

Kruschke, J. K. & Denton, S. E. (2010). Backward blocking of relevance-indicating cues: evidence for locally Bayesian learning. In: C. J. Mitchell & M. E. Le Pelley (Eds.), *Attention and Associative Learning: From Brain to Behaviour*, pp. 273-304. Oxford, UK: Oxford University Press. ISBN 9780199550531

## Chapter 11

# Backward blocking of relevance-indicating cues: Evidence for locally Bayesian learning

John K. Kruschke and Stephen E. Denton

The phenomenon in associative learning called ‘blocking’ has been a central target for theories that aim to explain learning and attention. The training phases of the blocking procedure, described in detail below, can be run in backward order, yet still produce the blocking effect (e.g., Shanks, 1985; Dickinson and Burke, 1996; Kruschke and Blair, 2000). Explaining forward and backward blocking, along with other associative learning phenomena, has proven to be challenging. Some accounts of forward blocking include an attentional mechanism that modulates the influence of the cues on learning and/or responding (e.g., Kruschke, 2001; Mackintosh, 1975). Some accounts of backward blocking employ a Bayesian framework in which different combinations of associative weights are considered simultaneously, with more belief allocated to the combination that is most consistent with training items (e.g., Dayan and Kakade, 2001; Tenenbaum and Griffiths, 2003).

A theoretical framework that is able to combine the attentional and Bayesian approaches is called ‘locally Bayesian learning’ (Kruschke, 2006b). The framework is based on the idea that a learning system may consist of a sequence of subsystems in a feed-forward chain, each of which is a locally Bayesian learner. The argument for locally-learning layers was as follows. First, Bayesian learning is very attractive for explaining retrospective revaluation effects such as backward blocking, among many other phenomena (Chater, Tenenbaum, and Yuille, 2006). Second, globally Bayesian learning may also be unattractive for a number of reasons. In a large learning system, there are too many combinations of parameters to keep track of in a monolithic joint parameter space. Furthermore, many globally Bayesian models do not explain learning phenomena (such as ‘highlighting’, Kruschke, 2010) that depend on training order, because the models treat all training items as equally representative of

the world to be learned, regardless of when the items occurred. Finally, the level of analysis for theories of learning is arbitrary: Learning occurs simultaneously at the levels of neurons, brain regions, functional components, individuals, committees, institutions, and societies, all of which may be modeled (in principle, if not accurately) as Bayesian learners. Therefore, a system of locally Bayesian learning components may retain some attractions of Bayesian models while also implementing Bayesian learning in smaller, tractable parameter spaces.

The general framework for locally Bayesian learning has been instantiated in a particular two-layer model, wherein one layer learns how to allocate attention to cues, and a second layer learns how to associate attended cues with outcomes (Kruschke, 2006b). The model showed retrospective revaluation effects such as backward blocking while also showing the order-sensitive phenomenon of highlighting. The two-layer model had locally Bayesian learning in both of its layers. Bayesian learning was needed in the upper, outcome layer to produce effects such as backward blocking. Attentional shifting was needed to produce effects such as highlighting. But there was no phenomenon that demanded Bayesian learning in the lower, attentional layer. Bayesian learning was conducted in the lower layer merely for mechanistic consistency and as a demonstration of the more general framework for locally Bayesian learning.

The purpose of the present chapter is to report a novel learning design, the results of which suggest the need for Bayesian learning in the lower, attentional layer. In essence, the results show backward blocking of cues to attentional allocation, as distinct from traditional designs that show backward blocking of cues to response allocation. This pattern of results can be accommodated by a model that has layers of locally Bayesian learning.

## Background

### Forward and backward blocking

In the standard forward blocking procedure, a learner is initially trained with cases of cue A leading to outcome X, denoted  $A \rightarrow X$ . Subsequently, training continues with cases in which two cues, A and B, lead to the same outcome X, denoted  $A \cdot B \rightarrow X$ . It turns out that the strength of association from cue B to outcome X is weaker than if the initial training with A had not occurred. In other words, the learning of the association from B to X has been ‘blocked’, or attenuated, by the previous learning of the association from A to X. The phenomenon of blocking, first reported by Kamin (1968), challenges theories that base strength of association on merely the number of co-occurrences of cue and outcome (but cf. Miller and Matzel, 1988).

1 There are many types of explanations of blocking, but one family of explanations  
 2 posits a role for attention. The intuition is that when learning cases of  
 3  $A \cdot B \rightarrow X$ , the learner re-allocates attention away from the redundant cue B  
 4 because it is distracting resources away from the cue A that is already known to  
 5 generate the correct outcome. In other words, the learner learns to ignore cue  
 6 B. Some evidence for the attentional explanation comes from studies of learning  
 7 about cue B after blocking. If learners have learned to ignore B, then subsequent  
 8 learning about it should be retarded. This prediction has been  
 9 confirmed (e.g., Kruschke and Blair, 2000; Kruschke, 2005; LePelley, Beesley,  
 10 and Suret, 2007; Mackintosh and Turner, 1971; Mitchell, Harris, Westbrook,  
 11 and Griffiths, 2008). Other evidence for reduced attention to the blocked cue  
 12 comes from eye tracking experiments, in which it has been shown that gaze  
 13 duration is reduced for blocked cues (Kruschke, Kappenman, and Hetrick,  
 14 2005; Wills, Lavric, Croft, and Hodgson, 2007).

15 Backward blocking, as a training procedure, simply reverses the phases of  
 16 training in standard forward blocking. In other words, learners are first trained  
 17 with cases of  $A \cdot B \rightarrow X$ , and subsequently trained with cases of  $A \rightarrow X$ . It turns  
 18 out that the strength of association from B to X is again weakened by the  $A \rightarrow X$   
 19 training, even though it happened after the training with B, and even though B  
 20 never appeared in the subsequent training (Shanks, 1985). Thus, cue B seems  
 21 to have been retrospectively revalued in its absence. Backward blocking and  
 22 other retrospective revaluations are especially challenging to theories of associative  
 23 learning (for reviews see DeHouwer and Beckers, 2002; Dickinson,  
 24 2001).

25 The attentional theories that account for blocking do not account for backward  
 26 blocking. This failure is caused by the fact that the theories rely on the  
 27 presence of a cue to learn how much to attend to it. If a cue is absent, the  
 28 models do not change its attention strength.

29 One class of theories that accommodates backward blocking is Bayesian  
 30 models of association (e.g., Dayan and Kakade, 2001; Sobel, Tenenbaum, and  
 31 Gopnik, 2004; Tenenbaum and Griffiths, 2003).<sup>1</sup> These Bayesian theories posit  
 32 a set of associative hypotheses simultaneously entertained by the learner.  
 33 For illustration, suppose that the learner considers three hypotheses: (1) A  
 34 indicates X and B is irrelevant, (2) B indicates X and A is irrelevant, and

---

<sup>1</sup> A class of non-Bayesian theories of backward blocking asserts that absent-but-expected cues acquire reduced associations when the expected outcome is present. For absent-but-expected cues, either the learning rate or encoding of the absent-cue is negative, resulting in a loss of association (Dickinson and Burke, 1996; Markman, 1989; Tassoni, 1995; Van Hamme and Wasserman, 1994).

(3) either A or B indicate X. The three hypotheses are mutually exclusive, and exhaust the space of possibilities for this particular learner. After seeing the initial cases of  $A \cdot B \rightarrow X$ , all three hypotheses have some credibility. But after seeing cases of  $A \rightarrow X$ , the first and third hypotheses gain credibility, because they are both consistent with the additional training. Because the set of hypotheses are mutually exclusive and exhaustive, when the first and third hypotheses gain credibility, the second hypothesis loses credibility. Therefore, across all the hypotheses, there is reduced strength of belief that B indicates X. For a detailed tutorial, see the discussion of Bayesian associative models by Kruschke (2008).

Bayesian approaches to learning are attractive for a variety of other reasons. Bayesian models are not limited to associative formalisms, but can instead incorporate complex structural representations into the hypothesis space. Bayesian models merely assume that the learner executes normatively correct learning (i.e., uses Bayes' rule) on whatever formal hypothesis space is posited. This representational flexibility allows Bayesian models to be applied to situations from learning by neurons (Deneve, 2008) to learning of language (Xu and Tenenbaum, 2007).

## Locally Bayesian learning

Bayesian formalisms are very attractive as theories of learning, but implementing them can be difficult because of their computational complexity. In principle, Bayesian systems need to keep track of the credibility of every possible hypothesis. In hypothesis spaces with many parameters, such as numerous associative weights, the system needs to keep track of the credibility of every possible combination of parameter values in a high-dimensionality joint parameter space. In some simple models, this can be done exactly and easily because the entire distribution across beliefs can be summarized by a simple function such as a multivariate normal distribution. In more complicated models, the infinite space of hypotheses can be represented by a large but finite random sample of representative hypotheses. In these 'particle filter' approximations, as new cue outcome cases are experienced, the representative hypotheses are resampled to reflect the new experience. These methods are not trivial, however, and Bayesian computation can be quite challenging in large hypothesis spaces.

One way to attack the problem is to recognize that learning happens at many levels of analysis simultaneously. Neurons learn, brain regions learn, functional components of mind learn, individual people learn, teams of people learn, and so on. Indeed, even more microscopic and macroscopic systems may learn. Any of those levels of analysis may be amenable to description as Bayesian learning.

1 Therefore, it is at least plausible that component processes of the mind may be  
 2 describable as Bayesian learners, because the components have a tractable  
 3 hypothesis space. Whether or not the system as a whole is Bayesian depends on  
 4 how the components interact. A scheme for interaction of hierarchically organ-  
 5 ized, locally Bayesian learners was described by Kruschke (2006b).

6 The general framework for locally Bayesian learning assumes that there are  
 7 component processes in a hierarchy from stimulus encoding to response gen-  
 8 eration. Each functional component in the hierarchical chain takes its local  
 9 input representation, transforms it, and delivers its local output representation  
 10 to the next component in the chain. The transformation is parameterized,  
 11 meaning that the exact quantitative behaviour of the local transformation  
 12 depends on its parameter values. As a simple example, consider a linear trans-  
 13 formation  $y = mx + b$ , where  $x$  is the input and  $y$  is the output. The slope  $m$  and  
 14 the intercept  $b$  are parameters that govern the quantitative value of the output  
 15 for any given input. The parameters are learned as exemplary  $\langle x, y \rangle$  values are  
 16 experienced. The learning of parameters is Bayesian reallocation of credibility  
 17 to combinations of parameter values that are most consistent with the incom-  
 18 ing stimuli and the target response. For example, if the component system has  
 19 experienced  $\langle x, y \rangle$  pairs such as  $\langle 1, 2.01 \rangle$ ,  $\langle 2, 3.99 \rangle$ , and  $\langle 3, 6.01 \rangle$ , then it will  
 20 allocate strong credibility to  $m = 2$  and  $b = 0$ , i.e.,  $y = 2x + 0$ , and weaker credi-  
 21 bility to other parameter values.

