

479

## 480 **Methods**

### 481 **Data used in meta-analysis**

482         To test the descriptive performance of the extended confidence model, we performed the meta-  
483 analysis using the same data as [13]. The analysis was conducted using eight types of experiment data, that  
484 is, AS95, BCC03.1, BCC03.3, HO07.1, HO07.2, LS00, W03.2, and W03.6.

485         In AS95, forty graduate and undergraduate students were recruited. They were given co-occurrence  
486 information about the presence or absence of drug treatment and the presence or absence of the side effects.  
487 The subjects were asked to judge a number of problems, and each problem involved a sequence of instances  
488 of these four information types. The frequencies of each information type varied from problem to problem.  
489 At the end of a problem, the subjects were asked to enter a number from 0 to 100 that best reflected their  
490 judgment of the drugs causing the side effects.

491 In BCC03.1, 109 undergraduate students were recruited and divided into two groups (preventive  
492 group and generative group). In preventive group, they evaluated how effectively each vaccine prevented the  
493 corresponding disease by giving a rating on a scale from 0 (the vaccine does not prevent the disease at all) to  
494 100 (the vaccine prevents the disease every time). They also evaluated the influence of ray exposure on the  
495 mutation of viruses in generative group.

496 Thirty-one undergraduate students participated in BCC03.3. With regard to the side effects of drugs  
497 that reduce allergy, participants determined whether there were side effects of headache and, if so, assessed  
498 the causal strength between the drug and the headache.

499 In HO07.1, participants were 39 undergraduate students. They were asked to assess the strength of  
500 the causal relation between a particular type of fertiliser and plants blooming. They only observed a sequence  
501 of scenes in which fertiliser and plant blooming were either present or absent. After observing a series of  
502 situations, participants rated the subjective strength of the causal relation with a value between 0 (completely  
503 unrelated) and 100 (completely related).

504 In HO07.2, participants were 50 undergraduate students. In this experiment the cause was ‘drinking  
505 milk’ and the effect was ‘stomach-ache’. They judged the causal strength between drinking milk and  
506 stomach-ache according to given co-occurrence information.

507 In LS00, the participants of experiments 1, 2 and 3 were 27, 16, and 24 students, respectively. They  
508 assessed the extent to which a certain chemical causes a mutation in animals’ DNA using a number from  
509 between 0 and 100, where 0 indicates that the chemical does not cause mutations at all and 100 indicates that  
510 the chemical causes a mutation.

511 In W03.2, the participants were 40 undergraduate students. They were given information on the  
512 additives (manganese trioxide) contained in the foods a patient has eaten, and information on whether the  
513 patient has developed an allergic reaction. They were asked to judge the extent to which the statement  
514 ‘Manganese trioxide causes the allergic reaction in this patient’ was right for that patient and to write a



515 number from 0 (zero) to 100, where 0 (zero) means that the statement is definitely not right, and 100 means  
516 that the statement is definitely right.

517 In W03.6. the participants were 43 first-year undergraduate students. Most features of method,  
518 including initial written instructions; format of stimulus presentations; and procedure, were the same as in  
519 W03.2. The studies differed in design, however.

520

## 521 **Extended Bayesian inference**

522 First, we can replace  $p$  and  $q$  in formula (1) with  $h_i$  and  $d$ .

$$c(d|h_i) = [(1 - \alpha)P(d|h_i)^m + \alpha P(h_i|d)^m]^{1/m} \quad (25)$$

523 Then we apply the Bayes' theorem to the right side of the formula (25), and we can obtain the conversion as  
524 follows.

$$\begin{aligned} c(d|h_i) &= \left[ (1 - \alpha) \left( \frac{P(h_i)P(d|h_i)}{P(h_i)} \right)^m + \alpha \left( \frac{P(h_i)P(d|h_i)}{P(d)} \right)^m \right]^{1/m} = \left[ (1 - \alpha) \left( \frac{1}{P(h_i)} \right)^m + \alpha \left( \frac{1}{P(d)} \right)^m \right]^{1/m} P(h_i)P(d|h_i) \\ &= [(1 - \alpha)(P(h_i))^{-m} + \alpha(P(d))^{-m}]^{1/m} P(h_i)P(d|h_i) = \frac{P(h_i)P(d|h_i)}{[(1 - \alpha)(P(h_i))^{-m} + \alpha(P(d))^{-m}]^{-1/m}} \end{aligned} \quad (26)$$

525 By replacing  $h_i$  and  $d$  in formula (25), to perform the same conversion, we get

$$c(h_i|d) = \frac{P(d)P(h_i|d)}{[(1 - \alpha)(P(d))^{-m} + \alpha(P(h_i))^{-m}]^{-1/m}} = \frac{P(h_i)P(d|h_i)}{[(1 - \alpha)(P(d))^{-m} + \alpha(P(h_i))^{-m}]^{-1/m}} \quad (27)$$

526 In the next step, we replace the conditional probability  $P$  on the right side of formulas (26) and (27)  
527 with the extended confidence  $C$  to make the formulas recursive, and then we replace the equation with the  
528 update formula.

$$c(d|h_i) \leftarrow \frac{C(h_i)C(d|h_i)}{[(1 - \alpha)(C(h_i))^{-m} + \alpha(C(d))^{-m}]^{-1/m}} \quad (28)$$

$$c(h_i|d) \leftarrow \frac{C(h_i)C(d|h_i)}{[(1 - \alpha)(C(d))^{-m} + \alpha(C(h_i))^{-m}]^{-1/m}} \quad (29)$$

