

## Cognition with neurons: A large-scale, biologically realistic model of the Wason task

Chris Eliasmith

Dept. of Philosophy and Dept. of Systems Design Engineering  
University of Waterloo, Waterloo, Ontario N2L 3G1 Canada

### Abstract

I present a cognitive model, dubbed BioSLIE, that integrates and extends recent advances in: 1) distributed, structure-sensitive representation; 2) neurocomputational modeling; and 3) our understanding of the neuroanatomy of inference. As a result, BioSLIE is biologically detailed, learns different behaviors in different contexts, and exhibits systematic, structure-sensitive generalization. Here, BioSLIE is applied to the Wason card selection task. Its performance meets Cosmides' (1989) challenge to mechanistically define domain-general procedures that can use induction to produce the observed domain specific performance on the Wason task. As well, it demonstrates the relevance of neural computation to understanding cognition, despite claims to the contrary by Fodor and Pylyshyn (1988) and Jackendoff (2002).

### Introduction

Fodor and Pylyshyn (1988), and more recently Jackendoff (2002), have argued that understanding neural computation is not relevant for understanding cognitive function. They have suggested that neurally plausible architectures do not naturally support structure-sensitive computation, and that such computation is essential for explaining cognition. They argue that this leads to the conclusion that neurons merely implement a classical system, and that characterizing the implementation itself is irrelevant for understanding the cognitive properties of the system. In this paper, I present a large-scale, biologically realistic model which demonstrates that this view may be incorrect. This model is also, I believe, the first demonstration of a structure-sensitive (i.e., language-based) cognitive function being exhibited by a biologically detailed model. In particular, this model captures the context sensitive inference exhibited by human subjects in the Wason card task using massively interconnected spiking single neurons. To do so, the model learns the relevant structural transformations appropriate for a given context, and is able to generalize them. Given these salient properties of the model, I refer to it as BioSLIE (BIOlogically-plausible Structure-sensitive Learning Inference Engine).

Beyond presenting this specific model, another purpose of this paper is to introduce a modeling methodology that can help researchers build similar models for other cognitive functions. In general, there remains a large difference between the kinds of models offered by cognitive neuroscientists or psychologists on the one hand, and those offered by systems neuroscientists on the other: the former tend to be high-level, where components of the model are large portions of cortex, while the latter tend to be low-level, where each

component is a single cell. This is true despite the fact that researchers in these areas share a similar interest in brain-based explanations of behavioral phenomena. Here I apply the neural engineering framework (NEF) methodology described in Eliasmith & Anderson (2003), which enables the construction of a model that is both high-level and low-level in this sense.

As mentioned, BioSLIE is an application of this methodology to the Wason card selection task (Wason, 1966). In the Wason task, subjects are given a conditional rule of the form "if P, then Q". They are then shown four cards. Each card expresses the satisfaction (or not) of condition P on one side and the satisfaction (or not) of condition Q on the other. The four visible card faces show representations of 'P', 'Q', 'not-P', and 'not-Q'. Subjects are instructed to select all cards which must be turned over in order to determine whether the conditional rule is true. A vast majority of subjects (greater than 90%) do not give the logically correct response (i.e., P and not-Q). Instead, the most common answer is to select the P and Q cards, or just the P card (Oaksford and Chater, 1994). However, it became apparent that performance on the task could be greatly facilitated by changing the content of the task to be more realistic or thematic, often by making the rule a permissive one (Sperber, 1995). To distinguish these two versions of the task, I refer to them as the 'abstract' and 'permissive' versions of the task respectively. Human performance on the Wason task is an ideal target for providing a neural model of cognition because it is generally considered a phenomena that can only be explained by invoking structure-sensitive processing. As a result, the task allows BioSLIE to demonstrate its ability to generalize across structures, i.e. to be systematic – a hallmark of cognitive systems (Fodor and Pylyshyn, 1988; Jackendoff, 2002). In addition, the context-sensitive performance of humans has not yet been adequately explained.

Indeed, there has been much debate over the nature of the mechanism which causes the marked difference in performance on the abstract and permissive versions of the task. Cosmides (1989) suggests that completely different, uniquely evolved modules are invoked under the two contexts. Cheng and Holyoak (1985), in contrast, suggest that the difference between deontic and non-deontic contexts results in the application of different, multi-step reasoning schemas. BioSLIE relies on a weaker claim than either of these. Specifically, in light of data that the deontic/non-deontic distinction does not adequately capture human performance variation (Oaksford and Chater, 1996), it embodies the assumption that the infer-

ence to be performed in a specific context has merely been learned in a similar context (regardless of its status as deontic or not). As well, it is assumed that the necessary inference need not be spelled out as a complex, multi-step schema, but is rather a direct mapping between the presented rule and appropriate responses in that context.