22 The challenge for such a framework is determining the target response  
 23 for interior components, because the outside world only indicates the  
 24 target response for the final component that generates an overt response.  
 25 Formally, the world only supplies the exterior stimulus  $x$  and the exterior  
 26 target response  $y$ . For a component buried in the interior of the hierarchy, the  
 27 component's input can be computed by propagating the exterior stimulus  $x$   
 28 up through the transformations leading into the component. The compo-  
 29 nent's target output, however, is not obvious, because the transformations  
 30 feed forward from  $x$  to  $y$ , not the opposite direction.

31 A heuristic for determining interior targets is as follows. The target for a  
 32 component should be whatever is the input to the next component that maxi-  
 33 mizes the probability of achieving the target of that next component. Start at  
 34 the final component and determine the input to that component that would  
 35 maximize the probability of achieving the known exterior target. Use that  
 36 input as the target for the penultimate component. Given that target for the  
 37 penultimate component, determine the input to the penultimate component  
 38 that would maximize the probability of achieving its target. Use that input as  
 39 the target for the preceding component. Continue this process down the hier-  
 40 archy of components until every component has a target.

1 Notice that the targets selected in this manner are the inputs that are most  
 2 consistent with the component's current beliefs. In other words, at any given  
 3 time, the component gives strongest credibility to some particular combina-  
 4 tions of parameter values. That component finds the candidate input that  
 5 would be most consistent with the parameter values that the component cur-  
 6 rently believes. That best candidate input is declared to be the target for the  
 7 preceding layer. In this way, each component is teaching the preceding com-  
 8 ponent to 'tell me what I want to hear' based on its current beliefs.

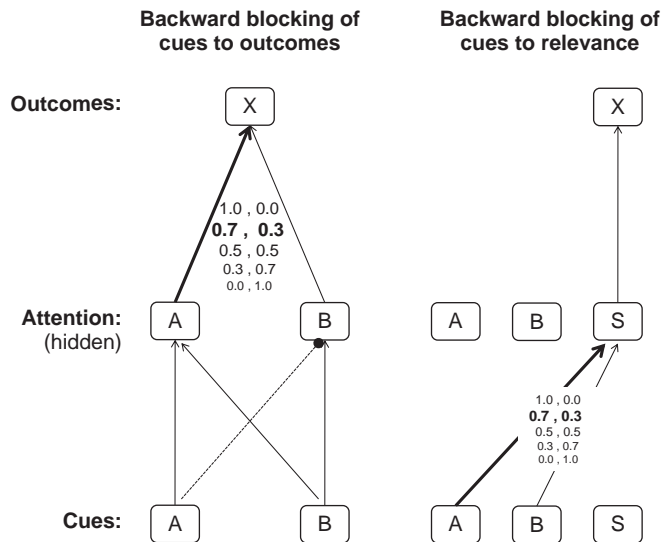
9 A crucial consequence of choosing interior targets this way is that the tem-  
 10 poral order of data has an effect on what is learned internally. The reason is  
 11 that the interior targets are determined by the components' current beliefs,  
 12 which depend on the set of data experienced so far. The interior targets gener-  
 13 ated for any particular exterior input-output pair depend on what has been  
 14 previously learned.

15 When a component transformation has a target, along with its input, then  
 16 the parameters of the component transformation are adjusted via Bayesian  
 17 learning. Parameter values that are consistent with the input and target are  
 18 deemed more credible. There are various temporal dynamics that could be  
 19 used for interleaving parameter learning and target determination (Kruschke,  
 20 2006b, p. 683, footnote 2).

## 21 **Attention in locally Bayesian learning**

22 Locally Bayesian learning is a general framework for models. One particular  
 23 application is to attention in associative learning. In this application, the first  
 24 component learns how to allocate attention across cues. Attention is an inte-  
 25 rior value, and is not explicitly provided by exterior data. The second compo-  
 26 nent learns how to associate attentionally-filtered cues to outcomes. Both the  
 27 attentional allocations and the outcome associations are acquired via locally  
 28 Bayesian learning.

29 Previous work demonstrated the usefulness of this architecture (Kruschke,  
 30 2006a, 2006b). In particular, locally Bayesian learning of the output associa-  
 31 tions allows the system to accommodate backward blocking and other retro-  
 32 spective revaluation phenomena. Figure 11.1 shows the basic elements of  
 33 locally Bayesian learning applied to attentional learning. The lowest nodes  
 34 encode the cues, the middle nodes represent the attentionally filtered cues,  
 35 and the upper nodes represent the (anticipated) outcome. The left side of  
 36 Figure 11.1 illustrates what happens when the system is trained in the back-  
 37 ward blocking procedure. In this procedure, the model first experiences cases  
 38 of  $A \cdot B \rightarrow X$ , and then experiences cases of  $A \rightarrow X$ . The upper-left part of the  
 39 figure shows a variety of output-weight combinations that the model finds



**Fig. 11.1** Schematic of locally Bayesian learning applied to backward blocking. *Left:* Backward blocking of cues to outcomes. *Right:* Backward blocking of context cues to relevance. The pairs of numbers suggest some of the weight combinations that the model finds credible.

1 credible after training. The size of the font suggests how strongly the model  
 2 believes in the weight combination. Thus, the model gives greatest credibility  
 3 to an associative weight of 0.7 from A and an associative weight of 0.3 from B.  
 4 But the model also gives some modest credibility to other weight combina-  
 5 tions that are reasonably consistent with the training items. The layer of weights  
 6 leading into the attentional gates also has a distribution of credibility across  
 7 possible weight combinations.

8 Attentional shifting, and the temporal dependency of the scheme for select-  
 9 ing interior targets, also allows the system to accommodate the highlighting  
 10 phenomenon (Kruschke, 2010). Highlighting is a trial-order effect that is par-  
 11 ticularly vexing for Bayesian models that treat all trials as equally representa-  
 12 tive of the world. The temporal dynamics of locally Bayesian learning also let  
 13 the system show other trial-order effects, such as stronger forward blocking  
 14 than backward blocking. A variety of other phenomena are addressed by the  
 15 model.

16 The ability of the model to show backward blocking (and other effects)  
 17 relied on Bayesian learning in the upper associative layer. But none of the  
 18 phenomena considered by Kruschke (2006b) demanded Bayesian learning  
 19 in the lower, attentional layer. In principle, the attentional layer could

1 have been a non-Bayesian learner, and all the same effects could have been  
2 produced.

3 The primary goal of the present article is to report a phenomenon that does  
4 suggest the need for Bayesian learning of attentional allocation. The argument  
5 goes like this: Backward blocking is naturally accounted for by Bayesian learn-  
6 ing. If it can be shown that there is backward blocking in the learning of atten-  
7 tional allocation, then one candidate explanation is that there is Bayesian  
8 learning of attentional allocation.

9 The right side of Figure 11.1 illustrates what is meant by backward blocking  
10 of attentional allocation. The complete design of the cue–outcome combina-  
11 tions will be described later, but the gist is provided here. Contrary to the  
12 standard design, cues A and B are not themselves diagnostic of the outcome:  
13 Across trials, the outcome occurs just as often when cues A and B are absent as  
14 when they are present. A different cue, labeled S in Figure 11.1, is perfectly  
15 predictive of the outcome, but only when cue A or cue B is present. When cues  
16 A and B are absent, cue S is not predictive. Thus, cues A and B do not indicate  
17 what outcome to anticipate, but they do indicate what other cue is relevant. In  
18 essence, cues A and B indicate what other cues should be attended.

19 The goal of the experiment reported below is demonstrate backward block-  
20 ing of such cues to relevance. In the first part of training, cues A and B are both  
21 present whenever S is diagnostic. Therefore the model should learn to attend  
22 to S in the presence of A and B. Later in training, cue A is present without cue  
23 B whenever S is diagnostic. If there is backward blocking of B, then the asso-  
24 ciation from B to S should be weakened. The weight pairs in the lower-right  
25 side of Figure 11.1 are intended to suggest the credible weight combinations  
26 after backward blocking of B as a cue to relevance.

27 The purpose of the experiment and modelling presented here is not to  
28 rule out other explanations or disconfirm other models. Instead, the goal is  
29 to bolster an assumption of locally Bayesian learning applied to attentional  
30 learning, a model which was already shown to address a spectrum of phenom-  
31 ena (Kruschke, 2006b). There are surely other models that can accommodate  
32 the data from the one new experiment presented here, but for a model to com-  
33 pete with locally Bayesian learning, the candidate model should also accom-  
34 modate the spectrum of other phenomena addressed by locally Bayesian  
35 learning.

36 The new experiment presented here is offered merely as a suggestive proof of  
37 concept. Future experiments will be necessary to generalize the conditions  
38 under which the effects are observed, and to rule out alternative explanations.  
39 Nevertheless, it is hoped that the experiment and modeling may provoke inter-  
40 esting new ideas and research. In particular, the experiment presented here



- 1 might not have been invented were it not for the implications of locally
- 2 Bayesian learning applied to attentional learning.

### 3 **Experiment: Blocking and backward blocking of** 4 **cues to relevance**

5 Different cues can be relevant in different contexts. For example, when driving  
6 an automobile, the colour of the stoplight is relevant to the decision to stop or  
7 go, but when walking, the colour of the pedestrian signal is relevant to the deci-  
8 sion to stop or go. In general, the cues in an environment can suggest which  
9 sources of information are relevant for determining a response. For example,  
10 the cue of having a steering wheel in your hands does not tell you whether to  
11 stop or go, but it does indicate that you should attend to stoplights, which will  
12 indicate whether you should stop or go.

13 Cues that are indirectly informative, such as the steering wheel in the previ-  
14 ous example, are a type of context cue. The term 'context' has no generally  
15 accepted technical definition. Sometimes context refers to stimulus attributes  
16 that are spatially ambient (not focal) or temporally extended (i.e. tonic not  
17 phasic). Other times, context refers to cues that can be focal and phasic but  
18 that are not directly correlated with outcomes. This latter character is empha-  
19 sized here. There has not been a vast amount of previous research into the role  
20 of context in learned attentional allocation, but several lines of work have indi-  
21 cated that people can learn about 'irrelevant' context as a cue to Attention  
22 (e.g., Chun, 2000; Nelson, 2002; Rosas, Callejas-Aguilera, Ramos-Alvarez, and  
23 Abad, 2006; Yang and Lewandowsky, 2003).

24 The present experiment is aimed at demonstrating forward and backward  
25 blocking of contextual cues to relevance. Continuing the example from driving  
26 discussed earlier, the idea is that people first experience steering wheels as a cue  
27 to attend to stoplights, and later people experience steering wheels along with  
28 a newly installed car stereo as a compound cue to attend to stoplights. The  
29 association from car stereo to attention may be blocked because of the previ-  
30 ously learned association from steering wheel to attention.