529

530 As seen in formula (28), in inverse Bayesian inference, the amount of modification to the model of  
 531 each hypothesis increases as  $\alpha$  becomes larger. However, in this article, not all hypothetical models are  
 532 uniformly modified, and the amount of modification changes according to confidence levels as follows.

$$\alpha_i = \alpha \cdot \pi(h_i) \quad (30)$$

533 However,

$$\pi(h_i) = \frac{\exp(C(h_i)/\tau)}{\sum_i \exp(C(h_i)/\tau)} \quad (31)$$

534  
 535 Here, formula (31) is a procedure from the field of machine learning called Softmax [29], and  $\tau$   
 536 ( $> 0$ ) is a parameter termed temperature.  $\pi$  remains the same value for all hypotheses if the temperature  
 537 is high with the limit  $\tau \rightarrow \infty$ . On the other hand, if the temperature is low,  $\pi$  becomes greater for hypotheses  
 538 with a higher confidence level.  $\pi$  takes the value 1 in the limit  $\tau \rightarrow 0$  for hypotheses with the highest  
 539 confidence level, and takes value 0 for all the other hypotheses.

540 In inverse Bayesian inference, the hypothetical model is modified when  $\mathbf{d}$  is observed using  $\alpha_i$  as  
 541 follows.

$$C(\mathbf{d}|h_i) \leftarrow \frac{C(h_i)C(\mathbf{d}|h_i)}{[(1 - \alpha_i)(C(h_i))^{-m} + \alpha_i(C(\mathbf{d}))^{-m}]^{-1/m}} \quad (32)$$

542 Here, there are reasons why the degree of modification for each hypothetical model changes  
 543 according to the level of confidence. First, this process is a modification of the hypothetical model, which  
 544 can be understood as a learning procedure rather than inference. Second, it is more likely that the currently  
 545 observed data  $\mathbf{d}$  derives from a hypothesis, if that hypothesis has a higher degree of confidence. Therefore,  
 546 when modifying the model for each hypothesis based on observed data  $\mathbf{d}$ , the hypothesis with a higher  
 547 degree of confidence requires a greater modification of its model. Of course, when  $\tau \rightarrow \infty$ , all hypothetical  
 548 models can equally be modified. On the other hand, when  $\tau \rightarrow 0$ , only the hypothesis model with the highest

549 degree of confidence is modified. Moreover, supposing  $\alpha = 0$ , in formula (30),  $\alpha_i = 0$  obtains for all  
550 hypotheses.

551 In the simulation,  $\tau \rightarrow 0$  was set. In other words, the inverse Bayesian inference was applied to only  
552 the hypothesis  $h$  that has the highest confidence value.

553

## 554 **Acknowledgements**

555 This research is supported by the Center of Innovation Program from the Japan Science and Technology  
556 Agency, JST and by JSPS KAKENHI Grant Numbers JP16K01408 and JP15H03002. We would like to  
557 thank Editage [<http://www.editage.com>] for editing and reviewing this manuscript for English language.

558

559

## 560 **References**

- 561 1. Sidman M, Rauzin R, Lazar R, Cunningham S, Tailby W, Carrigan P. A search for symmetry in the  
562 conditional discriminations of rhesus monkeys, baboons, and children. *J Exp Anal Behav.* 1982;37: 23-  
563 44.
- 564 2. Sidman M, Tailby, W. Conditional discrimination vs. matching-to-sample: an expansion of the testing  
565 paradigm. *J Exp Anal Behav.* 1982;37: 5-22.
- 566 3. Yamazaki Y. Logical and illogical behavior in nonhuman animals. *Jap Psychol Res.* 2004;46: 195-206.  
567
- 568 4. Markman E, Wachtel G. Children's use of mutual exclusivity to constrain the meanings of words. *Cogn*  
569 *Psychol.* 1988;20: 121-157.
- 570 5. Davidson D, Tell, D. Monolingual and bilingual children's use of mutual exclusivity in the naming of  
571 whole objects. *J Exp Child Psychol.* 2005;92: 25-45.  
572
- 573 6. Halberda, J. The development of a word-learning strategy. *Cognition.* 2003;87: B23-B34.  
574
- 575 7. Lipkens R, Hayes S, Hayes L. Longitudinal study of the development of derived relations in an infant. *J*  
576 *Exp Child Psychol.* 1993;56: 201-239.
- 577 8. O'Donnell J, Saunders K. Equivalence relations in individuals with language limitations and mental  
578 retardation. *J Exp Anal Behav.* 2003;80: 131-157.  
579  
580  
581  
582  
583