Notably, Cosmides (1989) challenges any theory based on induction, like that underlying BioSLIE, to lay out the mechanistically defined domain-general procedures that can take modern human experience as statistically encountered as input, and produce the observed domain specific performance in the selection task as output. This is a challenge that BioSLIE meets.

## Model description

BioSLIE integrates advances in structured vector representations, relevant physiological and anatomical data from frontal cortices (Wharton & Grafman, 1998), and the NEF, to explain human performance on the Wason task.

Since the early 1990s, there have been a series of suggestions as to how to incorporate structure-sensitive processing in models employing distributed representations (including Spatter Codes (Kanerva 1994); Holographic Reduced Representations (HRRs; Plate 1991); and Tensor Products (Smolensky 1990)). Few of these approaches have been used to build models of cognitive phenomena (although see Eliasmith & Thagard 2001). However, none of these methods have been employed in a biologically plausible computational setting. As described in the next section, I extend the NEF to incorporate HRRs, in order to integrate the structure sensitivity of the latter with the biological plausibility of the former.

Of course, to use this characterization of structure-sensitive processing in an explanatorily useful model, it is essential to suggest which anatomical structures may be performing the relevant functions. Only then is it possible to bring to bear the additional constraints of (and make predictions relating to) single cell physiology and functional imaging data. Figure 1 shows how BioSLIE is mapped to functional anatomy. Specifically, the network consists of: a) input from ventromedial prefrontal cortex (VMPFC) which provides familiarity, or context, information that is used to select the appropriate transformation (Adolphs et al. 1995); b) left language areas which provide representations of the rule to be examined (Parsons, Osherson, & Martinez 1999); and c) anterior cingulate cortex (ACC) which gives an error signal consisting of either the correct answer, or an indication that the response was correct or not (Holroyd & Coles 2002). The neural populations that make up BioSLIE itself model right inferior frontal cortex, where VMPFC and linguistic information is combined to select and apply the appropriate transformation to solve the Wason task (Parsons & Osherson 2001). It is during the application of the transformation that learning is also presumed to occur.

## Model derivation

### HRRs in spiking networks

Following Plate (1991) BioSLIE encodes structure in a distributed vector representation using circular convolution ( $\otimes$ ), which implements a kind of vector binding. In order to decode the structure, circular correlation ( $\oplus$ ) is used. These

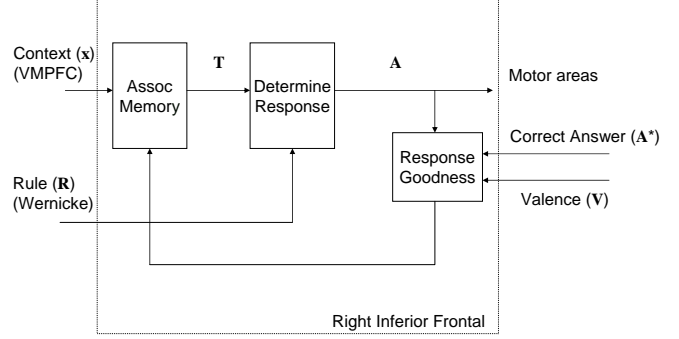


Figure 1: Functional decomposition and anatomical mapping of the model. The letters in bold indicate the vector signals in the model associated with the area.

operations are defined as:

$$\mathbf{C} = \mathbf{A} \otimes \mathbf{B} \quad \text{and} \quad \mathbf{B} \approx \mathbf{A} \oplus \mathbf{C}$$

$$c_j = \sum_{k=0}^{n-1} a_k b_{j-k} \quad b_j = \sum_{k=0}^{n-1} a_k c_{j+k}$$

where subscripts are modulo  $n$ . Conveniently, correlation can be defined in terms of convolution:  $\mathbf{A} \oplus \mathbf{C} = \mathbf{A}' \otimes \mathbf{C}$ , where  $'$  indicates an approximate inverse.