31 To test this idea, we conducted a series of experiments in which all the cues  
32 were simple words on a computer screen, such as 'radio' and 'ocean'. Some  
33 words were perfectly correlated with the correct key to press. Other words had  
34 zero correlation with the correct response key, but were perfect indicators of  
35 which other words on the screen were relevant to the choice of response key.  
36 For example, suppose that people have learned that radio indicates key X and  
37 ocean indicates key Y. Subsequently, both 'radio' and 'ocean' appear simulta-  
38 neously. Should the response be X or Y? The conflict is resolved by a third

1 word on the screen, e.g., ‘queen’, which indicates to attend to ‘radio’. Our  
 2 experiments revealed that it was difficult for people to learn this sort of contex-  
 3 tual dependency in a brief (< 20 min) experiment when the cues had no struc-  
 4 tural or semantic indicators of which were context cues and which were  
 5 response cues. For the few people who could learn such structures quickly, we  
 6 observed signs of blocking and backward blocking of cues to relevance. But it  
 7 was unsatisfying to base conclusions on a small subset of participants.  
 8 Presumably, all people could learn such structures if given enough practice,  
 9 but before subjecting people to endurance training, we explored other stimu-  
 10 lus arrangements.

11 In order to facilitate learning, and for purposes of an experiment that can act  
 12 as a proof of concept, we set up a cue arrangement in which it was natural to  
 13 think of some cues being indicators of responses, and other cues being indica-  
 14 tors of which response cues to attend to. The learners were instructed that they  
 15 were to diagnose the fictitious disease associated with symptoms, but only  
 16 indirectly, by learning which medical specialist knew about which symptoms.  
 17 For example, a patient might have the symptoms heartburn and myalgia, for  
 18 which Specialist 1 says the patient has disease F, but Specialist 2 says the patient  
 19 has disease J. After the learner makes a guess about the disease, corrective feed-  
 20 back indicates that it was disease J, thereby implying that Specialist 2 knows  
 21 about these symptoms. The learner should therefore learn that when symp-  
 22 toms heartburn and myalgia occur, s/he should attend to Specialist 2, and give  
 23 the response stated by Specialist 2. Importantly, there is no correlation between  
 24 symptoms and diseases across trials. For example, half the time that heartburn  
 25 and myalgia occur, the correct disease is F, but half the time the correct disease  
 26 is J. Specialist 2 always indicates the correct diagnosis, however, for these  
 27 symptoms.

28 The cues that indicate the correct overt response are here called ‘response  
 29 cues’. In the present scenario, the medical specialists are the response cues.  
 30 Other cues that indicate which response cues to attend to are here called ‘con-  
 31 text cues’. In the present scenario, the symptoms are the context cues. These  
 32 appellations, i.e., response cue versus context cue, are potentially misleading.  
 33 On the one hand, the so-called context cues do indicate a response, but that  
 34 response is an essentially covert re-allocation of attention (which may or may  
 35 not have overt signatures such as eye movements or other orienting responses).  
 36 On the other hand, the so-called response cues need not be known in advance  
 37 to be indicators of overt responses; the response cues might serve as context to  
 38 other cues. Despite these infelicities of nomenclature, a key aspect of intuiti-  
 39 vely contextual information *is* captured by the ‘context’ cues: They are not

- 1 directly predictive of the correct overt response. The context cues only indicate
- 2 which response cues to attend to, and the response cues, in turn, indicate
- 3 which overt response to make.

## 4 Method

### 5 Design

6 Table 11.1 shows the design components of the experiment. Each phase has a  
 7 different arrangement of context cues. In the Single Context phase, a single  
 8 context cue is present with two Specialists who give conflicting diagnoses. For  
 9 example, a trial might consist of  $A_{S1}S1_F S2_J \rightarrow F$ , which means that context cue  
 10  $A_{S1}$  occurred with specialists  $S1_F$  and  $S2_J$ , with correct diagnosis F. The sub-  
 11 scripts on the cues denote what the cue is intended to indicate. The notation  
 12  $A_{S1}$ , for example, means that cue A indicates specialist S1. This correspond-  
 13 ence had to be learned, however. Notice that across the eight cases of the Single  
 14 Context phase, the context cue A occurs twice with outcome F and twice with  
 15 outcome J. Hence the context cue is uncorrelated with the correct diagnosis. In  
 16 all cases of the Single Context phase, only a single context cue occurs.

17 In the Redundant Context phase, some cases had two context cues. In  
 18 particular, context cues A and B occurred together, and context cues C and D  
 19 occurred together. As these context cues never occurred separately, they

**Table 11.1** Components of the experiment design

Phase	Items	
Single Context	$A_{S1} S1_F S2_J \rightarrow F$	$A_{S1} S1_F S3_J \rightarrow F$
	$A_{S1} S1_J S2_F \rightarrow J$	$A_{S1} S1_J S3_F \rightarrow J$
	$E_{S3} S3_F S2_J \rightarrow F$	$E_{S3} S3_F S1_J \rightarrow F$
	$E_{S3} S3_J S2_F \rightarrow J$	$E_{S3} S3_J S1_F \rightarrow J$
Redundant Context	$A_{S1} B_{S1} S1_F S2_J \rightarrow F$	$A_{S1} B_{S1} S1_F S3_J \rightarrow F$
	$A_{S1} B_{S1} S1_J S2_F \rightarrow J$	$A_{S1} B_{S1} S1_J S3_F \rightarrow J$
	$C_{S2} D_{S2} S2_F S1_J \rightarrow F$	$C_{S2} D_{S2} S2_F S3_J \rightarrow F$
	$C_{S2} D_{S2} S2_J S1_F \rightarrow J$	$C_{S2} D_{S2} S2_J S3_F \rightarrow J$
	$E_{S3} S3_F S2_J \rightarrow F$	$E_{S3} S3_F S1_J \rightarrow F$
	$E_{S3} S3_J S2_F \rightarrow J$	$E_{S3} S3_J S1_F \rightarrow J$
Test: Conflicting Context	$A_{S1} C_{S2} S1_F S2_J \rightarrow ?$	$A_{S1} C_{S2} S1_J S2_F \rightarrow ?$
	$A_{S1} D_{S2} S1_F S2_J \rightarrow ?$	$A_{S1} D_{S2} S1_J S2_F \rightarrow ?$
	$B_{S1} C_{S2} S1_F S2_J \rightarrow ?$	$B_{S1} C_{S2} S1_J S2_F \rightarrow ?$
	$B_{S1} D_{S2} S1_F S2_J \rightarrow ?$	$B_{S1} D_{S2} S1_J S2_F \rightarrow ?$

Note: An item is shown in the format, Cues  $\rightarrow$  Correct Response. The subscripts on the cues indicate the design's intended correspondence from that cue. Cues  $A_{S1}$ ,  $B_{S1}$ ,  $C_{S2}$ ,  $D_{S2}$ , and  $E_{S3}$  are symptoms. Cues  $S1_F$ ,  $S1_J$ ,  $S2_F$ ,  $S2_J$ ,  $S3_F$ , and  $S3_J$  are specialists. Responses F and J are disease labels.

1 are called ‘redundant’. As in the single-context phase, all context cues are  
 2 uncorrelated with the correct diagnosis.

3 Notice that when the single-context phase occurs before the redundant-  
 4 context phase, the context cues instantiate a standard blocking sequence.  
 5 People learn first that context cue A indicates specialist 1, and then people see  
 6 that context cue A with context cue B also indicate specialist 1. There may be  
 7 blocking of learning about the redundant cue B. When the redundant-context  
 8 phase occurs before the single-context phase, this constitutes a backward  
 9 blocking design for the context cues.

10 Blocking is assessed in the test phase, wherein conflicting cues appear  
 11 together. In particular, cue  $A_{S1}$  appears with either  $C_{S2}$  or  $D_{S2}$ , and cue  $B_{S1}$   
 12 appears with either  $C_{S2}$  or  $D_{S2}$ . If context cue  $B_{S1}$  is blocked, then when it is  
 13 paired with either  $C_{S2}$  or  $D_{S2}$ , the response appropriate to Specialist 2 should  
 14 be given, in preference over the response appropriate to Specialist 1. As a par-  
 15 tial test that Specialist 2 is not favored generically whenever conflict occurs, the  
 16 cases in which  $A_{S1}$  appears with either  $C_{S2}$  or  $D_{S2}$  should favor Specialist 1, or  
 17 at least not Specialist 2. This issue is addressed further in the discussion after  
 18 the results are reported.

19 There were three different ‘epochs’ of the experiment. Each epoch had the  
 20 phases shown in Table 11.1, in different orderings or with different types of  
 21 feedback. The first epoch was a backward blocking sequence (redundant-  
 22 context phase before single-context phase), but with the correct specialist  
 23 explicitly and directly indicated by corrective feedback. This design consti-  
 24 tuted a replication of the backward blocking design of Kruschke and Blair  
 25 (2000, Experiment 3), except that the present procedure used a shorter-  
 26 duration training to criterion(details of which are described below). The  
 27 hope was that by replicating a previous design known to produce backward  
 28 blocking, we would observe the effect here too.

29 The second epoch was a forward blocking sequence(single-context phase  
 30 before redundant-context phase), with the correct specialist indicated only  
 31 indirectly and implicitly via the correct diagnosis. In other words, this epoch  
 32 was trained as shown in Table 11.1. New symptoms were used in all epochs, so  
 33 that novel learning was involved. As forward blocking tends to be more robust  
 34 than backward blocking (e.g., Kruschke and Blair, 2000; Shanks, 1985), it was  
 35 hoped that we would be able to observe forward blocking even in this complex,  
 36 indirect-feedback situation. This sort of procedure has never been reported  
 37 before, as far as we know.

38 The third epoch was a backward blocking sequence (redundant context  
 39 before single context), with the correct specialist indicated indirectly via  
 40 the correct diagnosis. In other words, it was just like the second epoch, but

1 with the phases of training reversed. This epoch constituted the main focus of  
 2 the experiment. We hoped to observe backward blocking of context cues.  
 3 Notice that the stimulus–outcome structure of Epoch 3 is identical to that of  
 4 Epoch 1; only the response and feedback are different between the two  
 5 epochs.

## 6 Procedure

7 Table 11.1 shows one ‘block’ of trials in each phase. Within each block, the  
 8 cases were presented in a random order. Training continued in each phase  
 9 until accuracy in a block exceeded 87% correct, meaning at least 7 correct in  
 10 8-trial blocks or at least 11 correct in 12-trial blocks. The maximum number of  
 11 blocks allowed in each phase was 8, at which point training progressed seam-  
 12 lessly to the next phase. Each test block was repeated twice. The entire experi-  
 13 ment took approximately 20 minutes or less.

14 All trials progressed in a seamless series. There were no pauses or markers  
 15 between phases. When a new epoch began, the novel symptoms appeared. The  
 16 response prompt for each trial indicated whether the learner was to guess the  
 17 correct specialist (in epoch 1) or the correct diagnosis (in epochs 2 and 3).