- 584 9. Takahashi T, Nakano M, Shinohara S. Cognitive symmetry: illogical but rational biases. *Symmetry: Cult*  
585 *Sci.* 2011;21: 275-294.
- 586
- 587 10. Kastak C, Schusterman R, Kastak D. Equivalence classification by California sea lions using class-  
588 specific reinforcers. *J Exp Anal Behav.* 2001;76: 131-158.
- 589
- 590 11. Tomonaga M, Matsuzawa T, Fujita K, Yamamoto J. Emergence of symmetry in a visual conditional  
591 discrimination by chimpanzees (*Pan troglodytes*). *Psycholo Rep.* 1991;68: 51-60.
- 592
- 593 12. Evans J, Handley S, Over D. Conditionals and conditional probability. *J Exp Psychol: Learn mem cogn.*  
594 2003;29: 331-335.
- 595
- 596 13. Hattori M, Oaksford M. Adaptive non-interventional heuristics for covariation detection in causal  
597 induction: model comparison and rational analysis. *Cogn Sci.* 2007;31: 765-814.
- 598
- 599 14. Jenkins H, Ward W. Judgment of contingency between responses and outcomes. *Psychol Monogr.*  
600 1965;79: 1-17.
- 601
- 602 15. Takahashi T, Oyo K, Tamatsukuri A. Correlation detection with and without the theory of conditionals:  
603 a model update of Hattori & Oaksford (2007). doi: 10.1101/247742.
- 604
- 605 16. Dehaene S. *Consciousness and the brain: Deciphering how the brain codes our thoughts.* New York:  
606 Viking; 2014.
- 607
- 608 17. Chater N, Oaksford M. *The probabilistic mind: Prospects for Bayesian cognitive science.* Oxford:  
609 Oxford University Press; 2008.
- 610
- 611 18. Arecchi F. Phenomenology of consciousness from apprehension to judgment: nonlinear dynamics.  
612 *Psychol Life Sci.* 2011;15: 359-375.
- 613
- 614 19. Gunji Y, Shinohara S, Haruna T, Basios V. Inverse Bayesian inference as a key of consciousness  
615 featuring a macroscopic quantum logical structure. *Biosystems.* 2016;152: 44-65.
- 616
- 617 20. Gunji Y, Murakami H, Tomaru T, Basios V. Inverse Bayesian inference in swarming behaviour of  
618 soldier crabs. *Philos Trans A Math Phys Eng Sci.* 2018;376: 20170370.
- 619
- 620 21. Horry Y, Yoshinari A, Nakamoto Y, Gunji Y. Modeling of decision-making process for moving  
621 straight using inverse Bayesian inference. *Biosystems.* 2018;163: 70-8.
- 622
- 623 22. Anderson J, Sheu C. Causal inferences as perceptual judgments. *Mem Cogn.* 1995;23: 510-524.
- 624
- 625 23. Buehner M, Cheng P, Clifford D. From covariation to causation: a test of the assumption of causal  
626 power. *J Exp Psychol Learn Mem Cogn.* 2003;29: 1119-1140.
- 627
- 628 24. Lober K, Shanks D. Is causal induction based on causal power? Critique of Cheng (1997). *Psychol Rev.*  
629 2000;107: 195-212.
- 630
- 631 25. White P. Making causal judgments from the proportion of confirming instances: the pCI rule. *J Exp*  
632 *Psychol Learn Mem Cogn.* 2003;29: 710-727.

633

634 26. Gauffroy C, Barrouillet P. Heuristic and analytic processes in mental models for conditionals: an  
635 integrative developmental theory. *Dev Rev.* 2009;29: 249-282.

636

637 27. Von Domarus E. The specific laws of logic in schizophrenia. In: Kasanin JS, editor. *Language and*  
638 *thought in schizophrenia*. Berkeley: University of California Press; 1944. pp. 104-113.

639

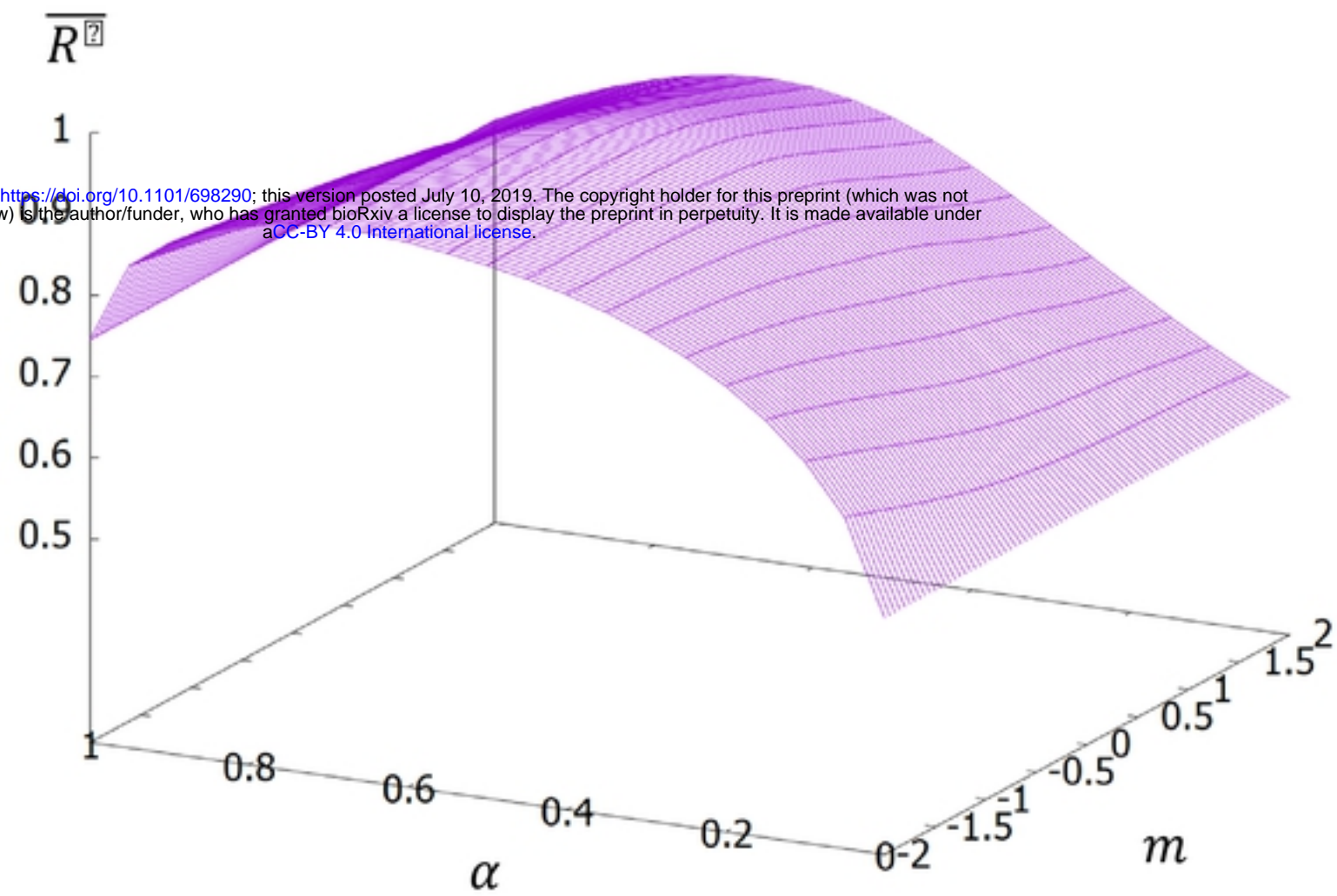
640 28. Arieti S. *Interpretation of schizophrenia*. New York: Basic Books; 1957.

