

Received STI

JAN 08 1991

Los Alamos National Laboratory is operated by the University of California for the United States Department of Energy under contract W-7405-ENG-36

LA-UR--90-4251

DE91 005878

TITLE: GENETIC ALGORITHMS AND THE IMMUNE SYSTEM

AUTHOR(S): STEPHANIE FORREST
ALAN S. PERELSON

SUBMITTED TO: For Proceedings, Workshop on Parallel Problem Solving
From Nature, 1990. Computer Society of IEEE et al., Dortmund, Germany
October, 1990. Springer-Verlag, New York

DISCLAIMER

This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency thereof, nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.

By acceptance of this article, the publisher recognizes that the U.S. Government retains a nonexclusive, royalty-free license to publish or reproduce the published form of this contribution, or to allow others to do so, for U.S. Government purposes.

The Los Alamos National Laboratory requests that the publisher identify this article as work performed under the auspices of the U.S. Department of Energy.

MASTER

Los Alamos Los Alamos National Laboratory
Los Alamos, New Mexico 87545

JKS

Genetic Algorithms and the Immune System *

Stephanie Forrest

forrest@unmvax.cs.unm.edu
Dept. of Computer Science
University of New Mexico
Albuquerque, N.M. 87131 USA

Alan S. Perelson

asp@receptor.lanl.gov
Theoretical Division
Los Alamos National Laboratory
Los Alamos, N.M. 87545 USA

Abstract

Using genetic algorithm techniques we introduce a model to examine the hypothesis that antibody and T cell receptor genes evolved so as to encode the information needed to recognize schemas that characterize common pathogens. We have implemented the algorithm on the Connection Machine for 16,384 64-bit antigens and 512 64-bit antibodies.

1 Introduction

The immune system is our basic defense system against bacteria, viruses and other disease-causing organisms. In order to provide its defense functions efficiently, the immune system must perform pattern recognition tasks to distinguish self molecules and cells from foreign ones (antigens). The number of foreign molecules that the immune system can recognize is unknown but it has been estimated to be greater than 10^{16} . In practical terms, essentially any foreign molecule presented to the immune system, even those created in the laboratory and thus never having appeared before in all of evolutionary time, are recognized as being foreign. Besides this immense recognition capacity, the other feature that distinguishes the vertebrate immune system from the defense systems of lower organisms is that the immune system learns and exhibits memory. Thus the response to the second exposure of the same antigen occurs more quickly and vigorously than the first response.

In this paper we introduce a model that addresses long-term learning and pattern recognition in immune systems. By "long-term" we mean evolutionary time scales rather than the time scale of an individual's response to antigen. In particular, we are interested in understanding what types of information the immune system needs to store and process in order to defend us against a wide range of pathogens.

Recognition in the immune system occurs via receptor molecules on the surfaces of a class of white blood cells known as lymphocytes. These receptors are very diverse, so that with high probability each lymphocyte has receptors with different specificity. For

*To appear in the Proceedings of the 1990 Workshop on Parallel Problem Solving from Nature

B lymphocytes, which secrete antibody molecules, the receptors are a membrane-bound form of the antibody the cell will secrete if stimulated to do so. For T lymphocytes, the receptor is simply called the T cell receptor. Both types of receptors are proteins, and as such are coded for in the DNA of the organism. A mouse is thought to be able to make on the order of 10^{11} different receptor molecules (Berek and Milstein, 1988), even though its entire genome probably contains fewer than 10^5 genes! The solution that the immune system discovered for genetically encoding the information to make 10^{11} receptors is combinatorics. The variable portion of each receptor is coded for by five gene segments, each segment being chosen apparently randomly from a different gene library.

A question that puzzles immunologists is how evolution selects appropriate gene segments, since each gene segment by itself does not code for a receptor. Further, there are thousands of gene segments that can randomly combine to form the gene that ultimately codes for the receptor. Thus changing one segment should have little effect on the survival of an individual or even a species. Similar questions arise in artificial systems built on evolutionary principles, such as classifier systems (Holland et al., 1986). In classifier systems, the individual components (rules) are each intended to play a unique and complementary role with respect to other components in the system. Yet, the success of the system is measured collectively and evolutionary pressure is applied at the global level. The question of how the appropriate components can evolve from a global selection process is thus common to both classifier and immune systems.

2 The Model

We have developed a model directed at understanding the genetic evolution of the gene segment libraries that control the production of antibodies and T cell receptors. The hypothesis is that over evolutionary time scales these libraries evolved biases towards recognizing common pathogens. Assuming that this type of learning takes place through natural selection, the genetic algorithm provides a natural model for studying its behavior.

The model is based on a universe in which both antigens and receptors on B cells and T cells are represented by binary strings (cf. Farmer et al., 1986). This is certainly a simplification from the real biology in which genes are specified by a four-letter alphabet and recognition between receptors and antigens is based on their three-dimensional shapes and physical properties. However, this abstract universe is rich enough to allow us to study how a relatively small number of building blocks (the entries in the gene libraries) can be combined to recognize large classes of composite patterns. Our experiments and calculations have been based on genotypes of length 64, although it is in principle possible to represent any length genotype in the model.