To implement convolution in a spiking network using the NEF, we must first define the encoding and decoding for a vector  $\mathbf{x}$  in a population  $a$  of neurons  $a_i$ . The encoding describes the biophysical processes that result in a series of rapid neural voltage changes (i.e., neural spikes). The decoding determines how much information those spikes carry about the original signal  $\mathbf{x}$ , by determining an estimate of that signal,  $\hat{\mathbf{x}}$ :

*Encoding*

$$a_i(t) = \sum_{n=1}^N \delta(t - t_{in}) = G_i \left[ \alpha_i \left\langle \mathbf{x} \cdot \tilde{\phi}_{im} \right\rangle J_i^{bias} \right]$$

*Decoding*

$$\hat{\mathbf{x}} = \sum_{i=1, n=1}^{N_i, N} h_i(t - t_n) \phi_i^{\mathbf{x}}$$

where  $\delta_i(\cdot)$  are the  $N_i$  spikes at times  $t_n$  for neuron  $a_i$ , generated by the spiking nonlinearity  $G_i$  in the population with  $N$  neurons. The neuron parameters  $\alpha_i$ ,  $\tilde{\phi}_i$ , and  $J_i^{bias}$  are the gain, preferred direction vector in stimulus space, and bias current respectively, which are chosen to reflect the heterogeneity of neuron responses observed in cortex (Eliasmith and Anderson, 2003). For the decoding,  $h_i(t)$  are the linear decoding filters, which for reasons of biological plausibility, are taken to be the post-synaptic currents (PSCs) generated in the subsequent neuron's dendrites, and the decoding vectors,  $\phi_i^{\mathbf{x}}$ , determine the importance of that neuron's response to the estimate of  $\mathbf{x}$ . Notably, the neural nonlinearity  $G_i$  can be as complex (i.e. biologically realistic) as desired. In BioSLIE we use a standard leaky integrate-and-fire (LIF) model.

Assuming this kind of vector representation in four populations,  $a$ ,  $b$ ,  $c$ , and  $d$ , it is possible to implement circular convolution. Using the convolution theorem, we know that any convolution in a domain  $\mathbf{x}$  is equivalent to multiplication in its Fourier domain. Thus, the two vectors to be convolved are projected through the Fourier matrix into a middle layer (using the encoding defined earlier):



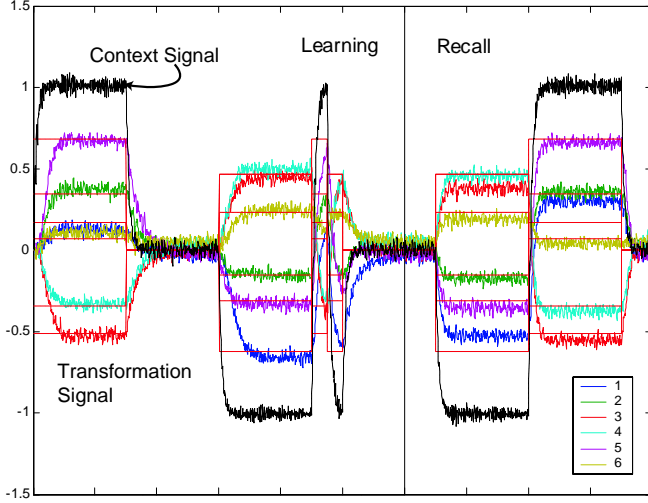


Figure 4: Learning and retrieval of a 6-dimensional vector in a spiking network. During the first two-thirds of the simulation, the context signal is changed while the input from  $z$  is changed to associate the context with the vector represented by  $z$ . In the last third, learning is turned off, and successful retrieval of the vectors is displayed given a context signal.

## Results on the Wason task

### Performance in different contexts

BioSLIE combines these subnetworks as shown in figure 5, resulting in a model that consists of ten interconnected neural populations, for a total of approximately twenty thousand neurons. The representations in this network have been scaled up to 100 dimensions, in order to encoded the vectors needed to perform the task.

The model is able to reproduce the typical results from the Wason task under both the abstract and permissive contexts, as shown in figure 6. In order to classify the results produced by the model, the resulting vectors must be ‘cleaned-up’. That is, they are compared to all possible labeled answers by taking a similarity measure (dot product) between the resulting vector and items in the ‘clean-up’ memory (i.e., all labeled vectors used in any simulation presented; 19 vectors). The labels on the graph indicate the similarity measures (maximum of 1). The top three responses are displayed to demonstrate the large difference in similarity between the provided answers and the next most similar vector. Simple thresholding can thus be used to determine what counts as an answer and what does not.