641

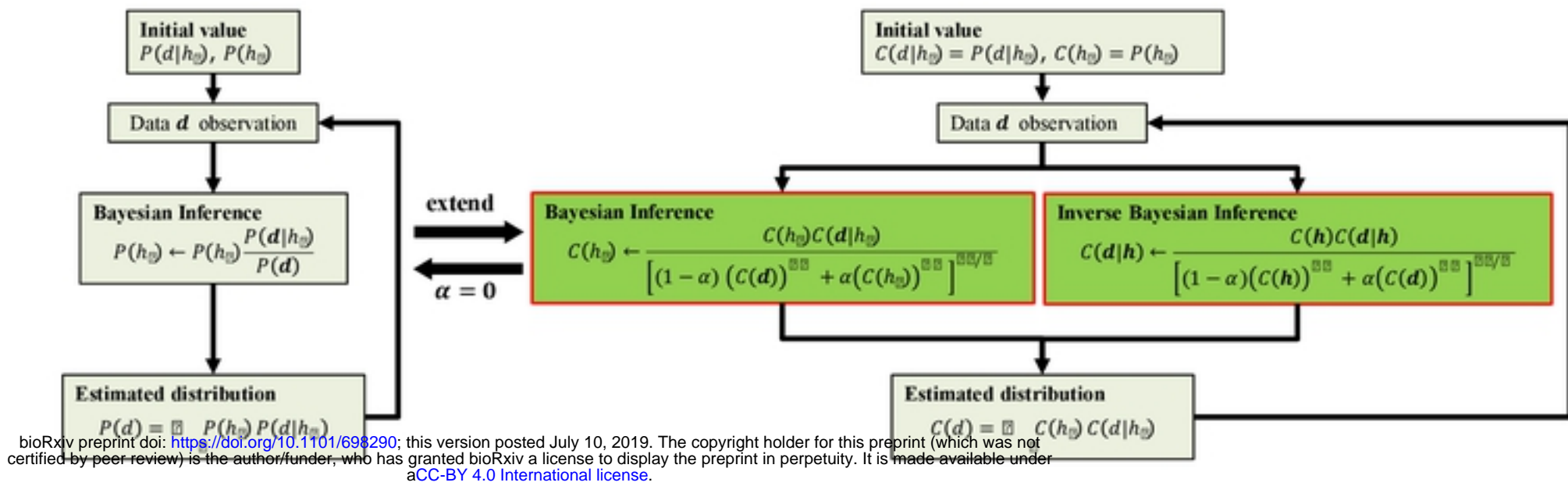
642 29. Sutton R, Barto A. *Reinforcement learning: An introduction*. Cambridge, MA: MIT press; 1998.

643

bioRxiv preprint doi: <https://doi.org/10.1101/698290>; this version posted July 10, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



**Fig. 1.** Mean coefficient of determination for each pair of parameter values. The mean value for the determination coefficient peaked when  $(\alpha, m) = (0.35, -1.15)$  and the value was  $R^2 = 0.93$ .