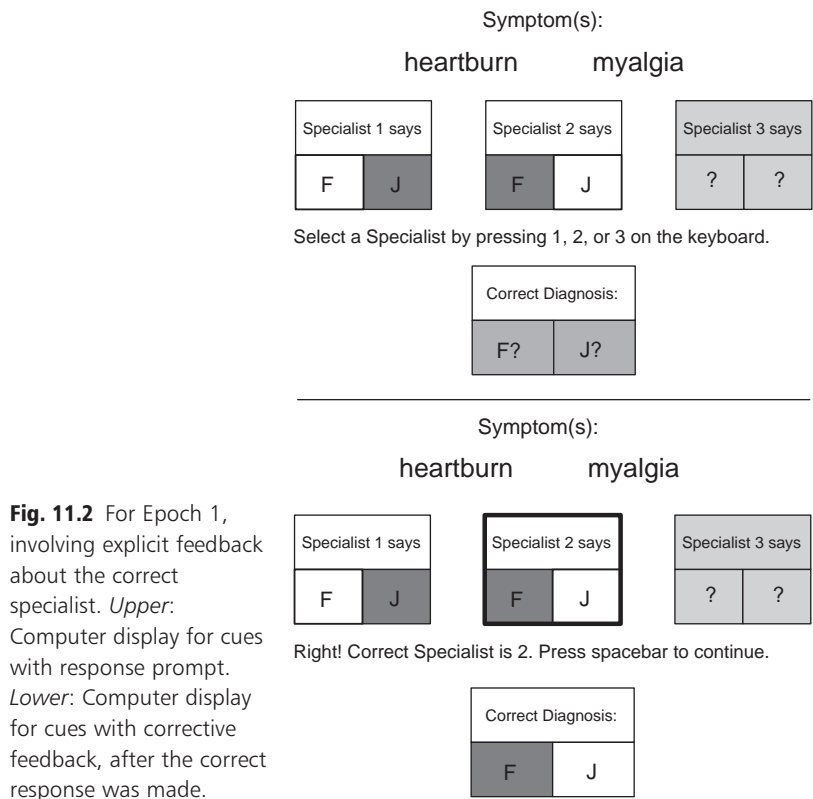
## 18 Stimuli

19 Figures 11.2 and 11.3 show screen shots of the stimuli. Figure 11.2 shows a  
 20 prompt and feedback screen from the first epoch, in which the correct special-  
 21 ist is explicitly and directly trained. The prompt asks the learner to press one of  
 22 the keys 1, 2, or 3, and the feedback states the correct specialist and highlights  
 23 the correct specialist with a heavy outline. The correct diagnosis is also indi-  
 24 cated by a white (instead of grey) background. Figure 11.3 shows a prompt and  
 25 feedback screen from the second and third epochs, in which the correct spe-  
 26 cialist is only indirectly trained via the diagnosis. Notice that the correct spe-  
 27 cialist is not indicated; only the correct diagnosis is shown. (It is only a random  
 28 coincidence that Figures 2 and 3 both show Specialist 3 without a diagnosis.  
 29 Across trials, the ‘missing’ specialist was counterbalanced, as indicated in the  
 30 design of Table 11.1.)

## 31 Results

### 32 Participants

33 Participants volunteered for partial credit in introductory psychology courses  
 34 at Indiana University. This subject pool has a median age of approximately 19  
 35 years, and is about 50–60% female. Procedures for protection of human sub-  
 36 jects were approved by the local Institutional Review Board. There were 188  
 37 participants.



1 Learning criterion

2 For purposes of data analysis, epochs were excluded if accuracy did not achieve  
3 at least 58% in both training phases by the final block or training. The criterion  
4 was selected arbitrarily as a compromise between excluding too many subjects  
5 and including too many poor learners. Results did not change in any qualita-  
6 tive way with different criteria. The criterion resulted in 180, 122, and 131  
7 subjects (out of 188) included in each of epochs 1, 2, and 3, respectively. Many  
8 or most of the excluded subjects appeared to have been unmotivated to learn,  
9 as their response times were on the order of 200 ms. or less, which indicates  
10 pressing a key as quickly as a stimulus appears without processing its  
11 attributes.

12 Choice in the test phase

13 Table 11.2 shows the choice preferences in the test phase, collapsed across  
14 participants. All three epochs show evidence of blocking, i.e., for the BC or  
15 BD tests, the response tends to be consistent with C or D more than with B.

Symptom(s):  
insomnia    bloating

Specialist 1 says		Specialist 2 says		Specialist 3 says	
F	J	F	J	?	?

Select a Diagnosis by pressing F or J on the keyboard.

Correct Diagnosis:	
F?	J?

Symptom(s):  
insomnia    bloating

Specialist 1 says		Specialist 2 says		Specialist 3 says	
F	J	F	J	?	?

Right! Correct Diagnosis is F. Press spacebar to continue.

Correct Diagnosis:	
F	J

**Fig. 11.3** For Epochs 2 and 3, involving feedback only about the correct diagnosis. *Upper:* Computer display for cues with response prompt. *Lower:* Computer display for cues with corrective feedback, after the correct response was made.

**Table 11.2** Test phase response percentages, collapsed across subjects

Test cues	Response consistent with A/B or C/D					
	Epoch 1 Backward Blocking Explicit Specialist		Epoch 2 Forward Blocking Indirect Feedback		Epoch 3 Backward Blocking Indirect Feedback	
	A/B	C/D	A/B	C/D	A/B	C/D
AC or AD	54.5	45.5	59.8	40.2	52.3	47.8
BC or BD	42.1	57.9	44.5	55.5	45.3	54.7

- 1 The magnitude of the preference is weak, however. For example, the magni-
- 2 tude of backward blocking in Epoch 1 is notably weaker than that reported by
- 3 Kruschke and Blair (2000).
- 4 We can only speculate as to why the blocking effect is so weak, but presum-
- 5 ably it is because of the complex stimulus display and distraction by the disease
- 6 diagnoses that were irrelevant in Epoch 1. Subsequent epochs also show weak

1 magnitudes of blocking, presumably because of the difficulty of inferring the  
 2 correct specialist indirectly from the diagnosis. We will say more about the  
 3 magnitude of blocking in the discussion after the statistical analysis.

#### 4 Bayesian statistical analysis

5 The data were analyzed using Bayesian methods. In a Bayesian analysis, a  
 6 descriptive model of the data is defined, and the parameter values of the model  
 7 are estimated. The Bayesian analysis yields an entire posterior distribution of  
 8 parameter values, not merely a single best-fitting parameter value. One reason  
 9 to prefer Bayesian methods over traditional null hypothesis significance testing  
 10 (NHST) is that the Bayesian analysis yields an explicit distribution regarding  
 11 the believability of various underlying choice probabilities, given the experi-  
 12 ment data. Another reason to prefer a Bayesian approach is that individual  
 13 differences are explicitly modelled and taken into account. In the following  
 14 paragraphs, the model is first defined, followed by a description of how the  
 15 posterior distribution was generated, followed, finally, by a description of the  
 16 posterior distribution itself.

17 In the test phase of any epoch, each participant saw the two context types  
 18 (i.e., either AC/AD or BC/BD) eight times, because there were two repetitions  
 19 of the test block in Table 1. For each test type, the  $i^{th}$  participant's 8 responses  
 20 to that type were modeled as a random sample from a binomial distribution  
 21 having underlying probability  $\theta_{A/B,i}$  of selecting the response consistent with  
 22 the A or B cue, and probability  $\theta_{C/D,i} = 1 - \theta_{A/B,i}$  of selecting the response con-  
 23 sistent with the C or D cue.

24 The individuals' probabilities,  $\theta_{C/D,i}$ , were modeled as a random selection  
 25 from an overarching beta distribution that had (1) a parameter  $\mu_{C/D}$  that  
 26 specifies the central tendency of the group, and (2) a parameter  $\kappa$  that specifies  
 27 how tightly the individuals were clustered around that central tendency.  
 28 (In detail, the two 'shape parameters' of the beta distribution for  $\mu_{C/D}$  were  
 29  $a = \mu_{C/D}\kappa + 1$  and  $b = (1 - \mu_{C/D})\kappa + 1$ .)

30 The primary goal of the analysis is to generate a posterior estimate of the  
 31  $\mu_{C/D}$  parameter for each test type. The  $\mu_{C/D}$  parameter represents the overall  
 32 response propensity for the type of context. The value of  $\mu_{C/D}$  can range  
 33 between 0 and 1, and can be thought of as the underlying probability of  
 34 responding consistently with the C/D cues. When  $\mu_{C/D} > .5$ , the C/D cues are  
 35 dominating the competing cues. When  $\mu_{C/D} < .5$ , the C/D cues are being dom-  
 36 inated by the competing cues. If there is blocking, then the analysis should  
 37 show that the credible values of  $\mu_{C/D}$  are greater than 0.5, when the competing  
 38 cue is B.



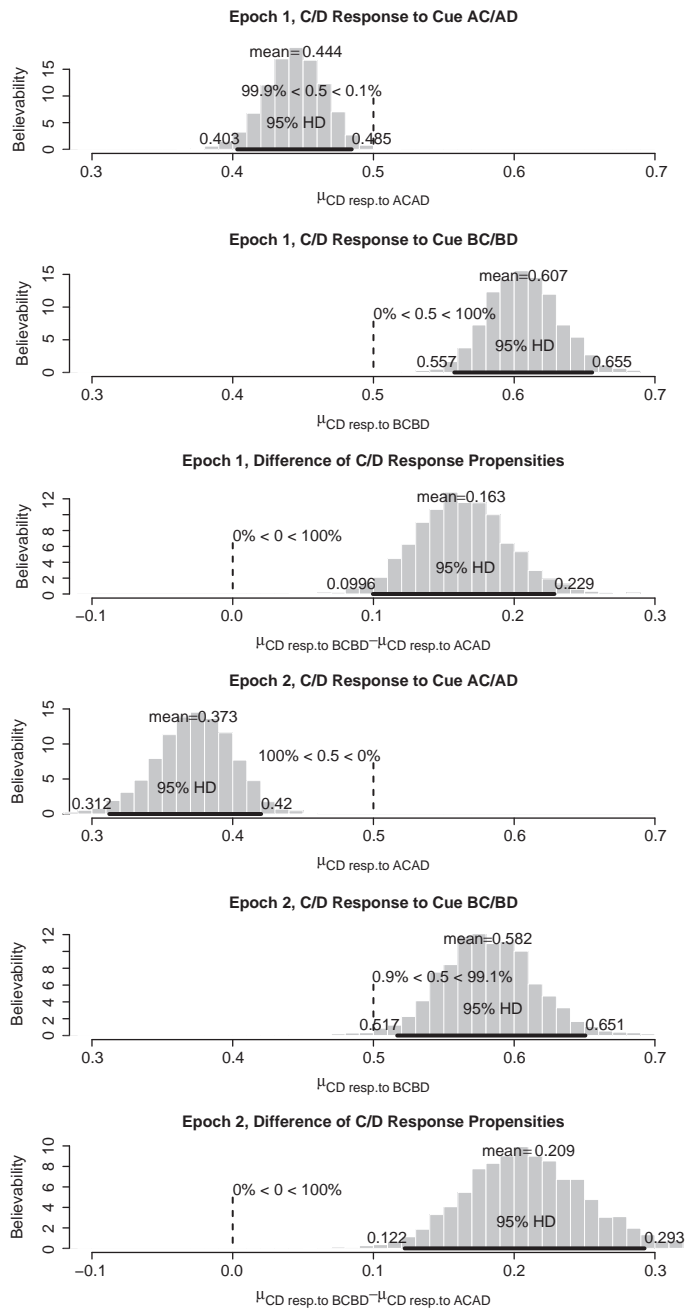
1 The hyperprior on  $\mu_{C/D}$  was a uniform on the interval (0,1). This means that  
 2 the prior was very noncommittal and gave the full range of  $\mu_{C/D}$  values equal  
 3 credibility. The hyperprior on  $\kappa$  was a gamma density with shape and rate  
 4 parameter values of 0.01 (censored at 0.3 so that the random samples in the  
 5 MCMC chain did not cause overflow errors in the beta density). This again  
 6 means that the prior on  $\kappa$  was very noncommittal, allowing a huge range of  
 7 possibilities, but emphasizing small values of  $\kappa$  that reflect large individual dif-  
 8 ferences and a conservative estimate of  $\mu_{C/D}$ . The posterior distributions were  
 9 robust to reasonable changes in the diffuse hyperpriors.