When the run in figure 6 begins, learning is initiated, the context is set to ‘abstract’ (i.e. 1), and the correct result in that context is present to the network. The network then learns to infer (in this context) the expected (incorrect) result (i.e.,  $a$  and  $b$ ). Both the context and expected result is then changed for the second phase, and the network learns a different transformation in the new ‘permissive’ context (i.e., -1; resulting in  $a$  and  $not-b$ ). Learning is then turned off, and the expected result is no longer presented to the network. Only the context signal is changed. As expected, in the abstract context the network’s answer is  $a$  and  $b$ , and in the permissive context

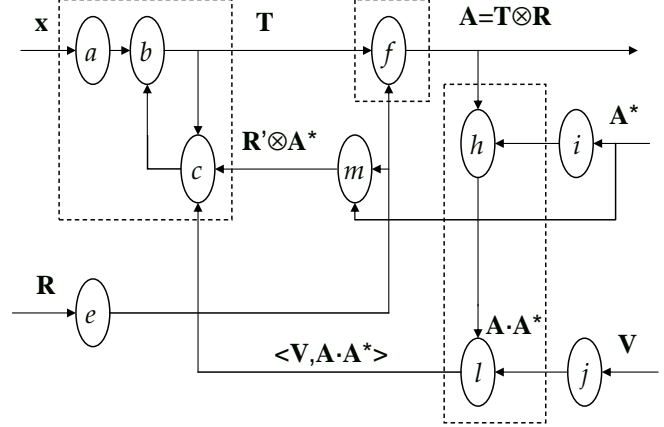


Figure 5: The complete network at the population level. The lower case letters indicate populations of approximately 2000 neurons each. Upper case letters indicate the signals being sent along the relevant projections. The dotted boxes indicate how this diagram relates to figure 3, and hence the anatomical mapping discussed earlier.

the network’s answer is  $a$  and  $not-b$ . BioSLIE has learned to perform different inferences in different contexts, resulting in similar performance to human subjects on the Wason task.

### Generalization within a context

To demonstrate that the network is truly learning a language-like transformation in a context, figure 7 shows that it generalizes learned, structure-sensitive transformations to new representations. This demonstrates that the system has learned a systematic regularity. That is, it can transform the structured representation based solely on the syntax of that representation.

This simulation is similar to that presented previously, except the context signal is kept constant and there are three separate rules that are presented to BioSLIE. During the learning ‘on’ phase, the rules *Implies(a,b)* and *Implies(c,d)* along with their expected answers are presented to the network. The learning is then turned off, and it is presented with *Implies(e,f)*. As expected, since the context is the same as the previous examples, the same transformation is applied, and BioSLIE infers that  $e$  and  $f$  are the expected answer. In the last quarter of the simulation, no rule is presented and thus no answer is produced (i.e., all similarity measures are very low). The similarity measures of the top three most similar vectors in the clean-up memory are displayed to demonstrate that the top two responses are the appropriate answers in this context.

## Conclusion

I have successfully built and simulated BioSLIE, a low-level, spiking neuron model of a high-level cognitive behavior, the Wason task. Compared to past models of the Wason task, this model has a number of advantages. In contrast to Cosmides’ explanation, I have presented a detailed computational model that demonstrates that a domain general mechanism can indeed account for the observed phenomena. As well, this