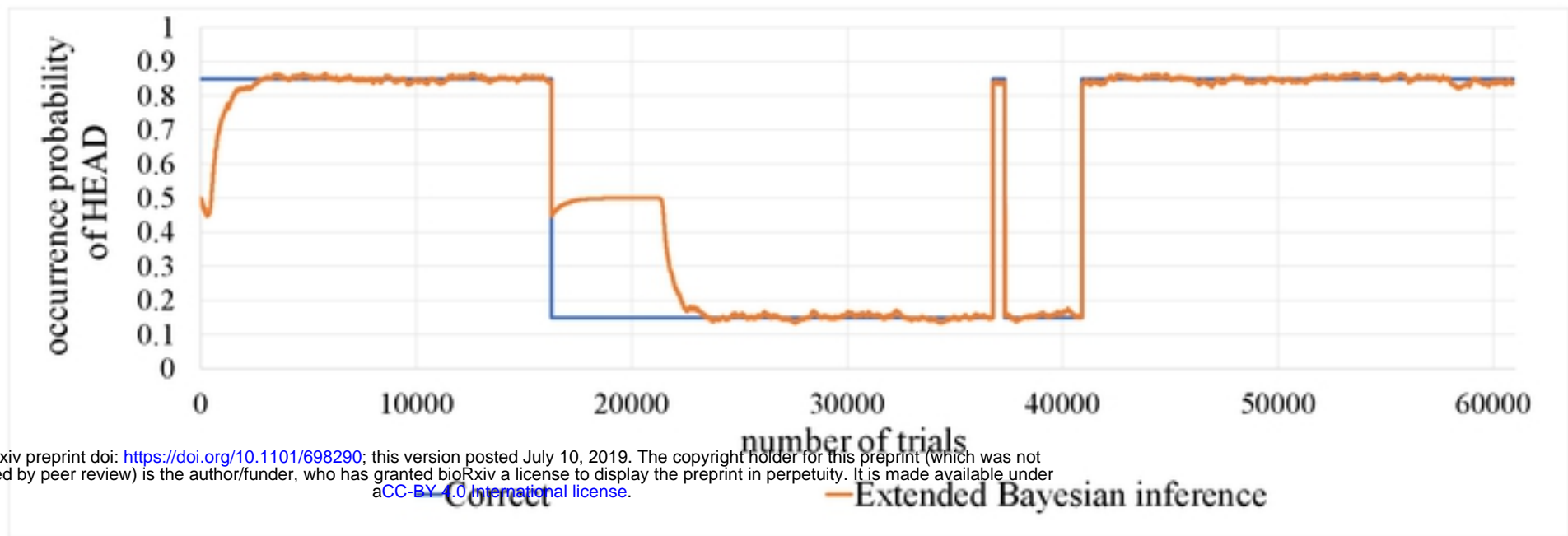
(a) Bayesian Inference

(b) Extended Bayesian Inference

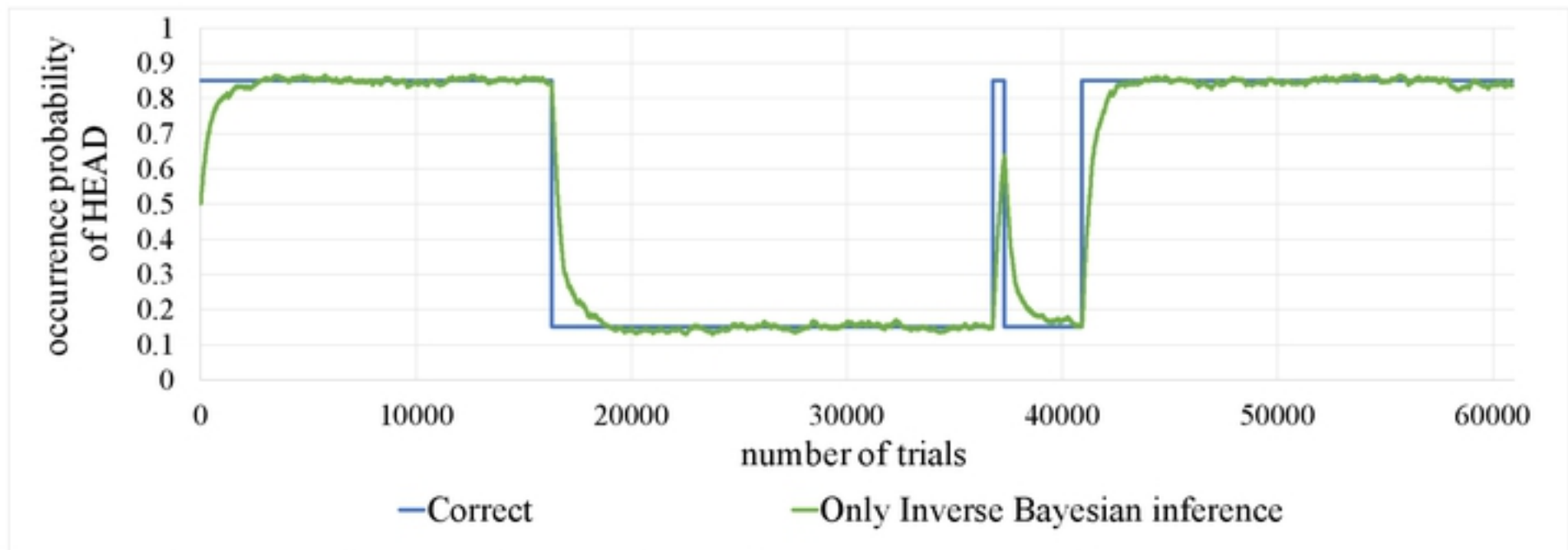
**Fig. 2.** A flow chart comparing Bayesian Inference to Extended Bayesian inference. (a) An overview of the processing flow of Bayesian inference; and (b) An overview of the processing flow of extended Bayesian inference. In (b), for simplicity, the portion that belongs to the normalisation process is omitted. In (b), if we suppose  $\alpha = 0$ , the portion of inverse Bayesian inference disappears, corresponding to the Bayesian inference in (a).