10 This hierarchical model allowed individual differences to be captured by  
 11 variation in participant-level binomial probabilities, which in turn were mutu-  
 12 ally informed by being modeled as representative samples from the same  
 13 higher level beta distribution. The higher level beta distribution captures  
 14 across-subject response tendencies for each context type. The posterior  
 15 certainty in the beta parameters depends on the consistency of response ten-  
 16 dencies across subjects.

17 There is no general analytical solution for deriving the forms of the posterior  
 18 distributions in hierarchical models. Nevertheless, the posteriors can be  
 19 accurately estimated by generating large representative samples. The large  
 20 samples include parameter values that are consistent with the data and the  
 21 prior. The samples are generated by taking a random walk through the high-  
 22 dimensional parameter space. Each step in the walk lands on a point for which  
 23 the combination of parameter values is credible, given the data. Thus, after a  
 24 large number of steps in the random walk, the sampled points provide an  
 25 accurate reflection of the underlying continuous posterior distribution. The  
 26 distribution of points also inherently reveals any correlations among credible  
 27 parameter values.

28 The posterior distribution was determined by Markov chain Monte Carlo  
 29 (MCMC) approximation. The program for generating the sample was written  
 30 in the R language (Ihaka and Gentleman, 1996), using the BRugs interface  
 31 (Thomas, 2004) to the OpenBUGS version (Thomas, O'Hara, Ligges, and  
 32 Sturtz, 2006) of BUGS (Gilks, Thomas, and Spiegelhalter, 1994). Three parallel  
 33 MCMC chains were simulated, using a burn-in of 10,000 steps and thinning of  
 34 200 steps. This burn-in and thinning produced well-mixed chains with small  
 35 auto-correlation, so the posterior sample is very trustworthy. From each of the  
 36 three chains, 1,000 steps were retained to represent the posterior, yielding  
 37 3,000 representative parameter values.

38 Figure 11.4 shows histograms of the believable values of the  $\mu_{C/D}$  param-  
 39 eters. The upper three histograms indicate results from Epoch 1, i.e. backward



**Fig. 11.4** Posterior distribution of C/D preference for the group. *Upper three histograms:* Epoch 1; backward blocking with direct feedback. *Middle three histograms:* Epoch 2; forward blocking with indirect feedback. *Lower three histograms (next page):* Epoch 3; backward blocking with indirect feedback. The dark bar labeled '95% HD' spans the 95% highest density interval, such that all parameter values within the interval have higher believability than values outside the interval, and the interval covers 95% of the believable values.

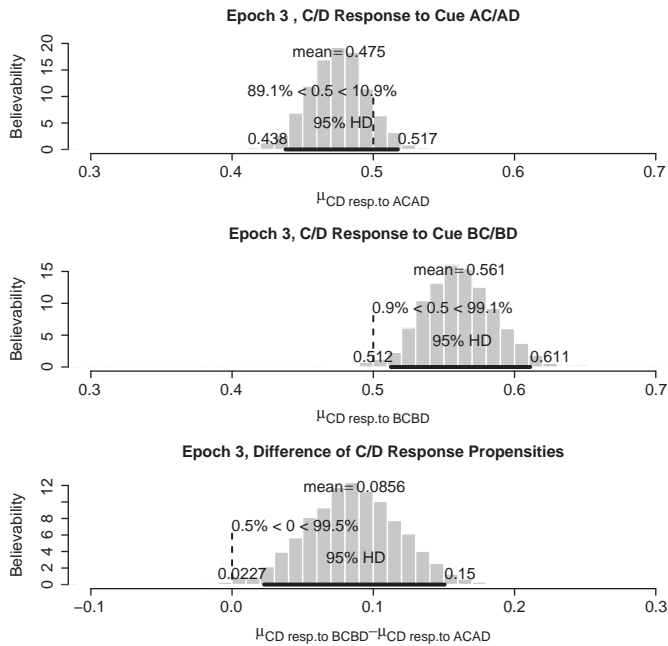


Fig. 11.4 (continued)

1 blocking with direct feedback regarding the relevant specialist. The upper  
 2 histogram indicates that for the test probes AC or AD, the credible values of  
 3 the response propensity are virtually all less than .5, i.e., people prefer the  
 4 response consistent with cue A. The second histogram of the set indicates that  
 5 for test probes BC or BD, the credible response propensities are all well above  
 6 0.5, indicating a robust backward blocking effect. The third histogram in the  
 7 set indicates that the difference between the two types of probes is credibly  
 8 different from zero.

9 The histograms in Figure 11.4 for Epoch 2 indicate that there was credible  
 10 forward blocking of cues to relevance, even when there was only indirect feed-  
 11 back regarding the correspondence of context cues to response cues. In par-  
 12 ticular, virtually all the believable response propensities to BC or BD cues are  
 13 in favor of the C/D consistent response.

14 Most important for our present purposes, the lower set of histograms in  
 15 Figure 11.4 for Epoch 3 indicates that there was credible *backward* blocking of  
 16 cues to relevance, even when there was only indirect feedback regarding the  
 17 correspondence of context cues to response cues. Specifically, the distribution  
 18 that estimates the C/D propensities for cues BC or BD falls mostly (99.1%)

1 above 0.5. The lowest histogram shows that the C/D propensity for BC or BD  
 2 tests is larger than the C/D propensity for AC or AD tests, with nearly all the  
 3 distribution falling above zero.

#### 4 **Summary and discussion of experiment results**

5 The Bayesian analysis of the data incorporated a model of individual differ-  
 6 ences and yielded an explicit representation of credible response propensities.  
 7 The analysis revealed that it is highly credible that there is forward and back-  
 8 ward blocking of cues to relevance.

9 For devotees of the twentieth century ritual of null hypothesis significance  
 10 testing, a chi-square analysis is hereby provided. In Table 11.2, for each epoch's  
 11  $2 \times 2$  table, a chi-square test of independence was conducted on the raw fre-  
 12 quencies. These three tests correspond to the lowest histogram in each of the  
 13 three panels of Figure 11.4. For Epoch 1,  $\chi^2(df = 1, N = 2880) = 44.55, p < .001$ ;  
 14 for Epoch 2,  $\chi^2(df = 1, N = 1952) = 46.19, p < .001$ ; for Epoch 3,  $\chi^2(df = 1,$   
 15  $N = 2096) = 9.90, p < .002$ . A chi-square test can also be conducted on the the  
 16 BD trials alone, i.e., the lower row of Table 11.2. These three tests correspond  
 17 to the middle histogram in each of the three panels of Figure 11.4. For Epoch 1,  
 18  $\chi^2(df = 1, N = 1440) = 36.10, p < .001$ ; for Epoch 2,  $\chi^2(df = 1, N = 976) = 11.95,$   
 19  $p < .001$ ; for Epoch 3,  $\chi^2(df = 1, N = 1048) = 9.16, p < .003$ . The implication  
 20 from these tests is that there was highly significant backward blocking in  
 21 all three Epochs. These standard analyses assume that all individuals have  
 22 the same underlying magnitude of response preference. This is not a reasona-  
 23 ble assumption; the Bayesian analysis does not make it. Furthermore, the  
 24 standard computation of  $p$  values assumes that  $N$  was fixed in advance, and  
 25 data collection was stopped when that  $N$  was reached. In fact, the data were  
 26 collected for a set number of weeks during which participants volunteered,  
 27 and data were excluded if learners did not reach the accuracy criterion. Bayesian  
 28 analysis has no reliance on why data collection stopped. Notice also that the  
 29 Bayesian analysis provides a complete distribution for credible values of the  
 30 parameters (in Figure 11.4), while the chi-square analysis provides no interval  
 31 estimate.

32 As mentioned earlier, the magnitudes of the forward and backward blocking  
 33 effects were weak. We speculate that the small effects occurred because process-  
 34 ing the corrective feedback was effortful and interfered with attentional  
 35 re-allocation. For example, when the corrective feedback indicates disease F,  
 36 the learner must first determine which specialist is consistent with that out-  
 37 come, before then being able to determine which symptoms are relevant to  
 38 selecting that specialist. If processing of the feedback interferes with allocation

1 of attention to symptoms, and if blocking of symptoms depends in part on  
 2 attentional allocation, then the extra processing of the feedback may impair  
 3 learned inattention to the blocked cue.

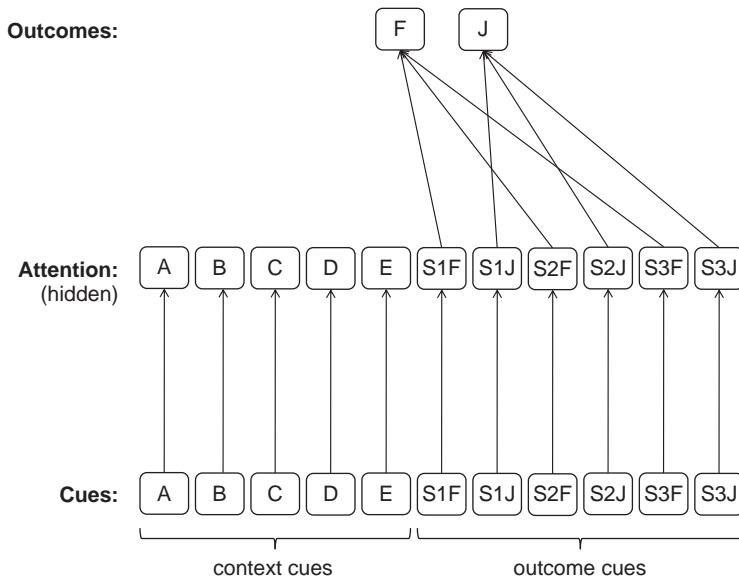
4 The limited variety of test trials in this particular experiment admits a differ-  
 5 ent explanation of the results. In this alternative explanation, there is no block-  
 6 ing during learning, but there is instead a response bias at test. On tests  
 7 involving BC or BD cues, the preference for C/D-consistent responses is a  
 8 result of a response bias to choose specialist S2 whenever two context cues  
 9 appear. The response bias is an inference from the training phases, for which it  
 10 is the case that only two-cue contexts occurred when specialist S2 was relevant.  
 11 This general bias is overcome on tests involving AC or AD cues, because the  
 12 association from A to S1 is very strong. Only future experiments will be able to  
 13 distinguish the two explanations, perhaps by including tests involving single-  
 14 cue contexts. For example, the two test items  $B_{S1} S1_F S2_J$  and  $D_{S2} S1_F S2_J$  involve  
 15 a single context cue, and therefore would not suffer the hypothesized two-  
 16 context-cue bias toward specialist S2, but blocking would predict a difference  
 17 in response preferences across the two items. Until follow-up experiments are  
 18 conducted, we must rest with the argument that structurally analogous previ-  
 19 ous experiments, involving blocking of cues to outcomes instead of cues to  
 20 relevance, have shown less ambiguous blocking phenomena (e.g., Kruschke  
 21 and Blair, 2000).