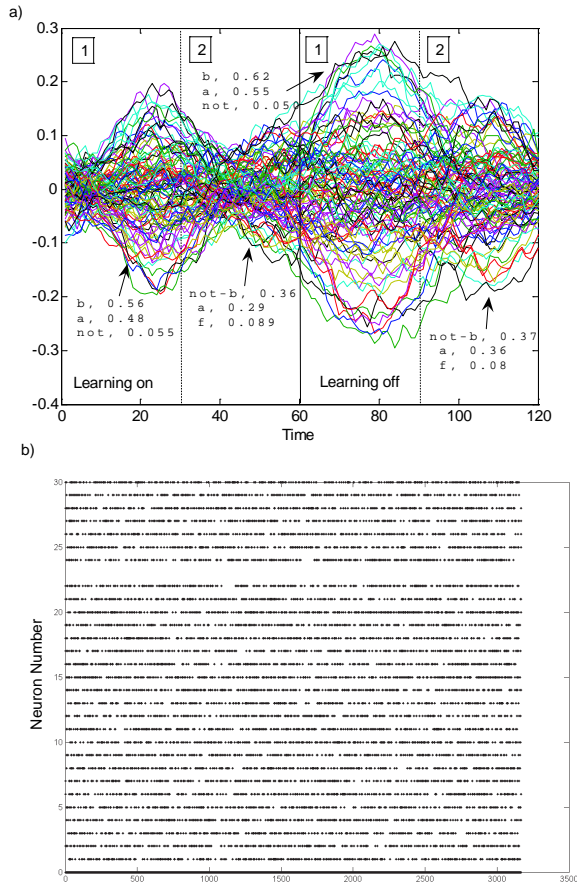


Figure 6: Results of the Wason task for the complete network. a) This is the decoded neural output representing the 100-dimensional vector. Results have been smoothed for legibility. The similarity measures for labeled vectors are written above the decoded neural output. These similarity measures indicate which answers are returned (and hence which inferences are performed) in those contexts. The numbers in boxes indicate the context (1 for abstract, 2 for permissive). This diagram shows that in the first two contexts, two different transformations are learned. Then, these transformations are appropriately applied when those contexts are later encountered. b) The neural spikes produced by every 100th neuron in the population representing these results.

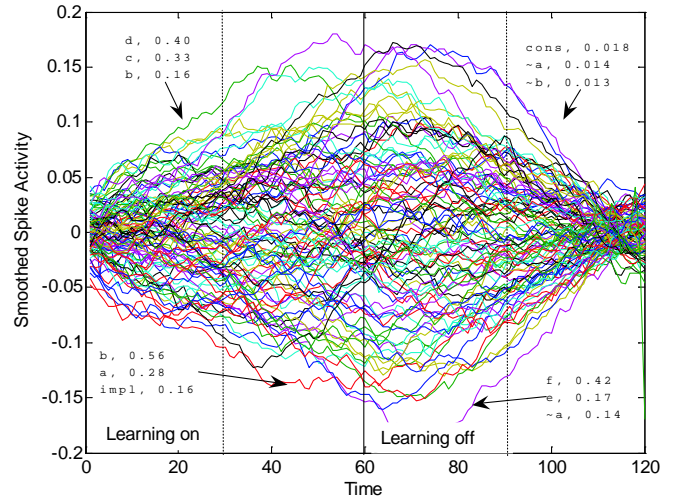


Figure 7: Generalization across different rules in the same context. See text for discussion.

model offers an advantage over pragmatic reasoning schemas, in that the transformations used to solve the problem are not limited to a pre-specified, discrete group of “schemas” which are identical for all subjects. Instead, each individual can solve the problem using her own estimation of the correct transformation in the given context, as determined by her idiosyncratic learning history in that context. BioSLIE is thus not restricted to the binary “deontic” and “non-deontic” distinctions made by Cheng and Holyoak.

More generally, because BioSLIE spans what are often considered disparate levels of description of cognitive phenomena, it is also able to support predictions at those various levels. I believe it is the first model to do so.

At the single neuron level, BioSLIE helps clarify the kinds of properties neurons involved in these computations need to have. For instance, in order to implement the high-dimensional nonlinear vector transformation necessary to capture structure-sensitive behavior, neurons in the simulation need to respond to at least two dimensions at the same time (one from each vector being convolved). However, despite the size of the vectors being convolved, *no more* than two dimensions need to be represented either (suggesting that the model will scale very well). As well, the learning rule derived to implement the network has implications for neuron connectivity. Specifically, it suggests that information carried by projections from one of the associative populations serves to direct the modification of synaptic weights between the memory population, and the other associative population. Thus the biophysical mechanisms (e.g. NO transport) that can support this kind of learning should be prevalent in these areas. Finally, examining the spike trains produced by the model show that despite the model being deterministic, the nonlinearities in the model serve to generate very random-looking spike trains, like those observed in cortex. This suggests that perhaps nonlinearities, not noise, are largely responsible for spike train variability in frontal areas.

At the behavioral level, the model not only meets Cosmides’ challenge of specifying an inductive, domain general

inference mechanism, it also makes it possible to predict behavioral variations on the task given a learning history. For instance, it should be possible to predict the effects of varying the kind of feedback that a subject receives in similar and dissimilar contexts. As well, I have not discussed the differing effects of explicit (i.e., the answer) versus implicit (i.e. ‘right’ or ‘wrong’) feedback on learning transformations. However, the model includes a valence signal that can be used to examine these differences. Finally, the ability of the model to generalize helps explain why those trained in logic do better on the content-independent tasks (Rinella, Bringsjord, & Yang 2001).