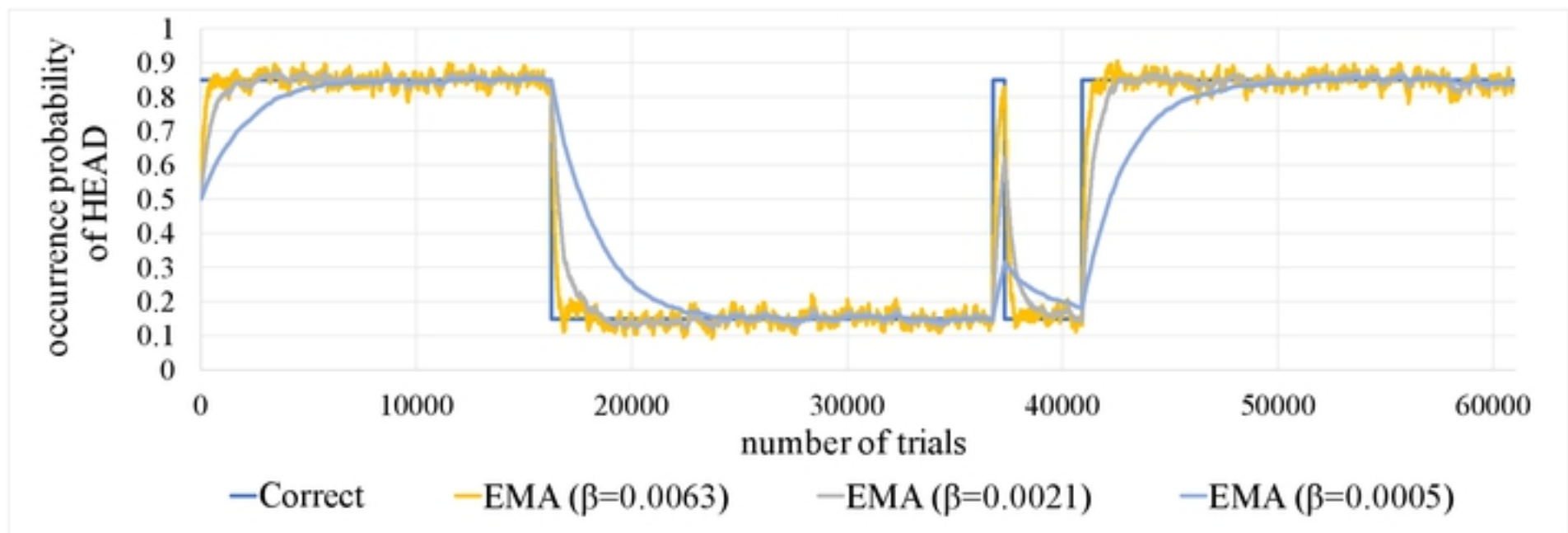
bioRxiv preprint doi: <https://doi.org/10.1101/698290>; this version posted July 10, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



(a) Estimated values by extended Bayesian inference  $(\alpha, m, \tau) = (0.01, -1.0, 0.0)$



(b) Estimated values by only inverse Bayesian inference  $(\alpha, m, \tau) = (0.01, -1.0, 0.0)$