The initial model makes the important simplification that a bitstring represents both the genes that code for a receptor and the phenotypic expression of the receptor molecule. As a further simplification, the model does not encode the concept of a library. Thus, all of the bits (genes) in a bitstring are used simultaneously to determine the bitstring's fitness.

The model includes only recognition of our idealized antigens by receptors and does not consider how the immune system neutralizes an antigen once it is recognized. A receptor

Antibody: 11001001000100100001001010101010
 Antigen: 01111100111001011110110101110100
 Complement: 10110101111101111111111111011110

Length of contiguous substrings: 1, 2, 1, 5, 13, 4

Figure 1: Scoring complementary matches between antigens and antibodies

or “antibody” is said to *match* an antigen if their bitstrings are complementary. Since each antibody must match against several different antigens simultaneously, we do not require perfect bit-wise matching. There are many possible match rules that make physiological sense (Perelson, 1989). Here we quantify the degree of match by a matching function $M : Antibody \times Antigen \rightarrow \mathfrak{R}$. M identifies contiguous regions of complementary bitwise matches within the string. M computes the lengths of the regions (l_i), and combines them such that long regions are rewarded more than short ones. Figure 1 illustrates the matching procedure. Using this basic idea, many different specific functions can be defined that are nonlinear in l_i . For our initial work we used an exponential function introduced by Stadnyk (1987), $\sum_i \frac{2^{l_i}}{(l - l_i + 1)}$, where l is the total length of the bitstrings.

Using the bitstring representation for antibodies and a fitness function M to score matches between antigens and antibodies, we then construct one population of antigens and one of antibodies. The antibodies are matched randomly against a sample of antigens, scored according to M , and replicated using a conventional genetic algorithm. To approximate the constraints of the real immune system, we use antigen populations of size 16,384 (2^{14}) and antibody populations of 512 (2^9). However, it will be interesting to vary this to understand the “recognition capacity” of different size antibody populations. Figure 2 illustrates the basic immune model.

The model is implemented on a Connection Machine. In the implementation, each processor holds one antigen and one antibody, so that all the matches are performed in parallel. Each of the basic genetic algorithm operations (fitness evaluation, selection and reproduction, mutation, and cross-over) are parallelized so that we can vary the size of antigen and antibody populations without affecting the running time of the model (up to the physical limits of the machine).

Since there are many fewer antibodies than antigens, the antibody population will recognize the most antigens if it can find common patterns among antigens (that is, if it can generalize across the antigen population). These patterns are analogous to the “schemas” described in the genetic algorithms literature (Goldberg, 1989). If we consider the probability that an arbitrary k -bit schema will appear in a significant fraction of the population, it is straightforward to show that this probability diminishes quickly as k and the size of the population approach reasonable values. To be more precise, consider the case in which the k -bit pattern is the first k bits. If 0 and 1 are chosen with equal probability in generating the antibody and antigen bitstrings, then p , the probability that the first k bits in an antibody will match the first k bits of an antigen, is given by $p = 2^{-k}$. In

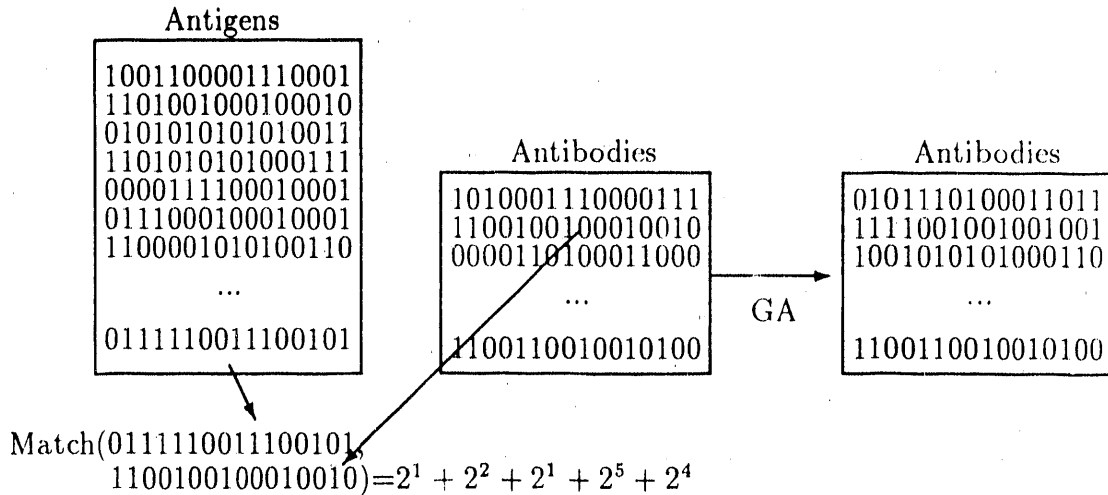


Figure 2: A schematic illustration of the immune model

an antigen population of size N , the probability that the antibody matches at least N_0 antigens in the first k bits is then given by

$$P(k; N_0) = \sum_{i=N_0}^N \binom{N}{i} p^i (1-p)^{N-i} = I_p(N_0, N - N_0 + 1),$$

where $I_p(N_0, N - N_0 + 1)$ is an incomplete beta function (Abramowitz and Stegun, 1964).