22 The current results at least indicate proof of concept: In principle, it is  
 23 possible for people to learn which context cues indicate which response cues  
 24 to attend to. And, most importantly, we have demonstrated data consistent  
 25 with backward blocking of cues to relevance. Future experiments will attempt  
 26 to use context and response cues that are not so blatantly distinct as symptoms  
 27 and medical specialists. With enough training, people should be able to learn  
 28 in those situations too, and presumably also show backward blocking.

## 29 **Modelling**

30 The experiment of the previous section showed results consistent with back-  
 31 ward blocking of contextual cues to relevance. In this section we show that the  
 32 behavior can be mimicked by a locally Bayesian learning model in which the  
 33 first layer learns to allocate attention and the second layer learns to generate  
 34 outcomes.

35 The architecture of the model is illustrated in Figure 11.5. Each node in  
 36 the lower layer represents a cue that can be present or absent. Notice that the  
 37 lower layer consists of 11 cues, including the symptoms and the specialist  
 38 information. The symptoms and specialists were presented in the experiment



**Fig. 11.5** Model architecture. Arrows denote the most credible associations at the beginning of training (i.e., the mean of the prior distribution). The cues are marked at the bottom as context or outcome cues, but this marking is purely for the benefit of the reader, as the model does not ‘know’ which cues are context cues and which cues are outcome cues.

- 1 as different types of cues, but in principle they are both just present/absent bits
- 2 of information. For example, when stimulus  $A_{S1} S1_F S2_J$  is presented, input
- 3 nodes A, S1F, and S2J are activated.
- 4 The middle layer, otherwise called the hidden layer, represents attentionally
- 5 gated cue activations. Each input cue has a corresponding node in the hidden
- 6 layer. By default, each cue calls attention to itself, indicated in Figure 11.5 by the
- 7 1-to-1 arrows. Learning can cause the magnitude of the 1-to-1 connection to
- 8 change, and also generate ‘lateral’ connections between input cues and atten-
- 9 tional nodes. In particular, when a cue is blocked, the 1-to-1 connection from
- 10 that cue to its own attentional gate may be reduced, and the lateral connection to
- 11 its attentional gate from the blocking cue may become negative. More accurately
- 12 described, one would say that lower values of the 1-to-1 connection become
- 13 more believable, and negative values from the blocking cue to the blocked
- 14 attention node become more believable. [Griffiths and LePelley (2009)
- 15 have shown that *forward* blocking of *response* cues is unlikely to produce strong
- 16 negative lateral connections, in at least some situations. It is not yet known

1 whether comparable experiments with backward blocking would lead to analo-  
 2 gous conclusions. The present model can accommodate these results, at least in  
 3 principle, by modulating the relative learnabilities of lateral and 1-to-1  
 4 weights.]

5 The upper layer represents outcomes, which are disease diagnoses in the  
 6 present application. The associations from attentionally gated cues to out-  
 7 comes must be learned. In principle, any of the cues could be indicative of any  
 8 outcome. In the present application, only the cues corresponding to medical  
 9 specialists happen to be correlated with the correct diagnoses. Moreover, the  
 10 experiment presented the specialist cues in such a way that the corresponding  
 11 diagnosis was explicitly indicated. In other words, the associative links from  
 12 specialists to diagnoses were already suggested, and therefore the model has  
 13 these links built in, as shown in Figure 11.5 by the arrows fanning into the  
 14 outcome nodes. Through learning, these links can be altered, and other links  
 15 from cues to outcomes can be established.

16 In locally Bayesian learning, multiple hypotheses about the associative  
 17 weights are maintained, for both layers of weights. Learning consists of shifting  
 18 credibility from unlikely hypothetical weights to likely hypothetical weights.  
 19 The arrows in Figure 5 indicate the hypothetical weights with highest credibil-  
 20 ity at the beginning of the experiment, before training.

21 Intuitively, locally Bayesian learning proceeds in Epoch 3 as follows. The  
 22 first phase, involving redundant-context cues, has cases of  $A_{S1} B_{S1} S1_F S2_J \rightarrow F$ ,  
 23 among others (see Table 11.1). This phase causes credibility to be enhanced  
 24 on positive connections from cues A and B to attention on  $S1_F$  and  $S1_J$ , and  
 25 perhaps also causes credibility to be enhanced on negative connections from  
 26 cues A and B to attention on  $S2_F$  and  $S2_J$ . In other words, the model learns  
 27 that when symptoms A or B are present, attend to specialist 1 and ignore spe-  
 28 cialist 2. Because the lower layer is a Bayesian system that keeps track of multi-  
 29 ple candidate weight combinations, the system ‘knows’ that the training cases  
 30 could be explained by either (1) A indicating  $S1$  and B being irrelevant, or  
 31 (2) B indicating  $S1$  and A being irrelevant, or (3) A or B indicating  $S1$ . All of  
 32 these plausible weight combinations retain some credibility. In the second  
 33 phase, involving single-context cues, there are cases of  $A_{S1} S1_F S2_J \rightarrow F$ , among  
 34 others (see Table 11.1). Of the weight combinations that remained credible  
 35 after the first phase, the ones involving cue A gain increased credibility, thereby  
 36 reducing the credibility of the ones involving cue B alone. In other words,  
 37 locally Bayesian learning on the lower layer can account for backward blocking  
 38 of contextual cues to relevance.

39 The remainder of this section is devoted to a detailed description of a  
 40 particular instantiation of this approach. In this particular instantiation, the

1 outcome and attention nodes are modeled as individual Kalman filters. Kalman  
 2 filters have been previously used by Dayan and Kakade (2001) to model back-  
 3 ward blocking, and the idea of using layers of locally learning Kalman filters  
 4 was mentioned by Kruschke (2006a, 2006b). We do not intend to argue that  
 5 Kalman filters are the best way to instantiate the components of locally Bayesian  
 6 learning. Indeed, Kalman filters can only represent linear mappings, and are  
 7 therefore quite limiting. We use Kalman filters merely because they are con-  
 8 venient for computational tractability. Presumably the model structure used  
 9 by Kruschke (2006b), which did not involve Kalman filters, would show simi-  
 10 lar qualitative behavior.

11 One goal of the detailed modeling is to demonstrate by computer simulation  
 12 that the intuitive argument provided in the previous paragraphs is actually  
 13 correct. A second goal of the reporting the mathematical details is to present  
 14 novel derivations that have not been previously presented elsewhere. There is  
 15 not sufficient space here to provide an extensive tutorial on Kalman filters, but  
 16 a tutorial regarding Kalman filters applied to associative learning has been  
 17 previously provided by Kruschke (2008). The reader, who is not concerned at  
 18 this time with the mathematical implementation of the model, is invited to  
 19 skip ahead to the next subsection where the model results are reported.

## 20 Formal description of layers of Kalman filters

21 The Kalman filter assumes that the output to be predicted is a scalar metric  
 22 value  $y$ . In our experiment, the outcome is dichotomous (present/absent),  
 23 which is represented as  $y = 1$  or  $y = 0$ . The Kalman filter assumes that  $y$  is a  
 24 linear function of the input vector  $\vec{x}$  (thought of as a column matrix), with  
 25 normally distributed noise in the output. The variance of the noise is the scalar  
 26 value  $v$ . The associative weights are denoted by the vector  $\vec{w}$  (thought of as a  
 27 column matrix). Formally, then, the probability of a value  $y$  is

$$28 \quad p(y | \vec{x}, \vec{w}, v) = \frac{1}{\sqrt{v}(2\pi)^{1/2}} \exp\left(-.5 \frac{(y - \vec{w}^T \vec{x})^2}{v}\right) \quad (11.1)$$

29 The value of  $v$  is considered to be a known constant, fixed by the modeler.  
 30 Equation 11.1 is the likelihood function that specifies the probabilistic  
 31 output of any node in the model. The normal distribution in Equation 1 is  
 32 also denoted  $N(y | \vec{w}^T \vec{x}, v)$ .

33 Each outcome and attention node is a distinct Kalman filter. For example,  
 34 the node for outcome  $F$  is a Kalman filter for which  $y = 1$  means  $F$  is present,  
 35 and  $y = 0$  means  $F$  is absent. The input to that outcome node is the vector of  
 36 attentional values from the hidden layer. The attentional values are the (means



1 of the) output values of the attentional Kalman filter nodes. Each attentional  
 2 node is also a Kalman filter. The input to each attentional node is the vector of  
 3 cue values.

4 In each Kalman filter, the value of  $\vec{w}$  has uncertainty, meaning that each  
 5 possible combination of weights has a degree of belief, and beliefs are spread  
 6 over a wide range of weight combinations. When the model learns, belief is  
 7 shifted toward weight combinations that are consistent with the training items.  
 8 The initial state of the network has beliefs spread out symmetrically and dif-  
 9 fusely over a wide range of weight combinations. The initial state is called the  
 10 prior distribution. It is assumed to be normal with mean  $\vec{\mu}$  (a column matrix)  
 11 and covariance matrix  $C$ . Thus, the prior on  $\vec{w}$  is

$$12 \quad p(\vec{w} | \vec{\mu}, C) = \frac{1}{\sqrt{|C|}(2\pi)^{d/2}} \exp\left(-.5(\vec{w} - \vec{\mu})^T C^{-1}(\vec{w} - \vec{\mu})\right) \quad (11.2)$$

13 where  $|C|$  is the determinant of  $C$ , and  $d$  is the dimensionality of  $\vec{w}$ , i.e., the  
 14 number of input nodes. Equation 11.2 is merely the well-known formula  
 15 for the multivariate normal distribution, and Equation 11.1 is merely the spe-  
 16 cial case of that formula when  $d = 1$  and  $\mu = \vec{w}^T \vec{x}$ . The normal distribution in  
 17 Equation 11.2 is also denoted  $N(\vec{w} | \vec{\mu}, C)$ .