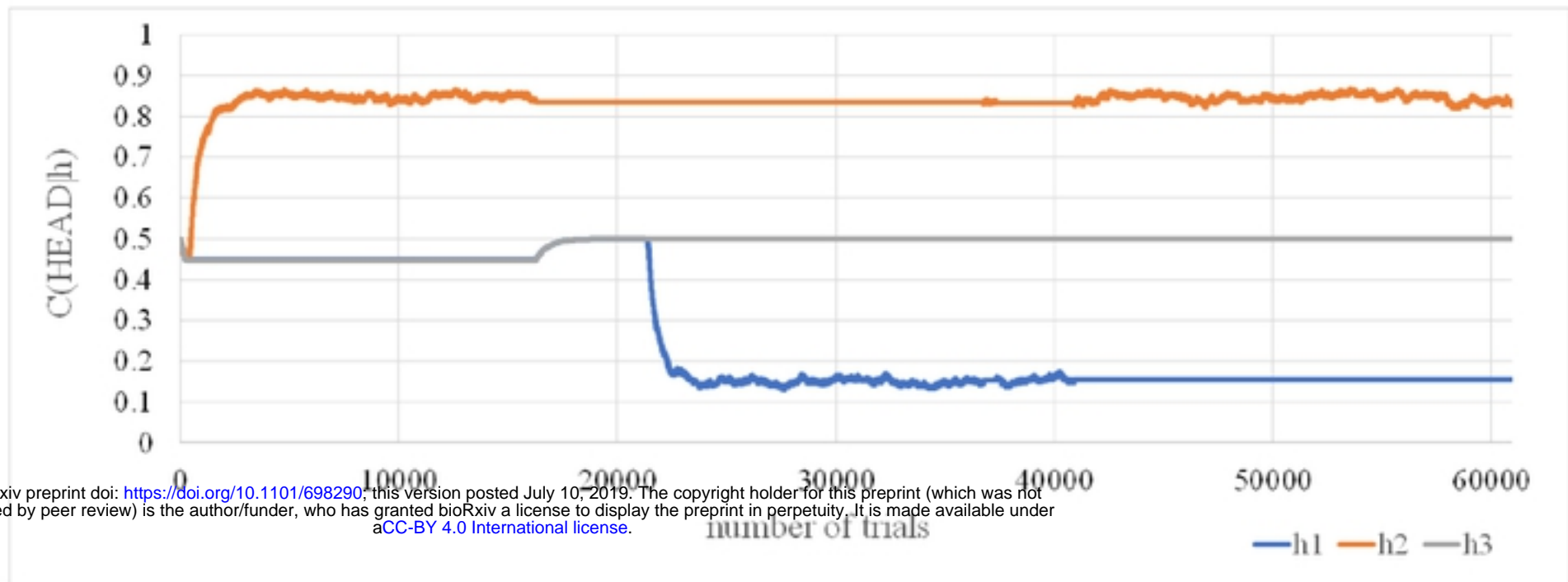


(c) Estimated values by EMA

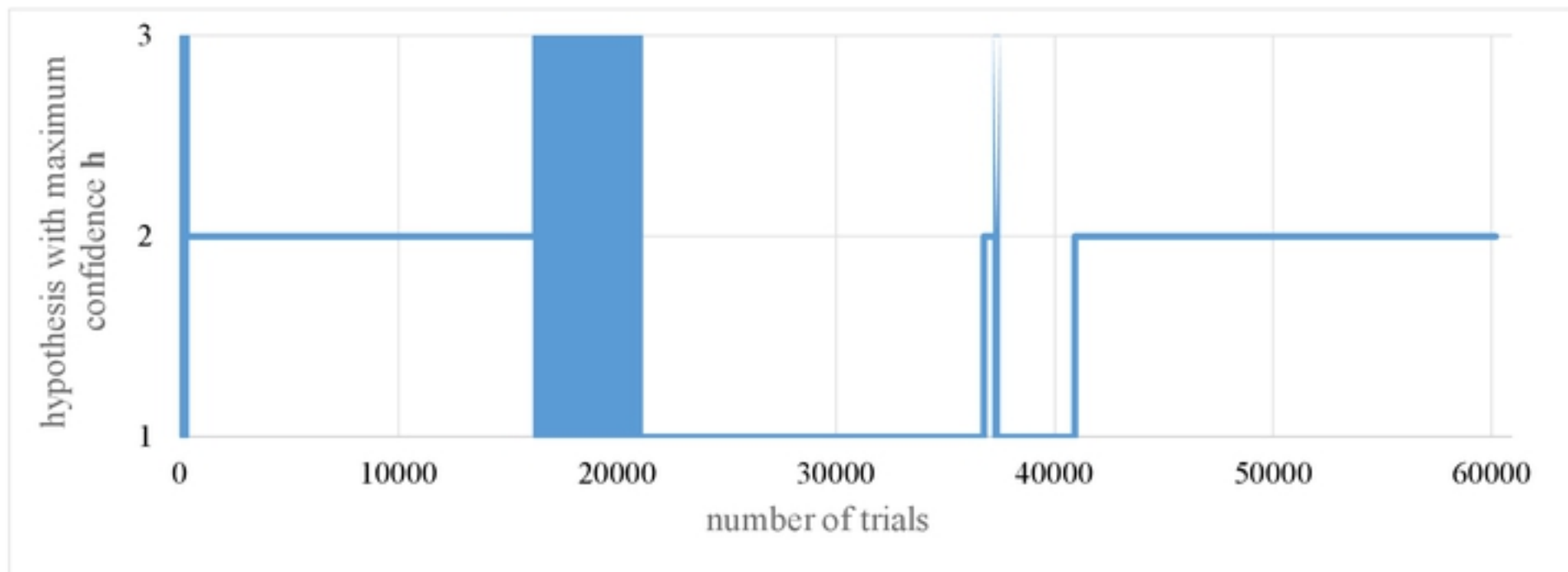
**Fig. 3** Time progress of the estimated values for the probability of head landing. The figure includes the correct probability.



bioRxiv preprint doi: <https://doi.org/10.1101/698290>; this version posted July 10, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



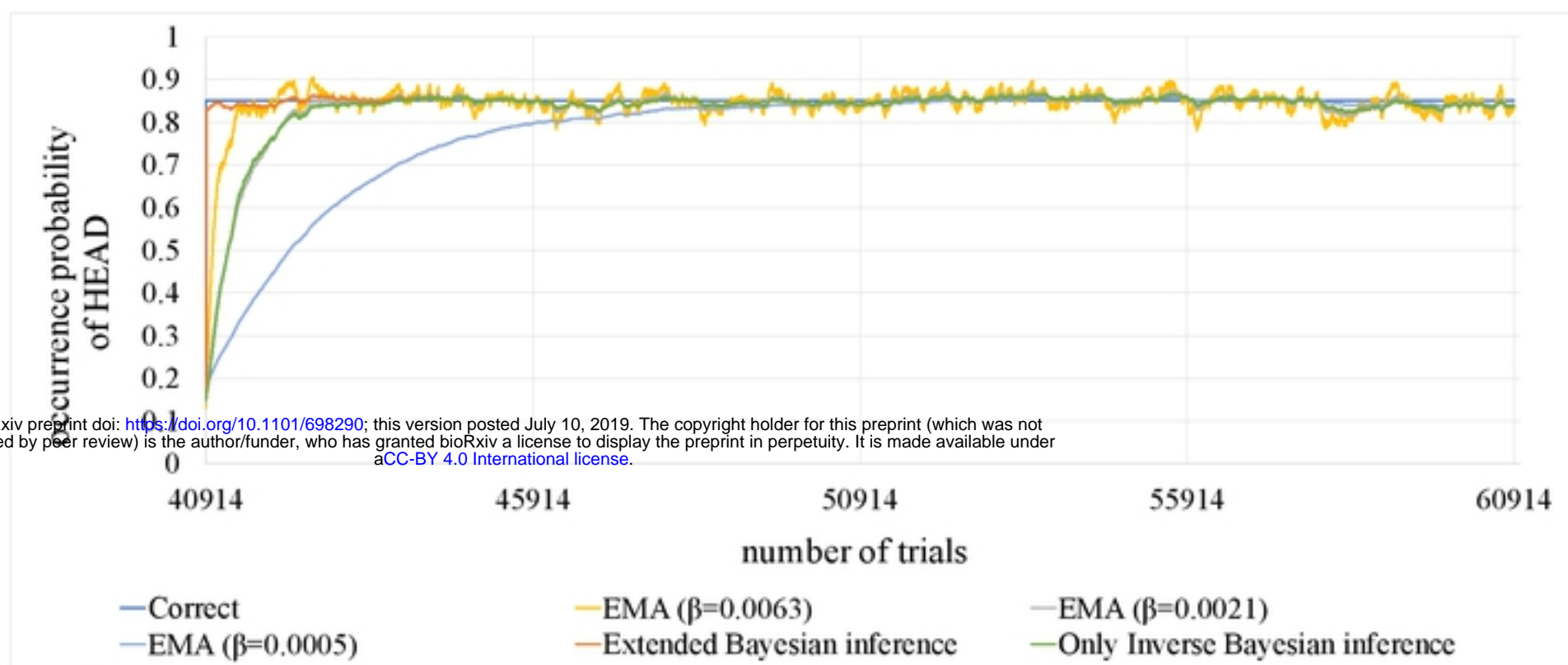
(a) Time progress of the probability of head landing for each hypothesis