For the genetic algorithm to succeed in identifying the first k bits as a schema, the pattern needs to appear in a significant fraction ($\frac{N_0}{N}$) of the antigens. Typically each antibody is matched against a sample of the antigens and thus the schema will need to occur frequently enough to be detected by the sampling process. For $N_0 = 0.1N$, $P(k) \approx 1$ for $k \leq 3$. For larger values of k , $P(k) \leq 10^{-75}$. Thus schemas of length greater than 3 will be impossible to discover by random antigen sampling. Patterns of length 3 are so common that they would also be present on self molecules and thus the immune system could not use them as a marker for antigen recognition.

We are thus led to conclude that in order for the immune system to solve the pattern recognition problem posed above, it must operate on slightly different principles from what we have formulated. For example, the population of antigens may not be random. A bias could exist in the antigen population if pathogens have particular structures on their surfaces that are different from the animals they infect. This is the case for bacteria with cell walls made of polysaccharides that are not found in mammals. These types of biases are incorporated in the model in two ways: (1) choosing 0 more frequently than 1 when constructing antigens, (2) by prespecifying certain schemas when the antigens are constructed. Alternatively, biases could be created implicitly if the antigens are evolving to evade immune detection. A second possibility is that the immune system has evolved to recognize antigen and not recognize self. Modeling self as another set of strings, the problem the immune system would now have to solve is recognizing patterns that occur in the antigen population but not the self population. There is some evidence that this may be the case (cf. Claverie et al., 1988).

A third alternative is that each antibody need not be compared with the entire antigen population when computing its fitness. Recall, we are trying to solve a "covering problem" in which each antibody recognizes a subset of antigens, and the union of the subsets of recognized antigens include all antigens. In our simulations the number of antigens is 32 times the number of antibodies. Thus, in the best possible solution each antibody need only match 32 antigens. Evaluating $P(k)$ with $N_0 = 32$, we find $P(k) \approx 1$ for $k \leq 8$, $P(9) = 0.52$, $P(10) = 0.00027$, and $P(11) = 1.3 \times 10^{-10}$. Thus, an algorithm in which each antibody is evolved to recognize 32 antigens may be able to discover schemas of length 10, which should be sufficiently unique to identify antigens and even distinguish them from self molecules. We are currently using our model to study these alternatives.

3 Acknowledgements

Portions of this work were done under the auspices of the U. S. Department of Energy and the Santa Fe Institute, and supported by a grant from the National Institutes of Health (AI28433). We are also grateful to John Holland and Rob De Boer for helpful comments.

4 References

- Abromowitz, M. and Stegun, I. A. (1964). Handbook of Mathematical Functions. National Bureau of Standards, Washington, D. C., p.263.
- Berek, C. & Milstein, C. (1988). The dynamics nature of the antibody repertoire. *Immunol. Rev.* **105**, 5-26.
- Claverie, J.-M., Kourilsky, P., Langlade-Demoyen, P., Chalufour-Prochnicka, A., Dadaglio, G., Plata, F. and Bougueleret, L. (1988). T-immunogenic peptides are constituted of rare sequence patterns. Use in the identification of T epitopes in the human immunodeficiency virus *gag* protein. *Eur. J. Immunol.* **18**, 1547-1553.
- Farmer, J. D., Packard, N. H. and Perelson, A. S. (1986). The immune system, adaptation, and machine learning. *Physica D* **22**: 187-204.
- Goldberg, D. E. (1989). *Genetic Algorithms in Search, Optimization, and Machine Learning*, Addison-Wesley, Reading, MA.
- Holland, J. H., Holyoak, K.J., Nisbett, R.E., and Thagard, P. (1986) *Induction: Processes of Inference, Learning, and Discovery* MIT Press, Cambridge, MA.
- Perelson, A. S. (1990). Theoretical immunology. In *1989 Lectures in Complex Systems*, SFI Studies in the Sciences of Complexity, Lect. Vol. II, E. Jen, ed., Addison-Wesley, Redwood City, CA, pp. 465-499.
- Stadnyk, I. (1987). Schema recombination in pattern recognition problems. *Proc. 2nd International Conference on Genetic Algorithms and their Applications*, Lawrence Erlbaum Assoc., Hillsdale, NJ.

END

DATE FILMED

02 / 13 / 91