18 The mean vector for the weights was set to all zeros except for specific com-  
 19 ponents that represented initial correspondences. The mean vector on weights  
 20 fanning in to the outcome nodes was initialized to represent the explicit  
 21 mappings indicated by the specialists. Therefore the mean weight to F from  
 22 S1F was set to 1, as was the weight to F from S2F, and to J from S1J, and so  
 23 forth. These mean connection weights of 1 are represented by the arrows in  
 24 Figure 11.5. The mean vector on weights fanning in to the attention nodes was  
 25 initialized to represent the 1-to-1 correspondence of cues to attention, but also  
 26 to take into account only partial attention to any given cue. Therefore the  
 27 1-to-1 connections were initialized arbitrarily at a mean value of 0.1. These  
 28 1-to-1 connections are represented by the arrows in Figure 11.5.

29 The prior covariance matrix for each node specifies the prior certainty in the  
 30 mean values. Because the specialists gave explicit diagnoses, the mean weights  
 31 from specialists to outcome nodes had very small variances, i.e., very high cer-  
 32 tainties. But the weights to the outcome nodes from the context attention  
 33 nodes had prior variances that were large, to reflect great initial uncertainty.  
 34 The prior covariance matrices for the attention nodes were likewise set with  
 35 small variances on the specialist cues and large variances on the context cues.  
 36 This was because of the prior assumption that specialists do not inform each  
 37 other, but the symptoms handled by each specialist must be learned.

### 1 Prediction by a Kalman-filter node

2 In propagating activation up a succession of Kalman filter nodes, the input to  
 3 the upper layer is the mean output of the lower layer. The mean output of a  
 4 Kalman filter is simply

$$\begin{aligned}\bar{y} &= \int d\bar{w} p(\bar{w} | \bar{\mu}, C) \int dy y p(y | \bar{x}, \bar{w}, v) \\ &= \int d\bar{w} p(\bar{w} | \bar{\mu}, C) \bar{w}^T \bar{x} \\ &= \left( \int d\bar{w} p(\bar{w} | \bar{\mu}, C) \bar{w} \right)^T \bar{x} \\ &= \bar{\mu}^T \bar{x}\end{aligned}\tag{11.3}$$

6 In the architecture we use for expressing attentional gating, the attentional  
 7 Kalman filter acts as a multiplier on the input cue. Formally, the hidden layer  
 8 activation that acts as the input to the outcome nodes is the mean activation of  
 9 the attentional Kalman filters times the corresponding input cue activations:  
 10  $x_i^{hid} = \bar{y}_i x_i^{cue}$  where  $\bar{y}_i$  is the mean output of the  $i$ -th attentional Kalman filter,  
 11 as determined from Equation 11.3.

### 12 Learning by a Kalman-filter node

13 The values of  $\bar{\mu}$  and  $C$  serve as the prior distribution for the trial's Bayesian  
 14 updating. On a given trial, the input vector is  $\bar{x}$  and the correct output  
 15 value, a.k.a. the target, is denoted  $t$  (a scalar). It turns out (Meinhold  
 16 and Singpurwalla, 1983) that the Bayesian updating formula simplifies to the  
 17 following expressions:

$$\bar{\mu}' = \bar{\mu} + C \bar{x} \left[ v + \bar{x}^T C \bar{x} \right]^{-1} (t - \bar{x}^T \bar{\mu})\tag{11.4}$$

$$C' = C - C \bar{x} \left[ v + \bar{x}^T C \bar{x} \right]^{-1} \bar{x}^T C\tag{11.5}$$

20 Thus, for any given outcome node or attention node, Equations 11.4 and 11.5  
 21 are applied to update its beliefs. The computational simplicity of these updat-  
 22 ing equations is what makes the Kalman filter appealing as an implementation  
 23 of Bayesian learning. The target for the outcome nodes is explicitly indicated  
 24 by corrective feedback, and the input to the outcome nodes is the pattern of  
 25 attentional activation as defined in the previous paragraph. For the attentional  
 26 nodes, the input is the cue activation, but the target values for the attentional  
 27 nodes will be defined in the next subsection.

28 In summary, at the beginning of a trial, the weight vector is distributed as  
 29 in Equation 11.2. Then target and input for the trial are provided, and beliefs

1 regarding credible weight combinations are adjusted by Bayes' rule. The pos-  
 2 terior distribution of the weights conveniently turns out to be again normal,  
 3 with mean and covariance given by Equations 11.4 and 11.5.<sup>2</sup>

#### 4 Finding the attention that maximizes a desired output

5 To find a target for the attention nodes, we want to find the attentional input  
 6 to the outcome nodes that would maximize the probability of the correct out-  
 7 come. In other words, given the target outcome values,  $t_k$  of outcome  $k$ , we  
 8 want to find the attention-node values  $\vec{x}_t$  that maximize the joint probability  
 9 of the outcome values:

$$\begin{aligned}\vec{x}_t &= \arg \max_{\vec{x}} \prod_k p(t_k | \vec{x}) \\ &= \arg \max_{\vec{x}} \prod_k \int d\vec{w} p(t_k | \vec{x}, \vec{w}, \nu) p(\vec{w} | \vec{\mu}, \mathbf{C})\end{aligned}\quad (11.6)$$

10

11 A formal identity regarding the product of Gaussians states that

$$\begin{aligned}N(\vec{w} | \vec{\mu}, \mathbf{C}) N(\vec{x}^T \vec{w} | t, \nu) \\ = N(t | \vec{x}^T \vec{\mu}, \nu + \vec{x}^T \mathbf{C} \vec{x}) N(\vec{w} | \vec{m}, \mathbf{V})\end{aligned}$$

12

13 where

$$\begin{aligned}\vec{m} &= \mathbf{V}(\mathbf{C}^{-1} \vec{\mu} + \vec{x} t / \nu) \\ \mathbf{V} &= (\mathbf{C}^{-1} + \vec{x} \vec{x}^T / \nu)^{-1}\end{aligned}$$

14

<sup>2</sup> The full Kalman filter also assumes that the distribution of  $\vec{w}$  changes in time, separately from and before any updating in beliefs inferred from observed data. This dynamic aspect of the weights is assumed to be linear at any given trial, so that the mean weights are dynamically changed into some linear transformation of the current mean weights. Dayan et al. (2000) used this mechanism to model unbiased diffusion of weights through time. Because aficionados of the Kalman filter may wonder what we did with the dynamic mechanism, it is summarized in this footnote. Denote the dynamic linear transformation by  $\mathbf{D}$ . There is also assumed to be some additive change in the covariances of the weights; typically this is thought of as a constant increment in the uncertainty of the weights as time progresses. Denote the added uncertainty by  $\mathbf{U}$ . Formally, then, at the beginning of each trial, the weight distribution is dynamically changed as follows:  $\vec{\mu}^* = \mathbf{D} \vec{\mu}$  and  $\mathbf{C}^* = \mathbf{D} \mathbf{C} \mathbf{D}^T + \mathbf{U}$ . In all of our applications, we assume that  $\mathbf{D}$  is the identity matrix and that  $\mathbf{U} = 0$ . That is, we assume that the weights are not systematically changing through time. This restriction implies that all the behavior of the model comes from Bayesian learning, not from additional dynamic assumption.

1 Hence the integral in Equation 11.6 can be rewritten as

$$\begin{aligned}
 & \int d\bar{w} p(t_k | \bar{x}, \bar{w}, \nu) p(\bar{w} | \bar{\mu}, C) \\
 &= \int d\bar{w} N(t_k | \bar{x}^T \bar{w}, \nu) N(\bar{w} | \bar{\mu}, C) \\
 &= \int d\bar{w} N(t_k | \bar{x}^T \bar{\mu}, \nu + \bar{x}^T C \bar{x}) N(\bar{w} | \bar{m}, V) \\
 &= N(t_k | \bar{x}^T \bar{\mu}, \nu + \bar{x}^T C \bar{x}) \int d\bar{w} N(\bar{w} | \bar{m}, V) \\
 &= N(t_k | \bar{x}^T \bar{\mu}, \nu + \bar{x}^T C \bar{x})
 \end{aligned} \tag{11.7}$$

2  
3 We used numerical approximation (specifically a Newton-Raphson method  
4 applied to the derivative) to find the  $\bar{x}$  that maximizes Equations 11.6  
5 and 11.7.

## 6 **Model result: Single-layer Bayesian model fails**

7 The experiment was designed such that no context cue is individually corre-  
8 lated with the correct diagnosis. Inspection of Table 11.1 reveals that every  
9 context cue occurs as often with outcome  $F$  and with outcome  $J$ . Therefore, a  
10 single-layer Kalman filter is unable to learn the correct responses in either the  
11 single-context phase or the redundant-context phase.

12 Because complex models have an uncanny way of confounding intuitions,  
13 especially when applied to complex designs, we simulated a single-layer  
14 Kalman filter to be sure that such a model truly was unable to learn the map-  
15 ping. The results verified that the model could only produce 50–50 (i.e.,  
16 chance) responding for Epochs 2 and 3.

## 17 **Model result: Locally Bayesian learning succeeds**

18 To fit the output of the model to the human response proportions in Table 11.2,  
19 we had to map the model output, expressed as  $\bar{y}_k$ , to response proportions.  
20 We did this via the often used softmax rule:  $p(K) = \exp(\gamma \bar{y}_K) / \sum_k \exp(\gamma \bar{y}_k)$ ,  
21 where  $\gamma > 0$  is a parameter called the ‘decisiveness’ of the choice. When  $\gamma$  is large,  
22 small differences in the outcome activations lead to large differences in choice  
23 proportion, but when  $\gamma$  is small, outcome activations must be very different to  
24 produce much difference in choice proportion. We arbitrarily set  $\gamma = 1$ .

25 The only free parameter was the noise variable  $\nu$  in Equation 11.1. The noise  
26 parameter acts much like a learning rate, as can be seen by its appearance in the  
27 update Equations 11.4 and 11.5. When  $\nu$  is larger, learning is slower.

28 The model was trained in each phase on the median number of blocks that  
29 human learners took. For successive phases in the three consecutive epochs,  
30 the median number of blocks was 3, 2, 4, 2, 3, and 2, respectively. The average  
31 of 40 simulated subjects, each with a different permutation of trials within

**Table 11.3** Locally Bayesian model behavior

Test cues	Response consistent with A/B or C/D					
	Epoch 1 Backward Blocking Explicit Specialist		Epoch 2 Forward Blocking Indirect Feedback		Epoch 3 Backward Blocking Indirect Feedback	
	A/B	C/D	A/B	C/D	A/B	C/D
AC or AD	56.1	43.9	56.2	43.8	56.0	44.0
BC or BD	44.2	55.8	44.4	55.6	44.3	55.7

blocks, was used as the model prediction. The fit of the model was measured simply as sum-squared-deviation between the human response percentages in Table 11.2 and the corresponding model response percentages.

The best fit used  $\nu = 0.047$ , and the average model output is shown in Table 11.3. The model shows forward and backward blocking very similar to the human preferences in Table 11.2. The qualitative trend of the model is quite robust against changes in parameter values. The main point of the simulation is to demonstrate that the locally Bayesian model does indeed produce backward blocking of cues to relevance, as shown by the model behavior in Epoch 3. The magnitude of this backward blocking can be made larger with other parameter values. The magnitude is small in this simulation specifically to match the small magnitude shown by human learners.