(b) Time progress of the hypothesis with the greatest degree of confidence

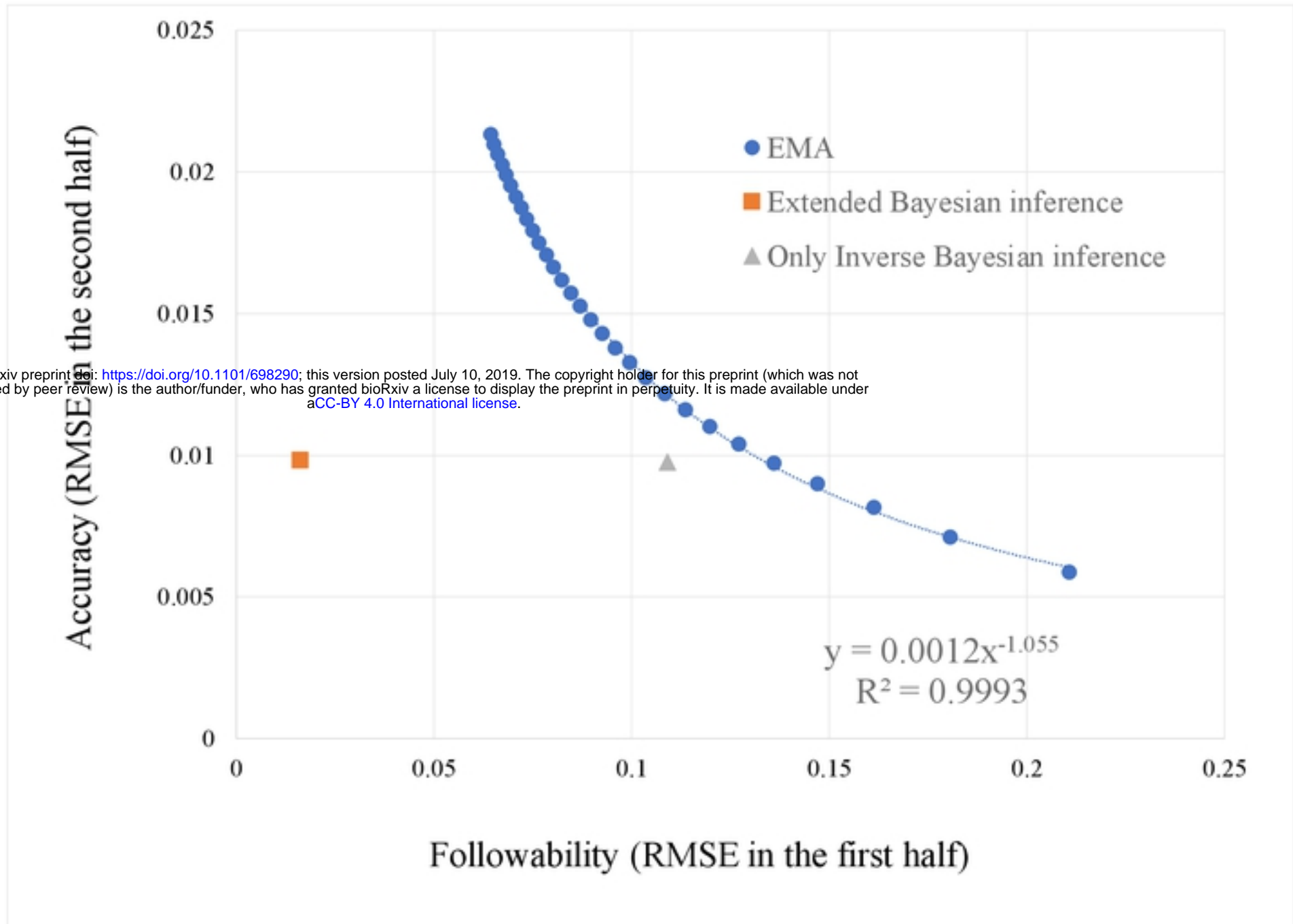
**Fig. 4** Internal state of extended Bayesian inference

bioRxiv preprint doi: <https://doi.org/10.1101/698290>; this version posted July 10, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



**Fig. 5** Time progress of the estimated values for the probability of head landing. The figure includes the correct probability.

bioRxiv preprint doi: <https://doi.org/10.1101/698290>; this version posted July 10, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



**Fig. 6** Relationship between the followability and the estimation accuracy in extended Bayesian inference, inverse Bayesian inference, and EMA estimations.