More theoretically, the model is at the very least an existence proof that understanding neural computation might have important implications for understanding cognitive behavior, *contra* Fodor, Pylyshyn, and Jackendoff. This is because HRRs, unlike classical symbols, are noisy representations. Thus, there are serious limitations on memory size, depth of structure, etc., that can be encoded by BioSLIE, just as there are for people. Understanding how well such noisy representations can be processed by a realistic neural system helps pave the way to better understanding these limitations. These same properties show how this model is also not a ‘mere’ implementation of a classical system. Classical systems are perfectly compositional and systematic. However, BioSLIE clearly is not, since encoding essentially blurs the represented constituents. Nevertheless, BioSLIE has enough compositionality and systematicity to model human cognitive performance. Thus, neurocomputational models, like BioSLIE, can help us understand the *degrees* of systematicity and compositionality possessed by real cognitive systems in ways that classical models cannot.

## References

- [1] Adolphs, R., A. Bechara, D. Tranel, H. Damasio, and A. Damasio, 1995. “Neuropsychological approaches to reasoning and decision-making.” In A. Damasio, H. Damasio and Y. Christen., eds., *Neurobiology of Decision-Making*. New York: Springer Verlag.
- [2] Cheng, P. W., and Holyoak, K. J. 1985. “Pragmatic reasoning schemas.” *Cognitive Psychology* 17:391–416.
- [3] Cosmides, 1989. “The logic of social exchange: Has natural selection shaped how humans reason? Studies with the Wason selection task.” *Cognition*, 31:187–276,
- [4] Eliasmith, C., and Anderson, C. H. 2003. *Neural engineering: Computation, representation, and dynamics in neurobiological systems*. Cambridge, MA: MIT Press.
- [5] Eliasmith, C., and Thagard, P. 2001. “Integrating structure and meaning: A distributed model of analogical mapping.” *Cognitive Science* 25:245–286.
- [6] Fodor, J., and Pylyshyn, Z. 1988. “Connectionism and cognitive science: A critical analysis.” *Behavioral and Brain Sciences* 28:3–71.
- [7] Holroyd, C., and Coles, M. 2002. “The neural basis of human error processing: Reinforcement learning, dopamine, and the error-related negativity.” *Psychological Review* 109:679–709.
- [8] Jackendoff, R. 2002. *Foundations of Language: Brain, Meaning, Grammar, Evolution*. Oxford University Press.
- [9] Kanerva, P. 1994. “The spatter code for encoding concepts at many levels.” *Proceedings of International Conference on Artificial Neural Networks* 46:226–229.
- [10] Neumann, J. 2001. *Holistic Processing of Hierarchical Structures in Connectionist Networks*. PhD dissertation, University of Edinburgh, Department of Computer Science.
- [11] Oaksford and Chater, 1994, “A rational analysis of the selection task as optimal data selection,” *Psychological Review*, 101(4), 608–631.
- [12] Oaksford and Chater, 1996. “Rational explanation of the selection task” *Psychological Review* 103(2):381–391
- [13] Parsons, L., and D. Osherson. 2001. “New evidence for distinct right and left brain systems for deductive versus probabilistic reasoning.” *Cerebral Cortex* 11:954–965.
- [14] Parsons, L., D. Osherson and M. Martinez. 1999. “Distinct neural mechanisms for propositional logic and probabilistic reasoning.” *Proceedings of the Psychonomic Society Meeting* 61–62.
- [15] Plate, A. 1991. “Holographic reduced representations: Convolution algebra for compositional distributed representations.” In Mylopoulos, J., and Reiter, R., eds., *Proceedings of the 12th International Joint Conference on Artificial Intelligence*. San Mateo, CA: Morgan Kaufmann.
- [16] Rinella, K., S. Bringsjord, Y. and Yang. 2001. “Efficient logic instruction: People are not irremediably poor deductive reasoners.” In Moore, J., and Stenning, K., eds., *Proceedings of the 23rd Annual Conference of the Cognitive Science Society*, Mahwah, NJ: Lawrence Erlbaum Associates. 851–856.
- [17] Smolensky, P. 1990. “Tensor product variable binding and the representation of symbolic structures in connectionist systems.” *Artificial Intelligence* 46:159–217.
- [18] Sperber, Cara, and Girotto, 1995. “Relevance theory explains the selection task.” *Cognition* 57: 31–95.
- [19] Wason, P. C. 1966. “Reasoning.” In Foss, B. M., ed., *New horizons in psychology*. Harmondsworth: Penguin.
- [20] Wharton, C., and Grafman, J. 1998. “Deductive reasoning and the brain.” *Trends in Cognitive Sciences* 2:54–59.