The fit by the model is not ‘spot on’ the data, but the model uses only one free parameter and there are many aspects of the human procedure not directly imitated in the model procedure. For instance, all simulated subjects used exactly the median number of training blocks, but many human subjects took more blocks to learn. All simulated subjects used the same value of  $V$ , but presumably different human learners had different learning rates. The simulated model had a constant value of  $\nu$  throughout all three epochs, but human learners may have ‘learned to learn’ as they became familiar with the stimulus arrangement and task during training. And, perhaps most importantly, the simulation started the 1-to-1 attention links at the same neutral prior for every epoch, but people may have had progressively more certain priors in successive epochs. The purpose of the model simulation is to demonstrate that even with these simplifications, the locally Bayesian model shows robust backward blocking of cues to relevance.

## Conclusion

We have provided some data that suggest that people exhibit backward blocking of cues to relevance. One way of modeling this behavior is with layers of locally Bayesian learning. The lower layer learns to allocate attention to cues,

1 and the upper layer learns to produce outcomes based on the attended cues.  
 2 Because there is Bayesian learning within each layer, the model can exhibit  
 3 backward blocking of associations within each layer.

4 The novel contributions of this chapter are (1) the empirical suggestion of  
 5 backward blocking of cues to relevance, and (2) the implementation of locally  
 6 Bayesian learning as layers of Kalman filters. Backward blocking of cues to  
 7 relevance may justify locally Bayesian learning in the attentional layer of  
 8 Kruschke's (2006) model. The results of the new experiment reported here  
 9 may have alternative interpretations, but we hope that this work may provoke  
 10 interesting follow-up research to distinguish alternative accounts. Any com-  
 11 peting model should also account for the spectrum of phenomena addressed  
 12 by locally Bayesian learning applied to attentional learning (Kruschke,  
 13 2006b).

## 14 Acknowledgements

15 For help administering the experiments, thanks go to Kelsey Buckingham,  
 16 Kevin Clemens, Alyssa Heggen, Kaitlyn Smith, Kari Vann, and Phaedra  
 17 Willson. For providing pilot data and discussion, we thank Rima Hanania and  
 18 Richard Hullinger. The editors, Chris Mitchell and Mike LePelley, provided  
 19 helpful suggestions for clarification and discussion. Correspondence can be  
 20 addressed to John K. Kruschke, Department of Psychological and Brain  
 21 Sciences, Indiana University, 1101 E. 10th St., Bloomington IN 47405-7007, or  
 22 via electronic mail to [kruschke@indiana.edu](mailto:kruschke@indiana.edu). Supplementary information can  
 23 be found at <http://www.indiana.edu/~kruschke/>

## 24 References

- 25 Chater, N., Tenenbaum, J.B., & Yuille, A. (Eds.). (2006, July). Special issue: Probabilistic  
 26 models of cognition. *Trends in Cognitive Sciences*, 10(7), 287–344.  
 27 Chun, M. M. (2000). Contextual cueing of visual attention. *Trends in Cognitive Sciences*,  
 28 4(5), 170–178.  
 29 Dayan, P., & Kakade, S. (2001). Explaining away in weight space. In T. Leen, T. Dietterich,  
 30 & V. Tresp (Eds.), *Advances in neural information processing systems* (Vol. 13,  
 31 pp. 451–457). Cambridge, MA: MIT Press. Dayan, P., Kakade, S., & Montague, P. R.  
 32 (2000). Learning and selective attention. *Nature Neuroscience*, 3, 1218–1223.  
 33 De Houwer, J. & Beckers, T. (2002). A review of recent developments in research and  
 34 theories on human contingency learning. *The Quarterly Journal of Experimental*  
 35 *Psychology*, 55B, 289–310.  
 36 Deneve, S. (2008). Bayesian spiking neurons II: Learning. *Neural Computation*, 20,  
 37 118–145.  
 38 Dickinson, A. (2001). Causal learning: Association versus computation. *Current Directions*  
 39 *in Psychological Science*, 10, 127–132.

- 1 Dickinson, A. & Burke, J. (1996). Within-compound associations mediate the retrospective  
2 revaluation of causality judgements. *Quarterly Journal of Experimental Psychology:*  
3 *Comparative and Physiological Psychology*, 49B, 60–80.
- 4 Gilks, W. R., Thomas, A., & Spiegelhalter, D. J. (1994). A language and program for  
5 complex Bayesian modelling. *The Statistician*, 43(1), 169–177.
- 6 Griffiths, O., & Le Pelley, M. E. (2009). Attentional changes in blocking are not a  
7 consequence of lateral inhibition. *Learning and Behavior*, 37(1), 27–41.
- 8 Ihaka, R. & Gentleman, R. (1996). R: a language for data analysis and graphics. *Journal of*  
9 *Computational and Graphical Statistics*, 299–314. (With other contributors listed  
10 at <http://www.r-project.org/>)
- 11 Kamin, L. J. (1968). 'Attention-like' processes in classical conditioning. In M. R. Jones  
12 (Ed.), *Miami symposium on the prediction of behavior: Aversive stimulation*  
13 (pp. 9–33). Coral Gables, FL: University of Miami Press.
- 14 Kruschke, J. K. (2001). Toward a unified model of attention in associative learning. *Journal*  
15 *of Mathematical Psychology*, 45, 812–863.
- 16 Kruschke, J. K. (2005). Learning involves attention. In G. Houghton (Ed.), *Connectionist*  
17 *models in cognitive psychology* (Ch 4, pp. 113–140). Hove, East Sussex: Psychology Press.
- 18 Kruschke, J. K. (2006a). Locally Bayesian learning. In Proceedings of the 28th Annual  
19 Conference of the Cognitive Science Society (pp. 453–458).
- 20 Kruschke, J. K. (2006b). Locally Bayesian learning with applications to retrospective  
21 revaluation and highlighting. *Psychological Review*, 113(4), 677–699.
- 22 Kruschke, J. K. (2008). Bayesian approaches to associative learning: From passive to active  
23 learning. *Learning and Behavior*, 36(3), 210–226.
- 24 Kruschke, J. K. (2009). Highlighting: A canonical experiment. In B. Ross (Ed.), *The*  
25 *psychology of learning and motivation* (Vol. 51, pp. 153–185). Academic Press.
- 26 Kruschke, J. K. & Blair, N. J. (2000). Blocking and backward blocking involve learned  
27 inattention. *Psychonomic Bulletin and Review*, 7, 636–645.
- 28 Kruschke, J. K., Kappenman, E. S., & Hetrick, W. P. (2005). Eye gaze and individual  
29 differences consistent with learned attention in associative blocking and highlighting.  
30 *Journal of Experimental Psychology: Learning, Memory and Cognition*, 31, 830–845.
- 31 Le Pelley, M. E., Beesley, T., & Suret, M. B. (2007). Blocking of human causal learning  
32 involves learned changes in stimulus processing. *The Quarterly Journal of Experimental*  
33 *Psychology*, 60(11), 1468–1476.
- 34 Mackintosh, N. J. (1975). A theory of attention: Variations in the associability of stimuli  
35 with reinforcement. *Psychological Review*, 82, 276–298.
- 36 Mackintosh, N. J., & Turner, C. (1971). Blocking as a function of novelty of CS and  
37 predictability of UCS. *Quarterly Journal of Experimental Psychology*, 23, 359–366.
- 38 Markman, A. B. (1989). LMS rules and the inverse base-rate effect: Comment on Gluck  
39 and Bower (1988). *Journal of Experimental Psychology: General*, 118, 417–421.
- 40 Meinhold, R. J., & Singpurwalla, N. D. (1983). Understanding the Kalman filter. *American*  
41 *Statistician*, 37(2), 123–127.
- 42 Miller, R. R. & Matzel, L. D. (1988). The comparator hypothesis: A response rule for the  
43 expression of associations. In G.H. Bower (Ed.), *The psychology of learning and*  
44 *motivation: Advances in research and theory* (Vol. 22, pp. 51–92). San Diego, CA:  
45 Academic Press.

- 1 Mitchell, C. J., Harris, J. A., Westbrook, R. F., & Griffiths, O. (2008). Changes in cue  
2 associability across training in human causal learning. *Journal of Experimental*  
3 *Psychology: Animal Behavior Processes*, 34(4), 423–436.
- 4 Nelson, J. B. (2002). Context specificity of excitation and inhibition in ambiguous stimuli.  
5 *Learning and Motivation*, 33, 284–310.
- 6 Rosas, J. M., Callejas-Aguilera, J. E., Ramos-Alvarez, M. M., & Abad, M. J. F. (2006).  
7 Revision of retrieval theory of forgetting: What does make information context-  
8 specific? *International Journal of Psychology and Psychological Therapy*, 6(2), 147–166.
- 9 Shanks, D.R. (1985). Forward and backward blocking in human contingency judgement.  
10 *Quarterly Journal of Experimental Psychology: Comparative and Physiological Psychology*,  
11 37B, 1–21.
- 12 Sobel, D. M., Tenenbaum, J. B., & Gopnik, A. (2004). Children's causal inferences from  
13 indirect evidence: Backwards blocking and Bayesian reasoning in preschoolers.  
14 *Cognitive Science*, 28, 303–333.
- 15 Tassoni, C. J. (1995). The least mean squares network with information coding: A model of  
16 cue learning. *Journal of Experimental Psychology: Learning, Memory, and Cognition*,  
17 21(1), 193–204.
- 18 Tenenbaum, J. B. & Griffiths, T. L. (2003). Theory-based causal inference. In S. Becker,  
19 S. Thrun, and K. Obermayer (Eds.), *Advances in neural information processing systems*  
20 (Vol. 15, pp. 35–42). Cambridge, MA: MIT Press.
- 21 Thomas, A. (2004). *BRugs user manual (the R interface to BUGS)*. [http://mathstat.](http://mathstat.helsinki.fi/openbugs/data/Docu/BRugs%20Manual.html)  
22 [helsinki.fi/openbugs/data/Docu/ BRugs%20Manual.html](http://mathstat.helsinki.fi/openbugs/data/Docu/BRugs%20Manual.html).
- 23 Thomas, A., O'Hara, B., Ligges, U., & Sturtz, S. (2006, March). Making BUGS open.  
24 *R News*, 6(1), 12–17.
- 25 Van Hamme, L.J. & Wasserman, E.A. (1994). Cue competition in causality judgments: The  
26 role of nonpresentation of compound stimulus elements. *Learning and Motivation*,  
27 25(3), 127–151.
- 28 Wills, A. J., Lavric, A., Croft, G. S., & Hodgson, T. L. (2007). Predictive learning, prediction  
29 errors, and attention: Evidence from event-related potentials and eye tracking. *Journal*  
30 *of Cognitive Neuroscience*, 19(5), 843–854.
- 31 Xu, F. & Tenenbaum, J. B. (2007). Word learning as Bayesian inference. *Psychological*  
32 *Review*, 114(2), 245–272.
- 33 Yang, L.-X., & Lewandowsky, S. (2003). Context-gated knowledge partitioning in  
34 categorization. *Learning, Memory*, 29(4), 663–679.